



Universidad de  
Oviedo



# ESCUELA POLITÉCNICA DE INGENIERÍA DE GIJÓN

MÁSTER UNIVERSITARIO EN INGENIERÍA INFORMÁTICA  
ÁREA DE CIENCIA DE LA COMPUTACIÓN E INTELIGENCIA  
ARTIFICIAL

## ANEXO I: CÓDIGO

Aprendiendo perfiles de usuario: distinción entre  
intereses consolidados e intereses recientes

Pérez Núñez, Pablo

TUTOR: D. Jorge Díez Peláez

FECHA: 11 de julio de 2018

# Índice

<b>1. Introducción</b>	<b>1</b>
<b>2. Clase LastFM</b>	<b>2</b>
<b>3. Métodos</b>	<b>13</b>
3.1. Word2Vec . . . . .	13
3.2. Doc2Vec-DM . . . . .	18
3.3. Doc2Vec-DBOW . . . . .	24
<b>4. Tests</b>	<b>29</b>
4.1. TEST 01 . . . . .	29
4.2. TEST 02 . . . . .	37
4.3. TEST 03 . . . . .	45
4.4. TEST 04 . . . . .	53
4.5. TEST 05 . . . . .	62
4.6. TEST 06 . . . . .	70
4.7. TEST 07 . . . . .	78
4.8. TEST 08 . . . . .	86
4.9. TEST 09 . . . . .	95
4.10. TEST 10 . . . . .	103
4.11. TEST 11 . . . . .	111
4.12. TEST 12 . . . . .	119
4.13. TEST 13 . . . . .	127
4.14. TEST 14 . . . . .	135
4.15. TEST 15 . . . . .	143
4.16. TEST 16 . . . . .	151
4.17. TEST 17 . . . . .	159
4.18. TEST 18 . . . . .	167
4.19. TEST 19 . . . . .	175
4.20. TEST 20 . . . . .	183
4.21. TEST 21 . . . . .	191
4.22. TEST 22 . . . . .	199
4.23. TEST 23 . . . . .	207
4.24. TEST 24 . . . . .	215

# 1. Introducción

A continuación se incluirán cada uno de los ficheros de código realizados para llevar a cabo la experimentación. Para mayor comodidad, se han subido todos estos ficheros a un repositorio público:

<https://github.com/uo232368/TFM>.

Destacar la ausencia de los datos dentro del repositorio debido al gran tamaño de los mismos.

## 2. Clase LastFM

Esta clase contiene procedimientos para llevar a cabo la limpieza y gestión de datos además de otros métodos genéricos utilizados a lo largo de la investigación.

```

1  # -*- coding: utf-8 -*-
2  import pandas as pd
3  import numpy as np
4  import scipy.sparse as sp
5
6
7  import sys, time, os
8  import sqlite3 as lite
9  import fileinput, csv
10 from ggplot import *
11 import collections
12 from difflib import SequenceMatcher
13 import pickle
14 from datetime import datetime, timedelta
15 import re
16
17 import tensorflow as tf
18
19 class bcolors:
20     HEADER = '\033[95m'
21     OKBLUE = '\033[94m'
22     OKGREEN = '\033[92m'
23     WARNING = '\033[93m'
24     FAIL = '\033[91m'
25     ENDC = '\033[0m'
26     BOLD = '\033[1m'
27     UNDERLINE = '\033[4m'
28
29 class LastFM:
30
31     DATA_PATH = "/Users/pablo/Desktop/W2VEC/Datasets/lastfm"
32     DATABASE_NAME = "LastFM.sqlite"
33
34     #####
35     # Metodos
36     #####
37
38     def __init__(self):
39         self.DATA_PATH = "datos/"
40
41         #reload(sys)
42         #sys.setdefaultencoding('utf8')
43
44         pd.set_option('display.max_rows', 10)
45         pd.set_option('display.max_columns', 500)
46         pd.set_option('display.width', 1000)
47
48     def doQuery(self, SQL):
49
50         con = lite.connect(self.DATA_PATH+"/"+self.DATABASE_NAME)
51         con.text_factory = lambda x: str(x, 'utf-8', 'ignore')
52         dft = pd.read_sql_query(SQL, con)

```

```

53     return dft
54
55 def createDict(self, pickleName="datos/DUP_ID_DICT", similarity=0.95):
56     '''
57     Este metodo genera un diccionario donde, para cada nombre de cancion diferente, se le
58     asigna un ID.
59     Si existen repetidos , les asigna a todos el mismo id, de forma que (a,b,b',b',b',c)
60     tendrian ids (0,1,1,1,1,2)
61
62     Funcionamiento basico.
63     Para cada artista, se obtiene su lista de canciones.
64     Se compara cada cancion con el resto y, si existen similares se les asigna a todas
65     (las similares y la que se compara) el mismo id.
66
67     Posee alguna optimizacion para reducir el orden de complejidad a logaritmico.
68
69     Los resultados se almacenan en dos pickles:
70     DUP_ID_DICT => Diccionario con clave = Nombre de cancion y valor= ID
71     DUP_ID_DICT_INV => Diccionario anterior inevertido (clave=valor,valor=clave)
72
73     Genera tambien los pickles:
74     ALL_SONGS => Todas las canciones (incluidas repeticiones)
75     ALL_SONGS_UNIQUE => Las canciones sin repeticion
76     ALL_SONGS_ID_COUNT => Todas las canciones con los id generados y el numero de
77     veces que aparece cada cancion
78
79     :param: Grado de similitud deseado.
80     :return: Diccionario con, nombre de cancion => Id asignado. Ademas se guarda en un
81     pickle
82     '''
83
84     #####
85
86 def similar(a, b):
87     # Metodo para ver si 2 strings son parecidos. Retorna porcentaje en tanto por uno
88     return SequenceMatcher(None, str(a), str(b)).ratio()
89
90 def checkDuplicates(CURR_ID, SONG_LIST, ARTIST):
91     # A partir de la lista de canciones, se encarga de detectar duplicados e
92     incrementar el id actual
93     # el cual retorna para continuar en la funcion principal.
94
95     AVOID_LIST = []
96     RES = {}
97
98     TEMP_LIST = SONG_LIST
99     i = 0
100
101     regex = r"(\d)"
102
103     for s in SONG_LIST:
104         TEMP_LIST = TEMP_LIST[1:]
105
106         SIMILARS = []

```

```

104 # Si ya se visito en una iteracion anterior nada
105 if s in AVOID_LIST:
106     continue
107
108 # Buscar items similares del actual
109 for s2 in TEMP_LIST:
110     SML = similar(s, s2)
111
112     # Si ya se visito en una iteracion anterior nada
113     if s2 in AVOID_LIST:
114         continue
115
116     # Si son muy similares
117     if (SML > similarity):
118
119         m1 = re.finditer(regex, s)
120         m2 = re.finditer(regex, s2)
121         m1_items = []
122         m2_items = []
123         arequal = True
124
125         for matchNum, match in enumerate(m1):
126             m1_items.append(match)
127         for matchNum, match in enumerate(m2):
128             m2_items.append(match)
129
130         if (len(m1_items) == len(m2_items)):
131             for i in range(len(m2_items)):
132                 arequal = m1_items[i].group() == m2_items[i].group()
133
134                 if not arequal: break
135         else:
136             arequal = False
137
138         # Si son iguales
139         if (arequal):
140             SIMILARS.append(s2)
141             AVOID_LIST.append(s2)
142
143     i += 1
144
145 # Si el item posee repetidos, asignar el mismo id a todos
146 if (len(SIMILARS) > 0):
147     SIMILARS.append(s)
148     for i in range(len(SIMILARS)):
149         RES.update({str(ARTIST) + " - " + str(SIMILARS[i]): CURR_ID});
150     CURR_ID += 1;
151
152 else:
153     SIMILARS.append(s)
154     RES.update({str(ARTIST) + " - " + str(s): CURR_ID})
155     CURR_ID += 1;
156
157 if (len(RES) != len(SONG_LIST)):
158     print("Algo va mal...")
159
160 return RES, CURR_ID

```

```

161
162 def getID(s, DICT):
163     id = DICT.get(s.longname, None)
164     if (id is None):
165         return (-1)
166     else:
167         return (id)
168
169 #####
170
171 # Si no existe el diccionario, generar
172 if os.path.isfile(pickleName + ".pkl"):
173     print("Ya existe un diccionario previo...")
174     return()
175
176 if not os.path.isfile("datos/ALL_SONGS.pkl") and not
177 os.path.isfile("datos/ALL_SONGS_UNIQUE.pkl"):
178     print("No existen los pickles:\n\t · ALL_SONGS.pkl\n\t · ALL_SONGS_UNIQUE.pkl")
179
180     ALL_SONGS = self.doQuery("select userid,timestamp,artname,traname from
181                             user_track")
182     ALL_SONGS.artname = ALL_SONGS.artname.astype(str)
183     ALL_SONGS.traname = ALL_SONGS.traname.astype(str)
184     ALL_SONGS["longname"] = ALL_SONGS.artname + " - " + ALL_SONGS.traname
185     ALL_SONGS.to_pickle("datos/ALL_SONGS.pkl");
186
187     ALL_SONGS_UNIQUE = ALL_SONGS.drop_duplicates("longname")
188     ALL_SONGS_UNIQUE = ALL_SONGS_UNIQUE.sort_values("longname")
189     ALL_SONGS_UNIQUE.to_pickle("datos/ALL_SONGS_UNIQUE.pkl");
190
191 else:
192     ALL_SONGS_UNIQUE = pd.read_pickle("datos/ALL_SONGS_UNIQUE.pkl")
193     ALL_SONGS = pd.read_pickle("datos/ALL_SONGS.pkl")
194
195 #####
196
197 CURR_ID = 0
198
199 RES = {}
200
201 ARTISTS = ALL_SONGS_UNIQUE.groupby("artname")
202
203 for i, g in ARTISTS:
204     ARTIST = i
205     # Si hay mas de una cancion para el artista, hay que ver si se repiten.
206     if (len(g) > 1):
207         # Comparar los nombres de las canciones buscando repetidos
208         PARTIAL_RES, CURR_ID = checkDuplicates(CURR_ID, g.traname.values, ARTIST)
209         RES.update(PARTIAL_RES)
210
211     # Si solo hay una cancion para el artista, se le asigna el numero.
212     else:
213         RES.update({str(g.longname.values[0]): CURR_ID})
214         CURR_ID += 1
215     print(len(RES))
216
217 #-----

```

```

216 # Crear diccionario con "Para cada string, el id al que se refiere" y el inverso
217 #-----
218
219 self.save_pickle(RES, pickleName)
220
221 # Crear el diccionario inverso
222 '''ESTE DICCCIONARIO TENDRA MENOS ITEMS DADO QUE SE ELIMINAN LAS CLAVES REPETIDAS'''
223
224 INV_DICT = dict(zip(RES.values(), RES.keys()))
225 self.save_pickle(INV_DICT, pickleName + "_INV")
226
227 #-----
228 # Anadir al dataframe el ID y el numero de reproducciones
229 #-----
230 if not os.path.isfile("datos/ALL_SONGS_ID_COUNT.pkl"):
231
232     # Anadir a cada cancion su id
233     ID_COL = ALL_SONGS.apply(lambda row: getID(row, RES), axis=1)
234     ALL_SONGS_ID = pd.concat([ALL_SONGS, ID_COL], axis=1)
235     ALL_SONGS_ID.columns = ["user", "timestamp", 'artname', 'traname', 'longname',
236                             'ID']
236
237     ID_COUNT = ALL_SONGS_ID.groupby(['ID']).size().reset_index(name='counts')
238     RES = pd.merge(ALL_SONGS_ID, ID_COUNT, how='left', on=['ID'])
239     RES.timestamp = pd.to_datetime(RES.timestamp, format="%Y-%m-%dT%H:%M:%SZ")
240     RES.to_pickle("datos/ALL_SONGS_ID_COUNT.pkl");
241 else:
242     print("Fichero ya existente...")
243
244 #-----
245
246 def generaFicherosUsuario(self, filename):
247     '''
248     A partir de los datos (YA FILTRADOS) , genera para cada usuario un fichero con las
249     canciones escuchadas,
250     en orden y representadas por su id
251     :param: Diccionario con "Titulo"=>ID
252     '''
253
254     DATA = pd.read_pickle(filename)
255
256     # Obtener los id de usuarios
257     USER_LIST = DATA.groupby("user")
258
259     #Para cada usuario
260     for u,r in USER_LIST:
261
262         #Imprimir ID
263         self.printB(u)
264
265         # Obtener todas las canciones del usuario
266         USER_SONG_LIST = r
267
268         # Ordenar por fecha
269         USER_SONG_LIST.sort_values("timestamp", ascending=True, inplace=True)
270         USER_SONG_LIST = USER_SONG_LIST.reset_index()
271

```



```

271     # Lista con los id
272     ID_LIST = USER_SONG_LIST.ID.values
273
274     # Crear fichero con canciones para cada uno de los usuarios (Separado por ;)
275     file = open("datos/usuarios/" + str(u) + ".txt", "w")
276     file.write(";".join(str(x) for x in ID_LIST))
277     file.close()
278
279 def
getCleanData(self, songMin=10, userMin=100, years=[2005, 2006, 2007, 2008, 2009], saveAs="MAIN_DATA"):
280
281     if os.path.isfile("datos/ALL_SONGS_ID_COUNT.pkl"):
282         ALL_SONGS_ID_COUNT = pd.read_pickle("datos/ALL_SONGS_ID_COUNT.pkl")
283     else:
284         print("No existen el pickle:\n\t · ALL_SONGS_ID_COUNT.pkl [generado en
                createDict()]")
285         return ()
286
287     INIT_LEN = len(ALL_SONGS_ID_COUNT)
288
289     ALL_SONGS_ID_COUNT["year"] = ALL_SONGS_ID_COUNT.timestamp.dt.year.astype(int)
290
291     print("-"*70)
292     #-----
293     print(" · Eliminar datos que NO sean de los anos " + (',' .join(str(x) for x in
                years))+".")
294     ALL_SONGS_ID_COUNT = ALL_SONGS_ID_COUNT.loc[ALL_SONGS_ID_COUNT.year.isin(years)]
295
296     #-----
297     print(" · Eliminar canciones con menos de " + str(songMin) + " rerproducciones.")
298
299     #print(len(ALL_SONGS_ID_COUNT.loc[ALL_SONGS_ID_COUNT.counts < songMin]))
300     ALL_SONGS_ID_COUNT = ALL_SONGS_ID_COUNT.loc[ALL_SONGS_ID_COUNT.counts >=songMin]
301
302     #-----
303     print(" · Eliminar usuarios con menos de " + str(userMin) + " rerproducciones.")
304
305     USERS = ALL_SONGS_ID_COUNT.groupby(["user"])
306     BAD_USERS=[]
307
308     for u,d in USERS:
309         if len(d)<userMin: BAD_USERS.append(u)
310
311
312     ALL_SONGS_ID_COUNT = ALL_SONGS_ID_COUNT.loc[~ALL_SONGS_ID_COUNT.user.isin(BAD_USERS)]
313
314
315     #-----
316     #-----
317
318     '''
319     Es necesario anadir unos nuevos ID para facilitar el OneHot y eliminar los no
                utilizados
320     '''
321
322     #Se crea un diccionario con los IDs actuales (con saltos) y el correspondiente ID
                nuevo

```

```

323 key = ALL_SONGS_ID_COUNT.ID.unique().tolist()
324 value = range(0, len(key))
325 dict = dict(zip(key, value))
326
327 #Eliminar vieja columna y anadir la nueva
328 ID_NEW = ALL_SONGS_ID_COUNT.ID.apply(lambda row: dict.get(row))
329 ALL_SONGS_ID_COUNT = ALL_SONGS_ID_COUNT.drop(['ID'], axis=1)
330 ALL_SONGS_ID_COUNT = pd.concat([ALL_SONGS_ID_COUNT, ID_NEW], axis=1)
331
332 #-----
333
334
335 print("-"*70)
336
337 print("Reproducciones iniciales: "+str(INIT_LEN))
338 print("Reproducciones eliminadas: "+str(INIT_LEN-len(ALL_SONGS_ID_COUNT)))
339 print("Usuarios eliminados: "+str(len(BAD_USERS))+"\n\t" + ', '.join(BAD_USERS))
340
341 ALL_SONGS_ID_COUNT.to_pickle(saveAs+".pkl");
342
343
344
345
346 return ALL_SONGS_ID_COUNT
347
348 #####
349 # Metodos W2V
350 #####
351
352 def getSongDataArrays(self, train=75):
353     '''
354     Retorna un array de arrays [[cancion1,cancion2... ],[cancion1, cancion2...]]...
355     :param: Array de arrays
356     '''
357
358     filename = "datos/CLEAN_DATA_[10,100].pkl"
359
360     DATA = pd.read_pickle(filename)
361
362     # Obtener los id de usuarios
363     USER_LIST = DATA.groupby("user")
364
365     RES_TRAIN=[]
366     RES_REC=[]
367
368     RES_USERS=[]
369
370     #Para cada usuario
371     for u,r in USER_LIST:
372
373         # Obtener todas las canciones del usuario
374         USER_SONG_LIST = r
375
376         # Ordenar por fecha
377         USER_SONG_LIST.sort_values("timestamp", ascending=True, inplace=True)
378         USER_SONG_LIST = USER_SONG_LIST.reset_index()
379

```

```

380     # Lista con los id
381     ID_LIST = USER_SONG_LIST.ID.values.tolist()
382     # Lista con las fechas
383
384     # Separar para embeddigns y para recomendacion
385     TRAIN_INDX = int(len(ID_LIST)*(train/100))
386
387     TRAIN_LIST = ID_LIST[:TRAIN_INDX]
388     REC_LIST = ID_LIST[TRAIN_INDX:]
389
390     #Pasar a string las listas
391     TRAIN_LIST = list(map(str, TRAIN_LIST))
392     REC_LIST = list(map(str, REC_LIST))
393
394     RES_TRAIN.extend(TRAIN_LIST)
395     RES_REC.extend(REC_LIST)
396
397     RES_USERS.append(u)
398
399     self.save_pickle(RES_TRAIN, "datos/Extend/ARRAY_ARRAYS_TRAIN")
400     self.save_pickle(RES_REC, "datos/Extend/ARRAY_ARRAYS_REC")
401
402     return(RES_TRAIN, RES_REC, RES_USERS)
403
404 def getUserDataArrays(self, train=75 , days=30, minSongs=5):
405     '''
406     Genera los datos para el metodo D2V.
407
408     :param train: Porcentaje de datos para TRAIN
409     :param days: Numero de dias para el perfil reciente
410     :param minSongs: Numero minimo de canciones que ha de tener un usuario en su perfil
411     reciente para no ser eliminado
412     :return: Retorna datos de TRAIN COMPLETOS y RECIENTES y datos de DEV y TEST
413     '''
414     filename = "datos/CLEAN_DATA_[10,100].pkl"
415
416     DATA = pd.read_pickle(filename)
417
418     # Obtener los id de usuarios
419     USER_LIST = DATA.groupby("user")
420
421     RES_TRAIN_PRESENT=[]
422     RES_TRAIN_COMPLETE=[]
423     RES_REC=[]
424
425     RES_USERS=[]
426
427     #Para cada usuario
428     for u,r in USER_LIST:
429
430         # Obtener todas las canciones del usuario
431         USER_SONG_LIST = r
432
433         # Ordenar por fecha
434         USER_SONG_LIST.sort_values("timestamp", ascending=True, inplace=True)
435         USER_SONG_LIST = USER_SONG_LIST.reset_index()

```

```

436
437     # Lista con los id y las fechas
438     ID_LIST = USER_SONG_LIST.ID.values.tolist()
439     DATES_LIST = USER_SONG_LIST.timestamp.values.tolist()
440     DATES_LIST = np.array(list(map(lambda x: datetime.datetime.fromtimestamp(x /
441                                     1000000000), DATES_LIST)))
442
443     # Separar TRAIN DEV y TEST
444     TRAIN_INDX = int(len(ID_LIST)*(train/100))
445
446     TRAIN_LIST = ID_LIST[:TRAIN_INDX]
447     REC_LIST = ID_LIST[TRAIN_INDX:]
448     DATES_TRAIN = DATES_LIST[:TRAIN_INDX]
449
450     # Obtener X dias antes
451     PRESENT_START = DATES_TRAIN[-1] - timedelta(days=days)
452     PAST_LENGTH = sum((DATES_TRAIN >= PRESENT_START) == False)
453     PRESENT_ITEMS = TRAIN_LIST[PAST_LENGTH:]
454
455     #Si el usuario no posee suficientes canciones en el perfil presente, no se anade
456     if(len(PRESENT_ITEMS)<=minSongs):
457         print("Skipping "+str(u))
458         continue
459
460     #Pasas a string las listas
461     PRESENT_LIST = list(map(str,PRESENT_ITEMS))
462     TRAIN_LIST = list(map(str, TRAIN_LIST))
463     REC_LIST = list(map(str, REC_LIST))
464
465     RES_TRAIN_COMPLETE.append(TRAIN_LIST)
466     RES_TRAIN_PRESENT.append(PRESENT_LIST)
467     RES_REC.append(REC_LIST)
468
469     RES_USERS.append(u)
470
471     print(len(RES_TRAIN_COMPLETE),len(RES_USERS))
472
473     self.save_pickle(RES_TRAIN_COMPLETE,"datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE")
474     self.save_pickle(RES_TRAIN_PRESENT,"datos/Append/ARRAY_ARRAYS_TRAIN_PRESENT")
475     self.save_pickle(RES_REC,"datos/Append/ARRAY_ARRAYS_REC")
476
477     return(RES_TRAIN_COMPLETE,RES_TRAIN_PRESENT,RES_REC,RES_USERS)
478
479 def getSkipGramData(self, window=2):
480
481     filename = "datos/ARRAY_ARRAYS_EXTEND"
482
483     if os.path.isfile(filename+".pkl"):
484         DATA = self.load_pickle(filename)
485         DATA = self.save_pickle(filename)
486
487     else:
488         print("No existe el pickle, creando...")
489         DATA = self.getDataArrays();
490
491     DATA = list(map(int, DATA))

```

```

492     LENGTH = max(DATA)+1
493
494     RES_LEN = (len(DATA)-(window*2))*(window*2)
495
496     batch_inputs = np.ndarray(shape=(RES_LEN), dtype=np.int32)
497     batch_context = np.ndarray(shape=(RES_LEN, 1), dtype=np.int32)
498
499     j=0
500
501     for i in range(window,len(DATA)-window):
502         for w in range(1,window+1):
503
504             batch_inputs[j] = DATA[i]
505             batch_context[j,0] = DATA[i-w]
506             j+=1
507
508             batch_inputs[j] = DATA[i]
509             batch_context[j, 0] = DATA[i+w]
510             j+=1
511
512     return(batch_inputs,batch_context,LENGTH)
513
514     #####
515     # Estadísticas e info acerca de los datos
516     #####
517
518     def getPlaysByUser(self, times = list(range(10,110,10)) ):
519
520         if os.path.isfile("datos/ALL_SONGS_ID_COUNT.pkl"):
521             ALL_SONGS_ID_COUNT = pd.read_pickle("datos/ALL_SONGS_ID_COUNT.pkl")
522
523         else:
524             print("No existen el pickle:\n\t · ALL_SONGS_ID_COUNT.pkl [generado en
525                 createDict()]")
526             return()
527
528         #Obtener los id de usuarios
529         USER_LIST = ALL_SONGS_ID_COUNT.groupby("user")
530
531         print('\t'.join(str(x) for x in ["user",1]+times))
532
533         for u, s in USER_LIST:
534             COUNTS = [len(s)]
535             for i in times: COUNTS.append(len(s.loc[s.counts>=i]))
536
537             COUNTS = [u]+COUNTS
538             print('\t'.join(str(x) for x in COUNTS))
539
540     def getPlaysByYear(self):
541
542         if os.path.isfile("datos/ALL_SONGS_ID_COUNT.pkl"):
543             ALL_SONGS_ID_COUNT = pd.read_pickle("datos/ALL_SONGS_ID_COUNT.pkl")
544
545         else:
546             print("No existen el pickle:\n\t · ALL_SONGS_ID_COUNT.pkl [generado en
547                 createDict()]")

```



## 3. Métodos

En los siguientes ficheros se crea, entrena y evalúan los métodos de aprendizaje de embeddings empleados durante la experimentación (W2V, D2V-DM y D2V-DBOW).

### 3.1. Word2Vec

```

1  # -*- coding: utf-8 -*-
2  import os.path
3  import tensorflow as tf
4  import numpy as np
5  import math
6  import time
7  import pickle
8
9  #####
10 # eMtodos
11 #####
12
13 def getSkipGramData( window=2):
14
15     def load_pickle(name):
16         with open(name + '.pkl', 'rb') as f:
17             return pickle.load(f)
18
19     filename = "datos/Extend/ARRAY_ARRAYS_TRAIN"
20
21     if os.path.isfile(filename + ".pkl"):
22         DATA = load_pickle(filename)
23     else:
24         print("No existe el pickle...")
25         exit()
26
27     DATA = list(map(int, DATA))
28     LENGTH = max(DATA) + 1
29
30     RES_LEN = (len(DATA) - (window * 2)) * (window * 2)
31
32     batch_inputs = np.ndarray(shape=(RES_LEN), dtype=np.int32)
33     batch_context = np.ndarray(shape=(RES_LEN, 1), dtype=np.int32)
34
35     j = 0
36
37     for i in range(window, len(DATA) - window):
38         for w in range(1, window + 1):
39             batch_inputs[j] = DATA[i]
40             batch_context[j, 0] = DATA[i - w]
41             j += 1
42
43             batch_inputs[j] = DATA[i]
44             batch_context[j, 0] = DATA[i + w]
45             j += 1
46
47     return (batch_inputs, batch_context, LENGTH)
48
49 #####

```

```

50 # Llamadas
51 #####
52
53 batch_size = 512
54 embedding_size = 64 # tamaño del embedding (capa oculta)
55 window = 2 # tamaño de la ventana del skipgram
56 num_sampled = 64 # Ejemplos negativos (Para el NSE)
57 seed = 100
58 learning_rate = 0.1
59 epochs = 1000
60
61 # Carpeta donde se almacena el modelo de la red
62 #model_path = 'models/w2v'
63 model_path = 'models/w2v/TITANXP'
64 model_name = 'model'
65 complete_path = model_path+"/"+model_name
66 emb_path = 'embeddings/w2v'
67
68 # Obtener todos los datos
69 print("Obteniendo datos...")
70 ALL_X, ALL_Y, vocabulary_size = getSkipGramData(window=window)
71 print("Datos obtenidos!")
72
73 #Mezclar datos
74 TOTAL_LENGTH = len(ALL_X)
75 indx = np.arange(TOTAL_LENGTH)
76 np.random.seed(seed)
77 np.random.shuffle(indx)
78
79 ALL_X=ALL_X[indx]
80 ALL_Y=ALL_Y[indx]
81
82 #Separar datos en TRAIN, DEV y TEST
83 TRAIN_LENGTH = int(TOTAL_LENGTH*0.98)
84 DEV_LENGTH= int(TOTAL_LENGTH*0.01)
85 TEST_LENGTH = int(TOTAL_LENGTH*0.01)
86 #Train
87 X= ALL_X[:TRAIN_LENGTH]
88 Y= ALL_Y[:TRAIN_LENGTH]
89 #Dev
90 DEV_X = ALL_X[TRAIN_LENGTH:TRAIN_LENGTH+DEV_LENGTH]
91 DEV_Y = ALL_Y[TRAIN_LENGTH:TRAIN_LENGTH+DEV_LENGTH]
92 #Test
93 TEST_X = ALL_X[TRAIN_LENGTH+DEV_LENGTH:]
94 TEST_Y = ALL_Y[TRAIN_LENGTH+DEV_LENGTH:]
95
96 #Creación del grafo de TF.
97 graph = tf.Graph()
98
99 with graph.as_default():
100
101     #Número global de iteraciones
102     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
103
104     #Datos de entrada
105     #Array del tamaño del batch con las X

```



```

106 train_inputs = tf.placeholder(tf.int32, shape=[None], name="train_inputs")
107 #Array del tamaño del batch con las Y asociadas
108 train_labels = tf.placeholder(tf.int32, shape=[None, 1], name="train_labels")
109
110 #Conexin primera y segunda capa


---


111 # Matriz de pesos entre la capa de entrada y la de embeddings
112 embeddings = tf.Variable(tf.random_uniform([vocabulary_size, embedding_size], -1.0, 1.0),
113                          name='embeddings')
114 #embeddings = tf.get_variable("embeddings", shape=[vocabulary_size,
115                          embedding_size], initializer=tf.contrib.layers.xavier_initializer(seed=seed))
116
117 # Como la relacion entre ambas es lineal, no hay bias ni función de activación
118 # Por tanto se puede obtener el embedding directamente de la matriz anterior ( Embedding
119 # de la palabra 10 == fila 10 de la matriz de pesos)
120 embed = tf.nn.embedding_lookup(embeddings, train_inputs, name="embed")
121
122 #Conexin segunda y última capa


---


123 #Pesos y biases
124 nce_weights = tf.Variable(tf.truncated_normal([vocabulary_size,
125                          embedding_size], stddev=1.0 / math.sqrt(embedding_size)), name="weights")
126 nce_biases = tf.Variable(tf.zeros([vocabulary_size]), name="biases")
127
128 #Cálculo de LOSS y optimizador


---


129 nce_loss = tf.reduce_mean(
130     tf.nn.nce_loss(weights=nce_weights,
131                   biases=nce_biases,
132                   labels=train_labels,
133                   inputs=embed,
134                   num_sampled=num_sampled,
135                   num_classes=vocabulary_size), name="nce_loss")
136
137 #optimizer =
138     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
139 optimizer =
140     tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
141
142 init = tf.global_variables_initializer()
143
144 #Crear objeto encargado de almacenar la red
145 saver = tf.train.Saver(max_to_keep=1)
146
147 gpu_options = tf.GPUOptions(per_process_gpu_memory_fraction=0.333)
148
149 with tf.Session(graph=graph) as session:
150
151     # We must initialize all variables before we use them.
152     init.run()
153     print('Initialized')
154
155     #Obtener el checkpoint
156     ckpt = tf.train.get_checkpoint_state(model_path)

```

```

154 #Si existe el checkpoint restaurar
155 if ckpt and ckpt.model_checkpoint_path:
156     saver.restore(session, ckpt.model_checkpoint_path)
157     print('Modelo cargado...')
158 else:
159     print('No existe modelo...')
160
161 #####
162 #Obtener la matriz de embeddings
163 print("Obteniendo embeddings...")
164
165 E = session.run(embeddings)
166 E.dump(emb_path+"/EMB_MATRIX")
167
168 print("Embedings obtenidos!")
169 #####
170
171 exit()
172
173 #Obtener el unmero de iteraciones necesario y los datos restantes
174 ITRS= int(len(X)/batch_size)
175 RES_ITMS = len(X) - (ITRS*batch_size)
176
177 #Si no es divisin exacta, se aade una iteracin
178 if(RES_ITMS>0):ITRS+=1
179
180 print("Iteraciones necesarias:"+str(ITRS))
181
182 for e in range(epochs):
183     print("Epoch "+str(e))
184
185     tss = time.time()
186
187     for step in range(ITRS):
188
189         #Si se aadi una iteracin, introducir los datos restantes
190         if(step==ITRS-1 and RES_ITMS>0):
191             batch_inputs = X[(ITRS-1)*batch_size:]
192             batch_context = Y[(ITRS-1)*batch_size:]
193         #Si no, es un batch normal
194         else:
195             st = step*batch_size
196             nd = (step+1)*batch_size
197
198             batch_inputs = X[st:nd]
199             batch_context = Y[st:nd]
200
201         feed_dict = {train_inputs: batch_inputs, train_labels: batch_context}
202
203         # We perform one update step by evaluating the optimizer op (including it
204         # in the list of returned values for session.run()
205         -, loss_val, gs = session.run([optimizer, nce_loss, global_step],
206                                     feed_dict=feed_dict)
207
208         if (gs % 100 == 0 and gs>0):
209             # The average loss is an estimate of the loss over the last 100 batches.

```

```

210     print('\tLast batch loss at global-step ', gs, ': ', loss_val)
211
212     # Now, save the graph
213     #saver.save(session, complete_path, global_step=global_step)
214
215     print(time.time()-tss)
216     saver.save(session, complete_path, global_step=global_step)
217
218     # Probar con DEV
219     feed_dict = {train_inputs: DEV_X, train_labels: DEV_Y}
220     loss_val = session.run([nce_loss], feed_dict=feed_dict)
221     print('DEV batch loss at global-step ' + str(gs) + ': ' + str(loss_val))
222     with open('train_results.txt', 'a') as the_file: the_file.write('DEV\t' + str(gs) +
223         '\t' + str(loss_val) + '\t' + str(time.time() - tss)+'\n')
224
225     # Probar con TEST
226     feed_dict = {train_inputs: TEST_X, train_labels: TEST_Y}
227     loss_val = session.run([nce_loss], feed_dict=feed_dict)
228     print('TEST batch loss at global-step ' + str(gs) + ': ' + str(loss_val))
229     with open('train_results.txt', 'a') as the_file: the_file.write('TEST\t' + str(gs) +
230         '\t' + str(loss_val) + '\t' + str(time.time() - tss)+'\n')
231
232     #-----

```

## 3.2. Doc2Vec-DM

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import tensorflow as tf
5  from tensorflow.python.client import device_lib
6  import numpy as np
7  import math
8  import time
9  import pickle
10
11 #####
12 # eMtodos
13 #####
14
15 def getData(window=2, seed = 100, shuffle=True, file=None):
16
17     #Cargar datos del tipo [[canciones user1],[canciones user2],...]
18     if(file==None):DATA = np.load("datos/Append/ARRAY_ARRAYS_TRAIN.pkl")
19     else:DATA = np.load(file)
20     #print("ELIMINAR FILTRO——>");DATA = DATA[:5]
21
22     #Obtener el unmero total de usuarios
23     NUM_USERS = len(DATA)
24
25     #Matriz resultado
26     RES_COLS = window + 1 + 1 # USER ID, [C1,C2], C3
27     RES = np.empty([0, RES_COLS], dtype=np.int32)
28
29     #Para cada usuario, crear ejemplos
30     for u in range(NUM_USERS):
31         tss = time.time()
32
33         USR_SONGS = DATA[u]
34         USR_SONGS_SIZE = len(USR_SONGS)
35
36         if(USR_SONGS_SIZE<window+1):print(str(u)+" : "+str(USR_SONGS_SIZE));continue
37
38         USR_RES_ROWS = USR_SONGS_SIZE - window
39
40         #Crear matriz de retorno
41         USR_RES = np.empty([USR_RES_ROWS, RES_COLS], dtype=np.int32)
42
43         #Cambiar id de usuario
44         USR_RES[:,0] = u
45
46         #Obtener los datos (uno ams para la clase (ocancin siguiente))
47         for i in range(window+1):
48             T0 = USR_RES_ROWS+i
49             USR_RES[:,1+i]= USR_SONGS[i:T0]
50
51         #Finalmente naadir estos datos al conjunto de todos los usuarios
52         RES = np.concatenate([RES,USR_RES])
53
54     #Mezclar
55     if(shuffle):

```

```

56     np.random.seed(seed)
57     np.random.shuffle(RES)
58
59     #Dividir en TRAIN, DEV y TEST
60     #98/1/1 TODO: VALE MAS ???????????
61
62     TOTAL_LENGTH = len(RES)
63     DEV_LENGTH = int(TOTAL_LENGTH * 0.01)
64     TEST_LENGTH = int(TOTAL_LENGTH * 0.01)
65     TRAIN_LENGTH = TOTAL_LENGTH - DEV_LENGTH - TEST_LENGTH
66
67     # Train
68     TRAIN = RES[:TRAIN_LENGTH,:]
69     X = TRAIN[:, :window + 1]
70     Y = np.matrix(TRAIN[:, -1]).T # Hay que pasarla a matriz y transponer
71     # Dev
72     DEV = RES[TRAIN_LENGTH:TRAIN_LENGTH+DEV_LENGTH,:]
73     DEV_X = DEV[:, :window + 1]
74     DEV_Y = np.matrix(DEV[:, -1]).T # Hay que pasarla a matriz y transponer
75     # Test
76     TEST = RES[TRAIN_LENGTH+DEV_LENGTH:,:]
77     TEST_X = TEST[:, :window + 1]
78     TEST_Y = np.matrix(TEST[:, -1]).T # Hay que pasarla a matriz y transponer
79
80     return(X,Y,DEV_X,DEV_Y,TEST_X,TEST_Y, NUM_USERS)
81
82     #####
83     # Llamadas
84     #####
85
86     #Cargar matriz de pesos W2V
87     W2V_emb_path = 'embeddings/W2V/EMB_MATRIX'
88     W2V_EMB = np.load(W2V_emb_path)
89
90     batch_size = 1024
91     vocabulary_size = W2V_EMB.shape[0]
92     embedding_size_w = W2V_EMB.shape[1] # tamaño del embedding (capa oculta)
93     embedding_size_d = embedding_size_w
94     window=2
95     seed = 100
96     num_sampled = 64 # Ejemplos negativos (Para el NSE)
97     learning_rate = 0.1
98     epochs = 1500
99
100    # Obtener todos lo datos
101    print("Obteniendo datos...")
102
103    model_name = 'complete'
104    model_path = 'models/d2v_dm/complete'
105    emb_path = 'embeddings/d2v_dm/complete'
106    X,Y,DEV_X,DEV_Y,TEST_X,TEST_Y, user_size = getData(window=window,seed=seed,
107        file="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl")
108
109    #model_name = 'present'
110    #model_path = 'models/d2v_dm/present'
111    #emb_path = 'embeddings/d2v_dm/present'
112    #X,Y,DEV_X,DEV_Y,TEST_X,TEST_Y, user_size = getData(window=window,seed=seed,

```

```

file="datos/Append/ARRAY_ARRAYS_TRAIN_PRESENT.pkl")
112
113 print("Datos obtenidos!")
114
115 # Carpeta donde se almacena el modelo de la red
116 complete_path = model_path+"/"+model_name
117
118 # Creación del grafo de TF.
119 graph = tf.Graph()
120
121 with graph.as_default():
122
123     # Número global de iteraciones
124     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
125
126     # Datos de entrada
127
128     # Array del tamaño del batch con las X
129     train_dataset = tf.placeholder(tf.int32, shape=[None, window + 1], name="train_inputs")
130
131     # Array del tamaño del batch con las Y asociadas
132     train_labels = tf.placeholder(tf.int32, shape=[None, 1], name="train_labels")
133
134     # Conexión primera y segunda capa
135
136     # Matriz W de embeddings de las canciones obtenida en el W2V previo
137     word_embeddings = tf.Variable(W2V_EMB, trainable=False, name="word_embeddings") #
138     FUNCIONA???
139     https://stackoverflow.com/questions/37326002/is-it-possible-to-make-a-trainable-variable-not-trainable
140
141     # Matriz D de embeddings para los usuarios (lo que se aprende ahora)
142     doc_embeddings = tf.Variable(tf.random_uniform([user_size, embedding_size_d], -1.0, 1.0),
143     name="doc_embeddings")
144
145     # Conexión segunda y última capa
146
147
148     combined_embed_vector_length = embedding_size_w + embedding_size_d
149     # softmax weights, W and D vectors should be concatenated before applying softmax
150     weights = tf.Variable(tf.truncated_normal([vocabulary_size,
151     combined_embed_vector_length], stddev=1.0 / math.sqrt(combined_embed_vector_length)),
152     name="weights")
153     # softmax biases
154     biases = tf.Variable(tf.zeros([vocabulary_size]), name="biases")
155
156     # Cálculo de LOSS y optimizador
157
158     embed = []
159     # Obtener el embedding del documento (Columna 0 del batch)
160     embed_d = tf.nn.embedding_lookup(doc_embeddings, train_dataset[:, 0])
161     embed.append(embed_d)
162
163     # Obtener la SUMA de los embeddings de las canciones (Columnas restantes del batch)
164     # El primero se hace fuera. (No se puede crear un vector de 0 inicial e ir añadiendo dado
165     que no se conoce el tamaño del batch)
166     embed_w = tf.nn.embedding_lookup(word_embeddings, train_dataset[:, 1])

```

```

158
159 for j in range(2,window+1):
160     embed_w += tf.nn.embedding_lookup(word_embeddings, train_dataset[:, j])
161 embed.append(embed_w)
162
163 #Se concatena el embedding del usuario con la SUMA de los de las canciones (No se suma el
164     embedding del doc) (hasta ahora era un array de 2 elementos)
165 embed = tf.concat(embed, 1)
166
167 nce_loss = tf.nn.nce_loss(weights=weights,
168                           biases=biases,
169                           labels=train_labels,
170                           inputs=embed,
171                           num_sampled=num_sampled,
172                           num_classes=vocabulary_size,
173                           name="nce_loss")
174
175 nce_loss = tf.reduce_mean(nce_loss)
176
177 optimizer =
178     tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
179
180 init = tf.global_variables_initializer()
181
182 #Crear objeto encargado de almacenar la red
183 saver = tf.train.Saver(max_to_keep=1)
184
185 gpu_options = tf.GPUOptions(per_process_gpu_memory_fraction=0.333)
186
187 with tf.Session(graph=graph,config=tf.ConfigProto(gpu_options=gpu_options)) as session:
188
189     # We must initialize all variables before we use them.
190     init.run()
191     print('Initialized')
192
193 #Obtener el checkpoint
194 ckpt = tf.train.get_checkpoint_state(model_path)
195
196 #Si existe el checkpoint restaurar
197 if ckpt and ckpt.model_checkpoint_path:
198     saver.restore(session, ckpt.model_checkpoint_path)
199     print('Modelo cargado...')
200 else:
201     print('No existe modelo...')
202
203 #####
204 #Obtener la matriz de embeddings
205 print("Obteniendo embeddings...")
206
207 E = session.run(doc_embeddings)
208 E.dump(emb_path+"/EMB_MATRIX")
209
210 print("Embedings obtenidos!")
211 #####
212

```

```

213     exit()
214
215
216
217     ITRS= int(len(X)/batch_size)
218     RES_ITMS = len(X) - (ITRS*batch_size)
219
220     #Si no es divisible exacta, se añade una iteración
221     if(RES_ITMS>0):ITRS+=1
222
223     print("Iteraciones necesarias:"+str(ITRS))
224
225     for e in range(epochs):
226         print("Epoch "+str(e))
227
228         tss = time.time()
229
230         for step in range(ITRS):
231
232             #Si se añade una iteración, introducir los datos restantes
233             if(step==ITRS-1 and RES_ITMS>0):
234                 batch_inputs = X[(ITRS-1)*batch_size:]
235                 batch_context = Y[(ITRS-1)*batch_size:]
236             #Si no, es un batch normal
237             else:
238                 st = step*batch_size
239                 nd = (step+1)*batch_size
240
241                 batch_inputs = X[st:nd]
242                 batch_context = Y[st:nd]
243
244             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
245             _, loss_val, gs = session.run([optimizer, nce_loss, global_step],
246                 feed_dict=feed_dict)
247
248             if (gs % 1000 == 0 and gs>0):
249
250                 # The average loss is an estimate of the loss over the last 100 batches.
251                 print('\tLast batch loss at global-step ', gs, ': ', loss_val)
252
253                 # Now, save the graph
254                 #saver.save(session, complete_path, global_step=global_step)
255
256             print(time.time()-tss)
257             saver.save(session, complete_path, global_step=global_step)
258
259             # Probar con DEV
260             feed_dict = {train_dataset: DEV_X, train_labels: DEV_Y}
261             loss_val = session.run([nce_loss], feed_dict=feed_dict)
262             print('DEV batch loss at global-step ' + str(gs) + ': ' + str(loss_val))
263             with open('train_results_D2V_DM_'+model_name+'.txt', 'a') as the_file:
264                 the_file.write('DEV\t' + str(gs) + '\t' + str(loss_val) + '\t' + str(time.time() -
265                     tss)+'\n')
266
267             # Probar con TEST
268             feed_dict = {train_dataset: TEST_X, train_labels: TEST_Y}

```



```
267 loss_val = session.run([nce_loss], feed_dict=feed_dict)
268 print('TEST batch loss at global-step ' + str(gs) + ': ' + str(loss_val))
269 with open('train_results_D2V_DM_'+model_name+'.txt', 'a') as the_file:
    the_file.write('TEST\t' + str(gs) + '\t' + str(loss_val) + '\t' + str(time.time()
    - tss)+'\n')
270
271
272 #-----
```

### 3.3. Doc2Vec-DBOW

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import tensorflow as tf
5  from tensorflow.python.client import device_lib
6  import numpy as np
7  import math
8  import time
9  import pickle
10
11  #####
12  # eMtodos
13  #####
14
15  def getData(seed = 100, shuffle=True, file=None):
16
17      #Cargar datos del tipo [[canciones user1],[canciones user2],...]
18      if(file==None):DATA = np.load("datos/Append/ARRAY_ARRAYS_TRAIN.pkl")
19      else:DATA = np.load(file)
20
21      #print("ELIMINAR FILTRO——>");DATA = DATA[:5]
22
23      #Obtener el unmero total de usuarios
24      NUM_USERS = len(DATA)
25
26      #Matriz resultado
27      RES_COLS = 1 + 1 # USER ID => C1
28      RES = np.empty([0, RES_COLS], dtype=np.int32)
29
30      #Para cada usuario, crear ejemplos
31      for u in range(NUM_USERS):
32
33          USR_SONGS = DATA[u]
34          USR_RES_ROWS = len(USR_SONGS)
35
36          #Crear una matriz "vaca de (0n de canciones) filas y 2 columnas (USER -> CANCION)
37          USR_RES = np.empty([USR_RES_ROWS, RES_COLS], dtype=np.int32)
38
39          #Rellenar
40          USR_RES[:,0] = u
41          USR_RES[:,1] = USR_SONGS
42
43          #Finalmente añadir estos datos al conjunto de todos los usuarios
44          RES = np.concatenate([RES,USR_RES])
45
46      #Mezclar
47      if(shuffle):
48          np.random.seed(seed)
49          np.random.shuffle(RES)
50
51      #Dividir en TRAIN, DEV y TEST
52      #98/1/1 TODO: VALE ¿AS ???????????
53
54      TOTAL_LENGTH = len(RES)
55      DEV_LENGTH = int(TOTAL_LENGTH * 0.01)

```

```

56     TEST_LENGTH = int(TOTAL_LENGTH * 0.01)
57     TRAIN_LENGTH = TOTAL_LENGTH - DEV_LENGTH - TEST_LENGTH
58
59     # Train
60     TRAIN = RES[:TRAIN_LENGTH,:]
61     X = TRAIN[:, 0]
62     Y = np.matrix(TRAIN[:, 1]).T # Hay que pasarla a matriz y transponer
63     # Dev
64     DEV = RES[TRAIN_LENGTH:TRAIN_LENGTH+DEV_LENGTH,:]
65     DEV_X = DEV[:, 0]
66     DEV_Y = np.matrix(DEV[:, 1]).T # Hay que pasarla a matriz y transponer
67     # Test
68     TEST = RES[TRAIN_LENGTH+DEV_LENGTH:,:]
69     TEST_X = TEST[:, 0]
70     TEST_Y = np.matrix(TEST[:,1]).T # Hay que pasarla a matriz y transponer
71
72     return(X,Y,DEV_X,DEV_Y,TEST_X,TEST_Y, NUM_USERS)
73
74     #####
75     # Llamadas
76     #####
77
78     #Cargar matriz de pesos W2V (NO SE NECESITA, SIMPLEMENTE ES PARA CONOCER EL NTAMAÑO DEL
79     VOCABULARIO)
80     W2V_emb_path = 'embeddings/W2V/EMB_MATRIX'
81     W2V_EMB = np.load(W2V_emb_path)
82
83     batch_size = 1024
84     vocabulary_size = W2V_EMB.shape[0]
85     embedding_size_d = 64
86     seed = 100
87     num_sampled = 64 # Ejemplos negativos (Para el NSE)
88     learning_rate = 0.1
89     epochs = 1500
90
91     # Obtener todos lo datos
92     print("Obteniendo datos...")
93
94     #model_name = 'complete'
95     #model_path = 'models/d2v.dbow/complete'
96     #emb_path = 'embeddings/d2v.dbow/complete'
97     #X,Y,DEV_X,DEV_Y,TEST_X,TEST_Y, user_size = getData(seed=seed,
98     file="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl")
99
100     model_name = 'present'
101     model_path = 'models/d2v.dbow/present'
102     emb_path = 'embeddings/d2v.dbow/present'
103     X,Y,DEV_X,DEV_Y,TEST_X,TEST_Y, user_size = getData(seed=seed,
104     file="datos/Append/ARRAY_ARRAYS_TRAIN_PRESENT.pkl")
105
106     print(user_size)
107
108     print("Datos obtenidos!")
109
110     # Carpeta donde se almacena el modelo de la red
111     complete_path = model_path+"/"+model_name

```

```

110
111 #Creacin del grafo de TF.
112 graph = tf.Graph()
113
114 with graph.as_default():
115
116     #Numero global de iteraciones
117     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
118
119     #Datos de entrada
120
121     #Array del tamaño del batch con las X
122     train_dataset = tf.placeholder(tf.int32, shape=[None], name="train_inputs")
123
124     #Array del tamaño del batch con las Y asociadas
125     train_labels = tf.placeholder(tf.int32, shape=[None, 1], name="train_labels")
126
127     #Conexin primera y segunda capa
128
129     # Matriz D de embeddigs para los usuarios (lo que se aprende ahora)
130     doc_embeddings = tf.Variable(tf.random_uniform([user_size, embedding_size_d], -1.0, 1.0),
131                                 name="doc_embeddings")
132
133     embed_d = tf.nn.embedding_lookup(doc_embeddings, train_dataset, name="embed")
134
135     #Conexin segunda y ultima capa
136
137     weights = tf.Variable(tf.truncated_normal([vocabulary_size, embedding_size_d], stddev=1.0
138                                             / math.sqrt(embedding_size_d)), name="weights")
139     # softmax biases
140     biases = tf.Variable(tf.zeros([vocabulary_size]), name="biases")
141
142     #Cálculo de LOSS y optimizador
143
144     nce_loss = tf.nn.nce_loss(weights=weights,
145                               biases=biases,
146                               labels=train_labels,
147                               inputs=embed_d,
148                               num_sampled=num_sampled,
149                               num_classes=vocabulary_size,
150                               name="nce_loss")
151
152     nce_loss = tf.reduce_mean(nce_loss)
153
154     optimizer =
155         tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
156
157     init = tf.global_variables_initializer()
158
159     #Crear objeto encargado de almacenar la red
160     saver = tf.train.Saver(max_to_keep=1)
161
162     gpu_options = tf.GPUOptions(per_process_gpu_memory_fraction=0.333)

```

```

160
161 with tf.Session(graph=graph,config=tf.ConfigProto(gpu_options=gpu_options)) as session:
162
163     # We must initialize all variables before we use them.
164     init.run()
165     print('Initialized')
166
167     #Obtener el checkpoint
168     ckpt = tf.train.get_checkpoint_state(model_path)
169
170     #Si existe el checkpoint restaurar
171     if ckpt and ckpt.model_checkpoint_path:
172         saver.restore(session, ckpt.model_checkpoint_path)
173         print('Modelo cargado...')
174     else:
175         print('No existe modelo...')
176
177     #####
178     #Obtener la matriz de embeddings
179     print("Obteniendo embeddings...")
180
181     E = session.run(doc_embeddings)
182     print(E.shape)
183     E.dump(emb_path+"/EMB_MATRIX")
184
185     print("Embedings obtenidos!")
186     #####
187
188     exit()
189
190
191
192     ITRS= int(len(X)/batch_size)
193     RES_ITMS = len(X) - (ITRS*batch_size)
194
195     #Si no es divisin exacta, se aade una iteracin
196     if(RES_ITMS>0):ITRS+=1
197
198     print("Iteraciones necesarias:"+str(ITRS))
199
200     for e in range(epochs):
201         print("Epoch "+str(e))
202
203         tss = time.time()
204
205         for step in range(ITRS):
206
207             #Si se aadi una iteracin, introducir los datos restantes
208             if(step==ITRS-1 and RES_ITMS>0):
209                 batch_inputs = X[(ITRS-1)*batch_size:]
210                 batch_context = Y[(ITRS-1)*batch_size:]
211             #Si no, es un batch normal
212             else:
213                 st = step*batch_size
214                 nd = (step+1)*batch_size
215
216                 batch_inputs = X[st:nd]

```

```

217         batch_context = Y[st:nd]
218
219     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
220     _, loss_val, gs = session.run([optimizer, nce_loss, global_step],
221                                   feed_dict=feed_dict)
222
223     if (gs % 100 == 0 and gs > 0):
224
225         # The average loss is an estimate of the loss over the last 100 batches.
226         print('\tLast batch loss at global-step ', gs, ': ', loss_val)
227
228         # Now, save the graph
229         #saver.save(session, complete_path, global_step=global_step)
230
231     saver.save(session, complete_path, global_step=global_step)
232
233
234     # Probar con DEV
235     feed_dict = {train_dataset: DEV_X, train_labels: DEV_Y}
236     loss_val = session.run([nce_loss], feed_dict=feed_dict)
237     print('DEV batch loss at global-step ' + str(gs) + ': ' + str(loss_val))
238     with open('train_results_D2V_DBOW_'+model_name+'.txt', 'a') as the_file:
239         the_file.write('DEV\t' + str(gs) + '\t' + str(loss_val) + '\t' + str(time.time() -
240                                     tss)+'\n')
241
242     # Probar con TEST
243     feed_dict = {train_dataset: TEST_X, train_labels: TEST_Y}
244     loss_val = session.run([nce_loss], feed_dict=feed_dict)
245     print('TEST batch loss at global-step ' + str(gs) + ': ' + str(loss_val))
246     with open('train_results_D2V_DBOW_'+model_name+'.txt', 'a') as the_file:
247         the_file.write('TEST\t' + str(gs) + '\t' + str(loss_val) + '\t' + str(time.time()
248                                     - tss)+'\n')

```

## 4. Tests

En esta sección final, se incluirán todos los ficheros referentes a los test llevados a cabo durante las tareas de recomendación.

### 4.1. TEST 01

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[], data
20            ="datos/Append/ARRAY_ARRAYS_REC.pkl", pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl", shuffl
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]

```

```

48
49 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
      ahora
50 USED_IDS = set(list(map(int, PAST_DATA[u])))
51 USED_IDS.update(TRAIN_DATA)
52 TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))
53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
      DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
      DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101

```



```

102
103     return
104         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106
107     # Si se noaadi una oiteracin, introducir los datos restantes
108     if (step == ITRS - 1 and RES_ITMS > 0):
109         batch_inputs = X[(ITRS - 1) * batch_size:]
110         batch_context = Y[(ITRS - 1) * batch_size:]
111     # Si no, es un batch normal
112     else:
113         st = step * batch_size
114         nd = (step + 1) * batch_size
115
116         batch_inputs = X[st:nd]
117         batch_context = Y[st:nd]
118     return(batch_inputs, batch_context)
119
120 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     song_emb_size = 64
135     hidden_size = 32
136
137     num_sampled = 64 # Ejemplos negativos (Para el NSE)
138
139     model_name = 'model'
140     model_path = 'models/' + TEST_NAME + '/model'
141     complete_path = model_path + "/" + model_name
142
143     # Se almacena la loss de dev de las 5 ultimas epochs
144     slope_size = 100
145     train_hist=[]
146     dev_hist=[]
147
148     # oCreacin del grafo de TF.
149     graph = tf.Graph()
150
151     with graph.as_default():
152
153         if(config["seed"]!=1):
154             tf.set_random_seed(config["seed"])
155
156         # uNumero global de iteraciones
157         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')

```

```

158
159 # Datos de entrada
160
161 # Array del tamaño del batch con las X
162 train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
163
164 # Array del tamaño del batch con las Y asociadas
165 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
166
167 # Embeddings
168
169 # Matriz W de embeddings de las canciones obtenida en el W2V previo
170
171 O1 = tf.Variable(tf.truncated_normal([user_size, user_emb_size], mean=0.0, stddev=1.0
172 / math.sqrt(user_emb_size)), name="O1")
173 B01 = tf.Variable(tf.zeros([user_emb_size]), name="B01")
174
175 W2V_EMB = tf.Variable(W2V, trainable=False, name="word_embeddings") # FUNCIONA???
176 https://stackoverflow.com/questions/37326002/is-it-possible-to-make-a-trainable-variable-no
177
178 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
179 stddev=1.0 / math.sqrt(hidden_size)), name="T1")
180 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
181
182 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
183 stddev=1.0 / math.sqrt(hidden_size)), name="T2")
184 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
185
186 # Operaciones
187
188 # Embedding de documento
189 ed = tf.nn.embedding_lookup(O1, train_dataset[:, 0])
190 ed = ed + B01
191
192 # Embedding de canción
193 ec = tf.nn.embedding_lookup(W2V_EMB, train_dataset[:, 1])
194
195 # Transformar a 32 documento
196 hd = tf.matmul(ed, T1) + B1
197
198 # Transformar a 32 canción
199 hw = tf.matmul(ec, T2) + B2
200
201 # Obtener el producto escalar
202 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
203
204 # Cálculo de LOSS y optimizador
205
206 # Obtener la loss
207 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
208 loss_softplus = tf.reduce_mean(softplus_batch)
209
210 # Minimizar la loss

```

```

206     train_step =
           tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus ,global_step
207
208     # Obtener la probabilidad
209     prob = tf.sigmoid(dot_prod)
210
211     # Inicializar variables
212     init = tf.global_variables_initializer()
213
214     # Crear objeto encargado de almacenar la red
215     saver = tf.train.Saver(max_to_keep=1)
216
217     config = tf.ConfigProto()
218     config.gpu_options.allow_growth = True
219     with tf.Session(graph=graph, config=config) as session:
220
221         # We must initialize all variables before we use them.
222         init.run()
223
224         if(save):
225             # Obtener el checkpoint
226             ckpt = tf.train.get_checkpoint_state(model_path)
227
228             # Si existe el checkpoint restaurar
229             if ckpt and ckpt.model_checkpoint_path:
230                 saver.restore(session, ckpt.model_checkpoint_path)
231
232         #Si se desea entrenar
233         if(train):
234             #-----
235             # TRAIN
236             #-----
237
238             ITRS= int(len(X)/batch_size)
239             RES_ITMS = len(X) - (ITRS*batch_size)
240             if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
241
242             for e in range(epochs):
243                 TRAIN_LOSS = 0.0
244
245                 print("Epoch "+str(e))
246
247                 for step in range(ITRS):
248
249                     batch_inputs, batch_context =
250                         getBatchData(X,Y,batch_size ,step,ITRS,RES_ITMS)
251
252                     feed_dict = {train.dataset: batch_inputs, train.labels: batch_context}
253
254                     ., loss_val,softplus, gs = session.run([train_step,
255                         loss_softplus,softplus_batch,global_step], feed_dict=feed_dict)
256
257                     TRAIN_LOSS += np.sum(softplus[:, 0])
258
259                 if (save): saver.save(session, complete_path, global_step=global_step)
260
261             #-----

```

```

260 # DEV
261 #-----
262
263 if(len(DX)>0):
264     ITRS_TEST = int(len(DX) / batch_size)
265     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
266     TEST_LOSS = 0.0
267     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisible exacta, se
        añade una iteración
268
269     for step in range(ITRS_TEST):
270
271         batch_inputs, batch_context =
            getBatchData(DX,DY, batch_size, step, ITRS_TEST, RES_ITMS_TEST)
272
273         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
274         loss_val, softplus = session.run([loss_softplus, softplus_batch],
            feed_dict=feed_dict)
275
276         TEST_LOSS+=np.sum(softplus[:,0])
277
278         dev_hist.append(TEST_LOSS/(len(DX)*1.0))
279         train_hist.append(TRAIN_LOSS/(len(X)*1.0))
280
281     if (len(dev_hist) >= slope_size):
282         # Calcular la pendiente
283         print(getSlope(dev_hist[-slope_size:]))
284         print(dev_hist[-slope_size:])
285
286         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
287
288             def printArray(data, data2):
289                 for i in range(len(data)):
290                     a=data[i]
291                     b=data2[i]
292                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
293
294             print("-" * 50)
295             print(e)
296             printArray(dev_hist, train_hist)
297             return
298
299 else:
300     #
301
302 #-----
303 # TEST FINAL
304 #
305
306 TRUE_POSITIVE = 0
307 FALSE_NEGATIVE = 0
308 FALSE_POSITIVE = 0
309 TRUE_NEGATIVE = 0
310
311 ITRS_TEST = int(len(TX) / batch_size)
312 RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
313 TEST_LOSS = 0.0

```

```

312 # Si no es divisin exacta, se aade una iteracin
313 if (RES_ITMS_TEST > 0): ITRS_TEST += 1
314
315 for step in range(ITRS_TEST):
316     # Si se aadi una iteracin, introducir los datos restantes
317     if (step == ITRS_TEST - 1 and RES_ITMS_TEST > 0):
318         batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
319         batch_context = TY[(ITRS_TEST - 1) * batch_size:]
320     # Si no, es un batch normal
321     else:
322         st = step * batch_size
323         nd = (step + 1) * batch_size
324
325         batch_inputs = TX[st:nd]
326         batch_context = TY[st:nd]
327
328         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
329         loss_val, loss_batch, p = session.run([loss_softplus, softplus_batch, prob],
330                                             feed_dict=feed_dict)
331
332         TEST_LOSS += np.sum(loss_batch)
333
334     # aClculo de TP/TN/FN/FP para cada batch
335
336     for i in range(len(p)):
337         real = batch_context[i, 0]
338         pred = p[i, 0]
339
340         if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
341         if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
342         if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
343         if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
344
345     PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
346     RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
347
348     F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
349
350     print("-" * 50)
351     print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
352     print("PRECISION:\t" + str(PRECISION))
353     print("RECALL:\t" + str(RECALL))
354     print("F1:\t" + str(F1))
355     print("-" * 50)
356     print("")
357
358 #####
359 # Llamadas
360 #####
361
362 TEST_NAME="test_021"
363
364 os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
365
366 #Para conocer los embeddings de las canciones
367 W2V = np.load('embeddings/w2v/EMB_MATRIX')

```

```

368
369 #Se cargan los datos de train de las canciones para conocer los IDs reales
370 DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
371
372 #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
373 ALL = set(list(range(W2V.shape[0])))
374 REAL = set(list(set(list(map(int, DATA))))))
375
376 #####
377 # TENSORFLOW
378 #####
379
380 #
381 # Grid Search + Early Stopping
382 #
383 '''
384 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
385 rates = [0.000001, 0.01, 0.0001]
386
387 for r in rates:
388     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
389
390     print("Learning Rate: "+str(r))
391     print("-" * 50)
392     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
393     print("-"*50)
394
395 exit()
396 '''
397 #
398 # Train
399 #
400 '''
401 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
402
403 #Para esta fase, se une TRAIN y DEV
404 X = np.concatenate((X,DX),axis=0)
405 Y = np.concatenate((Y,DY),axis=0)
406
407 DX=[]
408 DY=[]
409
410 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 102, "seed": 100}
411 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
412
413 exit()
414 '''
415 #
416 # Test
417 #
418
419
420 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
421
422 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 102, "seed": 100}
423 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```

## 4.2. TEST 02

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data
20             ="datos/Append/ARRAY_ARRAYS_REC.pkl",pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl",shuffl
21             seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
50
51         # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
52         ahora
53         USED_IDS = set(list(map(int, PAST_DATA[u])))
54         USED_IDS.update(TRAIN_DATA)
55         TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))

```

```

53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101
102
103    return
104    TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105    def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106

```



```

107 # Si se ha aadi una iteracin, introducir los datos restantes
108 if (step == ITRS - 1 and RES_ITMS > 0):
109     batch_inputs = X[(ITRS - 1) * batch_size:]
110     batch_context = Y[(ITRS - 1) * batch_size:]
111 # Si no, es un batch normal
112 else:
113     st = step * batch_size
114     nd = (step + 1) * batch_size
115
116     batch_inputs = X[st:nd]
117     batch_context = Y[st:nd]
118 return(batch_inputs, batch_context)
119
120 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     song_emb_size = 64
135     hidden_size = 32
136
137     num_sampled = 64 # Ejemplos negativos (Para el NSE)
138
139     model_name = 'model'
140     model_path = 'models/' + TEST_NAME + '/model'
141     complete_path = model_path + "/" + model_name
142
143     # Se almacena la loss de dev de las 5 ultimas epochs
144     slope_size = 100
145     train_hist = []
146     dev_hist = []
147
148     # Creacin del grafo de TF.
149     graph = tf.Graph()
150
151     with graph.as_default():
152
153         if(config["seed"]!=1):
154             tf.set_random_seed(config["seed"])
155
156         # nmero global de iteraciones
157         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
158
159         # Datos de entrada
160
161         # Array del tamao del batch con las X
162         train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")

```

```

163 # Array del tamaño del batch con las Y asociadas
164 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
165
166 # Embeddings


---


167 D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
168
169 W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
170
171 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
172     stddev=1.0 / math.sqrt(hidden_size)), name="T1")
173 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
174
175 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
176     stddev=1.0 / math.sqrt(hidden_size)), name="T2")
177 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
178
179 # Operaciones


---


180 # Embedding de documento
181 ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
182
183 # Embedding de canción
184 ec = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
185
186 # Transformar a 32 documento
187 hd = tf.matmul(ed, T1) + B1
188
189 # Transformar a 32 canción
190 hw = tf.matmul(ec, T2) + B2
191
192 # Obtener el producto escalar
193 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
194
195 # Cálculo de LOSS y optimizador


---


196
197 # Obtener la loss
198 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
199 loss_softplus = tf.reduce_mean(softplus_batch)
200
201 # Minimizar la loss
202 train_step =
203     tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus, global_step=
204
205 # Obtener la probabilidad
206 prob = tf.sigmoid(dot_prod)
207
208 # Inicializar variables
209 init = tf.global_variables_initializer()
210
211 # Crear objeto encargado de almacenar la red
212 saver = tf.train.Saver(max_to_keep=1)
213
214 config = tf.ConfigProto()

```

```

214 config.gpu_options.allow_growth = True
215 with tf.Session(graph=graph, config=config) as session:
216
217     # We must initialize all variables before we use them.
218     init.run()
219
220     if(save):
221         # Obtener el checkpoint
222         ckpt = tf.train.get_checkpoint_state(model_path)
223
224         # Si existe el checkpoint restaurar
225         if ckpt and ckpt.model_checkpoint_path:
226             saver.restore(session, ckpt.model_checkpoint_path)
227
228     #Si se desea entrenar
229     if(train):
230         #-----
231         # TRAIN
232         #-----
233
234         ITRS= int(len(X)/batch_size)
235         RES_ITMS = len(X) - (ITRS*batch_size)
236         if(RES_ITMS>0):ITRS+=1 #Si no es @divisin exacta, se naade una @iteracin
237
238         for e in range(epochs):
239             TRAIN_LOSS = 0.0
240
241             print("Epoch "+str(e))
242
243             for step in range(ITRS):
244
245                 batch_inputs, batch_context =
246                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
247
248                 feed_dict = {train.dataset: batch_inputs, train.labels: batch_context}
249
250                 ., loss_val,softplus, gs = session.run([train_step,
251                     loss_softplus,softplus_batch,global_step], feed_dict=feed_dict)
252
253                 TRAIN_LOSS += np.sum(softplus[:, 0])
254
255             if (save): saver.save(session, complete_path, global_step=global_step)
256
257             #-----
258             # DEV
259             #-----
260
261             if(len(DX)>0):
262                 ITRS_TEST = int(len(DX) / batch_size)
263                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
264                 TEST_LOSS = 0.0
265                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es @divisin exacta, se
266                     naade una @iteracin
267
268             for step in range(ITRS_TEST):
269
270                 batch_inputs, batch_context =

```

```

268         getBatchData(DX,DY, batch_size , step, ITRS_TEST , RES_ITMS_TEST)
269
270         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
271         loss_val, softplus = session.run([loss_softplus, softplus_batch],
272         feed_dict=feed_dict)
273
274         TEST_LOSS+=np.sum(softplus[:,0])
275
276         dev_hist.append(TEST_LOSS/(len(DX)*1.0))
277         train_hist.append(TRAIN_LOSS/(len(X)*1.0))
278
279         if (len(dev_hist) >= slope_size):
280             # Calcular la pendiente
281             print(getSlope(dev_hist[-slope_size:]))
282             print(dev_hist[-slope_size:])
283
284             if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
285
286                 def printArray(data, data2):
287                     for i in range(len(data)):
288                         a=data[i]
289                         b=data2[i]
290                         print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
291
292                 print("-" * 50)
293                 print(e)
294                 printArray(dev_hist, train_hist)
295                 return
296
297     else:
298         #
299
300     # TEST FINAL
301     #
302
303     TRUE_POSITIVE = 0
304     FALSE_NEGATIVE = 0
305     FALSE_POSITIVE = 0
306     TRUE_NEGATIVE = 0
307
308     ITRS_TEST = int(len(TX) / batch_size)
309     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
310     TEST_LOSS = 0.0
311     # Si no es divisin exacta, se aade una iteracin
312     if (RES_ITMS_TEST > 0): ITRS_TEST += 1
313
314     for step in range(ITRS_TEST):
315         # Si se aadi una iteracin, introducir los datos restantes
316         if (step == ITRS_TEST - 1 and RES_ITMS_TEST > 0):
317             batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
318             batch_context = TY[(ITRS_TEST - 1) * batch_size:]
319         # Si no, es un batch normal
320         else:
321             st = step * batch_size
322             nd = (step + 1) * batch_size

```

```

321         batch_inputs = TX[st:nd]
322         batch_context = TY[st:nd]
323
324         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
325         loss_val, loss_batch, p = session.run([loss_softplus, softplus_batch, prob],
326                                               feed_dict=feed_dict)
327
328         TEST_LOSS += np.sum(loss_batch)
329
330         # aClculo de TP/TN/FN/FP para cada batch
331
332         for i in range(len(p)):
333             real = batch_context[i, 0]
334             pred = p[i, 0]
335
336             if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
337             if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
338             if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
339             if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
340
341         PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
342         RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
343
344         F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
345
346         print("-" * 50)
347         print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
348         print("PRECISION:\t" + str(PRECISION))
349         print("RECALL:\t" + str(RECALL))
350         print("F1:\t" + str(F1))
351         print("-" * 50)
352         print("")
353
354         #-----
355         # Llamadas
356         #-----
357
358         TEST_NAME="test_031"
359
360         os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
361
362         #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
363         D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
364         D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
365         W2V = np.load('embeddings/w2v/EMB_MATRIX')
366
367         #Se concatenan los embeddings de los usuarios para generar el perfil
368         D2V_EMB = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
369         Consolidado DM + Consolidado DBOW
370
371         #Se cargan los datos de train de las canciones para conocer los IDs reales
372         DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
373
374         #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
375         ALL = set(list(range(W2V.shape[0])))

```

```

376 REAL = set(list(set(list(map(int, DATA))))))
377 NOT_EXISTS = list(ALL.difference-REAL))
378
379 #####
380 # TENSORFLOW
381 #####
382
383 #-----
384 # Grid Search + Early Stopping
385 #-----
386 '''
387 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
388 rates = [0.000001, 0.01, 0.0001]
389
390 for r in rates:
391     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
392
393     print("Learning Rate: "+str(r))
394     print("-" * 50)
395     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
396     print("-"*50)
397
398 exit()
399 '''
400 #-----
401 # Train
402 #-----
403 '''
404 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
405
406 #Para esta fase, se une TRAIN y DEV
407 X = np.concatenate((X,DX),axis=0)
408 Y = np.concatenate((Y,DY),axis=0)
409
410 DX=[]
411 DY=[]
412
413 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 199, "seed": 100}
414 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
415
416 exit()
417 '''
418 #-----
419 # Test
420 #-----
421
422
423 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
424
425 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 199, "seed": 100}
426 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```

### 4.3. TEST 03

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data
20             ="datos/Append/ARRAY_ARRAYS_REC.pkl",pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl",shuffl
21             seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
50
51         # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
52         ahora
53         USED_IDS = set(list(map(int, PAST_DATA[u])))
54         USED_IDS.update(TRAIN_DATA)
55         TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))

```

```

53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101
102
103    return
104    TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105    def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106

```



```

107 # Si se ha aadi una iteracin, introducir los datos restantes
108 if (step == ITRS - 1 and RES_ITMS > 0):
109     batch_inputs = X[(ITRS - 1) * batch_size:]
110     batch_context = Y[(ITRS - 1) * batch_size:]
111 # Si no, es un batch normal
112 else:
113     st = step * batch_size
114     nd = (step + 1) * batch_size
115
116     batch_inputs = X[st:nd]
117     batch_context = Y[st:nd]
118 return(batch_inputs, batch_context)
119
120 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     song_emb_size = 64
135     hidden_size = 32
136
137     num_sampled = 64 # Ejemplos negativos (Para el NSE)
138
139     model_name = 'model'
140     model_path = 'models/' + TEST_NAME + '/model'
141     complete_path = model_path + "/" + model_name
142
143     # Se almacena la loss de dev de las 5 ultimas epochs
144     slope_size = 100
145     train_hist=[]
146     dev_hist=[]
147
148     # Creacin del grafo de TF.
149     graph = tf.Graph()
150
151     with graph.as_default():
152
153         if(config["seed"]!=1):
154             tf.set_random_seed(config["seed"])
155
156         # nmero global de iteraciones
157         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
158
159         # Datos de entrada
160
161         # Array del tamao del batch con las X
162         train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")

```

```

163 # Array del tamaño del batch con las Y asociadas
164 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
165
166 # Embeddings


---


167 D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
168
169 W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
170
171 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
172     stddev=1.0 / math.sqrt(hidden_size)), name="T1")
173 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
174
175 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
176     stddev=1.0 / math.sqrt(hidden_size)), name="T2")
177 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
178
179 # Operaciones


---


180 # Embedding de documento
181 ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
182
183 # Embedding de canción
184 ec = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
185
186 # Transformar a 32 documento
187 hd = tf.matmul(ed, T1) + B1
188
189 # Transformar a 32 canción
190 hw = tf.matmul(ec, T2) + B2
191
192 # Obtener el producto escalar
193 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
194
195 # Cálculo de LOSS y optimizador


---


196
197 # Obtener la loss
198 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
199 loss_softplus = tf.reduce_mean(softplus_batch)
200
201 # Minimizar la loss
202 train_step =
203     tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus, global_step=
204
205 # Obtener la probabilidad
206 prob = tf.sigmoid(dot_prod)
207
208 # Inicializar variables
209 init = tf.global_variables_initializer()
210
211 # Crear objeto encargado de almacenar la red
212 saver = tf.train.Saver(max_to_keep=1)
213
214 config = tf.ConfigProto()

```

```

214 config.gpu_options.allow_growth = True
215 with tf.Session(graph=graph, config=config) as session:
216
217     # We must initialize all variables before we use them.
218     init.run()
219
220     if(save):
221         # Obtener el checkpoint
222         ckpt = tf.train.get_checkpoint_state(model_path)
223
224         # Si existe el checkpoint restaurar
225         if ckpt and ckpt.model_checkpoint_path:
226             saver.restore(session, ckpt.model_checkpoint_path)
227
228     #Si se desea entrenar
229     if(train):
230         #-----
231         # TRAIN
232         #-----
233
234         ITRS= int(len(X)/batch_size)
235         RES_ITMS = len(X) - (ITRS*batch_size)
236         if(RES_ITMS>0):ITRS+=1 #Si no es @divisin exacta, se naade una @iteracin
237
238         for e in range(epochs):
239             TRAIN_LOSS = 0.0
240
241             print("Epoch "+str(e))
242
243             for step in range(ITRS):
244
245                 batch_inputs, batch_context =
246                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
247
248                 feed_dict = {train.dataset: batch_inputs, train.labels: batch_context}
249
250                 ., loss_val,softplus, gs = session.run([train_step,
251                     loss_softplus,softplus_batch,global_step], feed_dict=feed_dict)
252
253                 TRAIN_LOSS += np.sum(softplus[:, 0])
254
255             if (save): saver.save(session, complete_path, global_step=global_step)
256
257             #-----
258             # DEV
259             #-----
260
261             if(len(DX)>0):
262                 ITRS_TEST = int(len(DX) / batch_size)
263                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
264                 TEST_LOSS = 0.0
265                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es @divisin exacta, se
266                     naade una @iteracin
267
268             for step in range(ITRS_TEST):
269
270                 batch_inputs, batch_context =

```

```

268         getBatchData(DX,DY, batch_size , step, ITRS_TEST , RES_ITMS_TEST)
269
270         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
271         loss_val, softplus = session.run([loss_softplus, softplus_batch],
272         feed_dict=feed_dict)
273
274         TEST_LOSS+=np.sum(softplus[:,0])
275
276         dev_hist.append(TEST_LOSS/(len(DX)*1.0))
277         train_hist.append(TRAIN_LOSS/(len(X)*1.0))
278
279         if (len(dev_hist) >= slope_size):
280             # Calcular la pendiente
281             print(getSlope(dev_hist[-slope_size:]))
282             print(dev_hist[-slope_size:])
283
284             if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
285
286                 def printArray(data, data2):
287                     for i in range(len(data)):
288                         a=data[i]
289                         b=data2[i]
290                         print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
291
292                 print("--" * 50)
293                 print(e)
294                 printArray(dev_hist, train_hist)
295                 return
296
297     else:
298         #
299
300     # TEST FINAL
301     #
302
303     TRUE_POSITIVE = 0
304     FALSE_NEGATIVE = 0
305     FALSE_POSITIVE = 0
306     TRUE_NEGATIVE = 0
307
308     ITRS_TEST = int(len(TX) / batch_size)
309     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
310     TEST_LOSS = 0.0
311     # Si no es divisin exacta, se aade una iteracin
312     if (RES_ITMS_TEST > 0): ITRS_TEST += 1
313
314     for step in range(ITRS_TEST):
315         # Si se aadi una iteracin, introducir los datos restantes
316         if (step == ITRS_TEST - 1 and RES_ITMS_TEST > 0):
317             batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
318             batch_context = TY[(ITRS_TEST - 1) * batch_size:]
319         # Si no, es un batch normal
320         else:
321             st = step * batch_size
322             nd = (step + 1) * batch_size

```

```

321         batch_inputs = TX[st:nd]
322         batch_context = TY[st:nd]
323
324         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
325         loss_val, loss_batch, p = session.run([loss_softplus, softplus_batch, prob],
326                                               feed_dict=feed_dict)
327
328         TEST_LOSS += np.sum(loss_batch)
329
330         # aClculo de TP/TN/FN/FP para cada batch
331
332         for i in range(len(p)):
333             real = batch_context[i, 0]
334             pred = p[i, 0]
335
336             if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
337             if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
338             if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
339             if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
340
341         PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
342         RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
343
344         F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
345
346         print("-" * 50)
347         print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
348         print("PRECISION:\t" + str(PRECISION))
349         print("RECALL:\t" + str(RECALL))
350         print("F1:\t" + str(F1))
351         print("-" * 50)
352         print("")
353
354         #-----
355         # Llamadas
356         #-----
357
358         TEST_NAME="test_041"
359
360         os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
361
362         #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
363         D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
364         D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
365         W2V = np.load('embeddings/w2v/EMB_MATRIX')
366
367         #Se concatenan los embeddings de los usuarios para generar el perfil
368         D2V_EMB = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El reciente = Reciente
369         DM + Reciente DBOW
370
371         #Se cargan los datos de train de las canciones para conocer los IDs reales
372         DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
373
374         #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
375         ALL = set(list(range(W2V.shape[0])))

```

```

376 REAL = set(list(set(list(map(int, DATA))))))
377 NOT_EXISTS = list(ALL.difference-REAL))
378
379 #####
380 # TENSORFLOW
381 #####
382
383 #-----
384 # Grid Search + Early Stopping
385 #-----
386 '''
387 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
388 rates = [0.000001, 0.01, 0.0001]
389
390 for r in rates:
391     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
392
393     print("Learning Rate: "+str(r))
394     print("-" * 50)
395     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
396     print("-"*50)
397
398 exit()
399 '''
400 #-----
401 # Train
402 #-----
403 '''
404 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
405
406 #Para esta fase, se une TRAIN y DEV
407 X = np.concatenate((X,DX),axis=0)
408 Y = np.concatenate((Y,DY),axis=0)
409
410 DX=[]
411 DY=[]
412
413 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 202, "seed": 100}
414 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
415
416 exit()
417 '''
418 #-----
419 # Test
420 #-----
421
422
423 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
424
425 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 202, "seed": 100}
426 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```

## 4.4. TEST 04

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data
20             ="datos/Append/ARRAY_ARRAYS_REC.pkl",pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl",shuffl
21             seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
50
51         # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
52         ahora
53         USED_IDS = set(list(map(int, PAST_DATA[u])))
54         USED_IDS.update(TRAIN_DATA)
55         TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))

```

```

53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101
102
103    return
104    TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105    def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106

```



```

107 # Si se ha aadi una iteracin, introducir los datos restantes
108 if (step == ITRS - 1 and RES_ITMS > 0):
109     batch_inputs = X[(ITRS - 1) * batch_size:]
110     batch_context = Y[(ITRS - 1) * batch_size:]
111 # Si no, es un batch normal
112 else:
113     st = step * batch_size
114     nd = (step + 1) * batch_size
115
116     batch_inputs = X[st:nd]
117     batch_context = Y[st:nd]
118 return(batch_inputs, batch_context)
119
120 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     user_concat_size = user_emb_size * 2
135     song_emb_size = 64
136     hidden_size = 32
137
138     num_sampled = 64 # Ejemplos negativos (Para el NSE)
139
140     model_name = 'model'
141     model_path = 'models/' + TEST_NAME + '/model'
142     complete_path = model_path + "/" + model_name
143
144     # Se almacena la loss de dev de las 5 ultimas epochs
145     slope_size = 100
146     train_hist=[]
147     dev_hist=[]
148
149     # Creacin del grafo de TF.
150     graph = tf.Graph()
151
152     with graph.as_default():
153
154         if(config["seed"]!=1):
155             tf.set_random_seed(config["seed"])
156
157         # nmero global de iteraciones
158         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
159
160         # Datos de entrada
161
162         # Array del tamao del batch con las X
163         train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")

```

```

163
164 # Array del tamaño del batch con las Y asociadas
165 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
166
167 # Embeddings


---


168
169 D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
170
171 T0 = tf.Variable(tf.truncated_normal([user_concat_size, user_emb_size], mean=0.0,
172     stddev=1.0 / math.sqrt(hidden_size)), name="T0")
173 B0 = tf.Variable(tf.zeros([user_emb_size]), name="B0")
174
175 W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
176
177 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
178     stddev=1.0 / math.sqrt(hidden_size)), name="T1")
179 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
180
181 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
182     stddev=1.0 / math.sqrt(hidden_size)), name="T2")
183 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
184
185 # Operaciones


---


186
187 # Embedding de documento
188 ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
189
190 #Transformar concat documento de 256 a 128
191 ed = tf.matmul(ed, T0)+B0
192
193 # Embedding de canción
194 ec = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
195
196 # Transformar a 32 documento
197 hd = tf.matmul(ed, T1)+B1
198
199 # Transformar a 32 canción
200 hw = tf.matmul(ec, T2)+B2
201
202 # Obtener el producto escalar
203 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
204
205 # aClculo de LOSS y optimizador—


---


206
207 # Obtener la loss
208 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
209 loss_softplus = tf.reduce_mean(softplus_batch)
210
211 # Minimizar la loss
212 train_step =
213     tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus, global_step
214
215 # Obtener la probabilidad
216 prob = tf.sigmoid(dot_prod)

```

```

213
214 # Inicializar variables
215 init = tf.global_variables_initializer()
216
217 # Crear objeto encargado de almacenar la red
218 saver = tf.train.Saver(max_to_keep=1)
219
220 #config = tf.ConfigProto()
221 #config.gpu_options.allow_growth = True
222 #with tf.Session(graph=graph,config=config) as session:
223 with tf.Session(graph=graph) as session:
224
225     # We must initialize all variables before we use them.
226     init.run()
227
228     if(save):
229         # Obtener el checkpoint
230         ckpt = tf.train.get_checkpoint_state(model_path)
231
232         # Si existe el checkpoint restaurar
233         if ckpt and ckpt.model_checkpoint_path:
234             saver.restore(session, ckpt.model_checkpoint_path)
235
236     #Si se desea entrenar
237     if(train):
238         #-----
239         # TRAIN
240         #-----
241
242         ITRS= int(len(X)/batch_size)
243         RES_ITMS = len(X) - (ITRS*batch_size)
244         if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
245
246         for e in range(epochs):
247             TRAIN_LOSS = 0.0
248
249             print("Epoch "+str(e))
250
251             for step in range(ITRS):
252
253                 batch_inputs, batch_context =
254                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
255
256                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
257
258                 _, loss_val,softplus, gs = session.run([train_step,
259                     loss_softplus,softplus_batch,global_step], feed_dict=feed_dict)
260
261                 TRAIN_LOSS += np.sum(softplus[:, 0])
262
263             if (save): saver.save(session, complete_path, global_step=global_step)
264
265             #-----
266             # DEV
267             #-----
268
269         if(len(DX)>0):

```

```

268     ITRS_TEST = int(len(DX) / batch_size)
269     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
270     TEST_LOSS = 0.0
271     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es ÷divisin exacta, se
        naade una ÷iteracin
272
273     for step in range(ITRS_TEST):
274
275         batch_inputs, batch_context =
            getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
276
277         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
278         loss_val, softplus = session.run([loss_softplus, softplus_batch],
            feed_dict=feed_dict)
279
280         TEST_LOSS+=np.sum(softplus[:,0])
281
282         dev_hist.append(TEST_LOSS/(len(DX)*1.0))
283         train_hist.append(TRAIN_LOSS/(len(X)*1.0))
284
285     if (len(dev_hist) >= slope_size):
286         # Calcular la pendiente
287         print(getSlope(dev_hist[-slope_size:]))
288         print(dev_hist[-slope_size:])
289
290     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
291
292         def printArray(data, data2):
293             for i in range(len(data)):
294                 a=data[i]
295                 b=data2[i]
296                 print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
297
298         print("-" * 50)
299         print(e)
300         printArray(dev_hist, train_hist)
301         return
302
303 else:
304     #
305
306     # TEST FINAL
307     #
308
309     TRUE_POSITIVE = 0
310     FALSE_NEGATIVE = 0
311     FALSE_POSITIVE = 0
312     TRUE_NEGATIVE = 0
313
314     ITRS_TEST = int(len(TX) / batch_size)
315     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
316     TEST_LOSS = 0.0
317     # Si no es ÷divisin exacta, se naade una ÷iteracin
318     if (RES_ITMS_TEST > 0): ITRS_TEST += 1
319
320     for step in range(ITRS_TEST):

```

```

320     # Si se reañadi una iteración, introducir los datos restantes
321     if (step == ITRS_TEST - 1 and RES_ITMS_TEST > 0):
322         batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
323         batch_context = TY[(ITRS_TEST - 1) * batch_size:]
324     # Si no, es un batch normal
325     else:
326         st = step * batch_size
327         nd = (step + 1) * batch_size
328
329         batch_inputs = TX[st:nd]
330         batch_context = TY[st:nd]
331
332     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
333     loss_val, loss_batch, p = session.run([loss_softplus, softplus_batch, prob],
334         feed_dict=feed_dict)
335
336     TEST_LOSS += np.sum(loss_batch)
337
338     # cálculo de TP/TN/FN/FP para cada batch
339
340     for i in range(len(p)):
341         real = batch_context[i, 0]
342         pred = p[i, 0]
343
344         if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
345         if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
346         if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
347         if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
348
349     PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
350     RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
351
352     F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
353
354     print("-" * 50)
355     print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
356     print("PRECISION:\t" + str(PRECISION))
357     print("RECALL:\t" + str(RECALL))
358     print("F1:\t" + str(F1))
359     print("-" * 50)
360     print("")
361
362     #####
363     # Llamadas
364     #####
365
366     TEST_NAME="test_051"
367
368     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
369
370     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
371     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
372     D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
373
374     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
375     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')

```

```

376 D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
377
378 W2V = np.load('embeddings/w2v/EMB_MATRIX')
379
380 #Se concatenan los embeddings de los usuarios para generar el perfil
381 D2V_CONS = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
    Consolidado DM + Consolidado DBOW
382 D2V_PRES = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El presente = Presente
    DM + Presente DBOW
383
384 D2V_EMB = np.concatenate((D2V_CONS, D2V_PRES), axis=1) # El perfil completo CONS + PRES
385
386 #Se cargan los datos de train de las canciones para conocer los IDs reales
387 DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
388
389 #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
390 ALL = set(list(range(W2V.shape[0])))
391 REAL = set(list(set(list(map(int, DATA)))))
392 NOT_EXISTS = list(ALL.difference(REAL))
393
394 #####
395 # TENSORFLOW
396 #####
397
398 #-----
399 # Grid Search + Early Stopping
400 #-----
401 '''
402 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
403 rates = [0.000001, 1, 0.01, 0.0001]
404
405 for r in rates:
406     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
407
408     print("Learning Rate: "+str(r))
409     print("-" * 50)
410     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
411     print("-"*50)
412
413 exit()
414 '''
415 #-----
416 # Train
417 #-----
418 '''
419 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
420
421 #Para esta fase, se une TRAIN y DEV
422 X = np.concatenate((X, DX), axis=0)
423 Y = np.concatenate((Y, DY), axis=0)
424
425 DX=[]
426 DY=[]
427
428 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 194, "seed": 100}
429 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
430

```

```
431 exit()
432 ', '
433 #
434 # Test
435 #
436
437
438 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
439
440 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs": 194, "seed": 100}
441 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)
```

## 4.5. TEST 05

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data
20             ="datos/Append/ARRAY_ARRAYS_REC.pkl",pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl",shuffl
21             seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
50
51         # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
52         ahora
53         USED_IDS = set(list(map(int, PAST_DATA[u])))
54         USED_IDS.update(TRAIN_DATA)
55         TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))

```



```

53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101
102
103    return
104    TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105    def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106

```

```

107 # Si se noaadi una iteracin, introducir los datos restantes
108 if (step == ITRS - 1 and RES_ITMS > 0):
109     batch_inputs = X[(ITRS - 1) * batch_size:]
110     batch_context = Y[(ITRS - 1) * batch_size:]
111 # Si no, es un batch normal
112 else:
113     st = step * batch_size
114     nd = (step + 1) * batch_size
115
116     batch_inputs = X[st:nd]
117     batch_context = Y[st:nd]
118 return(batch_inputs, batch_context)
119
120 def trainNet(X,Y, DX,DY, TX,TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     song_emb_size = 64
135     hidden_size = 32
136
137     num_sampled = 64 # Ejemplos negativos (Para el NSE)
138
139     model_name = 'model'
140     model_path = 'models/' + TEST_NAME + '/model'
141     complete_path = model_path + "/" + model_name
142
143     # Se almacena la loss de dev de las 5 ultimas epochs
144     slope_size = 100
145     train_hist=[]
146     dev_hist=[]
147
148     # Creacin del grafo de TF.
149     graph = tf.Graph()
150
151     with graph.as_default():
152
153         if(config["seed"]!=1):
154             tf.set_random_seed(config["seed"])
155
156         # nmero global de iteraciones
157         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
158
159         # Datos de entrada
160
161         # Array del tamaño del batch con las X
162         train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")

```

```

163 # Array del tamaño del batch con las Y asociadas
164 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
165
166 # Embeddings


---


167
168 # Matriz W de embeddings de las canciones obtenida en el W2V previo
169 #word_embeddings = tf.Variable(W2V, trainable=False, name="word_embeddings") #
    FUNCIONA???
```

<https://stackoverflow.com/questions/37326002/is-it-possible-to-make-a-trainable-variable-no>

```

170
171 O1 = tf.Variable(tf.truncated_normal([user_size, user_emb_size], mean=0.0, stddev=1.0
    / math.sqrt(user_emb_size)), name="O1")
172 B01 = tf.Variable(tf.zeros([user_emb_size]), name="B01")
173
174 O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
    / math.sqrt(song_emb_size)), name="O2")
175 B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
176
177 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
    stddev=1.0 / math.sqrt(hidden_size)), name="T1")
178 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
179
180 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
    stddev=1.0 / math.sqrt(hidden_size)), name="T2")
181 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
182
183 # Operaciones


---


184
185 # Embedding de documento
186 ed = tf.nn.embedding_lookup(O1, train_dataset[:, 0])
187 ed = ed + B01
188
189 # Embedding de canción
190 ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
191 ec = ec + B02
192
193 # Transformar a 32 documento
194 hd = tf.matmul(ed, T1) + B1
195
196 # Transformar a 32 canción
197 hw = tf.matmul(ec, T2) + B2
198
199 # Obtener el producto escalar
200 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
201
202 # Cálculo de LOSS y optimizador


---


203
204 # Obtener la loss
205 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
206 loss_softplus = tf.reduce_mean(softplus_batch)
207
208 # Minimizar la loss
209 train_step =
    tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus, global_step
```

```

210
211 # Obtener la probabilidad
212 prob = tf.sigmoid(dot_prod)
213
214 # Inicializar variables
215 init = tf.global_variables_initializer()
216
217 # Crear objeto encargado de almacenar la red
218 saver = tf.train.Saver(max_to_keep=1)
219
220 config = tf.ConfigProto()
221 config.gpu_options.allow_growth = True
222 with tf.Session(graph=graph, config=config) as session:
223
224     # We must initialize all variables before we use them.
225     init.run()
226
227     if(save):
228         # Obtener el checkpoint
229         ckpt = tf.train.get_checkpoint_state(model_path)
230
231         # Si existe el checkpoint restaurar
232         if ckpt and ckpt.model_checkpoint_path:
233             saver.restore(session, ckpt.model_checkpoint_path)
234
235     #Si se desea entrenar
236     if(train):
237         #-----
238         # TRAIN
239         #-----
240
241         ITRS= int(len(X)/batch_size)
242         RES_ITMS = len(X) - (ITRS*batch_size)
243         if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se añade una iteracin
244
245         for e in range(epochs):
246             TRAIN_LOSS = 0.0
247
248             print("Epoch "+str(e))
249
250             for step in range(ITRS):
251
252                 batch_inputs, batch_context =
253                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
254
255                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
256
257                 ., loss_val,softplus, gs = session.run([train_step,
258                     loss_softplus,softplus_batch,global_step], feed_dict=feed_dict)
259
260                 TRAIN_LOSS += np.sum(softplus[:, 0])
261
262             if (save): saver.save(session, complete_path, global_step=global_step)
263
264         #-----
265         # DEV
266         #-----

```

```

265
266     if(len(DX)>0):
267         ITRS_TEST = int(len(DX) / batch_size)
268         RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
269         TEST_LOSS = 0.0
270         if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es ÷divisin exacta, se
                ñaade una ÷iteracin
271
272         for step in range(ITRS_TEST):
273
274             batch_inputs, batch_context =
                getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
275
276             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
277             loss_val, softplus = session.run([loss_softplus,softplus_batch],
                feed_dict=feed_dict)
278
279             TEST_LOSS+=np.sum(softplus[:,0])
280
281             dev_hist.append(TEST_LOSS/(len(DX)*1.0))
282             train_hist.append(TRAIN_LOSS/(len(X)*1.0))
283
284
285         if (len(dev_hist) >= slope_size):
286             # Calcular la pendiente
287             print(getSlope(dev_hist[-slope_size:]))
288             print(dev_hist[-slope_size:])
289
290             if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
291
292                 def printArray(data,data2):
293                     for i in range(len(data)):
294                         a=data[i]
295                         b=data2[i]
296                         print(str(a).replace(".",",")+ "\t"+str(b).replace(".",","))
297
298                 print("-" * 50)
299                 print(e)
300                 printArray(dev_hist,train_hist)
301                 return
302
303     else:
304         #
305
306     # TEST FINAL
307     #
308
309     TRUE_POSITIVE = 0
310     FALSE_NEGATIVE = 0
311     FALSE_POSITIVE = 0
312     TRUE_NEGATIVE = 0
313
314     ITRS_TEST = int(len(TX) / batch_size)
315     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
316     TEST_LOSS = 0.0
317     # Si no es ÷divisin exacta, se ñaade una ÷iteracin

```

```

317     if (RES.ITMS.TEST > 0): ITRS_TEST += 1
318
319     for step in range(ITRS_TEST):
320         # Si se reañadi una iteración, introducir los datos restantes
321         if (step == ITRS_TEST - 1 and RES.ITMS.TEST > 0):
322             batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
323             batch_context = TY[(ITRS_TEST - 1) * batch_size:]
324         # Si no, es un batch normal
325         else:
326             st = step * batch_size
327             nd = (step + 1) * batch_size
328
329             batch_inputs = TX[st:nd]
330             batch_context = TY[st:nd]
331
332             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
333             loss_val, loss_batch, p = session.run([loss_softplus, softplus_batch, prob],
334                                                  feed_dict=feed_dict)
335
336             TEST_LOSS += np.sum(loss_batch)
337
338             # aClculo de TP/TN/FN/FP para cada batch
339
340             for i in range(len(p)):
341                 real = batch_context[i, 0]
342                 pred = p[i, 0]
343
344                 if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
345                 if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
346                 if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
347                 if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
348
349             PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
350             RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
351
352             F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
353
354             print("-" * 50)
355             print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
356             print("PRECISION:\t" + str(PRECISION))
357             print("RECALL:\t" + str(RECALL))
358             print("F1:\t" + str(F1))
359             print("-" * 50)
360             print("")
361
362     #####
363     # Llamadas
364     #####
365
366     TEST_NAME="test_011"
367
368     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
369
370     #Para conocer los embeddings de las canciones
371     W2V = np.load('embeddings/w2v/EMB_MATRIX')
372

```

```

373 #Se cargan los datos de train de las canciones para conocer los IDs reales
374 DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
375
376 #Obtener una lista con los ID de canciones que no aectn el la matriz de embeddings
377 ALL = set(list(range(W2V.shape[0])))
378 REAL = set(list(set(list(map(int, DATA)))))
379
380 #####
381 # TENSORFLOW
382 #####
383
384 # -----
385 # Grid Search + Early Stopping
386 # -----
387 '''
388 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
389 rates = [0.000001, 0.01, 0.0001]
390
391 for r in rates:
392     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
393
394     print("Learning Rate: "+str(r))
395     print("-" * 50)
396     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
397     print("-"*50)
398
399 exit()
400 '''
401 # -----
402 # Train
403 # -----
404 '''
405 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=REAL)
406
407 #Para esta fase, se une TRAIN y DEV
408 X = np.concatenate((X,DX),axis=0)
409 Y = np.concatenate((Y,DY),axis=0)
410
411 DX=[]
412 DY=[]
413
414 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 410, "seed": 100}
415 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
416
417 exit()
418 '''
419 # -----
420 # Test
421 # -----
422
423
424 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
425
426 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 410, "seed": 100}
427 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```

## 4.6. TEST 06

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data
20             ="datos/Append/ARRAY_ARRAYS_REC.pkl",pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl",shuffl
21             seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
50
51         # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
52         ahora
53         USED_IDS = set(list(map(int, PAST_DATA[u])))
54         USED_IDS.update(TRAIN_DATA)
55         TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))

```



```

53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101
102
103    return
104        TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105    def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106

```

```

107 # Si se ha aadi una iteracin, introducir los datos restantes
108 if (step == ITRS - 1 and RES_ITMS > 0):
109     batch_inputs = X[(ITRS - 1) * batch_size:]
110     batch_context = Y[(ITRS - 1) * batch_size:]
111 # Si no, es un batch normal
112 else:
113     st = step * batch_size
114     nd = (step + 1) * batch_size
115
116     batch_inputs = X[st:nd]
117     batch_context = Y[st:nd]
118 return(batch_inputs, batch_context)
119
120 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     song_emb_size = 64
135     hidden_size = 32
136
137     num_sampled = 64 # Ejemplos negativos (Para el NSE)
138
139     model_name = 'model'
140     model_path = 'models/' + TEST_NAME + '/model'
141     complete_path = model_path + "/" + model_name
142
143     # Se almacena la loss de dev de las 5 ultimas epochs
144     slope_size = 100
145     train_hist=[]
146     dev_hist=[]
147
148     # Creacin del grafo de TF.
149     graph = tf.Graph()
150
151     with graph.as_default():
152
153         if(config["seed"]!=1):
154             tf.set_random_seed(config["seed"])
155
156         # nmero global de iteraciones
157         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
158
159         # Datos de entrada
160
161         # Array del tamao del batch con las X
162         train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")

```

```

163 # Array del tamaño del batch con las Y asociadas
164 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
165
166 # Embeddings


---


167 D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
168
169 O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
170 / math.sqrt(song_emb_size)), name="O2")
171 B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
172
173 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
174 stddev=1.0 / math.sqrt(hidden_size)), name="T1")
175 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
176
177 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
178 stddev=1.0 / math.sqrt(hidden_size)), name="T2")
179 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
180
181 # Operaciones


---


182 # Embedding de documento
183 ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
184
185 # Embedding de canción
186 ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
187 ec = ec + B02
188
189 # Transformar a 32 documento
190 hd = tf.matmul(ed, T1) + B1
191
192 # Transformar a 32 canción
193 hw = tf.matmul(ec, T2) + B2
194
195 # Obtener el producto escalar
196 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
197
198 # Cálculo de LOSS y optimizador


---


199 # Obtener la loss
200 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
201 loss_softplus = tf.reduce_mean(softplus_batch)
202
203 # Minimizar la loss
204 train_step =
205     tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus, global_step=
206
207 # Obtener la probabilidad
208 prob = tf.sigmoid(dot_prod)
209
210 # Inicializar variables
211 init = tf.global_variables_initializer()
212
213 # Crear objeto encargado de almacenar la red

```

```

213     saver = tf.train.Saver(max_to_keep=1)
214
215     config = tf.ConfigProto()
216     config.gpu_options.allow_growth = True
217     with tf.Session(graph=graph, config=config) as session:
218
219         # We must initialize all variables before we use them.
220         init.run()
221
222         if(save):
223             # Obtener el checkpoint
224             ckpt = tf.train.get_checkpoint_state(model_path)
225
226             # Si existe el checkpoint restaurar
227             if ckpt and ckpt.model_checkpoint_path:
228                 saver.restore(session, ckpt.model_checkpoint_path)
229
230         #Si se desea entrenar
231         if(train):
232             #-----
233             # TRAIN
234             #-----
235
236             ITRS= int(len(X)/batch_size)
237             RES_ITMS = len(X) - (ITRS*batch_size)
238             if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
239
240             for e in range(epochs):
241                 TRAIN_LOSS = 0.0
242
243                 print("Epoch "+str(e))
244
245                 for step in range(ITRS):
246
247                     batch_inputs, batch_context =
248                         getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
249
250                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
251
252                     _, loss_val,softplus, gs = session.run([train_step,
253                         loss_softplus,softplus_batch,global_step], feed_dict=feed_dict)
254
255                     TRAIN_LOSS += np.sum(softplus[:, 0])
256
257                     if (save): saver.save(session, complete_path, global_step=global_step)
258
259                     #-----
260                     # DEV
261                     #-----
262
263                     if(len(DX)>0):
264                         ITRS_TEST = int(len(DX) / batch_size)
265                         RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
266                         TEST_LOSS = 0.0
267                         if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se
268                             aade una iteracin

```

```

267     for step in range(ITRS_TEST):
268
269         batch_inputs, batch_context =
                getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
270
271         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
272         loss_val, softplus = session.run([loss_softplus, softplus_batch],
                feed_dict=feed_dict)
273
274         TEST_LOSS+=np.sum(softplus[:,0])
275
276         dev_hist.append(TEST_LOSS/(len(DX)*1.0))
277         train_hist.append(TRAIN_LOSS/(len(X)*1.0))
278
279     if (len(dev_hist) >= slope_size):
280         # Calcular la pendiente
281         print(getSlope(dev_hist[-slope_size:]))
282         print(dev_hist[-slope_size:])
283
284     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
285
286         def printArray(data,data2):
287             for i in range(len(data)):
288                 a=data[i]
289                 b=data2[i]
290                 print(str(a).replace(".",",")+ "\t"+str(b).replace(".",","))
291
292         print("--" * 50)
293         print(e)
294         printArray(dev_hist,train_hist)
295         return
296
297 else:
298     #
299
300     # TEST FINAL
301     #
302
303     TRUE_POSITIVE = 0
304     FALSE_NEGATIVE = 0
305     FALSE_POSITIVE = 0
306     TRUE_NEGATIVE = 0
307
308     ITRS_TEST = int(len(TX) / batch_size)
309     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
310     TEST_LOSS = 0.0
311     # Si no es divisible exacta, se añade una iteración
312     if (RES_ITMS_TEST > 0): ITRS_TEST += 1
313
314     for step in range(ITRS_TEST):
315         # Si se añade una iteración, introducir los datos restantes
316         if (step == ITRS_TEST - 1 and RES_ITMS_TEST > 0):
317             batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
318             batch_context = TY[(ITRS_TEST - 1) * batch_size:]
319         # Si no, es un batch normal
320         else:

```

```

320         st = step * batch_size
321         nd = (step + 1) * batch_size
322
323         batch_inputs = TX[st:nd]
324         batch_context = TY[st:nd]
325
326         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
327         loss_val, loss_batch, p = session.run([loss_softplus, softplus_batch, prob],
328             feed_dict=feed_dict)
329
330         TEST_LOSS += np.sum(loss_batch)
331
332         # aClculo de TP/TN/FN/FP para cada batch
333
334         for i in range(len(p)):
335             real = batch_context[i, 0]
336             pred = p[i, 0]
337
338             if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
339             if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
340             if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
341             if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
342
343         PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
344         RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
345
346         F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
347
348         print("-" * 50)
349         print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
350         print("PRECISION:\t" + str(PRECISION))
351         print("RECALL:\t" + str(RECALL))
352         print("F1:\t" + str(F1))
353         print("-" * 50)
354         print("")
355
356         #-----
357         # Llamadas
358         #-----
359
360         TEST_NAME="test_061"
361
362         os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
363
364         #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
365         D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
366         D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
367         W2V = np.load('embeddings/w2v/EMB_MATRIX')
368
369         #Se concatenan los embeddings de los usuarios para generar el perfil
370         D2V_EMB = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
371             Consolidado DM + Consolidado DBOW
372
373         #Se cargan los datos de train de las canciones para conocer los IDs reales
374         DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")

```

```

375 #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
376 ALL = set(list(range(W2V.shape[0])))
377 REAL = set(list(set(list(map(int, DATA))))))
378 NOT_EXISTS = list(ALL.difference(REAL))
379
380 #####
381 # TENSORFLOW
382 #####
383
384 # -----
385 # Grid Search + Early Stopping
386 # -----
387 '''
388 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
389 rates = [0.000001, 0.01, 0.0001]
390
391 for r in rates:
392     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
393
394     print("Learning Rate: "+str(r))
395     print("-" * 50)
396     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
397     print("-"*50)
398
399 exit()
400 '''
401 # -----
402 # Train
403 # -----
404 '''
405 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
406
407 #Para esta fase, se une TRAIN y DEV
408 X = np.concatenate((X,DX),axis=0)
409 Y = np.concatenate((Y,DY),axis=0)
410
411 DX=[]
412 DY=[]
413
414 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 464, "seed": 100}
415 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
416
417 exit()
418 '''
419 # -----
420 # Test
421 # -----
422
423
424 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
425
426 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 464, "seed": 100}
427 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```

## 4.7. TEST 07

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data
20             ="datos/Append/ARRAY_ARRAYS_REC.pkl",pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl",shuffl
21             seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
50
51         # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
52         ahora
53         USED_IDS = set(list(map(int, PAST_DATA[u])))
54         USED_IDS.update(TRAIN_DATA)
55         TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))

```



```

53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101
102
103    return
104    TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105    def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106

```

```

107 # Si se noaadi una iteracin, introducir los datos restantes
108 if (step == ITRS - 1 and RES_ITMS > 0):
109     batch_inputs = X[(ITRS - 1) * batch_size:]
110     batch_context = Y[(ITRS - 1) * batch_size:]
111 # Si no, es un batch normal
112 else:
113     st = step * batch_size
114     nd = (step + 1) * batch_size
115
116     batch_inputs = X[st:nd]
117     batch_context = Y[st:nd]
118 return(batch_inputs, batch_context)
119
120 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     song_emb_size = 64
135     hidden_size = 32
136
137     num_sampled = 64 # Ejemplos negativos (Para el NSE)
138
139     model_name = 'model'
140     model_path = 'models/' + TEST_NAME + '/model'
141     complete_path = model_path + "/" + model_name
142
143     # Se almacena la loss de dev de las 5 ultimas epochs
144     slope_size = 100
145     train_hist=[]
146     dev_hist=[]
147
148     # Creacin del grafo de TF.
149     graph = tf.Graph()
150
151     with graph.as_default():
152
153         if(config["seed"]!=1):
154             tf.set_random_seed(config["seed"])
155
156         # nmero global de iteraciones
157         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
158
159         # Datos de entrada
160
161         # Array del tamaño del batch con las X
162         train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")

```

```

163 # Array del tamaño del batch con las Y asociadas
164 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
165
166 # Embeddings


---


167 D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
168
169 O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
170 / math.sqrt(song_emb_size)), name="O2")
171 B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
172
173 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
174 stddev=1.0 / math.sqrt(hidden_size)), name="T1")
175 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
176
177 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
178 stddev=1.0 / math.sqrt(hidden_size)), name="T2")
179 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
180
181 # Operaciones


---


182 # Embedding de documento
183 ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
184
185 # Embedding de canción
186 ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
187 ec = ec + B02
188
189 # Transformar a 32 documento
190 hd = tf.matmul(ed, T1) + B1
191
192 # Transformar a 32 canción
193 hw = tf.matmul(ec, T2) + B2
194
195 # Obtener el producto escalar
196 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
197
198 # Cálculo de LOSS y optimizador


---


199 # Obtener la loss
200 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
201 loss_softplus = tf.reduce_mean(softplus_batch)
202
203 # Minimizar la loss
204 train_step =
205     tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus, global_step=
206
207 # Obtener la probabilidad
208 prob = tf.sigmoid(dot_prod)
209
210 # Inicializar variables
211 init = tf.global_variables_initializer()
212
213 # Crear objeto encargado de almacenar la red

```

```

213     saver = tf.train.Saver(max_to_keep=1)
214
215 config = tf.ConfigProto()
216 config.gpu_options.allow_growth = True
217 with tf.Session(graph=graph, config=config) as session:
218
219     # We must initialize all variables before we use them.
220     init.run()
221
222     if(save):
223         # Obtener el checkpoint
224         ckpt = tf.train.get_checkpoint_state(model_path)
225
226         # Si existe el checkpoint restaurar
227         if ckpt and ckpt.model_checkpoint_path:
228             saver.restore(session, ckpt.model_checkpoint_path)
229
230     #Si se desea entrenar
231     if(train):
232         #-----
233         # TRAIN
234         #-----
235
236         ITRS= int(len(X)/batch_size)
237         RES_ITMS = len(X) - (ITRS*batch_size)
238         if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
239
240         for e in range(epochs):
241             TRAIN_LOSS = 0.0
242
243             print("Epoch "+str(e))
244
245             for step in range(ITRS):
246
247                 batch_inputs, batch_context =
248                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
249
250                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
251
252                 _, loss_val,softplus, gs = session.run([train_step,
253                     loss_softplus,softplus_batch,global_step], feed_dict=feed_dict)
254
255                 TRAIN_LOSS += np.sum(softplus[:, 0])
256
257                 if (save): saver.save(session, complete_path, global_step=global_step)
258
259                 #-----
260                 # DEV
261                 #-----
262
263                 if(len(DX)>0):
264                     ITRS_TEST = int(len(DX) / batch_size)
265                     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
266                     TEST_LOSS = 0.0
267                     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se
268                         aade una iteracin

```

```

267     for step in range(ITRS_TEST):
268
269         batch_inputs, batch_context =
                getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
270
271         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
272         loss_val, softplus = session.run([loss_softplus, softplus_batch],
                feed_dict=feed_dict)
273
274         TEST_LOSS+=np.sum(softplus[:,0])
275
276         dev_hist.append(TEST_LOSS/(len(DX)*1.0))
277         train_hist.append(TRAIN_LOSS/(len(X)*1.0))
278
279     if (len(dev_hist) >= slope_size):
280         # Calcular la pendiente
281         print(getSlope(dev_hist[-slope_size:]))
282         print(dev_hist[-slope_size:])
283
284     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
285
286         def printArray(data,data2):
287             for i in range(len(data)):
288                 a=data[i]
289                 b=data2[i]
290                 print(str(a).replace(".",",")+ "\t"+str(b).replace(".",","))
291
292         print("--" * 50)
293         print(e)
294         printArray(dev_hist,train_hist)
295         return
296
297 else:
298     #
299
300     # TEST FINAL
301     #
302
303     TRUE_POSITIVE = 0
304     FALSE_NEGATIVE = 0
305     FALSE_POSITIVE = 0
306     TRUE_NEGATIVE = 0
307
308     ITRS_TEST = int(len(TX) / batch_size)
309     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
310     TEST_LOSS = 0.0
311     # Si no es divisible exacta, se añade una iteración
312     if (RES_ITMS_TEST > 0): ITRS_TEST += 1
313
314     for step in range(ITRS_TEST):
315         # Si se añade una iteración, introducir los datos restantes
316         if (step == ITRS_TEST - 1 and RES_ITMS_TEST > 0):
317             batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
318             batch_context = TY[(ITRS_TEST - 1) * batch_size:]
319         # Si no, es un batch normal
320         else:

```

```

320         st = step * batch_size
321         nd = (step + 1) * batch_size
322
323         batch_inputs = TX[st:nd]
324         batch_context = TY[st:nd]
325
326         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
327         loss_val, loss_batch, p = session.run([loss_softplus, softplus.batch, prob],
328             feed_dict=feed_dict)
329
330         TEST_LOSS += np.sum(loss_batch)
331
332         # aClculo de TP/TN/FN/FP para cada batch
333
334         for i in range(len(p)):
335             real = batch_context[i, 0]
336             pred = p[i, 0]
337
338             if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
339             if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
340             if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
341             if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
342
343         PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
344         RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
345
346         F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
347
348         print("-" * 50)
349         print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
350         print("PRECISION:\t" + str(PRECISION))
351         print("RECALL:\t" + str(RECALL))
352         print("F1:\t" + str(F1))
353         print("-" * 50)
354         print("")
355
356         #-----
357         #####
358         # Llamadas
359         #####
360         TEST_NAME="test_071"
361
362         os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
363
364         #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
365         D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
366         D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
367         W2V = np.load('embeddings/w2v/EMB_MATRIX')
368
369         #Se concatenan los embeddings de los usuarios para generar el perfil
370         D2V_EMB = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El reciente = Reciente
371         DM + Reciente DBOW
372
373         #Se cargan los datos de train de las canciones para conocer los IDs reales
374         DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")

```

```

375
376 #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
377 ALL = set(list(range(W2V.shape[0])))
378 REAL = set(list(set(list(map(int, DATA))))))
379 NOT_EXISTS = list(ALL.difference(REAL))
380
381 #####
382 # TENSORFLOW
383 #####
384
385 # -----
386 # Grid Search + Early Stopping
387 # -----
388 '''
389 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
390 rates = [0.000001, 0.01, 0.0001]
391
392 for r in rates:
393     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
394
395     print("Learning Rate: "+str(r))
396     print("-" * 50)
397     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
398     print("-"*50)
399
400 exit()
401 '''
402 # -----
403 # Train
404 # -----
405 '''
406 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
407
408 #Para esta fase, se une TRAIN y DEV
409 X = np.concatenate((X,DX),axis=0)
410 Y = np.concatenate((Y,DY),axis=0)
411
412 DX=[]
413 DY=[]
414
415 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 401, "seed": 100}
416 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
417
418 exit()
419 '''
420 # -----
421 # Test
422 # -----
423
424
425 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
426
427 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 401, "seed": 100}
428 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```

## 4.8. TEST 08

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from scipy.stats import linregress
13 from collections import Counter
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data
20             ="datos/Append/ARRAY_ARRAYS_REC.pkl",pastData="datos/Append/ARRAY_ARRAYS_TRAIN_COMPLETE.pkl",shuffl
21             seed=100, train=.7):
22
23     DATA = np.load(data)
24     PAST_DATA = np.load(pastData)
25     REAL_SET = set(REAL_SONG_IDS)
26
27     random.seed(seed)
28
29     #DATA = DATA[:20];print("ELIMINAR FILTRO")
30
31     #User, cancion, +1 o 0
32     TRAIN = np.empty([0, 3], dtype=np.int32)
33     DEV = np.empty([0, 3], dtype=np.int32)
34     TEST = np.empty([0, 3], dtype=np.int32)
35
36     for u in range(len(DATA)):
37
38         # Eliminar canciones escuchadas que no tienen embedding
39         USER_PLAYS = DATA[u]
40         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
41         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
42
43         #Separar en train dev test
44         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
45         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
46
47         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
48         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
49         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
50
51         # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE y en TRAIN De
52         ahora
53         USED_IDS = set(list(map(int, PAST_DATA[u])))
54         USED_IDS.update(TRAIN_DATA)
55         TRAIN_NOT_USED = list(REAL_SET.difference(USED_IDS))

```



```

53
54 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV
55 USED_IDS.update(DEV_DATA)
56 DEV_NOT_USED = list(REAL_SET.difference(USED_IDS))
57
58 # Obtener IDS de canciones que no han sido escuchadas en TRAIN_COMPLETE , en TRAIN y
    DEV y TEST
59 USED_IDS.update(TEST_DATA)
60 TEST_NOT_USED = list(REAL_SET.difference(USED_IDS))
61
62 # Crear casos de train
63 TRAIN_RES = np.empty([len(TRAIN_DATA) * 2, 3], dtype=np.int32)
64 TRAIN_RES[:, 0] = u
65 TRAIN_RES[:len(TRAIN_DATA), 1] = TRAIN_DATA
66 TRAIN_RES[:len(TRAIN_DATA), 2] = 1
67
68 TRAIN_RES[len(TRAIN_DATA):, 1] = random.sample(TRAIN_NOT_USED, len(TRAIN_DATA))
69 TRAIN_RES[len(TRAIN_DATA):, 2] = 0
70
71 TRAIN = np.concatenate([TRAIN, TRAIN_RES])
72
73 # Crear casos de dev
74 DEV_RES = np.empty([len(DEV_DATA) * 2, 3], dtype=np.int32)
75 DEV_RES[:, 0] = u
76 DEV_RES[:len(DEV_DATA), 1] = DEV_DATA
77 DEV_RES[:len(DEV_DATA), 2] = 1
78
79 DEV_RES[len(DEV_DATA):, 1] = random.sample(DEV_NOT_USED, len(DEV_DATA))
80 DEV_RES[len(DEV_DATA):, 2] = 0
81
82 DEV = np.concatenate([DEV, DEV_RES])
83
84 # Crear casos de test
85 TEST_RES = np.empty([len(TEST_DATA) * 2, 3], dtype=np.int32)
86 TEST_RES[:, 0] = u
87 TEST_RES[:len(TEST_DATA), 1] = TEST_DATA
88 TEST_RES[:len(TEST_DATA), 2] = 1
89
90 TEST_RES[len(TEST_DATA):, 1] = random.sample(TEST_NOT_USED, len(TEST_DATA))
91 TEST_RES[len(TEST_DATA):, 2] = 0
92
93 TEST = np.concatenate([TEST, TEST_RES])
94
95 #Mezclar
96 if(shuffle):
97     np.random.seed(seed)
98     np.random.shuffle(TRAIN)
99     np.random.shuffle(DEV)
100    np.random.shuffle(TEST)
101
102
103    return
104    TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
105    def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
106

```

```

107 # Si se ha aadi una iteracin, introducir los datos restantes
108 if (step == ITRS - 1 and RES_ITMS > 0):
109     batch_inputs = X[(ITRS - 1) * batch_size:]
110     batch_context = Y[(ITRS - 1) * batch_size:]
111 # Si no, es un batch normal
112 else:
113     st = step * batch_size
114     nd = (step + 1) * batch_size
115
116     batch_inputs = X[st:nd]
117     batch_context = Y[st:nd]
118 return(batch_inputs, batch_context)
119
120 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
121
122     def getSlope(data):
123         r = linregress(range(len(data)), data)
124         return r.slope
125
126     batch_size = config['batch_size']
127     learning_rate = config['learning_rate']
128     epochs = config['epochs']
129
130     user_size = user_size
131     song_size = W2V.shape[0]
132
133     user_emb_size = 128
134     user_concat_size = user_emb_size * 2
135     song_emb_size = 64
136     hidden_size = 32
137
138     num_sampled = 64 # Ejemplos negativos (Para el NSE)
139
140     model_name = 'model'
141     model_path = 'models/' + TEST_NAME + '/model'
142     complete_path = model_path + "/" + model_name
143
144     # Se almacena la loss de dev de las 5 ultimas epochs
145     slope_size = 100
146     train_hist=[]
147     dev_hist=[]
148
149     # Creacin del grafo de TF.
150     graph = tf.Graph()
151
152     with graph.as_default():
153
154         if(config["seed"]!=1):
155             tf.set_random_seed(config["seed"])
156
157         # nmero global de iteraciones
158         global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
159
160         # Datos de entrada
161
162         # Array del tamao del batch con las X
163         train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")

```

```

163
164 # Array del tamaño del batch con las Y asociadas
165 train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
166
167 # Embeddings


---


168
169 D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
170
171 T0 = tf.Variable(tf.truncated_normal([user_concat_size, user_emb_size], mean=0.0,
172     stddev=1.0 / math.sqrt(hidden_size)), name="T0")
173 B0 = tf.Variable(tf.zeros([user_emb_size]), name="B0")
174
175 O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
176     / math.sqrt(song_emb_size)), name="O2")
177 B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
178
179 T1 = tf.Variable(tf.truncated_normal([user_emb_size, hidden_size], mean=0.0,
180     stddev=1.0 / math.sqrt(hidden_size)), name="T1")
181 B1 = tf.Variable(tf.zeros([hidden_size]), name="B1")
182
183 T2 = tf.Variable(tf.truncated_normal([song_emb_size, hidden_size], mean=0.0,
184     stddev=1.0 / math.sqrt(hidden_size)), name="T2")
185 B2 = tf.Variable(tf.zeros([hidden_size]), name="B2")
186
187 # Operaciones


---


188
189 # Embedding de documento
190 ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
191
192 #Transformar concat documento de 256 a 128
193 ed = tf.matmul(ed, T0)+B0
194
195 # Embedding de canción
196 ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
197 ec = ec + B02
198
199 # Transformar a 32 documento
200 hd = tf.matmul(ed, T1)+B1
201
202 # Transformar a 32 canción
203 hw = tf.matmul(ec, T2)+B2
204
205 # Obtener el producto escalar
206 dot_prod = tf.reduce_sum(tf.multiply(hd, hw), 1, keep_dims=True)
207
208 # aClculo de LOSS y optimizador—


---


209
210 # Obtener la loss
211 softplus_batch = tf.nn.softplus((1 - 2 * train_labels) * dot_prod)
212 loss_softplus = tf.reduce_mean(softplus_batch)
213
214 # Minimizar la loss
215 train_step =
216     tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss=loss_softplus, global_step

```

```

212
213     # Obtener la probabilidad
214     prob = tf.sigmoid(dot_prod)
215
216     # Inicializar variables
217     init = tf.global_variables_initializer()
218
219     # Crear objeto encargado de almacenar la red
220     saver = tf.train.Saver(max_to_keep=1)
221
222     config = tf.ConfigProto()
223     config.gpu_options.allow_growth = True
224     with tf.Session(graph=graph, config=config) as session:
225     #with tf.Session(graph=graph) as session:
226
227         # We must initialize all variables before we use them.
228         init.run()
229
230     if(save):
231         # Obtener el checkpoint
232         ckpt = tf.train.get_checkpoint_state(model_path)
233
234         # Si existe el checkpoint restaurar
235         if ckpt and ckpt.model_checkpoint_path:
236             saver.restore(session, ckpt.model_checkpoint_path)
237
238     #Si se desea entrenar
239     if(train):
240         #-----
241         # TRAIN
242         #-----
243
244         ITRS= int(len(X)/batch_size)
245         RES_ITMS = len(X) - (ITRS*batch_size)
246         if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
247
248         for e in range(epochs):
249             TRAIN_LOSS = 0.0
250
251             print("Epoch "+str(e))
252
253             for step in range(ITRS):
254
255                 batch_inputs, batch_context =
256                     getBatchData(X,Y, batch_size, step, ITRS, RES_ITMS)
257
258                 feed_dict = {train.dataset: batch_inputs, train.labels: batch_context}
259
260                 ., loss_val, softplus, gs = session.run([train_step,
261                     loss_softplus, softplus_batch, global_step], feed_dict=feed_dict)
262
263                 TRAIN_LOSS += np.sum(softplus[:, 0])
264
265             if (save): saver.save(session, complete_path, global_step=global_step)
266
267         #-----
268         # DEV

```

```

267 #
268
269 if(len(DX)>0):
270     ITRS_TEST = int(len(DX) / batch_size)
271     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
272     TEST_LOSS = 0.0
273     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es ÷divisin exacta, se
274         añade una ÷iteracin
275
276     for step in range(ITRS_TEST):
277
278         batch_inputs, batch_context =
279             getBatchData(DX,DY, batch_size, step, ITRS_TEST, RES_ITMS_TEST)
280
281         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
282         loss_val, softplus = session.run([loss_softplus, softplus_batch],
283             feed_dict=feed_dict)
284
285         TEST_LOSS+=np.sum(softplus[:,0])
286
287         dev_hist.append(TEST_LOSS/(len(DX)*1.0))
288         train_hist.append(TRAIN_LOSS/(len(X)*1.0))
289
290     if (len(dev_hist) >= slope_size):
291         # Calcular la pendiente
292         print(getSlope(dev_hist[-slope_size:]))
293         print(dev_hist[-slope_size:])
294
295         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
296
297             def printArray(data, data2):
298                 for i in range(len(data)):
299                     a=data[i]
300                     b=data2[i]
301                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
302
303             print("-" * 50)
304             print(e)
305             printArray(dev_hist, train_hist)
306             return
307
308 else:
309     #
310
311 # TEST FINAL
312 #
313
314
315 TRUE_POSITIVE = 0
316 FALSE_NEGATIVE = 0
317 FALSE_POSITIVE = 0
318 TRUE_NEGATIVE = 0
319
320 ITRS_TEST = int(len(TX) / batch_size)
321 RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
322 TEST_LOSS = 0.0
323 # Si no es ÷divisin exacta, se añade una ÷iteracin

```

```

319     if (RES.ITMS.TEST > 0): ITRS_TEST += 1
320
321     for step in range(ITRS_TEST):
322         # Si se reañadi una iteración, introducir los datos restantes
323         if (step == ITRS_TEST - 1 and RES.ITMS.TEST > 0):
324             batch_inputs = TX[(ITRS_TEST - 1) * batch_size:]
325             batch_context = TY[(ITRS_TEST - 1) * batch_size:]
326         # Si no, es un batch normal
327         else:
328             st = step * batch_size
329             nd = (step + 1) * batch_size
330
331             batch_inputs = TX[st:nd]
332             batch_context = TY[st:nd]
333
334             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
335             loss_val, loss_batch, p = session.run([loss_softplus, softplus_batch, prob],
336                 feed_dict=feed_dict)
337
338             TEST_LOSS += np.sum(loss_batch)
339
340             # aClculo de TP/TN/FN/FP para cada batch
341
342             for i in range(len(p)):
343                 real = batch_context[i, 0]
344                 pred = p[i, 0]
345
346                 if (real == 1) and (pred > 0.5): TRUE_POSITIVE += 1
347                 if (real == 1) and (pred <= 0.5): FALSE_NEGATIVE += 1
348                 if (real == 0) and (pred > 0.5): FALSE_POSITIVE += 1
349                 if (real == 0) and (pred <= 0.5): TRUE_NEGATIVE += 1
350
351             PRECISION = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_POSITIVE * 1.0)
352             RECALL = TRUE_POSITIVE / (TRUE_POSITIVE + FALSE_NEGATIVE * 1.0)
353
354             F1 = 2.0 * ((PRECISION * RECALL) / (PRECISION + RECALL))
355
356             print("-" * 50)
357             print("LOSS:\t" + str(TEST_LOSS / (len(TX) * 1.0)))
358             print("PRECISION:\t" + str(PRECISION))
359             print("RECALL:\t" + str(RECALL))
360             print("F1:\t" + str(F1))
361             print("-" * 50)
362             print("")
363
364             #-----
365             #####
366             # Llamadas
367             #####
368             TEST_NAME="test_081"
369
370             os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
371
372             #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
373             D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
374             D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')

```

```

375
376 #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
377 D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
378 D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
379
380 W2V = np.load('embeddings/w2v/EMB_MATRIX')
381
382 #Se concatenan los embeddings de los usuarios para generar el perfil
383 D2V_CONS = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
    Consolidado DM + Consolidado DBOW
384 D2V_PRES = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El presente = Presente
    DM + Presente DBOW
385
386 D2V_EMB = np.concatenate((D2V_CONS, D2V_PRES), axis=1) # El perfil completo CONS + PRES
387
388 #Se cargan los datos de train de las canciones para conocer los IDs reales
389 DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
390
391 #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
392 ALL = set(list(range(W2V.shape[0])))
393 REAL = set(list(set(list(map(int, DATA)))))
394 NOT_EXISTS = list(ALL.difference(REAL))
395
396 #####
397 # TENSORFLOW
398 #####
399
400 # -----
401 # Grid Search + Early Stopping
402 # -----
403 '''
404 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
405 rates = [0.000001, 1, 0.01, 0.0001]
406
407 for r in rates:
408     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
409
410     print("Learning Rate: "+str(r))
411     print("-" * 50)
412     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
413     print("-"*50)
414
415 exit()
416 '''
417 # -----
418 # Train
419 # -----
420 '''
421 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
422
423 #Para esta fase, se une TRAIN y DEV
424 X = np.concatenate((X, DX), axis=0)
425 Y = np.concatenate((Y, DY), axis=0)
426
427 DX=[]
428 DY=[]
429

```

```
430 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 435, "seed": 100}
431 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
432
433 exit()
434 '''
435 #
436 # Test
437 #
438
439
440 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
441
442 config = {"batch_size": 1024, "learning_rate": 0.000001, "epochs": 435, "seed": 100}
443 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)
```



## 4.9. TEST 09

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data ="datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
20            seed=100, train=.7):
21
22     DATA = np.load(data)
23
24     #DATA = DATA[:20];print("ELIMINAR FILTRO")
25
26     #User, cancion 1, cancion 2
27     TRAIN = np.empty([0, 3], dtype=np.int32)
28     DEV = np.empty([0, 3], dtype=np.int32)
29     TEST = np.empty([0, 3], dtype=np.int32)
30
31     for u in range(len(DATA)):
32
33         # Eliminar canciones escuchadas que no tienen embedding
34         USER_PLAYS = DATA[u]
35         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
36         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
37
38         #Separar en train dev test
39         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
40         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
41
42         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
43         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
44         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
45
46         # Crear ejemplos
47         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
48         USR_TRAIN[:,0] = u
49         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-1]
50         USR_TRAIN[:,2] = TRAIN_DATA[1:]
51
52         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
53         USR_DEV[:, 0] = u
54         USR_DEV[:, 1] = DEV_DATA[:len(DEV_DATA) - 1]
55         USR_DEV[:, 2] = DEV_DATA[1:]

```

```

55
56     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
57     USR_TEST[:, 0] = u
58     USR_TEST[:, 1] = TEST_DATA[:len(TEST_DATA) - 1]
59     USR_TEST[:, 2] = TEST_DATA[1:]
60
61
62     # Añadir al global
63     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
64     DEV = np.concatenate([DEV,USR_DEV])
65     TEST = np.concatenate([TEST,USR_TEST])
66
67
68     #Mezclar
69     if(shuffle):
70         np.random.seed(seed)
71         np.random.shuffle(TRAIN)
72         np.random.shuffle(DEV)
73         np.random.shuffle(TEST)
74
75
76     return
77         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
78
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se no aadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return(batch_inputs, batch_context)
93
94 def trainNet(X,Y, DX,DY, TX,TY, config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     song_emb_size = 64
109     hidden_size = 32
110
111     profile_size = user_emb_size + song_emb_size

```

```

111
112 num_sampled = 64 # Ejemplos negativos (Para el NSE)
113
114 model_name = 'model'
115 model_path = 'models/' + TEST_NAME + '/model'
116 complete_path = model_path + "/" + model_name
117
118 # Se almacena la loss de dev de las 5 ultimas epochs
119 slope_size = 100
120 train_hist = []
121 dev_hist = []
122
123 # Creacion del grafo de TF.
124 graph = tf.Graph()
125
126 with graph.as_default():
127
128     if(config["seed"]!=1):
129         tf.set_random_seed(config["seed"])
130
131     # Numero global de iteraciones
132     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
133
134     # Datos de entrada
135
136     # Array del tamaño del batch con las X
137     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
138
139     # Array del tamaño del batch con las Y asociadas
140     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
141
142     # Embeddings
143
144     O1 = tf.Variable(tf.truncated_normal([user_size, user_emb_size], mean=0.0, stddev=1.0
145         / math.sqrt(user_emb_size)), name="O1")
146     B01 = tf.Variable(tf.zeros([user_emb_size]), name="B01")
147
148     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
149
150     # Operaciones
151
152     # Embedding de documento
153     ed = tf.nn.embedding_lookup(O1, train_dataset[:, 0])
154     ed = ed + B01
155
156     # Embedding de cancion
157     ec = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
158
159     # Calculo de LOSS y optimizador
160
161     embed = []
162     embed.append(ed)
163     embed.append(ec)
164     embed = tf.concat(embed, 1)

```

```

163
164 # softmax weights, W and D vectors should be concatenated before applying softmax
165 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
166                                     math.sqrt(song_size)),
167                               name="weights")
168 # softmax biases
169 biases = tf.Variable(tf.zeros([song_size]), name="biases")
170
171 # aClculo de LOSS y optimizador—
172
173 nce_loss = tf.nn.nce_loss(weights=weights,
174                           biases=biases,
175                           labels=train_labels,
176                           inputs=embed,
177                           num_sampled=num_sampled,
178                           num_classes=song_size,
179                           name="nce_loss")
180
181 batch_loss = nce_loss
182 nce_loss = tf.reduce_mean(nce_loss)
183
184 logits = tf.matmul(embed, tf.transpose(weights))
185 logits = tf.nn.bias_add(logits, biases)
186 # logits = tf.nn.softmax(logits,axis=0)
187 top_v,top_i = tf.nn.top_k(logits,k=song_size)
188
189 # optimizer =
190     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
191 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
192                               global_step=global_step)
193
194 # Inicializar variables
195 init = tf.global_variables_initializer()
196
197 # Crear objeto encargado de almacenar la red
198 saver = tf.train.Saver(max_to_keep=1)
199
200 #config = tf.ConfigProto()
201 #config.gpu_options.allow_growth = True
202 #with tf.Session(graph=graph, config=config) as session:
203 with tf.Session(graph=graph) as session:
204
205     # We must initialize all variables before we use them.
206     init.run()
207
208     if(save):
209         # Obtener el checkpoint
210         ckpt = tf.train.get_checkpoint_state(model_path)
211
212         # Si existe el checkpoint restaurar
213         if ckpt and ckpt.model_checkpoint_path:
214             saver.restore(session, ckpt.model_checkpoint_path)
215
216     #Si se desea entrenar
217     if(train):
218         #

```

```

216 # TRAIN
217 #-----
218
219 ITRS= int(len(X)/batch_size)
220 RES_ITMS = len(X) - (ITRS*batch_size)
221 if(RES_ITMS>0):ITRS+=1 #Si no es divisible exacta, se añade una iteración
222
223 for e in range(epochs):
224     TRAIN_LOSS = 0.0
225
226     print("Epoch "+str(e))
227
228     for step in range(ITRS):
229
230         batch_inputs, batch_context =
231             getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
232
233         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
234
235         -, loss_val, batch_loss, gs = session.run([train_step,
236             nce_loss,batch_loss,global_step], feed_dict=feed_dict)
237
238         TRAIN_LOSS += np.sum(batch_loss)
239
240         #print('Epoch '+str(e)+', last batch loss: '+str(loss_val))
241
242         if (save): saver.save(session, complete_path, global_step=global_step)
243
244 #-----
245 # DEV
246 #-----
247
248 if(len(DX)>0):
249
250     ITRS_TEST = int(len(DX) / batch_size)
251     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
252     TEST_LOSS = 0.0
253     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisible exacta, se
254         añade una iteración
255
256     for step in range(ITRS_TEST):
257
258         batch_inputs, batch_context =
259             getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
260
261         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
262         loss_val,batch_loss = session.run([nce_loss,batch_loss],
263             feed_dict=feed_dict)
264
265         TEST_LOSS+=np.sum(batch_loss)
266
267         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
268         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
269
270     if (len(dev_hist) >= slope_size):
271         # Calcular la pendiente

```

```

268     print(getSlope(dev_hist[-slope_size:]))
269     print(dev_hist[-slope_size:])
270
271     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
272
273         def printArray(data, data2):
274             for i in range(len(data)):
275                 a=data[i]
276                 b=data2[i]
277                 print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
278
279         print("-" * 50)
280         print(e)
281         printArray(dev_hist, train_hist)
282         return
283 else:
284     #
285
286     # TEST FINAL
287     #
288
289     batch_size = 300
290
291     IN_TOP_5 = 0
292     IN_TOP_10 = 0
293     IN_TOP_20 = 0
294     IN_TOP_50 = 0
295     IN_TOP_100 = 0
296     IN_TOP_1000 = 0
297     real_positions = []
298
299     items = 0
300
301     ITRS_TEST = int(len(TX) / batch_size)
302     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
303     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se aade una
304     oiteracin
305
306     for step in range(ITRS_TEST):
307         # tss = time.time()
308
309         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
310             ITRS_TEST, RES_ITMS_TEST)
311         top_values, top_indx = session.run([top_v, top_i],
312             feed_dict={train.dataset: batch_inputs,
313                 train.labels: batch_context})
314
315         for j in range(len(batch_context)):
316             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
317             real_positions.append(real_position)
318
319             IN_TOP_5 += 1 if real_position <= 5 else 0
320             IN_TOP_10 += 1 if real_position <= 10 else 0
321             IN_TOP_20 += 1 if real_position <= 20 else 0
322             IN_TOP_50 += 1 if real_position <= 50 else 0
323             IN_TOP_100 += 1 if real_position <= 100 else 0
324             IN_TOP_1000 += 1 if real_position <= 1000 else 0

```

```

320
321     items += len(batch_inputs) * 1.0
322
323     PTF = IN_TOP_5 / items
324     PTT = IN_TOP_10 / items
325     PTTW = IN_TOP_20 / items
326     PTFT = IN_TOP_50 / items
327     PTCIEN = IN_TOP_100 / items
328     PTMIL = IN_TOP_1000 / items
329
330     print(np.mean(real_positions), np.std(real_positions),
331           np.median(real_positions))
332     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
333           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") +
334           "\t"
335           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
336
337     # print(time.time()-tss)
338
339     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
340
341     #-----
342     #####
343     # Llamadas
344     #####
345     TEST_NAME="test_022"
346
347     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
348
349     #Solamente se carga para conocer el unmero de canciones
350     W2V = np.load('embeddings/w2v/EMB_MATRIX')
351
352     #Se cargan los datos de train de las canciones para conocer los IDs reales
353     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
354
355     #Obtener una lista con los ID de canciones que no aestn el la matriz de embeddings
356     ALL = set(list(range(W2V.shape[0])))
357     REAL = set(list(set(list(map(int, DATA)))))
358
359     #####
360     # TENSORFLOW
361     #####
362     ' '
363     #-----
364     # Grid Search + Early Stopping
365     #-----
366
367     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
368     rates = [0.000001, 1, 0.01, 0.0001]
369
370     for r in rates:
371         config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
372
373         print("Learning Rate: "+str(r))
374         print("-" * 50)

```

```
375     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
376     print("-"*50)
377
378     exit()
379     '''
380     #-----
381     # Train
382     #-----
383     '''
384     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
385
386     #Para esta fase, se une TRAIN y DEV
387     X = np.concatenate((X,DX),axis=0)
388     Y = np.concatenate((Y,DY),axis=0)
389
390     DX=[]
391     DY=[]
392
393     config = {"batch_size": 1024, "learning_rate": 1, "epochs":172, "seed":100}
394     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
395
396     exit()
397     '''
398     #-----
399     # Test
400     #-----
401
402     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
403
404     config = {"batch_size": 1024, "learning_rate": 1, "epochs":172, "seed":100}
405     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)
```



## 4.10. TEST 10

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data ="datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
20            seed=100, train=.7):
21
22     DATA = np.load(data)
23
24     #DATA = DATA[:20];print("ELIMINAR FILTRO")
25
26     #User, cancion 1, cancion 2
27     TRAIN = np.empty([0, 3], dtype=np.int32)
28     DEV = np.empty([0, 3], dtype=np.int32)
29     TEST = np.empty([0, 3], dtype=np.int32)
30
31     for u in range(len(DATA)):
32
33         # Eliminar canciones escuchadas que no tienen embedding
34         USER_PLAYS = DATA[u]
35         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
36         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
37
38         #Separar en train dev test
39         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
40         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
41
42         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
43         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
44         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
45
46         # Crear ejemplos
47         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
48         USR_TRAIN[:,0] = u
49         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-1]
50         USR_TRAIN[:,2] = TRAIN_DATA[1:]
51
52         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
53         USR_DEV[:, 0] = u
54         USR_DEV[:, 1] = DEV_DATA[:len(DEV_DATA) - 1]
55         USR_DEV[:, 2] = DEV_DATA[1:]

```

```

55
56     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
57     USR_TEST[:, 0] = u
58     USR_TEST[:, 1] = TEST_DATA[:len(TEST_DATA) - 1]
59     USR_TEST[:, 2] = TEST_DATA[1:]
60
61
62     # Añadir al global
63     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
64     DEV = np.concatenate([DEV,USR_DEV])
65     TEST = np.concatenate([TEST,USR_TEST])
66
67
68     #Mezclar
69     if(shuffle):
70         np.random.seed(seed)
71         np.random.shuffle(TRAIN)
72         np.random.shuffle(DEV)
73         np.random.shuffle(TEST)
74
75
76     return
77         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
78
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se noaadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return(batch_inputs, batch_context)
93
94 def trainNet(X,Y, DX,DY, TX,TY,config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     song_emb_size = 64
109     hidden_size = 32
110
111     profile_size = user_emb_size + song_emb_size

```

```

111
112 num_sampled = 64 # Ejemplos negativos (Para el NSE)
113
114 model_name = 'model'
115 model_path = 'models/' + TEST_NAME + '/model'
116 complete_path = model_path + "/" + model_name
117
118 # Se almacena la loss de dev de las 5 ultimas epochs
119 slope_size = 100
120 train_hist = []
121 dev_hist = []
122
123 # Creacin del grafo de TF.
124 graph = tf.Graph()
125
126 with graph.as_default():
127
128     if(config["seed"]!=1):
129         tf.set_random_seed(config["seed"])
130
131     # nmero global de iteraciones
132     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
133
134     # Datos de entrada
135
136     # Array del tamaño del batch con las X
137     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
138
139     # Array del tamaño del batch con las Y asociadas
140     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
141
142     # Embeddings
143
144     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
145
146     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
147
148     # Operaciones
149
150     # Embedding de documento
151     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
152
153     # Embedding de cancion
154     ec = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
155
156     # Cálculo de LOSS y optimizador
157
158     embed = []
159     embed.append(ed)
160     embed.append(ec)
161     embed = tf.concat(embed, 1)
162
163     # softmax weights, W and D vectors should be concatenated before applying softmax
164     weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /

```

```

164         math.sqrt(song_size)),
165             name="weights")
166 # softmax biases
167 biases = tf.Variable(tf.zeros([song_size]), name="biases")
168
169 # aClculo de LOSS y optimizador—
170
171 nce_loss = tf.nn.nce_loss(weights=weights,
172                          biases=biases,
173                          labels=train.labels,
174                          inputs=embed,
175                          num_sampled=num_sampled,
176                          num_classes=song_size,
177                          name="nce_loss")
178
179 batch_loss = nce_loss
180 nce_loss = tf.reduce_mean(nce_loss)
181
182 logits = tf.matmul(embed, tf.transpose(weights))
183 logits = tf.nn.bias_add(logits, biases)
184 # logits = tf.nn.softmax(logits,axis=0)
185 top_v,top_i = tf.nn.top_k(logits,k=song_size)
186
187 # optimizer =
188     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
189 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
190                               global_step=global_step)
191
192 # Inicializar variables
193 init = tf.global_variables_initializer()
194
195 # Crear objeto encargado de almacenar la red
196 saver = tf.train.Saver(max_to_keep=1)
197
198 #config = tf.ConfigProto()
199 #config.gpu_options.allow_growth = True
200 #with tf.Session(graph=graph, config=config) as session:
201 with tf.Session(graph=graph) as session:
202
203     # We must initialize all variables before we use them.
204     init.run()
205
206     if(save):
207         # Obtener el checkpoint
208         ckpt = tf.train.get_checkpoint_state(model_path)
209
210         # Si existe el checkpoint restaurar
211         if ckpt and ckpt.model_checkpoint_path:
212             saver.restore(session, ckpt.model_checkpoint_path)
213
214     #Si se desea entrenar
215     if(train):
216         #-----
217         # TRAIN
218         #-----

```

```

217 ITRS= int(len(X)/batch_size)
218 RES_ITMS = len(X) - (ITRS*batch_size)
219 if(RES_ITMS>0):ITRS+=1 #Si no es odivisin exacta, se naade una oiteracin
220
221 for e in range(epochs):
222     TRAIN_LOSS = 0.0
223
224     print("Epoch "+str(e))
225
226     for step in range(ITRS):
227
228         batch_inputs, batch_context =
229             getBatchData(X,Y,batch_size ,step,ITRS,RES_ITMS)
230
231         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
232
233         _, loss_val, batch_lss, gs = session.run([train_step,
234             nce_loss,batch_loss,global_step], feed_dict=feed_dict)
235
236         TRAIN_LOSS += np.sum(batch_lss)
237
238         #print('Epoch '+str(e)+', last batch loss: '+str(loss_val))
239
240         if (save): saver.save(session, complete_path, global_step=global_step)
241
242         #-----
243         # DEV
244         #-----
245
246         if(len(DX)>0):
247
248             ITRS_TEST = int(len(DX) / batch_size)
249             RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
250             TEST_LOSS = 0.0
251             if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se
252                 naadin una oiteracin
253
254             for step in range(ITRS_TEST):
255
256                 batch_inputs, batch_context =
257                     getBatchData(DX,DY,batch_size ,step,ITRS_TEST,RES_ITMS_TEST)
258
259                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
260                 loss_val,batch_lss = session.run([nce_loss,batch_loss],
261                     feed_dict=feed_dict)
262
263                 TEST_LOSS+=np.sum(batch_lss)
264
265                 dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
266                 train_hist.append(TRAIN_LOSS / (len(X) *1.0))
267
268             if (len(dev_hist) >= slope_size):
269                 # Calcular la pendiente
270                 print(getSlope(dev_hist[-slope_size:]))
271                 print(dev_hist[-slope_size:])

```

```

269         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
270
271             def printArray(data, data2):
272                 for i in range(len(data)):
273                     a=data[i]
274                     b=data2[i]
275                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
276
277             print("-" * 50)
278             print(e)
279             printArray(dev_hist, train_hist)
280             return
281 else:
282     #
283
284     # TEST FINAL
285     #
286
287     batch_size = 300
288
289     IN_TOP_5 = 0
290     IN_TOP_10 = 0
291     IN_TOP_20 = 0
292     IN_TOP_50 = 0
293     IN_TOP_100 = 0
294     IN_TOP_1000 = 0
295     real_positions = []
296
297     items = 0
298
299     ITRS_TEST = int(len(TX) / batch_size)
300     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
301     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
302     iteracin
303
304     for step in range(ITRS_TEST):
305         #tss = time.time()
306
307         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
308             ITRS_TEST, RES_ITMS_TEST)
309         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
310             batch_inputs, train_labels: batch_context})
311
312         for j in range(len(batch_context)):
313             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
314             real_positions.append(real_position)
315
316             IN_TOP_5 += 1 if real_position <= 5 else 0
317             IN_TOP_10 += 1 if real_position <= 10 else 0
318             IN_TOP_20 += 1 if real_position <= 20 else 0
319             IN_TOP_50 += 1 if real_position <= 50 else 0
320             IN_TOP_100 += 1 if real_position <= 100 else 0
321             IN_TOP_1000 += 1 if real_position <= 1000 else 0
322
323         items += len(batch_inputs)*1.0
324
325     PTF = IN_TOP_5 / items

```

```

321     PTT = IN_TOP_10 / items
322     PTTW = IN_TOP_20 / items
323     PTFT = IN_TOP_50 / items
324     PTCIEN = IN_TOP_100 / items
325     PTMIL = IN_TOP_1000 / items
326
327     print(np.mean(real_positions), np.std(real_positions),
           np.median(real_positions))
328     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
329
330
331
332     #print(time.time()-tss)
333
334     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
335
336     #-----
337
338     #####
339     # Llamadas
340     #####
341
342     TEST_NAME="test_032"
343
344     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
345
346     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
347     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
348     D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
349     W2V = np.load('embeddings/w2v/EMB_MATRIX')
350
351     #Se concatenan los embeddings de los usuarios para generar el perfil
352     D2V_EMB = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
           Consolidado DM + Consolidado DBOW
353
354     #Se cargan los datos de train de las canciones para conocer los IDs reales
355     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
356
357     #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
358     ALL = set(list(range(W2V.shape[0])))
359     REAL = set(list(set(list(map(int, DATA)))))
360     NOT_EXISTS = list(ALL.difference(REAL))
361
362     #####
363     # TENSORFLOW
364     #####
365     ' '
366     #-----
367     # Grid Search + Early Stopping
368     #-----
369
370     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
371     rates = [0.000001, 1, 0.01, 0.0001]
372
373     for r in rates:
374         config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
375

```

```

376     print("Learning Rate: "+str(r))
377     print("-" * 50)
378     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
379     print("-"*50)
380
381     exit()
382     '''
383     #-----
384     # Train
385     #-----
386     '''
387     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
388
389     #Para esta fase, se une TRAIN y DEV
390     X = np.concatenate((X,DX),axis=0)
391     Y = np.concatenate((Y,DY),axis=0)
392
393     DX=[]
394     DY=[]
395
396     config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
397     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
398
399     exit()
400     '''
401     #-----
402     # Test
403     #-----
404
405     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
406
407     config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
408     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.11. TEST 11

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data ="datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
20            seed=100, train=.7):
21
22     DATA = np.load(data)
23
24     #DATA = DATA[:20];print("ELIMINAR FILTRO")
25
26     #User, cancion 1, cancion 2
27     TRAIN = np.empty([0, 3], dtype=np.int32)
28     DEV = np.empty([0, 3], dtype=np.int32)
29     TEST = np.empty([0, 3], dtype=np.int32)
30
31     for u in range(len(DATA)):
32
33         # Eliminar canciones escuchadas que no tienen embedding
34         USER_PLAYS = DATA[u]
35         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
36         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
37
38         #Separar en train dev test
39         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
40         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
41
42         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
43         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
44         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
45
46         # Crear ejemplos
47         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
48         USR_TRAIN[:,0] = u
49         USR_TRAIN[:,1] = TRAIN_DATA[:,len(TRAIN_DATA)-1]
50         USR_TRAIN[:,2] = TRAIN_DATA[1:]
51
52         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
53         USR_DEV[:, 0] = u
54         USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA) - 1]
55         USR_DEV[:, 2] = DEV_DATA[1:]

```

```

55
56     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
57     USR_TEST[:, 0] = u
58     USR_TEST[:, 1] = TEST_DATA[:len(TEST_DATA) - 1]
59     USR_TEST[:, 2] = TEST_DATA[1:]
60
61
62     # Añadir al global
63     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
64     DEV = np.concatenate([DEV,USR_DEV])
65     TEST = np.concatenate([TEST,USR_TEST])
66
67
68     #Mezclar
69     if(shuffle):
70         np.random.seed(seed)
71         np.random.shuffle(TRAIN)
72         np.random.shuffle(DEV)
73         np.random.shuffle(TEST)
74
75
76     return
77         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
78
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se no aadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return(batch_inputs, batch_context)
93
94 def trainNet(X,Y, DX,DY, TX,TY, config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     song_emb_size = 64
109     hidden_size = 32
110
111     profile_size = user_emb_size + song_emb_size

```

```

111
112 num_sampled = 64 # Ejemplos negativos (Para el NSE)
113
114 model_name = 'model'
115 model_path = 'models/' + TEST_NAME + '/model'
116 complete_path = model_path + "/" + model_name
117
118 # Se almacena la loss de dev de las 5 ultimas epochs
119 slope_size = 100
120 train_hist = []
121 dev_hist = []
122
123 # Creacin del grafo de TF.
124 graph = tf.Graph()
125
126 with graph.as_default():
127
128     if(config["seed"]!=1):
129         tf.set_random_seed(config["seed"])
130
131     # nmero global de iteraciones
132     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
133
134     # Datos de entrada
135
136     # Array del tamaño del batch con las X
137     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
138
139     # Array del tamaño del batch con las Y asociadas
140     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
141
142     # Embeddings
143
144     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
145
146     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
147
148     # Operaciones
149
150     # Embedding de documento
151     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
152
153     # Embedding de cancion
154     ec = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
155
156     # Cálculo de LOSS y optimizador
157
158     embed = []
159     embed.append(ed)
160     embed.append(ec)
161     embed = tf.concat(embed, 1)
162
163     # softmax weights, W and D vectors should be concatenated before applying softmax
164     weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /

```

```

164         math.sqrt(song_size)),
165             name="weights")
166 # softmax biases
167 biases = tf.Variable(tf.zeros([song_size]), name="biases")
168
169 # aClculo de LOSS y optimizador—
170
171 nce_loss = tf.nn.nce_loss(weights=weights,
172                           biases=biases,
173                           labels=train.labels,
174                           inputs=embed,
175                           num_sampled=num_sampled,
176                           num_classes=song_size,
177                           name="nce_loss")
178
179 batch_loss = nce_loss
180 nce_loss = tf.reduce_mean(nce_loss)
181
182 logits = tf.matmul(embed, tf.transpose(weights))
183 logits = tf.nn.bias_add(logits, biases)
184 # logits = tf.nn.softmax(logits,axis=0)
185 top_v,top_i = tf.nn.top_k(logits,k=song_size)
186
187 # optimizer =
188     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
189 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
190     global_step=global_step)
191
192 # Inicializar variables
193 init = tf.global_variables_initializer()
194
195 # Crear objeto encargado de almacenar la red
196 saver = tf.train.Saver(max_to_keep=1)
197
198 #config = tf.ConfigProto()
199 #config.gpu_options.allow_growth = True
200 #with tf.Session(graph=graph, config=config) as session:
201 with tf.Session(graph=graph) as session:
202
203     # We must initialize all variables before we use them.
204     init.run()
205
206     if(save):
207         # Obtener el checkpoint
208         ckpt = tf.train.get_checkpoint_state(model_path)
209
210         # Si existe el checkpoint restaurar
211         if ckpt and ckpt.model_checkpoint_path:
212             saver.restore(session, ckpt.model_checkpoint_path)
213
214     #Si se desea entrenar
215     if(train):
216         #-----
217         # TRAIN
218         #-----

```

```

217 ITRS= int(len(X)/batch_size)
218 RES_ITMS = len(X) - (ITRS*batch_size)
219 if(RES_ITMS>0):ITRS+=1 #Si no es divisible exacta, se añade una iteración
220
221 for e in range(epochs):
222     TRAIN_LOSS = 0.0
223
224     print("Epoch "+str(e))
225
226     for step in range(ITRS):
227
228         batch_inputs, batch_context =
229             getBatchData(X,Y,batch_size ,step,ITRS,RES_ITMS)
230
231         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
232
233         _, loss_val, batch_loss, gs = session.run([train_step,
234             nce_loss,batch_loss,global_step], feed_dict=feed_dict)
235
236         TRAIN_LOSS += np.sum(batch_loss)
237
238         #print('Epoch '+str(e)+', last batch loss: '+str(loss_val))
239
240         if (save): saver.save(session, complete_path, global_step=global_step)
241
242         #-----
243         # DEV
244         #-----
245
246         if(len(DX)>0):
247
248             ITRS_TEST = int(len(DX) / batch_size)
249             RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
250             TEST_LOSS = 0.0
251             if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisible exacta, se
252                 añade una iteración
253
254             for step in range(ITRS_TEST):
255
256                 batch_inputs, batch_context =
257                     getBatchData(DX,DY,batch_size ,step,ITRS_TEST,RES_ITMS_TEST)
258
259                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
260                 loss_val,batch_loss = session.run([nce_loss,batch_loss],
261                     feed_dict=feed_dict)
262
263                 TEST_LOSS+=np.sum(batch_loss)
264
265                 dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
266                 train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
267
268             if (len(dev_hist) >= slope_size):
269                 # Calcular la pendiente
270                 print(getSlope(dev_hist[-slope_size:]))
271                 print(dev_hist[-slope_size:])

```

```

269         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
270
271             def printArray(data, data2):
272                 for i in range(len(data)):
273                     a=data[i]
274                     b=data2[i]
275                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
276
277             print("-" * 50)
278             print(e)
279             printArray(dev_hist, train_hist)
280             return
281 else:
282     #
283
284     # TEST FINAL
285     #
286
287     batch_size = 300
288
289     IN_TOP_5 = 0
290     IN_TOP_10 = 0
291     IN_TOP_20 = 0
292     IN_TOP_50 = 0
293     IN_TOP_100 = 0
294     IN_TOP_1000 = 0
295     real_positions = []
296
297     items = 0
298
299     ITRS_TEST = int(len(TX) / batch_size)
300     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
301     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
302     iteracin
303
304     for step in range(ITRS_TEST):
305         #tss = time.time()
306
307         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
308             ITRS_TEST, RES_ITMS_TEST)
309         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
310             batch_inputs, train_labels: batch_context})
311
312         for j in range(len(batch_context)):
313             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
314             real_positions.append(real_position)
315
316             IN_TOP_5 += 1 if real_position <= 5 else 0
317             IN_TOP_10 += 1 if real_position <= 10 else 0
318             IN_TOP_20 += 1 if real_position <= 20 else 0
319             IN_TOP_50 += 1 if real_position <= 50 else 0
320             IN_TOP_100 += 1 if real_position <= 100 else 0
321             IN_TOP_1000 += 1 if real_position <= 1000 else 0
322
323         items += len(batch_inputs)*1.0
324
325     PTF = IN_TOP_5 / items

```

```

321     PTT = IN_TOP_10 / items
322     PTTW = IN_TOP_20 / items
323     PTFT = IN_TOP_50 / items
324     PTCIEN = IN_TOP_100 / items
325     PTMIL = IN_TOP_1000 / items
326
327     print(np.mean(real_positions), np.std(real_positions),
           np.median(real_positions))
328     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
329
330
331
332     #print(time.time()-tss)
333
334     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
335
336     #-----
337
338     #####
339     # Llamadas
340     #####
341
342     TEST_NAME="test_042"
343
344     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
345
346     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
347     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
348     D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
349     W2V = np.load('embeddings/w2v/EMB_MATRIX')
350
351     #Se concatenan los embeddings de los usuarios para generar el perfil
352     D2V_EMB = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El reciente = Reciente
           DM + Reciente DBOW
353
354     #Se cargan los datos de train de las canciones para conocer los IDs reales
355     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
356
357     #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
358     ALL = set(list(range(W2V.shape[0])))
359     REAL = set(list(set(list(map(int, DATA)))))
360     NOT_EXISTS = list(ALL.difference(REAL))
361
362     #####
363     # TENSORFLOW
364     #####
365
366     #-----
367     # Grid Search + Early Stopping
368     #-----
369     '''
370     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
371     rates = [0.000001, 1, 0.01, 0.0001]
372
373     for r in rates:
374         config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
375

```

```

376     print("Learning Rate: "+str(r))
377     print("-" * 50)
378     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
379     print("-"*50)
380
381     exit()
382     '''
383     #-----
384     # Train
385     #-----
386     '''
387     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
388
389     #Para esta fase, se une TRAIN y DEV
390     X = np.concatenate((X,DX),axis=0)
391     Y = np.concatenate((Y,DY),axis=0)
392
393     DX=[]
394     DY=[]
395
396     config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
397     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
398
399     exit()
400     '''
401     #-----
402     # Test
403     #-----
404
405     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
406
407     config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
408     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.12. TEST 12

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data ="datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
20            seed=100, train=.7):
21
22     DATA = np.load(data)
23
24     #DATA = DATA[:20];print("ELIMINAR FILTRO")
25
26     #User, cancion 1, cancion 2
27     TRAIN = np.empty([0, 3], dtype=np.int32)
28     DEV = np.empty([0, 3], dtype=np.int32)
29     TEST = np.empty([0, 3], dtype=np.int32)
30
31     for u in range(len(DATA)):
32
33         # Eliminar canciones escuchadas que no tienen embedding
34         USER_PLAYS = DATA[u]
35         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
36         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
37
38         #Separar en train dev test
39         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
40         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
41
42         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
43         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
44         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
45
46         # Crear ejemplos
47         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
48         USR_TRAIN[:,0] = u
49         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-1]
50         USR_TRAIN[:,2] = TRAIN_DATA[1:]
51
52         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
53         USR_DEV[:, 0] = u
54         USR_DEV[:, 1] = DEV_DATA[:len(DEV_DATA) - 1]
55         USR_DEV[:, 2] = DEV_DATA[1:]

```

```

55
56     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
57     USR_TEST[:, 0] = u
58     USR_TEST[:, 1] = TEST_DATA[:, len(TEST_DATA) - 1]
59     USR_TEST[:, 2] = TEST_DATA[1:]
60
61
62     # Añadir al global
63     TRAIN = np.concatenate([TRAIN, USR_TRAIN])
64     DEV = np.concatenate([DEV, USR_DEV])
65     TEST = np.concatenate([TEST, USR_TEST])
66
67
68     # Mezclar
69     if(shuffle):
70         np.random.seed(seed)
71         np.random.shuffle(TRAIN)
72         np.random.shuffle(DEV)
73         np.random.shuffle(TEST)
74
75
76     return
77         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
78
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se no aadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return(batch_inputs, batch_context)
93
94 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     user_concat_size = user_emb_size * 2
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
146
147     T0 = tf.Variable(tf.truncated_normal([user_concat_size, user_emb_size], mean=0.0,
148         stddev=1.0 / math.sqrt(hidden_size)), name="T0")
149     B0 = tf.Variable(tf.zeros([user_emb_size]), name="B0")
150
151     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
152
153     # Operaciones
154
155     # Embedding de documento
156     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
157
158     #Transformar concat documento de 256 a 128
159     ed = tf.matmul(ed,T0)+B0
160
161     # Embedding de cancion
162     ec = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
163
164     # Cálculo de LOSS y optimizador—

```

```

163
164 embed = []
165 embed.append(ed)
166 embed.append(ec)
167 embed = tf.concat(embed, 1)
168
169 # softmax weights, W and D vectors should be concatenated before applying softmax
170 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
171                                     math.sqrt(song_size)),
172                               name="weights")
173 # softmax biases
174 biases = tf.Variable(tf.zeros([song_size]), name="biases")
175
176 # aClculo de LOSS y optimizador—
177
178 nce_loss = tf.nn.nce_loss(weights=weights,
179                           biases=biases,
180                           labels=train_labels,
181                           inputs=embed,
182                           num_sampled=num_sampled,
183                           num_classes=song_size,
184                           name="nce_loss")
185
186 batch_loss = nce_loss
187 nce_loss = tf.reduce_mean(nce_loss)
188
189 logits = tf.matmul(embed, tf.transpose(weights))
190 logits = tf.nn.bias_add(logits, biases)
191 # logits = tf.nn.softmax(logits,axis=0)
192 top_v,top_i = tf.nn.top_k(logits,k=song_size)
193
194 # optimizer =
195     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
196 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
197                               global_step=global_step)
198
199 # Inicializar variables
200 init = tf.global_variables_initializer()
201
202 # Crear objeto encargado de almacenar la red
203 saver = tf.train.Saver(max_to_keep=1)
204
205 #config = tf.ConfigProto()
206 #config.gpu_options.allow_growth = True
207 #with tf.Session(graph=graph, config=config) as session:
208 with tf.Session(graph=graph) as session:
209
210     # We must initialize all variables before we use them.
211     init.run()
212
213 if(save):
214     # Obtener el checkpoint
215     ckpt = tf.train.get_checkpoint_state(model_path)
216
217     # Si existe el checkpoint restaurar
218     if ckpt and ckpt.model_checkpoint_path:

```

```

216         saver.restore(session, ckpt.model_checkpoint_path)
217
218     #Si se desea entrenar
219     if(train):
220         #-----
221         # TRAIN
222         #-----
223
224         ITRS= int(len(X)/batch_size)
225         RES_ITMS = len(X) - (ITRS*batch_size)
226         if(RES_ITMS>0):ITRS+=1 #Si no es òdivisin exacta, se aade una òiteracin
227
228         for e in range(epochs):
229             TRAIN_LOSS = 0.0
230
231             print("Epoch "+str(e))
232
233             for step in range(ITRS):
234
235                 batch_inputs, batch_context =
236                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
237
238                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
239
240                 _, loss_val, batch_lss, gs = session.run([train_step,
241                 nce_loss,batch_loss,global_step], feed_dict=feed_dict)
242
243                 TRAIN_LOSS += np.sum(batch_lss)
244
245                 #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
246
247                 if (save): saver.save(session, complete_path, global_step=global_step)
248
249                 #-----
250                 # DEV
251                 #-----
252
253                 if(len(DX)>0):
254
255                     ITRS_TEST = int(len(DX) / batch_size)
256                     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
257                     TEST_LOSS = 0.0
258                     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es òdivisin exacta, se
259                         aade una òiteracin
260
261                     for step in range(ITRS_TEST):
262
263                         batch_inputs, batch_context =
264                             getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
265
266                         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
267                         loss_val,batch_lss = session.run([nce_loss,batch_loss],
268                         feed_dict=feed_dict)
269
270                         TEST_LOSS+=np.sum(batch_lss)

```

```

268         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
269         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
270
271     if (len(dev_hist) >= slope_size):
272         # Calcular la pendiente
273         print(getSlope(dev_hist[-slope_size:]))
274         print(dev_hist[-slope_size:])
275
276     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
277
278         def printArray(data, data2):
279             for i in range(len(data)):
280                 a=data[i]
281                 b=data2[i]
282                 print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
283
284         print("-" * 50)
285         print(e)
286         printArray(dev_hist, train_hist)
287         return
288 else:
289     #
290
291     # TEST FINAL
292     #
293
294     batch_size = 300
295
296     IN_TOP_5 = 0
297     IN_TOP_10 = 0
298     IN_TOP_20 = 0
299     IN_TOP_50 = 0
300     IN_TOP_100 = 0
301     IN_TOP_1000 = 0
302     real_positions = []
303
304     items = 0
305
306     ITRS_TEST = int(len(TX) / batch_size)
307     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
308     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se aade una
309     iteracin
310
311     for step in range(ITRS_TEST):
312         #tss = time.time()
313
314         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
315             ITRS_TEST, RES_ITMS_TEST)
316         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
317             batch_inputs, train_labels: batch_context})
318
319         for j in range(len(batch_context)):
320             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
321             real_positions.append(real_position)
322
323         IN_TOP_5 += 1 if real_position <= 5 else 0
324         IN_TOP_10 += 1 if real_position <= 10 else 0

```

```

320         IN_TOP_20 += 1 if real_position <= 20 else 0
321         IN_TOP_50 += 1 if real_position <= 50 else 0
322         IN_TOP_100 += 1 if real_position <= 100 else 0
323         IN_TOP_1000 += 1 if real_position <= 1000 else 0
324
325         items += len(batch_inputs)*1.0
326
327         PTF = IN_TOP_5 / items
328         PTT = IN_TOP_10 / items
329         PTTW = IN_TOP_20 / items
330         PTFT = IN_TOP_50 / items
331         PTCIEN = IN_TOP_100 / items
332         PTMIL = IN_TOP_1000 / items
333
334         print(np.mean(real_positions), np.std(real_positions),
335               np.median(real_positions))
336         print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
337               + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
338               + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
339
340         #print(time.time()-tss)
341
342         print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
343
344     #-----
345     #####
346     # Llamadas
347     #####
348
349     TEST_NAME="test_052"
350
351     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
352
353     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
354     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
355     D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
356
357     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
358     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
359     D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
360
361     W2V = np.load('embeddings/w2v/EMB_MATRIX')
362
363     #Se concatenan los embeddings de los usuarios para generar el perfil
364     D2V_CONS = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
365     Consolidado DM + Consolidado DBOW
366
367     D2V_PRES = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El presente = Presente
368     DM + Presente DBOW
369
370     D2V_EMB = np.concatenate((D2V_CONS, D2V_PRES), axis=1) # El perfil completo CONS + PRES
371
372     #Se cargan los datos de train de las canciones para conocer los IDs reales
373     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
374
375     #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
376     ALL = set(list(range(W2V.shape[0])))

```

```

374 REAL = set(list(set(list(map(int, DATA))))))
375 NOT_EXISTS = list(ALL.difference-REAL))
376
377 #####
378 # TENSORFLOW
379 #####
380
381 #-----
382 # Grid Search + Early Stopping
383 #-----
384 '''
385 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
386 rates = [0.000001, 0.01, 0.0001]
387
388 for r in rates:
389     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
390
391     print("Learning Rate: "+str(r))
392     print("-" * 50)
393     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
394     print("-"*50)
395
396 exit()
397 '''
398 #-----
399 # Train
400 #-----
401 '''
402 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
403
404 #Para esta fase, se une TRAIN y DEV
405 X = np.concatenate((X,DX),axis=0)
406 Y = np.concatenate((Y,DY),axis=0)
407
408 DX=[]
409 DY=[]
410
411 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":164, "seed":100}
412 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
413
414 exit()
415 '''
416 #-----
417 # Test
418 #-----
419
420 X, Y, DX, DY, TX, TY, user_size = getData-REAL_IDS=list-REAL))
421
422 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":164, "seed":100}
423 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.13. TEST 13

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2
28     TRAIN = np.empty([0, 3], dtype=np.int32)
29     DEV = np.empty([0, 3], dtype=np.int32)
30     TEST = np.empty([0, 3], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-1]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:]
52
53         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
54         USR_DEV[:, 0] = u
55         USR_DEV[:, 1] = DEV_DATA[:len(DEV_DATA) - 1]

```

```

55     USR_DEV[:, 2] = DEV_DATA[1:]
56
57     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
58     USR_TEST[:, 0] = u
59     USR_TEST[:, 1] = TEST_DATA[:len(TEST_DATA) - 1]
60     USR_TEST[:, 2] = TEST_DATA[1:]
61
62
63     # Añadir al global
64     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
65     DEV = np.concatenate([DEV,USR_DEV])
66     TEST = np.concatenate([TEST,USR_TEST])
67
68
69     #Mezclar
70     if(shuffle):
71         np.random.seed(seed)
72         np.random.shuffle(TRAIN)
73         np.random.shuffle(DEV)
74         np.random.shuffle(TEST)
75
76
77     return
78         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se noaadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X,Y, DX,DY, TX,TY, config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     O1 = tf.Variable(tf.truncated_normal([user_size, user_emb_size], mean=0.0, stddev=1.0
146     / math.sqrt(user_emb_size)), name="O1")
147     B01 = tf.Variable(tf.zeros([user_emb_size]), name="B01")
148
149     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
150     / math.sqrt(song_emb_size)), name="O2")
151     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
152
153     # Operaciones
154
155     # Embedding de documento
156     ed = tf.nn.embedding_lookup(O1, train_dataset[:, 0])
157     ed = ed + B01
158
159     # Embedding de cancion
160     ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
161     ec = ec + B02
162
163     # Cálculo de LOSS y optimizador—

```

```

162     embed = []
163     embed.append(ed)
164     embed.append(ec)
165     embed = tf.concat(embed, 1)
166
167     # softmax weights, W and D vectors should be concatenated before applying softmax
168     weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
169                                           math.sqrt(song_size)),
170                           name="weights")
171     # softmax biases
172     biases = tf.Variable(tf.zeros([song_size]), name="biases")
173
174     # aClculo de LOSS y optimizador
175     nce_loss = tf.nn.nce_loss(weights=weights,
176                               biases=biases,
177                               labels=train_labels,
178                               inputs=embed,
179                               num_sampled=num_sampled,
180                               num_classes=song_size,
181                               name="nce_loss")
182
183     batch_loss = nce_loss
184     nce_loss = tf.reduce_mean(nce_loss)
185
186     logits = tf.matmul(embed, tf.transpose(weights))
187     logits = tf.nn.bias_add(logits, biases)
188
189     top_v, top_i = tf.nn.top_k(logits, k=song_size)
190
191     # logits = tf.nn.softmax(logits, axis=0)
192
193     # optimizer =
194     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
195     train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
196                                     global_step=global_step)
197
198     # Inicializar variables
199     init = tf.global_variables_initializer()
200
201     # Crear objeto encargado de almacenar la red
202     saver = tf.train.Saver(max_to_keep=1)
203
204     #config = tf.ConfigProto()
205     #config.gpu_options.allow_growth = True
206     #with tf.Session(graph=graph, config=config) as session:
207     with tf.Session(graph=graph) as session:
208
209         # We must initialize all variables before we use them.
210         init.run()
211
212         if(save):
213             # Obtener el checkpoint
214             ckpt = tf.train.get_checkpoint_state(model_path)
215
216             # Si existe el checkpoint restaurar

```

```

215     if ckpt and ckpt.model_checkpoint_path:
216         saver.restore(session, ckpt.model_checkpoint_path)
217
218 #Si se desea entrenar
219 if(train):
220     #
221     # TRAIN
222     #
223
224     ITRS= int(len(X)/batch_size)
225     RES_ITMS = len(X) - (ITRS*batch_size)
226     if(RES_ITMS>0):ITRS+=1 #Si no es divisible exacta, se añade una iteración
227
228     for e in range(epochs):
229         TRAIN_LOSS = 0.0
230
231         print("Epoch "+str(e))
232
233         for step in range(ITRS):
234
235             batch_inputs, batch_context =
236                 getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
237
238             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
239
240             -, loss_val, batch_loss, gs = session.run([train_step,
241                 nce_loss,batch_loss,global_step], feed_dict=feed_dict)
242
243             TRAIN_LOSS += np.sum(batch_loss)
244
245             #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
246
247             if (save): saver.save(session, complete_path, global_step=global_step)
248
249             #
250             # DEV
251             #
252
253             if(len(DX)>0):
254
255                 ITRS_TEST = int(len(DX) / batch_size)
256                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
257                 TEST_LOSS = 0.0
258                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisible exacta, se
259                     añade una iteración
260
261                 for step in range(ITRS_TEST):
262
263                     batch_inputs, batch_context =
264                         getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
265
266                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
267                     loss_val,batch_loss = session.run([nce_loss,batch_loss],
268                         feed_dict=feed_dict)
269
270                     TEST_LOSS+=np.sum(batch_loss)

```

```

267         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
268         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
269
270
271     if (len(dev_hist) >= slope_size):
272         # Calcular la pendiente
273         print(getSlope(dev_hist[-slope_size:]))
274         print(dev_hist[-slope_size:])
275
276         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
277
278             def printArray(data, data2):
279                 for i in range(len(data)):
280                     a=data[i]
281                     b=data2[i]
282                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
283
284             print("-" * 50)
285             print(e)
286             printArray(dev_hist, train_hist)
287             return
288 else:
289     #
290
291 # TEST FINAL
292 #
293
294 batch_size = 300
295
296 IN_TOP_5 = 0
297 IN_TOP_10 = 0
298 IN_TOP_20 = 0
299 IN_TOP_50 = 0
300 IN_TOP_100 = 0
301 IN_TOP_1000 = 0
302 real_positions = []
303
304 items = 0
305
306 ITRS_TEST = int(len(TX) / batch_size)
307 RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
308 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
309 iteracin
310
311 for step in range(ITRS_TEST):
312     #tss = time.time()
313
314     batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
315         ITRS_TEST, RES_ITMS_TEST)
316     top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
317         batch_inputs, train_labels: batch_context})
318
319     for j in range(len(batch_context)):
320         real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
321         real_positions.append(real_position)
322
323     IN_TOP_5 += 1 if real_position <= 5 else 0

```

```

319         IN_TOP_10 += 1 if real_position <= 10 else 0
320         IN_TOP_20 += 1 if real_position <= 20 else 0
321         IN_TOP_50 += 1 if real_position <= 50 else 0
322         IN_TOP_100 += 1 if real_position <= 100 else 0
323         IN_TOP_1000 += 1 if real_position <= 1000 else 0
324
325         items += len(batch_inputs)*1.0
326
327         PTF = IN_TOP_5 / items
328         PTT = IN_TOP_10 / items
329         PTTW = IN_TOP_20 / items
330         PTFT = IN_TOP_50 / items
331         PTCIEN = IN_TOP_100 / items
332         PTMIL = IN_TOP_1000 / items
333
334         print(np.mean(real_positions), np.std(real_positions),
335               np.median(real_positions))
336         print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
337               + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
338               + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
339
340         #print(time.time()-tss)
341
342         print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
343
344         #-----
345
346         #####
347         # Llamadas
348         #####
349
350         TEST_NAME="test_012"
351
352         os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
353
354         #Solamente se carga para conocer el unmero de canciones
355         W2V = np.load('embeddings/w2v/EMB_MATRIX')
356
357         #Se cargan los datos de train de las canciones para conocer los IDs reales
358         DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
359
360         #Obtener una lista con los ID de canciones que no aestn el la matriz de embeddings
361         ALL = set(list(range(W2V.shape[0])))
362         REAL = set(list(set(list(map(int, DATA)))))
363
364         #####
365         # TENSORFLOW
366         #####
367
368         #-----
369         # Grid Search + Early Stopping
370         #-----
371         '''
372         X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
373         rates = [0.000001,1, 0.01, 0.0001]
374

```

```

375 for r in rates:
376     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
377
378     print("Learning Rate: "+str(r))
379     print("-" * 50)
380     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
381     print("-"*50)
382
383 exit()
384 '''
385 #-----
386 # Train
387 #-----
388 '''
389 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL)
390
391 #Para esta fase, se une TRAIN y DEV
392 X = np.concatenate((X,DX),axis=0)
393 Y = np.concatenate((Y,DY),axis=0)
394
395 DX=[]
396 DY=[]
397
398 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
399 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
400
401 exit()
402 '''
403 #-----
404 # Test
405 #-----
406
407 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL)
408
409 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
410 trainNet(X, Y, DX, DY, TX, TY, config, user_size,train=False, save=True)

```



## 4.14. TEST 14

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data ="datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
20            seed=100, train=.7):
21
22     DATA = np.load(data)
23
24     #DATA = DATA[:20];print("ELIMINAR FILTRO")
25
26     #User, cancion 1, cancion 2
27     TRAIN = np.empty([0, 3], dtype=np.int32)
28     DEV = np.empty([0, 3], dtype=np.int32)
29     TEST = np.empty([0, 3], dtype=np.int32)
30
31     for u in range(len(DATA)):
32
33         # Eliminar canciones escuchadas que no tienen embedding
34         USER_PLAYS = DATA[u]
35         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
36         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
37
38         #Separar en train dev test
39         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
40         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
41
42         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
43         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
44         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
45
46         # Crear ejemplos
47         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
48         USR_TRAIN[:,0] = u
49         USR_TRAIN[:,1] = TRAIN_DATA[:,len(TRAIN_DATA)-1]
50         USR_TRAIN[:,2] = TRAIN_DATA[1:]
51
52         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
53         USR_DEV[:, 0] = u
54         USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA) - 1]
55         USR_DEV[:, 2] = DEV_DATA[1:]

```

```

55
56     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
57     USR_TEST[:, 0] = u
58     USR_TEST[:, 1] = TEST_DATA[:len(TEST_DATA) - 1]
59     USR_TEST[:, 2] = TEST_DATA[1:]
60
61
62     # Añadir al global
63     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
64     DEV = np.concatenate([DEV,USR_DEV])
65     TEST = np.concatenate([TEST,USR_TEST])
66
67
68     #Mezclar
69     if(shuffle):
70         np.random.seed(seed)
71         np.random.shuffle(TRAIN)
72         np.random.shuffle(DEV)
73         np.random.shuffle(TEST)
74
75
76     return
77         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
78
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se no aadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return(batch_inputs, batch_context)
93
94 def trainNet(X,Y, DX,DY, TX,TY, config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     song_emb_size = 64
109     hidden_size = 32
110
111     profile_size = user_emb_size + song_emb_size

```

```

111
112 num_sampled = 64 # Ejemplos negativos (Para el NSE)
113
114 model_name = 'model'
115 model_path = 'models/' + TEST_NAME + '/model'
116 complete_path = model_path + "/" + model_name
117
118 # Se almacena la loss de dev de las 5 ultimas epochs
119 slope_size = 100
120 train_hist = []
121 dev_hist = []
122
123 # Creacin del grafo de TF.
124 graph = tf.Graph()
125
126 with graph.as_default():
127
128     if(config["seed"]!=1):
129         tf.set_random_seed(config["seed"])
130
131     # nmero global de iteraciones
132     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
133
134     # Datos de entrada
135
136     # Array del tamaño del batch con las X
137     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
138
139     # Array del tamaño del batch con las Y asociadas
140     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
141
142     # Embeddings
143
144     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
145
146     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
147         / math.sqrt(song_emb_size)), name="O2")
148     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
149
150     # Operaciones
151
152     # Embedding de documento
153     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
154
155     # Embedding de cancion
156     ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
157     ec = ec + B02
158
159     # Cálculo de LOSS y optimizador
160
161     embed = []
162     embed.append(ed)
163     embed.append(ec)
164     embed = tf.concat(embed, 1)

```

```

163
164 # softmax weights, W and D vectors should be concatenated before applying softmax
165 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
166                                     math.sqrt(song_size)),
167                               name="weights")
168 # softmax biases
169 biases = tf.Variable(tf.zeros([song_size]), name="biases")
170
171 # aClculo de LOSS y optimizador—
172
173 nce_loss = tf.nn.nce_loss(weights=weights,
174                           biases=biases,
175                           labels=train_labels,
176                           inputs=embed,
177                           num_sampled=num_sampled,
178                           num_classes=song_size,
179                           name="nce_loss")
180
181 batch_loss = nce_loss
182 nce_loss = tf.reduce_mean(nce_loss)
183
184 logits = tf.matmul(embed, tf.transpose(weights))
185 logits = tf.nn.bias_add(logits, biases)
186 # logits = tf.nn.softmax(logits,axis=0)
187 top_v,top_i = tf.nn.top_k(logits,k=song_size)
188
189 # optimizer =
190     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
191 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
192     global_step=global_step)
193
194 # Inicializar variables
195 init = tf.global_variables_initializer()
196
197 # Crear objeto encargado de almacenar la red
198 saver = tf.train.Saver(max_to_keep=1)
199
200 config = tf.ConfigProto()
201 config.gpu_options.allow_growth = True
202 with tf.Session(graph=graph, config=config) as session:
203 #with tf.Session(graph=graph) as session:
204
205     # We must initialize all variables before we use them.
206     init.run()
207
208     if(save):
209         # Obtener el checkpoint
210         ckpt = tf.train.get_checkpoint_state(model_path)
211
212         # Si existe el checkpoint restaurar
213         if ckpt and ckpt.model_checkpoint_path:
214             saver.restore(session, ckpt.model_checkpoint_path)
215
216     #Si se desea entrenar
217     if(train):
218         #

```

```

216 # TRAIN
217 #-----
218
219 ITRS= int(len(X)/batch_size)
220 RES_ITMS = len(X) - (ITRS*batch_size)
221 if(RES_ITMS>0):ITRS+=1 #Si no es divisible exacta, se añade una iteración
222
223 for e in range(epochs):
224     TRAIN_LOSS = 0.0
225
226     print("Epoch "+str(e))
227
228     for step in range(ITRS):
229
230         batch_inputs, batch_context =
231             getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
232
233         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
234
235         -, loss_val, batch_loss, gs = session.run([train_step,
236             nce_loss,batch_loss,global_step], feed_dict=feed_dict)
237
238         TRAIN_LOSS += np.sum(batch_loss)
239
240         #print('Epoch '+str(e)+', last batch loss: '+str(loss_val))
241
242         if (save): saver.save(session, complete_path, global_step=global_step)
243
244 #-----
245 # DEV
246 #-----
247
248 if(len(DX)>0):
249
250     ITRS_TEST = int(len(DX) / batch_size)
251     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
252     TEST_LOSS = 0.0
253     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisible exacta, se
254         añade una iteración
255
256     for step in range(ITRS_TEST):
257
258         batch_inputs, batch_context =
259             getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
260
261         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
262         loss_val,batch_loss = session.run([nce_loss,batch_loss],
263             feed_dict=feed_dict)
264
265         TEST_LOSS+=np.sum(batch_loss)
266
267         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
268         train_hist.append(TRAIN_LOSS / (len(X) *1.0))
269
270     if (len(dev_hist) >= slope_size):
271         # Calcular la pendiente

```

```

268     print(getSlope(dev_hist[-slope_size:]))
269     print(dev_hist[-slope_size:])
270
271     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
272
273         def printArray(data, data2):
274             for i in range(len(data)):
275                 a=data[i]
276                 b=data2[i]
277                 print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
278
279         print("-" * 50)
280         print(e)
281         printArray(dev_hist, train_hist)
282         return
283 else:
284     #
285
286     # TEST FINAL
287     #
288
289     batch_size = 300
290
291     IN_TOP_5 = 0
292     IN_TOP_10 = 0
293     IN_TOP_20 = 0
294     IN_TOP_50 = 0
295     IN_TOP_100 = 0
296     IN_TOP_1000 = 0
297     real_positions = []
298
299     items = 0
300
301     ITRS_TEST = int(len(TX) / batch_size)
302     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
303     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se aade una
304     oiteracin
305
306     for step in range(ITRS_TEST):
307         #tss = time.time()
308
309         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
310             ITRS_TEST, RES_ITMS_TEST)
311         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
312             batch_inputs, train_labels: batch_context})
313
314         for j in range(len(batch_context)):
315             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
316             real_positions.append(real_position)
317
318             IN_TOP_5 += 1 if real_position <= 5 else 0
319             IN_TOP_10 += 1 if real_position <= 10 else 0
320             IN_TOP_20 += 1 if real_position <= 20 else 0
321             IN_TOP_50 += 1 if real_position <= 50 else 0
322             IN_TOP_100 += 1 if real_position <= 100 else 0
323             IN_TOP_1000 += 1 if real_position <= 1000 else 0

```

```

320         items += len(batch_inputs)*1.0
321
322         PTF = IN_TOP_5 / items
323         PT1 = IN_TOP_10 / items
324         PT2 = IN_TOP_20 / items
325         PT5 = IN_TOP_50 / items
326         PT100 = IN_TOP_100 / items
327         PT1000 = IN_TOP_1000 / items
328
329         print(np.mean(real_positions), np.std(real_positions),
330               np.median(real_positions))
331         print(str(PTF).replace(".", ",") + "\t" + str(PT1).replace(".", ",") + "\t"
332               + str(PT2).replace(".", ",") + "\t" + str(PT5).replace(".", ",") + "\t"
333               + str(PT100).replace(".", ",") + "\t" + str(PT1000).replace(".", ","))
334
335         #print(time.time()-tss)
336
337         print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
338
339     #-----
340     #####
341     # Llamadas
342     #####
343
344     TEST_NAME="test_062"
345
346     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
347
348     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
349     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
350     D2V_D2V_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
351     W2V = np.load('embeddings/w2v/EMB_MATRIX')
352
353     #Se concatenan los embeddings de los usuarios para generar el perfil
354     D2V_EMB = np.concatenate((D2V_DM_COMPLETE, D2V_D2V_COMPLETE), axis=1) # El consolidado =
355         Consolidado DM + Consolidado DBOW
356
357     #Se cargan los datos de train de las canciones para conocer los IDs reales
358     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
359
360     #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
361     ALL = set(list(range(W2V.shape[0])))
362     REAL = set(list(set(list(map(int, DATA)))))
363     NOT_EXISTS = list(ALL.difference(REAL))
364
365     #####
366     # TENSORFLOW
367     #####
368     #-----
369     # Grid Search + Early Stopping
370     #-----
371     '''
372     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
373     rates = [0.000001, 1, 0.01, 0.0001]
374

```

```

375 for r in rates:
376     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
377
378     print("Learning Rate: "+str(r))
379     print("-" * 50)
380     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
381     print("-"*50)
382
383 exit()
384 '''
385 #-----
386 # Train
387 #-----
388 '''
389 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
390
391 #Para esta fase, se une TRAIN y DEV
392 X = np.concatenate((X,DX),axis=0)
393 Y = np.concatenate((Y,DY),axis=0)
394
395 DX=[]
396 DY=[]
397
398 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":624, "seed":100}
399 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
400
401 exit()
402 '''
403 #-----
404 # Test
405 #-----
406
407 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
408
409 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":624, "seed":100}
410 trainNet(X, Y, DX, DY, TX, TY, config, user_size,train=False, save=True)

```



## 4.15. TEST 15

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data ="datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
20            seed=100, train=.7):
21
22     DATA = np.load(data)
23
24     #DATA = DATA[:20];print("ELIMINAR FILTRO")
25
26     #User, cancion 1, cancion 2
27     TRAIN = np.empty([0, 3], dtype=np.int32)
28     DEV = np.empty([0, 3], dtype=np.int32)
29     TEST = np.empty([0, 3], dtype=np.int32)
30
31     for u in range(len(DATA)):
32
33         # Eliminar canciones escuchadas que no tienen embedding
34         USER_PLAYS = DATA[u]
35         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
36         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
37
38         #Separar en train dev test
39         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
40         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
41
42         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
43         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
44         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
45
46         # Crear ejemplos
47         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
48         USR_TRAIN[:,0] = u
49         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-1]
50         USR_TRAIN[:,2] = TRAIN_DATA[1:]
51
52         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
53         USR_DEV[:, 0] = u
54         USR_DEV[:, 1] = DEV_DATA[:len(DEV_DATA) - 1]
55         USR_DEV[:, 2] = DEV_DATA[1:]

```

```

55
56     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
57     USR_TEST[:, 0] = u
58     USR_TEST[:, 1] = TEST_DATA[:len(TEST_DATA) - 1]
59     USR_TEST[:, 2] = TEST_DATA[1:]
60
61
62     # Añadir al global
63     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
64     DEV = np.concatenate([DEV,USR_DEV])
65     TEST = np.concatenate([TEST,USR_TEST])
66
67
68     #Mezclar
69     if(shuffle):
70         np.random.seed(seed)
71         np.random.shuffle(TRAIN)
72         np.random.shuffle(DEV)
73         np.random.shuffle(TEST)
74
75
76     return
77         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
78
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se no aadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return(batch_inputs, batch_context)
93
94 def trainNet(X,Y, DX,DY, TX,TY, config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     song_emb_size = 64
109     hidden_size = 32
110
111     profile_size = user_emb_size + song_emb_size

```

```

111
112 num_sampled = 64 # Ejemplos negativos (Para el NSE)
113
114 model_name = 'model'
115 model_path = 'models/' + TEST_NAME + '/model'
116 complete_path = model_path + "/" + model_name
117
118 # Se almacena la loss de dev de las 5 ultimas epochs
119 slope_size = 100
120 train_hist = []
121 dev_hist = []
122
123 # Creacin del grafo de TF.
124 graph = tf.Graph()
125
126 with graph.as_default():
127
128     if(config["seed"]!=1):
129         tf.set_random_seed(config["seed"])
130
131     # nmero global de iteraciones
132     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
133
134     # Datos de entrada
135
136     # Array del tamaño del batch con las X
137     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
138
139     # Array del tamaño del batch con las Y asociadas
140     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
141
142     # Embeddings
143
144     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
145
146     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
147         / math.sqrt(song_emb_size)), name="O2")
148     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
149
150     # Operaciones
151
152     # Embedding de documento
153     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
154
155     # Embedding de cancion
156     ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
157     ec = ec + B02
158
159     # Cálculo de LOSS y optimizador
160
161     embed = []
162     embed.append(ed)
163     embed.append(ec)
164     embed = tf.concat(embed, 1)

```

```

163
164 # softmax weights, W and D vectors should be concatenated before applying softmax
165 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
166                                     math.sqrt(song_size)),
167                               name="weights")
168 # softmax biases
169 biases = tf.Variable(tf.zeros([song_size]), name="biases")
170
171 # aClculo de LOSS y optimizador—
172
173 nce_loss = tf.nn.nce_loss(weights=weights,
174                           biases=biases,
175                           labels=train_labels,
176                           inputs=embed,
177                           num_sampled=num_sampled,
178                           num_classes=song_size,
179                           name="nce_loss")
180
181 batch_loss = nce_loss
182 nce_loss = tf.reduce_mean(nce_loss)
183
184 logits = tf.matmul(embed, tf.transpose(weights))
185 logits = tf.nn.bias_add(logits, biases)
186 # logits = tf.nn.softmax(logits,axis=0)
187 top_v,top_i = tf.nn.top_k(logits,k=song_size)
188
189 # optimizer =
190     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
191 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
192     global_step=global_step)
193
194 # Inicializar variables
195 init = tf.global_variables_initializer()
196
197 # Crear objeto encargado de almacenar la red
198 saver = tf.train.Saver(max_to_keep=1)
199
200 config = tf.ConfigProto()
201 config.gpu_options.allow_growth = True
202 with tf.Session(graph=graph, config=config) as session:
203 #with tf.Session(graph=graph) as session:
204
205     # We must initialize all variables before we use them.
206     init.run()
207
208     if(save):
209         # Obtener el checkpoint
210         ckpt = tf.train.get_checkpoint_state(model_path)
211
212         # Si existe el checkpoint restaurar
213         if ckpt and ckpt.model_checkpoint_path:
214             saver.restore(session, ckpt.model_checkpoint_path)
215
216     #Si se desea entrenar
217     if(train):
218         #

```

```

216 # TRAIN
217 #-----
218
219 ITRS= int(len(X)/batch_size)
220 RES_ITMS = len(X) - (ITRS*batch_size)
221 if(RES_ITMS>0):ITRS+=1 #Si no es divisible exacta, se añade una iteración
222
223 for e in range(epochs):
224     TRAIN_LOSS = 0.0
225
226     print("Epoch "+str(e))
227
228     for step in range(ITRS):
229
230         batch_inputs, batch_context =
231             getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
232
233         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
234
235         -, loss_val, batch_loss, gs = session.run([train_step,
236             nce_loss,batch_loss,global_step], feed_dict=feed_dict)
237
238         TRAIN_LOSS += np.sum(batch_loss)
239
240         #print('Epoch '+str(e)+', last batch loss: '+str(loss_val))
241
242         if (save): saver.save(session, complete_path, global_step=global_step)
243
244 #-----
245 # DEV
246 #-----
247
248 if(len(DX)>0):
249
250     ITRS_TEST = int(len(DX) / batch_size)
251     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
252     TEST_LOSS = 0.0
253     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisible exacta, se
254         añade una iteración
255
256     for step in range(ITRS_TEST):
257
258         batch_inputs, batch_context =
259             getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
260
261         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
262         loss_val,batch_loss = session.run([nce_loss,batch_loss],
263             feed_dict=feed_dict)
264
265         TEST_LOSS+=np.sum(batch_loss)
266
267         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
268         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
269
270     if (len(dev_hist) >= slope_size):
271         # Calcular la pendiente

```

```

268     print(getSlope(dev_hist[-slope_size:]))
269     print(dev_hist[-slope_size:])
270
271     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
272
273         def printArray(data, data2):
274             for i in range(len(data)):
275                 a=data[i]
276                 b=data2[i]
277                 print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
278
279         print("-" * 50)
280         print(e)
281         printArray(dev_hist, train_hist)
282         return
283 else:
284     #
285
286     # TEST FINAL
287     #
288
289     batch_size = 300
290
291     IN_TOP_5 = 0
292     IN_TOP_10 = 0
293     IN_TOP_20 = 0
294     IN_TOP_50 = 0
295     IN_TOP_100 = 0
296     IN_TOP_1000 = 0
297     real_positions = []
298
299     items = 0
300
301     ITRS_TEST = int(len(TX) / batch_size)
302     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
303     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se aade una
304     oiteracin
305
306     for step in range(ITRS_TEST):
307         #tss = time.time()
308
309         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
310             ITRS_TEST, RES_ITMS_TEST)
311         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
312             batch_inputs, train_labels: batch_context})
313
314         for j in range(len(batch_context)):
315             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
316             real_positions.append(real_position)
317
318             IN_TOP_5 += 1 if real_position <= 5 else 0
319             IN_TOP_10 += 1 if real_position <= 10 else 0
320             IN_TOP_20 += 1 if real_position <= 20 else 0
321             IN_TOP_50 += 1 if real_position <= 50 else 0
322             IN_TOP_100 += 1 if real_position <= 100 else 0
323             IN_TOP_1000 += 1 if real_position <= 1000 else 0

```

```

320         items += len(batch_inputs)*1.0
321
322         PTF = IN_TOP_5 / items
323         PT1 = IN_TOP_10 / items
324         PT2 = IN_TOP_20 / items
325         PT5 = IN_TOP_50 / items
326         PT100 = IN_TOP_100 / items
327         PT1000 = IN_TOP_1000 / items
328
329         print(np.mean(real_positions), np.std(real_positions),
330               np.median(real_positions))
331         print(str(PTF).replace(".", ",") + "\t" + str(PT1).replace(".", ",") + "\t"
332               + str(PT2).replace(".", ",") + "\t" + str(PT5).replace(".", ",") + "\t"
333               + str(PT100).replace(".", ",") + "\t" + str(PT1000).replace(".", ","))
334
335         #print(time.time()-tss)
336
337         print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
338
339     #-----
340     #####
341     # Llamadas
342     #####
343
344     TEST_NAME="test_072"
345
346     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
347
348     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
349     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
350     D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
351     W2V = np.load('embeddings/w2v/EMB_MATRIX')
352
353     #Se concatenan los embeddings de los usuarios para generar el perfil
354     D2V_EMB = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El reciente = Reciente
355     DM + Reciente DBOW
356
357     #Se cargan los datos de train de las canciones para conocer los IDs reales
358     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
359
360     #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
361     ALL = set(list(range(W2V.shape[0])))
362     REAL = set(list(set(list(map(int, DATA)))))
363     NOT_EXISTS = list(ALL.difference(REAL))
364
365     #####
366     # TENSORFLOW
367     #####
368     #-----
369     # Grid Search + Early Stopping
370     #-----
371     '''
372     X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
373     rates = [0.000001, 1, 0.01, 0.0001]
374

```

```

375 for r in rates:
376     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
377
378     print("Learning Rate: "+str(r))
379     print("-" * 50)
380     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
381     print("-"*50)
382
383 exit()
384 '''
385 #-----
386 # Train
387 #-----
388 '''
389 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
390
391 #Para esta fase, se une TRAIN y DEV
392 X = np.concatenate((X,DX),axis=0)
393 Y = np.concatenate((Y,DY),axis=0)
394
395 DX=[]
396 DY=[]
397
398 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
399 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
400
401 exit()
402 '''
403 #-----
404 # Test
405 #-----
406
407 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
408
409 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
410 trainNet(X, Y, DX, DY, TX, TY, config, user_size,train=False, save=True)

```



## 4.16. TEST 16

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14
15 #####
16 # eMtodos
17 #####
18
19 def getData(REAL_SONG_IDS=[],data ="datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
20            seed=100, train=.7):
21
22     DATA = np.load(data)
23
24     #DATA = DATA[:20];print("ELIMINAR FILTRO")
25
26     #User, cancion 1, cancion 2
27     TRAIN = np.empty([0, 3], dtype=np.int32)
28     DEV = np.empty([0, 3], dtype=np.int32)
29     TEST = np.empty([0, 3], dtype=np.int32)
30
31     for u in range(len(DATA)):
32
33         # Eliminar canciones escuchadas que no tienen embedding
34         USER_PLAYS = DATA[u]
35         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
36         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
37
38         #Separar en train dev test
39         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
40         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
41
42         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
43         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
44         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
45
46         # Crear ejemplos
47         USR_TRAIN = np.empty([len(TRAIN_DATA)-1, 3], dtype=np.int32)
48         USR_TRAIN[:,0] = u
49         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-1]
50         USR_TRAIN[:,2] = TRAIN_DATA[1:]
51
52         USR_DEV = np.empty([len(DEV_DATA) - 1, 3], dtype=np.int32)
53         USR_DEV[:, 0] = u
54         USR_DEV[:, 1] = DEV_DATA[:len(DEV_DATA) - 1]
55         USR_DEV[:, 2] = DEV_DATA[1:]

```

```

55
56     USR_TEST = np.empty([len(TEST_DATA) - 1, 3], dtype=np.int32)
57     USR_TEST[:, 0] = u
58     USR_TEST[:, 1] = TEST_DATA[:, len(TEST_DATA) - 1]
59     USR_TEST[:, 2] = TEST_DATA[1:]
60
61
62     # Añadir al global
63     TRAIN = np.concatenate([TRAIN, USR_TRAIN])
64     DEV = np.concatenate([DEV, USR_DEV])
65     TEST = np.concatenate([TEST, USR_TEST])
66
67
68     # Mezclar
69     if (shuffle):
70         np.random.seed(seed)
71         np.random.shuffle(TRAIN)
72         np.random.shuffle(DEV)
73         np.random.shuffle(TEST)
74
75
76     return
77         TRAIN[:, :2], np.matrix(TRAIN[:, -1]).T, DEV[:, :2], np.matrix(DEV[:, -1]).T, TEST[:, :2], np.matrix(TEST
78
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se no aadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return (batch_inputs, batch_context)
93
94 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     user_concat_size = user_emb_size * 2
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 2], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
146
147     T0 = tf.Variable(tf.truncated_normal([user_concat_size, user_emb_size], mean=0.0,
148         stddev=1.0 / math.sqrt(hidden_size)), name="T0")
149     B0 = tf.Variable(tf.zeros([user_emb_size]), name="B0")
150
151     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
152         / math.sqrt(song_emb_size)), name="O2")
153     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
154
155     # Operaciones
156
157     # Embedding de documento
158     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
159
160     #Transformar concat documento de 256 a 128
161     ed = tf.matmul(ed,T0)+B0
162
163     # Embedding de cancion
164     ec = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
165     ec = ec + B02

```

```

163
164 # aClculo de LOSS y optimizador—
165
166 embed = []
167 embed.append(ed)
168 embed.append(ec)
169 embed = tf.concat(embed, 1)
170
171 # softmax weights, W and D vectors should be concatenated before applying softmax
172 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
173                                     math.sqrt(song_size)),
174                       name="weights")
175 # softmax biases
176 biases = tf.Variable(tf.zeros([song_size]), name="biases")
177
178 # aClculo de LOSS y optimizador—
179
180 nce_loss = tf.nn.nce_loss(weights=weights,
181                          biases=biases,
182                          labels=train_labels,
183                          inputs=embed,
184                          num_sampled=num_sampled,
185                          num_classes=song_size,
186                          name="nce_loss")
187
188 batch_loss = nce_loss
189 nce_loss = tf.reduce_mean(nce_loss)
190
191 logits = tf.matmul(embed, tf.transpose(weights))
192 logits = tf.nn.bias_add(logits, biases)
193 # logits = tf.nn.softmax(logits,axis=0)
194 top_v,top_i = tf.nn.top_k(logits,k=song_size)
195
196 # optimizer =
197     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss,global_step=global_step)
198 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
199                               global_step=global_step)
200
201 # Inicializar variables
202 init = tf.global_variables_initializer()
203
204 # Crear objeto encargado de almacenar la red
205 saver = tf.train.Saver(max_to_keep=1)
206
207 config = tf.ConfigProto()
208 config.gpu_options.allow_growth = True
209 with tf.Session(graph=graph, config=config) as session:
210 #with tf.Session(graph=graph) as session:
211
212     # We must initialize all variables before we use them.
213     init.run()
214
215 if(save):
216     # Obtener el checkpoint
217     ckpt = tf.train.get_checkpoint_state(model_path)

```

```

215
216     # Si existe el checkpoint restaurar
217     if ckpt and ckpt.model_checkpoint_path:
218         saver.restore(session, ckpt.model_checkpoint_path)
219
220 #Si se desea entrenar
221 if(train):
222     #
223     # TRAIN
224     #
225
226     ITRS= int(len(X)/batch_size)
227     RES_ITMS = len(X) - (ITRS*batch_size)
228     if(RES_ITMS>0):ITRS+=1 #Si no es ÷divisin exacta, se ñaade una ÷iteracin
229
230     for e in range(epochs):
231         TRAIN_LOSS = 0.0
232
233         print("Epoch "+str(e))
234
235         for step in range(ITRS):
236
237             batch_inputs, batch_context =
238                 getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
239
240             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
241
242             _, loss_val, batch_loss, gs = session.run([train_step,
243                 nce_loss,batch_loss,global_step], feed_dict=feed_dict)
244
245             TRAIN_LOSS += np.sum(batch_loss)
246
247             #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
248
249             if (save): saver.save(session, complete_path, global_step=global_step)
250
251             #
252             # DEV
253             #
254
255             if(len(DX)>0):
256
257                 ITRS_TEST = int(len(DX) / batch_size)
258                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
259                 TEST_LOSS = 0.0
260                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es ÷divisin exacta, se
261                     ñaade una ÷iteracin
262
263                 for step in range(ITRS_TEST):
264
265                     batch_inputs, batch_context =
266                         getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
267
268                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
269                     loss_val,batch_loss = session.run([nce_loss,batch_loss],
270                         feed_dict=feed_dict)

```

```

267
268         TEST_LOSS+=np.sum(batch_loss)
269
270         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
271         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
272
273     if (len(dev_hist) >= slope_size):
274         # Calcular la pendiente
275         print(getSlope(dev_hist[-slope_size:]))
276         print(dev_hist[-slope_size:])
277
278         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
279
280             def printArray(data, data2):
281                 for i in range(len(data)):
282                     a=data[i]
283                     b=data2[i]
284                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
285
286             print("-" * 50)
287             print(e)
288             printArray(dev_hist, train_hist)
289             return
290 else:
291     #
292
293     # TEST FINAL
294     #
295
296     batch_size = 300
297
298     IN_TOP_5 = 0
299     IN_TOP_10 = 0
300     IN_TOP_20 = 0
301     IN_TOP_50 = 0
302     IN_TOP_100 = 0
303     IN_TOP_1000 = 0
304     real_positions = []
305
306     items = 0
307
308     ITRS_TEST = int(len(TX) / batch_size)
309     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
310     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
311     iteracin
312
313     for step in range(ITRS_TEST):
314         #tss = time.time()
315
316         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
317             ITRS_TEST, RES_ITMS_TEST)
318         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
319             batch_inputs, train_labels: batch_context})
320
321         for j in range(len(batch_context)):
322             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
323             real_positions.append(real_position)

```

```

319
320     IN_TOP_5 += 1 if real_position <= 5 else 0
321     IN_TOP_10 += 1 if real_position <= 10 else 0
322     IN_TOP_20 += 1 if real_position <= 20 else 0
323     IN_TOP_50 += 1 if real_position <= 50 else 0
324     IN_TOP_100 += 1 if real_position <= 100 else 0
325     IN_TOP_1000 += 1 if real_position <= 1000 else 0
326
327     items += len(batch_inputs)*1.0
328
329     PTF = IN_TOP_5 / items
330     PTT = IN_TOP_10 / items
331     PTTW = IN_TOP_20 / items
332     PTFT = IN_TOP_50 / items
333     PTCIEN = IN_TOP_100 / items
334     PTMIL = IN_TOP_1000 / items
335
336     print(np.mean(real_positions), np.std(real_positions),
337           np.median(real_positions))
338     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
339           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
340           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
341
342     #print(time.time()-tss)
343
344     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
345
346 #-----
347 #####
348 # Llamadas
349 #####
350
351 TEST_NAME="test_082"
352
353 os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
354
355 #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
356 D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
357 D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
358
359 #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
360 D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
361 D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
362
363 W2V = np.load('embeddings/w2v/EMB_MATRIX')
364
365 #Se concatenan los embeddings de los usuarios para generar el perfil
366 D2V_CONS = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
367     Consolidado DM + Consolidado DBOW
368
369 D2V_PRES = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El presente = Presente
370     DM + Presente DBOW
371
372 D2V_EMB = np.concatenate((D2V_CONS, D2V_PRES), axis=1) # El perfil completo CONS + PRES
373
374 #Se cargan los datos de train de las canciones para conocer los IDs reales
375 DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")

```

```

373
374 #Obtener una lista con los ID de canciones que no estén en la matriz de embeddings
375 ALL = set(list(range(W2V.shape[0])))
376 REAL = set(list(set(list(map(int, DATA))))))
377 NOT_EXISTS = list(ALL.difference(REAL))
378
379 #####
380 # TENSORFLOW
381 #####
382
383 # -----
384 # Grid Search + Early Stopping
385 # -----
386 '''
387 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
388 rates = [0.000001, 0.01, 0.0001]
389
390 for r in rates:
391     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
392
393     print("Learning Rate: "+str(r))
394     print("-" * 50)
395     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
396     print("-"*50)
397
398 exit()
399 '''
400 # -----
401 # Train
402 # -----
403 '''
404 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
405
406 #Para esta fase, se une TRAIN y DEV
407 X = np.concatenate((X,DX),axis=0)
408 Y = np.concatenate((Y,DY),axis=0)
409
410 DX=[]
411 DY=[]
412
413 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs":1000, "seed":100}
414 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
415
416 exit()
417 '''
418 # -----
419 # Test
420 # -----
421
422 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
423
424 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs":1000, "seed":100}
425 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.17. TEST 17

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:, len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:, len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN, USR_TRAIN])
67     DEV = np.concatenate([DEV, USR_DEV])
68     TEST = np.concatenate([TEST, USR_TEST])
69
70 #Mezclar
71 if shuffle:
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
79 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
80
81     # Si se noaadi una iteracin, introducir los datos restantes
82     if (step == ITRS - 1 and RES_ITMS > 0):
83         batch_inputs = X[(ITRS - 1) * batch_size:]
84         batch_context = Y[(ITRS - 1) * batch_size:]
85     # Si no, es un batch normal
86     else:
87         st = step * batch_size
88         nd = (step + 1) * batch_size
89
90         batch_inputs = X[st:nd]
91         batch_context = Y[st:nd]
92     return(batch_inputs, batch_context)
93
94 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
95
96     def getSlope(data):
97         r = linregress(range(len(data)), data)
98         return r.slope
99
100     batch_size = config['batch_size']
101     learning_rate = config['learning_rate']
102     epochs = config['epochs']
103
104     user_size = user_size
105     song_size = W2V.shape[0]
106
107     user_emb_size = 128
108     song_emb_size = 64
109     hidden_size = 32
110

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     O1 = tf.Variable(tf.truncated_normal([user_size, user_emb_size], mean=0.0, stddev=1.0
146         / math.sqrt(user_emb_size)), name="O1")
147     B01 = tf.Variable(tf.zeros([user_emb_size]), name="B01")
148
149     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
150
151     # Operaciones
152
153     # Embedding de documento
154     ed = tf.nn.embedding_lookup(O1, train_dataset[:, 0])
155     ed = ed + B01
156
157     # Embedding de cancion 1
158     ec1 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
159
160     # Embedding de cancion 2
161     ec2 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 2])
162
163     # Cálculo de LOSS y optimizador—

```

```

163 embed = []
164 embed.append(ed)
165 embed.append(ec1)
166 embed.append(ec2)
167 embed = tf.concat(embed, 1)
168
169 # softmax weights, W and D vectors should be concatenated before applying softmax
170 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
171     math.sqrt(song_size)), name="weights")
172 # softmax biases
173 biases = tf.Variable(tf.zeros([song_size]), name="biases")
174
175 # aClculo de LOSS y optimizador—
176
177 nce_loss = tf.nn.nce_loss(weights=weights,
178     biases=biases,
179     labels=train_labels,
180     inputs=embed,
181     num_sampled=num_sampled,
182     num_classes=song_size,
183     name="nce_loss")
184
185 batch_loss = nce_loss
186 nce_loss = tf.reduce_mean(nce_loss)
187
188 logits = tf.matmul(embed, tf.transpose(weights))
189 logits = tf.nn.bias_add(logits, biases)
190
191 top_v, top_i = tf.nn.top_k(logits, k=song_size)
192
193 # logits = tf.nn.softmax(logits, axis=0)
194
195 # optimizer =
196     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
197 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
198     global_step=global_step)
199
200 # Inicializar variables
201 init = tf.global_variables_initializer()
202
203 # Crear objeto encargado de almacenar la red
204 saver = tf.train.Saver(max_to_keep=1)
205
206 config = tf.ConfigProto()
207 config.gpu_options.allow_growth = True
208 with tf.Session(graph=graph, config=config) as session:
209     #with tf.Session(graph=graph) as session:
210
211     # We must initialize all variables before we use them.
212     init.run()
213
214 if(save):
215     # Obtener el checkpoint
216     ckpt = tf.train.get_checkpoint_state(model_path)
217
218     # Si existe el checkpoint restaurar

```

```

216     if ckpt and ckpt.model_checkpoint_path:
217         saver.restore(session, ckpt.model_checkpoint_path)
218
219 #Si se desea entrenar
220 if(train):
221     #
222     # TRAIN
223     #
224
225     ITRS= int(len(X)/batch_size)
226     RES_ITMS = len(X) - (ITRS*batch_size)
227     if(RES_ITMS>0):ITRS+=1 #Si no es òdivisin exacta, se ñaade una òiteracin
228
229     for e in range(epochs):
230         TRAIN_LOSS = 0.0
231
232         print("Epoch "+str(e))
233
234         for step in range(ITRS):
235
236             batch_inputs, batch_context =
237                 getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
238
239             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
240
241             -, loss_val, batch_lss, gs = session.run([train_step,
242                 nce_loss,batch_loss,global_step], feed_dict=feed_dict)
243
244             TRAIN_LOSS += np.sum(batch_lss)
245
246             #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
247
248             if (save): saver.save(session, complete_path, global_step=global_step)
249
250             #
251             # DEV
252             #
253
254             if(len(DX)>0):
255
256                 ITRS_TEST = int(len(DX) / batch_size)
257                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
258                 TEST_LOSS = 0.0
259                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es òdivisin exacta, se
260                     ñaade una òiteracin
261
262                 for step in range(ITRS_TEST):
263
264                     batch_inputs, batch_context =
265                         getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
266
267                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
268                     loss_val,batch_lss = session.run([nce_loss,batch_loss],
269                         feed_dict=feed_dict)
270
271                     TEST_LOSS+=np.sum(batch_lss)

```

```

268         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
269         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
270
271
272     if (len(dev_hist) >= slope_size):
273         # Calcular la pendiente
274         print(getSlope(dev_hist[-slope_size:]))
275         print(dev_hist[-slope_size:])
276
277     if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
278
279         def printArray(data, data2):
280             for i in range(len(data)):
281                 a=data[i]
282                 b=data2[i]
283                 print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
284
285         print("-" * 50)
286         print(e)
287         printArray(dev_hist, train_hist)
288         return
289     else:
290         #
291
292     # TEST FINAL
293     #
294
295     batch_size = 300
296
297     IN_TOP_5 = 0
298     IN_TOP_10 = 0
299     IN_TOP_20 = 0
300     IN_TOP_50 = 0
301     IN_TOP_100 = 0
302     IN_TOP_1000 = 0
303     real_positions = []
304
305     items = 0
306
307     ITRS_TEST = int(len(TX) / batch_size)
308     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
309     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
310     iteracin
311
312     for step in range(ITRS_TEST):
313         #tss = time.time()
314
315         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
316             ITRS_TEST, RES_ITMS_TEST)
317         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
318             batch_inputs, train_labels: batch_context})
319
320         for j in range(len(batch_context)):
321             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
322             real_positions.append(real_position)
323
324         IN_TOP_5 += 1 if real_position <= 5 else 0

```

```

320         IN_TOP_10 += 1 if real_position <= 10 else 0
321         IN_TOP_20 += 1 if real_position <= 20 else 0
322         IN_TOP_50 += 1 if real_position <= 50 else 0
323         IN_TOP_100 += 1 if real_position <= 100 else 0
324         IN_TOP_1000 += 1 if real_position <= 1000 else 0
325
326         items += len(batch_inputs)*1.0
327
328         PTF = IN_TOP_5 / items
329         PTT = IN_TOP_10 / items
330         PTTW = IN_TOP_20 / items
331         PTFT = IN_TOP_50 / items
332         PTCIEN = IN_TOP_100 / items
333         PTMIL = IN_TOP_1000 / items
334
335         print(np.mean(real_positions), np.std(real_positions),
336               np.median(real_positions))
337         print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
338               + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
339               + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
340
341         #print(time.time()-tss)
342
343         print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
344
345         #-----
346
347         #####
348         # Llamadas
349         #####
350
351         TEST_NAME="test_023"
352
353         os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
354
355         #Solamente se carga para conocer el unmero de canciones
356         W2V = np.load('embeddings/w2v/EMB_MATRIX')
357
358         #Se cargan los datos de train de las canciones para conocer los IDs reales
359         DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
360
361         #Obtener una lista con los ID de canciones que no aestn el la matriz de embeddings
362         ALL = set(list(range(W2V.shape[0])))
363         REAL = set(list(set(list(map(int, DATA)))))
364
365         #####
366         # TENSORFLOW
367         #####
368
369         #-----
370         # Grid Search + Early Stopping
371         #-----
372         '''
373         X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
374         rates = [0.000001,1, 0.01, 0.0001]
375

```

```

376
377 for r in rates:
378     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
379
380     print("Learning Rate: "+str(r))
381     print("-" * 50)
382     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
383     print("-"*50)
384
385 exit()
386 '''
387 #
388 # Train
389 #
390 '''
391 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
392
393 #Para esta fase, se une TRAIN y DEV
394 X = np.concatenate((X,DX),axis=0)
395 Y = np.concatenate((Y,DY),axis=0)
396
397 DX=[]
398 DY=[]
399
400 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
401 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
402
403 exit()
404 '''
405 #
406 # Test
407 #
408
409 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
410
411 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
412 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.18. TEST 18

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:,len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
67     DEV = np.concatenate([DEV,USR_DEV])
68     TEST = np.concatenate([TEST,USR_TEST])
69
70 #Mezclar
71 if(shuffle):
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se ha aadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X,Y, DX,DY, TX, TY, config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
146
147     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
148
149     # Operaciones
150
151     # Embedding de documento
152     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
153
154     # Embedding de cancion 1
155     ec1 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
156
157     # Embedding de cancion 2
158     ec2 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 2])
159
160     # Cálculo de LOSS y optimizador
161
162     embed = []
163     embed.append(ed)
164     embed.append(ec1)

```

```

164     embed.append(ec2)
165     embed = tf.concat(embed, 1)
166
167     # softmax weights, W and D vectors should be concatenated before applying softmax
168     weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
169         math.sqrt(song_size)), name="weights")
170     # softmax biases
171     biases = tf.Variable(tf.zeros([song_size]), name="biases")
172
173     # aClculo de LOSS y optimizador—
174     nce_loss = tf.nn.nce_loss(weights=weights,
175                             biases=biases,
176                             labels=train_labels,
177                             inputs=embed,
178                             num_sampled=num_sampled,
179                             num_classes=song_size,
180                             name="nce_loss")
181
182     batch_loss = nce_loss
183     nce_loss = tf.reduce_mean(nce_loss)
184
185     logits = tf.matmul(embed, tf.transpose(weights))
186     logits = tf.nn.bias_add(logits, biases)
187
188     top_v, top_i = tf.nn.top_k(logits, k=song_size)
189
190     # logits = tf.nn.softmax(logits, axis=0)
191
192     # optimizer =
193     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
194     train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
195         global_step=global_step)
196
197     # Inicializar variables
198     init = tf.global_variables_initializer()
199
200     # Crear objeto encargado de almacenar la red
201     saver = tf.train.Saver(max_to_keep=1)
202
203     config = tf.ConfigProto()
204     config.gpu_options.allow_growth = True
205     with tf.Session(graph=graph, config=config) as session:
206     #with tf.Session(graph=graph) as session:
207
208         # We must initialize all variables before we use them.
209         init.run()
210
211     if(save):
212         # Obtener el checkpoint
213         ckpt = tf.train.get_checkpoint_state(model_path)
214
215         # Si existe el checkpoint restaurar
216         if ckpt and ckpt.model_checkpoint_path:
217             saver.restore(session, ckpt.model_checkpoint_path)

```

```

217 #Si se desea entrenar
218 if(train):
219     #-----
220     # TRAIN
221     #-----
222
223     ITRS= int(len(X)/batch_size)
224     RES_ITMS = len(X) - (ITRS*batch_size)
225     if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
226
227     for e in range(epochs):
228         TRAIN_LOSS = 0.0
229
230         print("Epoch "+str(e))
231
232         for step in range(ITRS):
233
234             batch_inputs, batch_context =
235                 getBatchData(X,Y, batch_size ,step, ITRS, RES_ITMS)
236
237             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
238
239             ., loss_val, batch_loss, gs = session.run([train_step,
240                 nce_loss, batch_loss, global_step], feed_dict=feed_dict)
241
242             TRAIN_LOSS += np.sum(batch_loss)
243
244             #print('Epoch '+str(e)+', last batch loss: '+str(loss_val))
245
246             if (save): saver.save(session, complete_path, global_step=global_step)
247
248             #-----
249             # DEV
250             #-----
251
252             if(len(DX)>0):
253
254                 ITRS_TEST = int(len(DX) / batch_size)
255                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
256                 TEST_LOSS = 0.0
257                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se
258                     aade una iteracin
259
260                 for step in range(ITRS_TEST):
261
262                     batch_inputs, batch_context =
263                         getBatchData(DX,DY, batch_size ,step, ITRS_TEST, RES_ITMS_TEST)
264
265                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
266                     loss_val, batch_loss = session.run([nce_loss, batch_loss],
267                         feed_dict=feed_dict)
268
269                     TEST_LOSS+=np.sum(batch_loss)
270
271                 dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
272                 train_hist.append(TRAIN_LOSS / (len(X) * 1.0))

```

```

269
270     if (len(dev_hist) >= slope_size):
271         # Calcular la pendiente
272         print(getSlope(dev_hist[-slope_size:]))
273         print(dev_hist[-slope_size:])
274
275         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
276
277             def printArray(data, data2):
278                 for i in range(len(data)):
279                     a=data[i]
280                     b=data2[i]
281                     print(str(a).replace(".",",")+ "\t" + str(b).replace(".",","))
282
283             print("—" * 50)
284             print(e)
285             printArray(dev_hist, train_hist)
286             return
287 else:
288     #
289
290     # TEST FINAL
291     #
292
293     batch_size = 300
294
295     IN_TOP_5 = 0
296     IN_TOP_10 = 0
297     IN_TOP_20 = 0
298     IN_TOP_50 = 0
299     IN_TOP_100 = 0
300     IN_TOP_1000 = 0
301     real_positions = []
302
303     items = 0
304
305     ITRS_TEST = int(len(TX) / batch_size)
306     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
307     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se aaade una
308         oiteracin
309
310     for step in range(ITRS_TEST):
311         #tss = time.time()
312
313         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
314             ITRS_TEST, RES_ITMS_TEST)
315         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
316             batch_inputs, train_labels: batch_context})
317
318         for j in range(len(batch_context)):
319             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
320             real_positions.append(real_position)
321
322             IN_TOP_5 += 1 if real_position <= 5 else 0
323             IN_TOP_10 += 1 if real_position <= 10 else 0
324             IN_TOP_20 += 1 if real_position <= 20 else 0
325             IN_TOP_50 += 1 if real_position <= 50 else 0

```

```

321         IN_TOP_100 += 1 if real_position <= 100 else 0
322         IN_TOP_1000 += 1 if real_position <= 1000 else 0
323
324     items += len(batch_inputs)*1.0
325
326     PTF = IN_TOP_5 / items
327     PTT = IN_TOP_10 / items
328     PTTW = IN_TOP_20 / items
329     PTFT = IN_TOP_50 / items
330     PTCIEN = IN_TOP_100 / items
331     PTMIL = IN_TOP_1000 / items
332
333     print(np.mean(real_positions), np.std(real_positions),
334           np.median(real_positions))
335     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
336           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
337           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
338
339     #print(time.time()-tss)
340
341     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
342
343     #-----
344
345     #####
346     # Llamadas
347     #####
348
349     TEST_NAME="test_033"
350
351     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
352
353     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
354     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
355     D2V_DBWOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
356     W2V = np.load('embeddings/w2v/EMB_MATRIX')
357
358     #Se concatenan los embeddings de los usuarios para generar el perfil
359     D2V_EMB = np.concatenate((D2V_DM_COMPLETE, D2V_DBWOW_COMPLETE), axis=1) # El consolidado =
360     Consolidado DM + Consolidado DBOW
361
362     #Se cargan los datos de train de las canciones para conocer los IDs reales
363     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
364
365     #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
366     ALL = set(list(range(W2V.shape[0])))
367     REAL = set(list(set(list(map(int, DATA)))))
368
369     #####
370     # TENSORFLOW
371     #####
372     #-----
373     # Grid Search + Early Stopping
374     #-----
375     ' ' '

```

```

376 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
377 rates = [0.000001,1, 0.01, 0.0001]
378
379 for r in rates:
380     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
381
382     print("Learning Rate: "+str(r))
383     print("-" * 50)
384     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
385     print("-"*50)
386
387 exit()
388 '''
389 #
390 # Train
391 #
392 '''
393 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
394
395 #Para esta fase, se une TRAIN y DEV
396 X = np.concatenate((X,DX),axis=0)
397 Y = np.concatenate((Y,DY),axis=0)
398
399 DX=[]
400 DY=[]
401
402 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
403 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
404
405 exit()
406 '''
407 #
408 # Test
409 #
410
411 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
412
413 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
414 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.19. TEST 19

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:,len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
67     DEV = np.concatenate([DEV,USR_DEV])
68     TEST = np.concatenate([TEST,USR_TEST])
69
70 #Mezclar
71 if(shuffle):
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se noaadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
146
147     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
148
149     # Operaciones
150
151     # Embedding de documento
152     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
153
154     # Embedding de cancion 1
155     ec1 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])
156
157     # Embedding de cancion 2
158     ec2 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 2])
159
160     # Cálculo de LOSS y optimizador
161
162     embed = []
163     embed.append(ed)
164     embed.append(ec1)

```

```

164     embed.append(ec2)
165     embed = tf.concat(embed, 1)
166
167     # softmax weights, W and D vectors should be concatenated before applying softmax
168     weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
169         math.sqrt(song_size)), name="weights")
170     # softmax biases
171     biases = tf.Variable(tf.zeros([song_size]), name="biases")
172
173     # aClculo de LOSS y optimizador—
174     nce_loss = tf.nn.nce_loss(weights=weights,
175         biases=biases,
176         labels=train_labels,
177         inputs=embed,
178         num_sampled=num_sampled,
179         num_classes=song_size,
180         name="nce_loss")
181
182     batch_loss = nce_loss
183     nce_loss = tf.reduce_mean(nce_loss)
184
185     logits = tf.matmul(embed, tf.transpose(weights))
186     logits = tf.nn.bias_add(logits, biases)
187
188     top_v, top_i = tf.nn.top_k(logits, k=song_size)
189
190     # logits = tf.nn.softmax(logits, axis=0)
191
192     # optimizer =
193     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
194     train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
195         global_step=global_step)
196
197     # Inicializar variables
198     init = tf.global_variables_initializer()
199
200     # Crear objeto encargado de almacenar la red
201     saver = tf.train.Saver(max_to_keep=1)
202
203     config = tf.ConfigProto()
204     config.gpu_options.allow_growth = True
205     with tf.Session(graph=graph, config=config) as session:
206     #with tf.Session(graph=graph) as session:
207
208         # We must initialize all variables before we use them.
209         init.run()
210
211     if(save):
212         # Obtener el checkpoint
213         ckpt = tf.train.get_checkpoint_state(model_path)
214
215         # Si existe el checkpoint restaurar
216         if ckpt and ckpt.model_checkpoint_path:
217             saver.restore(session, ckpt.model_checkpoint_path)

```

```

217 #Si se desea entrenar
218 if(train):
219     #-----
220     # TRAIN
221     #-----
222
223     ITRS= int(len(X)/batch_size)
224     RES_ITMS = len(X) - (ITRS*batch_size)
225     if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
226
227     for e in range(epochs):
228         TRAIN_LOSS = 0.0
229
230         print("Epoch "+str(e))
231
232         for step in range(ITRS):
233
234             batch_inputs, batch_context =
235                 getBatchData(X,Y, batch_size ,step, ITRS, RES_ITMS)
236
237             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
238
239             ., loss_val, batch_lss, gs = session.run([train_step,
240                 nce_loss, batch_loss, global_step], feed_dict=feed_dict)
241
242             TRAIN_LOSS += np.sum(batch_lss)
243
244             #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
245
246             if (save): saver.save(session, complete_path, global_step=global_step)
247
248             #-----
249             # DEV
250             #-----
251
252             if(len(DX)>0):
253
254                 ITRS_TEST = int(len(DX) / batch_size)
255                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
256                 TEST_LOSS = 0.0
257                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se
258                     aade una iteracin
259
260                 for step in range(ITRS_TEST):
261
262                     batch_inputs, batch_context =
263                         getBatchData(DX,DY, batch_size ,step, ITRS_TEST, RES_ITMS_TEST)
264
265                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
266                     loss_val, batch_lss = session.run([nce_loss, batch_loss],
267                         feed_dict=feed_dict)
268
269                     TEST_LOSS+=np.sum(batch_lss)
270
271                 dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
272                 train_hist.append(TRAIN_LOSS / (len(X) * 1.0))

```

```

269
270     if (len(dev_hist) >= slope_size):
271         # Calcular la pendiente
272         print(getSlope(dev_hist[-slope_size:]))
273         print(dev_hist[-slope_size:])
274
275         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
276
277             def printArray(data, data2):
278                 for i in range(len(data)):
279                     a=data[i]
280                     b=data2[i]
281                     print(str(a).replace(".",",")+ "\t" + str(b).replace(".",","))
282
283             print("—" * 50)
284             print(e)
285             printArray(dev_hist, train_hist)
286             return
287 else:
288     #
289
290     # TEST FINAL
291     #
292
293     batch_size = 300
294
295     IN_TOP_5 = 0
296     IN_TOP_10 = 0
297     IN_TOP_20 = 0
298     IN_TOP_50 = 0
299     IN_TOP_100 = 0
300     IN_TOP_1000 = 0
301     real_positions = []
302
303     items = 0
304
305     ITRS_TEST = int(len(TX) / batch_size)
306     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
307     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se aaade una
308         oiteracin
309
310     for step in range(ITRS_TEST):
311         #tss = time.time()
312
313         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
314             ITRS_TEST, RES_ITMS_TEST)
315         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
316             batch_inputs, train_labels: batch_context})
317
318         for j in range(len(batch_context)):
319             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
320             real_positions.append(real_position)
321
322             IN_TOP_5 += 1 if real_position <= 5 else 0
323             IN_TOP_10 += 1 if real_position <= 10 else 0
324             IN_TOP_20 += 1 if real_position <= 20 else 0
325             IN_TOP_50 += 1 if real_position <= 50 else 0

```

```

321         IN_TOP_100 += 1 if real_position <= 100 else 0
322         IN_TOP_1000 += 1 if real_position <= 1000 else 0
323
324     items += len(batch_inputs)*1.0
325
326     PTF = IN_TOP_5 / items
327     PTT = IN_TOP_10 / items
328     PTW = IN_TOP_20 / items
329     PTFT = IN_TOP_50 / items
330     PTCIEN = IN_TOP_100 / items
331     PTMIL = IN_TOP_1000 / items
332
333     print(np.mean(real_positions), np.std(real_positions),
334           np.median(real_positions))
335     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
336           + str(PTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
337           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
338
339     #print(time.time()-tss)
340
341     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
342
343     #-----
344
345     #####
346     # Llamadas
347     #####
348
349     TEST_NAME="test_043"
350
351     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
352
353     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
354     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
355     D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
356     W2V = np.load('embeddings/w2v/EMB_MATRIX')
357
358     #Se concatenan los embeddings de los usuarios para generar el perfil
359     D2V_EMB = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El reciente = Reciente
360     DM + Reciente DBOW
361
362     #Se cargan los datos de train de las canciones para conocer los IDs reales
363     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
364
365     #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
366     ALL = set(list(range(W2V.shape[0])))
367     REAL = set(list(set(list(map(int, DATA)))))
368
369     #####
370     # TENSORFLOW
371     #####
372     #-----
373     # Grid Search + Early Stopping
374     #-----
375     ' '

```

```

376 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
377 rates = [0.000001,1, 0.01, 0.0001]
378
379 for r in rates:
380     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
381
382     print("Learning Rate: "+str(r))
383     print("-" * 50)
384     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
385     print("-"*50)
386
387 exit()
388 '''
389 #
390 # Train
391 #
392 '''
393 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
394
395 #Para esta fase, se une TRAIN y DEV
396 X = np.concatenate((X,DX),axis=0)
397 Y = np.concatenate((Y,DY),axis=0)
398
399 DX=[]
400 DY=[]
401
402 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
403 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
404
405 exit()
406 '''
407 #
408 # Test
409 #
410
411 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
412
413 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
414 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.20. TEST 20

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:,len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
67     DEV = np.concatenate([DEV,USR_DEV])
68     TEST = np.concatenate([TEST,USR_TEST])
69
70 #Mezclar
71 if(shuffle):
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se ha aadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X,Y, DX,DY, TX, TY, config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112 user_concat_size = user_emb_size * 2
113
114 num_sampled = 64 # Ejemplos negativos (Para el NSE)
115
116 model_name = 'model'
117 model_path = 'models/' + TEST_NAME + '/model'
118 complete_path = model_path + "/" + model_name
119
120 # Se almacena la loss de dev de las 5 ultimas epochs
121 slope_size = 100
122 train_hist = []
123 dev_hist = []
124
125 # Creacin del grafo de TF.
126 graph = tf.Graph()
127
128 with graph.as_default():
129
130     if(config["seed"]!=1):
131         tf.set_random_seed(config["seed"])
132
133     # Numero global de iteraciones
134     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
135
136     # Datos de entrada
137
138     # Array del tamaño del batch con las X
139     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
140
141     # Array del tamaño del batch con las Y asociadas
142     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
143
144     # Embeddings
145
146     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
147
148     T0 = tf.Variable(tf.truncated_normal([user_concat_size, user_emb_size], mean=0.0,
149         stddev=1.0 / math.sqrt(hidden_size)), name="T0")
150     B0 = tf.Variable(tf.zeros([user_emb_size]), name="B0")
151
152     W2V_EMB_VRB = tf.Variable(W2V, trainable=False, name="word_embeddings")
153
154     # Operaciones
155
156     # Embedding de documento
157     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
158
159     #Transformar concat documento de 256 a 128
160     ed = tf.matmul(ed, T0)+B0
161
162     # Embedding de cancion 1
163     ec1 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 1])

```

```

164 ec2 = tf.nn.embedding_lookup(W2V_EMB_VRB, train_dataset[:, 2])
165
166 # aClculo de LOSS y optimizador—
167
168 embed = []
169 embed.append(ed)
170 embed.append(ec1)
171 embed.append(ec2)
172 embed = tf.concat(embed, 1)
173
174 # softmax weights, W and D vectors should be concatenated before applying softmax
175 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
176     math.sqrt(song_size)), name="weights")
177 # softmax biases
178 biases = tf.Variable(tf.zeros([song_size]), name="biases")
179
180 # aClculo de LOSS y optimizador—
181
182 nce_loss = tf.nn.nce_loss(weights=weights,
183     biases=biases,
184     labels=train_labels,
185     inputs=embed,
186     num_sampled=num_sampled,
187     num_classes=song_size,
188     name="nce_loss")
189
190 batch_loss = nce_loss
191 nce_loss = tf.reduce_mean(nce_loss)
192
193 logits = tf.matmul(embed, tf.transpose(weights))
194 logits = tf.nn.bias_add(logits, biases)
195
196 top_v, top_i = tf.nn.top_k(logits, k=song_size)
197
198 # logits = tf.nn.softmax(logits, axis=0)
199
200 # optimizer =
201     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
202 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
203     global_step=global_step)
204
205 # Inicializar variables
206 init = tf.global_variables_initializer()
207
208 # Crear objeto encargado de almacenar la red
209 saver = tf.train.Saver(max_to_keep=1)
210
211 config = tf.ConfigProto()
212 config.gpu_options.allow_growth = True
213 with tf.Session(graph=graph, config=config) as session:
214     #with tf.Session(graph=graph) as session:
215
216     # We must initialize all variables before we use them.
217     init.run()

```

```

216     if(save):
217         # Obtener el checkpoint
218         ckpt = tf.train.get_checkpoint_state(model_path)
219
220         # Si existe el checkpoint restaurar
221         if ckpt and ckpt.model_checkpoint_path:
222             saver.restore(session, ckpt.model_checkpoint_path)
223
224     #Si se desea entrenar
225     if(train):
226         #-----
227         # TRAIN
228         #-----
229
230         ITRS= int(len(X)/batch_size)
231         RES_ITMS = len(X) - (ITRS*batch_size)
232         if(RES_ITMS>0):ITRS+=1 #Si no es divisible exacta, se añade una iteración
233
234         for e in range(epochs):
235             TRAIN_LOSS = 0.0
236
237             print("Epoch "+str(e))
238
239             for step in range(ITRS):
240
241                 batch_inputs, batch_context =
242                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
243
244                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
245
246                 _, loss_val, batch_loss, gs = session.run([train_step,
247                     nce_loss,batch_loss,global_step], feed_dict=feed_dict)
248
249                 TRAIN_LOSS += np.sum(batch_loss)
250
251                 #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
252
253             if (save): saver.save(session, complete_path, global_step=global_step)
254
255             #-----
256             # DEV
257             #-----
258
259             if(len(DX)>0):
260
261                 ITRS_TEST = int(len(DX) / batch_size)
262                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
263                 TEST_LOSS = 0.0
264                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisible exacta, se
265                     añade una iteración
266
267                 for step in range(ITRS_TEST):
268
269                     batch_inputs, batch_context =
270                         getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)

```

```

269         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
270         loss_val, batch_loss = session.run([nce_loss, batch_loss],
271                                           feed_dict=feed_dict)
272
273         TEST_LOSS += np.sum(batch_loss)
274
275         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
276         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
277
278     if (len(dev_hist) >= slope_size):
279         # Calcular la pendiente
280         print(getSlope(dev_hist[-slope_size:]))
281         print(dev_hist[-slope_size:])
282
283         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
284
285             def printArray(data, data2):
286                 for i in range(len(data)):
287                     a=data[i]
288                     b=data2[i]
289                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
290
291             print("-" * 50)
292             print(e)
293             printArray(dev_hist, train_hist)
294             return
295
296 else:
297     #
298
299     # TEST FINAL
300     #
301
302     batch_size = 300
303
304     IN_TOP_5 = 0
305     IN_TOP_10 = 0
306     IN_TOP_20 = 0
307     IN_TOP_50 = 0
308     IN_TOP_100 = 0
309     IN_TOP_1000 = 0
310     real_positions = []
311
312     items = 0
313
314     ITRS_TEST = int(len(TX) / batch_size)
315     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
316     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
317     iteracin
318
319     for step in range(ITRS_TEST):
320         #tss = time.time()
321
322         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
323                                                    ITRS_TEST, RES_ITMS_TEST)
324         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
325                                batch_inputs, train_labels: batch_context})

```

```

320     for j in range(len(batch_context)):
321         real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
322         real_positions.append(real_position)
323
324         IN_TOP_5 += 1 if real_position <= 5 else 0
325         IN_TOP_10 += 1 if real_position <= 10 else 0
326         IN_TOP_20 += 1 if real_position <= 20 else 0
327         IN_TOP_50 += 1 if real_position <= 50 else 0
328         IN_TOP_100 += 1 if real_position <= 100 else 0
329         IN_TOP_1000 += 1 if real_position <= 1000 else 0
330
331     items += len(batch_inputs)*1.0
332
333     PTF = IN_TOP_5 / items
334     PTT = IN_TOP_10 / items
335     PTTW = IN_TOP_20 / items
336     PTFT = IN_TOP_50 / items
337     PTCIEN = IN_TOP_100 / items
338     PTMIL = IN_TOP_1000 / items
339
340     print(np.mean(real_positions), np.std(real_positions),
341           np.median(real_positions))
342     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
343           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
344           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
345
346     #print(time.time()-tss)
347
348     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
349
350     #-----
351
352     #####
353     # Llamadas
354     #####
355
356     TEST_NAME="test_053"
357
358     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
359
360     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
361     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
362     D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
363
364     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
365     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
366     D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
367
368     W2V = np.load('embeddings/w2v/EMB_MATRIX')
369
370     #Se concatenan los embeddings de los usuarios para generar el perfil
371     D2V_CONS = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
372     Consolidado DM + Consolidado DBOW
373     D2V_PRES = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El presente = Presente
374     DM + Presente DBOW

```

```

374 D2V_EMB = np.concatenate((D2V_CONS,D2V_PRES),axis=1) # El perfil completo CONS + PRES
375
376 #Se cargan los datos de train de las canciones para conocer los IDs reales
377 DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
378
379 #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
380 ALL = set(list(range(W2V.shape[0])))
381 REAL = set(list(set(list(map(int, DATA)))))
382
383 #####
384 # TENSORFLOW
385 #####
386
387 # -----
388 # Grid Search + Early Stopping
389 # -----
390 '''
391 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
392 rates = [0.000001,1, 0.01, 0.0001]
393
394 for r in rates:
395     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
396
397     print("Learning Rate: "+str(r))
398     print("-" * 50)
399     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
400     print("-"*50)
401
402 exit()
403 '''
404 # -----
405 # Train
406 # -----
407 '''
408 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
409
410 #Para esta fase, se une TRAIN y DEV
411 X = np.concatenate((X,DX),axis=0)
412 Y = np.concatenate((Y,DY),axis=0)
413
414 DX=[]
415 DY=[]
416
417 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":167, "seed":100}
418 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
419
420 exit()
421 '''
422 # -----
423 # Test
424 # -----
425
426 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
427
428 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":167, "seed":100}
429 trainNet(X, Y, DX, DY, TX, TY, config, user_size,train=False, save=True)

```



## 4.21. TEST 21

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:,len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
67     DEV = np.concatenate([DEV,USR_DEV])
68     TEST = np.concatenate([TEST,USR_TEST])
69
70 #Mezclar
71 if(shuffle):
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se noaadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X,Y, DX,DY, TX, TY, config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     O1 = tf.Variable(tf.truncated_normal([user_size, user_emb_size], mean=0.0, stddev=1.0
146     / math.sqrt(user_emb_size)), name="O1")
147     B01 = tf.Variable(tf.zeros([user_emb_size]), name="B01")
148
149     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
150     / math.sqrt(song_emb_size)), name="O2")
151     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
152
153     # Operaciones
154
155     # Embedding de documento
156     ed = tf.nn.embedding_lookup(O1, train_dataset[:, 0])
157     ed = ed + B01
158
159     # Embedding de cancion 1
160     ec1 = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
161     ec1 = ec1 + B02
162
163     # Embedding de cancion 2
164     ec2 = tf.nn.embedding_lookup(O2, train_dataset[:, 2])
165     ec2 = ec2 + B02

```

```

163
164 # aClculo de LOSS y optimizador—
165
166 embed = []
167 embed.append(ed)
168 embed.append(ec1)
169 embed.append(ec2)
170 embed = tf.concat(embed, 1)
171
172 # softmax weights, W and D vectors should be concatenated before applying softmax
173 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
174     math.sqrt(song_size)), name="weights")
175 # softmax biases
176 biases = tf.Variable(tf.zeros([song_size]), name="biases")
177
178 # aClculo de LOSS y optimizador—
179
180 nce_loss = tf.nn.nce_loss(weights=weights,
181     biases=biases,
182     labels=train_labels,
183     inputs=embed,
184     num_sampled=num_sampled,
185     num_classes=song_size,
186     name="nce_loss")
187
188 batch_loss = nce_loss
189 nce_loss = tf.reduce_mean(nce_loss)
190
191 logits = tf.matmul(embed, tf.transpose(weights))
192 logits = tf.nn.bias_add(logits, biases)
193
194 top_v, top_i = tf.nn.top_k(logits, k=song_size)
195
196 # logits = tf.nn.softmax(logits, axis=0)
197
198 # optimizer =
199     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
200 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
201     global_step=global_step)
202
203 # Inicializar variables
204 init = tf.global_variables_initializer()
205
206 # Crear objeto encargado de almacenar la red
207 saver = tf.train.Saver(max_to_keep=1)
208
209 config = tf.ConfigProto()
210 config.gpu_options.allow_growth = True
211 with tf.Session(graph=graph, config=config) as session:
212 #with tf.Session(graph=graph) as session:
213
214     # We must initialize all variables before we use them.
215     init.run()
216
217 if(save):

```

```

215     # Obtener el checkpoint
216     ckpt = tf.train.get_checkpoint_state(model_path)
217
218     # Si existe el checkpoint restaurar
219     if ckpt and ckpt.model_checkpoint_path:
220         saver.restore(session, ckpt.model_checkpoint_path)
221
222     #Si se desea entrenar
223     if(train):
224         #-----
225         # TRAIN
226         #-----
227
228         ITRS= int(len(X)/batch_size)
229         RES_ITMS = len(X) - (ITRS*batch_size)
230         if(RES_ITMS>0):ITRS+=1 #Si no es odivisin exacta, se naade una oiteracin
231
232         for e in range(epochs):
233             TRAIN_LOSS = 0.0
234
235             print("Epoch "+str(e))
236
237             for step in range(ITRS):
238
239                 batch_inputs, batch_context =
240                     getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
241
242                 feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
243
244                 -, loss_val, batch_loss, gs = session.run([train_step,
245                     nce_loss,batch_loss,global_step], feed_dict=feed_dict)
246
247                 TRAIN_LOSS += np.sum(batch_loss)
248
249                 #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
250
251                 if (save): saver.save(session, complete_path, global_step=global_step)
252
253                 #-----
254                 # DEV
255                 #-----
256
257                 if(len(DX)>0):
258
259                     ITRS_TEST = int(len(DX) / batch_size)
260                     RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
261                     TEST_LOSS = 0.0
262                     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se
263                         naade una oiteracin
264
265                     for step in range(ITRS_TEST):
266
267                         batch_inputs, batch_context =
268                             getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
269
270                         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}

```

```

268         loss_val, batch_loss = session.run([nce_loss, batch_loss],
269             feed_dict=feed_dict)
270
271         TEST_LOSS += np.sum(batch_loss)
272
273         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
274         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
275
276     if (len(dev_hist) >= slope_size):
277         # Calcular la pendiente
278         print(getSlope(dev_hist[-slope_size:]))
279         print(dev_hist[-slope_size:])
280
281         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
282
283             def printArray(data, data2):
284                 for i in range(len(data)):
285                     a=data[i]
286                     b=data2[i]
287                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
288
289             print("-" * 50)
290             print(e)
291             printArray(dev_hist, train_hist)
292             return
293
294 else:
295     #
296
297     # TEST FINAL
298     #
299
300     batch_size = 300
301
302     IN_TOP_5 = 0
303     IN_TOP_10 = 0
304     IN_TOP_20 = 0
305     IN_TOP_50 = 0
306     IN_TOP_100 = 0
307     IN_TOP_1000 = 0
308     real_positions = []
309
310     items = 0
311
312     ITRS_TEST = int(len(TX) / batch_size)
313     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
314     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
315     iteracin
316
317     for step in range(ITRS_TEST):
318         #tss = time.time()
319
320         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
321             ITRS_TEST, RES_ITMS_TEST)
322         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
323             batch_inputs, train_labels: batch_context})
324
325     for j in range(len(batch_context)):

```

```

319         real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
320         real_positions.append(real_position)
321
322         IN_TOP_5 += 1 if real_position <= 5 else 0
323         IN_TOP_10 += 1 if real_position <= 10 else 0
324         IN_TOP_20 += 1 if real_position <= 20 else 0
325         IN_TOP_50 += 1 if real_position <= 50 else 0
326         IN_TOP_100 += 1 if real_position <= 100 else 0
327         IN_TOP_1000 += 1 if real_position <= 1000 else 0
328
329         items += len(batch_inputs)*1.0
330
331         PTF = IN_TOP_5 / items
332         PTT = IN_TOP_10 / items
333         PTTW = IN_TOP_20 / items
334         PTFT = IN_TOP_50 / items
335         PTCIEN = IN_TOP_100 / items
336         PTMIL = IN_TOP_1000 / items
337
338         print(np.mean(real_positions), np.std(real_positions),
339               np.median(real_positions))
340         print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
341               + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
342               + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
343
344         #print(time.time()-tss)
345
346         print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
347
348         #-----
349
350         #####
351         # Llamadas
352         #####
353
354         TEST_NAME="test_013"
355
356         os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
357
358         #Solamente se carga para conocer el unmero de canciones
359         W2V = np.load('embeddings/w2v/EMB_MATRIX')
360
361         #Se cargan los datos de train de las canciones para conocer los IDs reales
362         DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
363
364         #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
365         ALL = set(list(range(W2V.shape[0])))
366         REAL = set(list(set(list(map(int, DATA)))))
367
368         #####
369         # TENSORFLOW
370         #####
371
372         #-----
373         # Grid Search + Early Stopping
374         #-----

```

```

375 '''
376 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL))
377 rates = [0.000001,1, 0.01, 0.0001]
378
379
380 for r in rates:
381     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
382
383     print("Learning Rate: "+str(r))
384     print("-" * 50)
385     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
386     print("-"*50)
387
388 exit()
389 '''
390 #-----
391 # Train
392 #-----
393 '''
394 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL))
395
396 #Para esta fase, se une TRAIN y DEV
397 X = np.concatenate((X,DX),axis=0)
398 Y = np.concatenate((Y,DY),axis=0)
399
400 DX=[]
401 DY=[]
402
403 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
404 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
405
406 exit()
407 '''
408 #-----
409 # Test
410 #-----
411
412 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL))
413
414 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
415 trainNet(X, Y, DX, DY, TX, TY, config, user_size,train=False, save=True)

```



## 4.22. TEST 22

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:,len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
67     DEV = np.concatenate([DEV,USR_DEV])
68     TEST = np.concatenate([TEST,USR_TEST])
69
70 #Mezclar
71 if(shuffle):
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se ha aadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X,Y, DX,DY, TX, TY, config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
146
147     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
148         / math.sqrt(song_emb_size)), name="O2")
149     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
150
151     # Operaciones
152
153     # Embedding de documento
154     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
155
156     # Embedding de cancion 1
157     ec1 = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
158     ec1 = ec1 + B02
159
160     # Embedding de cancion 2
161     ec2 = tf.nn.embedding_lookup(O2, train_dataset[:, 2])
162     ec2 = ec2 + B02
163
164     # Cálculo de LOSS y optimizador—

```

```

163
164 embed = []
165 embed.append(ed)
166 embed.append(ec1)
167 embed.append(ec2)
168 embed = tf.concat(embed, 1)
169
170 # softmax weights, W and D vectors should be concatenated before applying softmax
171 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
172     math.sqrt(song_size)), name="weights")
173 # softmax biases
174 biases = tf.Variable(tf.zeros([song_size]), name="biases")
175
176 # aClculo de LOSS y optimizador—
177
178 nce_loss = tf.nn.nce_loss(weights=weights,
179     biases=biases,
180     labels=train_labels,
181     inputs=embed,
182     num_sampled=num_sampled,
183     num_classes=song_size,
184     name="nce_loss")
185
186 batch_loss = nce_loss
187 nce_loss = tf.reduce_mean(nce_loss)
188
189 logits = tf.matmul(embed, tf.transpose(weights))
190 logits = tf.nn.bias_add(logits, biases)
191
192 top_v, top_i = tf.nn.top_k(logits, k=song_size)
193
194 # logits = tf.nn.softmax(logits, axis=0)
195
196 # optimizer =
197     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
198 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
199     global_step=global_step)
200
201 # Inicializar variables
202 init = tf.global_variables_initializer()
203
204 # Crear objeto encargado de almacenar la red
205 saver = tf.train.Saver(max_to_keep=1)
206
207 config = tf.ConfigProto()
208 config.gpu_options.allow_growth = True
209 with tf.Session(graph=graph, config=config) as session:
210 #with tf.Session(graph=graph) as session:
211
212     # We must initialize all variables before we use them.
213     init.run()
214
215 if(save):
216     # Obtener el checkpoint
217     ckpt = tf.train.get_checkpoint_state(model_path)

```

```

216     # Si existe el checkpoint restaurar
217     if ckpt and ckpt.model_checkpoint_path:
218         saver.restore(session, ckpt.model_checkpoint_path)
219
220 #Si se desea entrenar
221 if(train):
222     #
223     # TRAIN
224     #
225
226     ITRS= int(len(X)/batch_size)
227     RES_ITMS = len(X) - (ITRS*batch_size)
228     if(RES_ITMS>0):ITRS+=1 #Si no es odivisin exacta, se naade una oiteracin
229
230     for e in range(epochs):
231         TRAIN_LOSS = 0.0
232
233         print("Epoch "+str(e))
234
235         for step in range(ITRS):
236
237             batch_inputs, batch_context =
238                 getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
239
240             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
241
242             ., loss_val, batch_lss, gs = session.run([train_step,
243                 nce_loss,batch_loss,global_step], feed_dict=feed_dict)
244
245             TRAIN_LOSS += np.sum(batch_lss)
246
247             #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
248
249             if (save): saver.save(session, complete_path, global_step=global_step)
250
251             #
252             # DEV
253             #
254
255             if(len(DX)>0):
256
257                 ITRS_TEST = int(len(DX) / batch_size)
258                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
259                 TEST_LOSS = 0.0
260                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se
261                     naade una oiteracin
262
263                 for step in range(ITRS_TEST):
264
265                     batch_inputs, batch_context =
266                         getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
267
268                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
269                     loss_val,batch_lss = session.run([nce_loss,batch_loss],
270                         feed_dict=feed_dict)

```

```

268         TEST_LOSS+=np.sum(batch_loss)
269
270     dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
271     train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
272
273     if (len(dev_hist) >= slope_size):
274         # Calcular la pendiente
275         print(getSlope(dev_hist[-slope_size:]))
276         print(dev_hist[-slope_size:])
277
278         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
279
280             def printArray(data, data2):
281                 for i in range(len(data)):
282                     a=data[i]
283                     b=data2[i]
284                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
285
286             print("-" * 50)
287             print(e)
288             printArray(dev_hist, train_hist)
289             return
290 else:
291     #
292
293     # TEST FINAL
294     #
295
296     batch_size = 300
297
298     IN_TOP_5 = 0
299     IN_TOP_10 = 0
300     IN_TOP_20 = 0
301     IN_TOP_50 = 0
302     IN_TOP_100 = 0
303     IN_TOP_1000 = 0
304     real_positions = []
305
306     items = 0
307
308     ITRS_TEST = int(len(TX) / batch_size)
309     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
310     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
311     iteracin
312
313     for step in range(ITRS_TEST):
314         #tss = time.time()
315
316         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
317             ITRS_TEST, RES_ITMS_TEST)
318         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
319             batch_inputs, train_labels: batch_context})
320
321         for j in range(len(batch_context)):
322             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
323             real_positions.append(real_position)

```

```

320     IN_TOP_5 += 1 if real_position <= 5 else 0
321     IN_TOP_10 += 1 if real_position <= 10 else 0
322     IN_TOP_20 += 1 if real_position <= 20 else 0
323     IN_TOP_50 += 1 if real_position <= 50 else 0
324     IN_TOP_100 += 1 if real_position <= 100 else 0
325     IN_TOP_1000 += 1 if real_position <= 1000 else 0
326
327     items += len(batch_inputs)*1.0
328
329     PTF = IN_TOP_5 / items
330     PTT = IN_TOP_10 / items
331     PTTW = IN_TOP_20 / items
332     PTFT = IN_TOP_50 / items
333     PTCIEN = IN_TOP_100 / items
334     PTMIL = IN_TOP_1000 / items
335
336     print(np.mean(real_positions), np.std(real_positions),
337           np.median(real_positions))
338     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
339           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
340           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
341
342     #print(time.time()-tss)
343
344     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
345
346     #-----
347
348     #####
349     # Llamadas
350     #####
351
352     TEST_NAME="test_063"
353
354     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
355
356     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
357     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
358     D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
359     W2V = np.load('embeddings/w2v/EMB_MATRIX')
360
361     #Se concatenan los embeddings de los usuarios para generar el perfil
362     D2V_EMB = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =
363     Consolidado DM + Consolidado DBOW
364
365     #Se cargan los datos de train de las canciones para conocer los IDs reales
366     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
367
368     #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
369     ALL = set(list(range(W2V.shape[0])))
370     REAL = set(list(set(list(map(int, DATA)))))
371
372     #####
373     # TENSORFLOW
374     #####

```

```

375 #
376 # Grid Search + Early Stopping
377 #
378 '''
379 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL
380 rates = [0.000001,1, 0.01, 0.0001]
381
382 for r in rates:
383     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
384
385     print("Learning Rate: "+str(r))
386     print("-" * 50)
387     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
388     print("-"*50)
389
390 exit()
391 '''
392 #
393 # Train
394 #
395 '''
396 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL
397
398 #Para esta fase, se une TRAIN y DEV
399 X = np.concatenate((X,DX),axis=0)
400 Y = np.concatenate((Y,DY),axis=0)
401
402 DX=[]
403 DY=[]
404
405 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":597, "seed":100}
406 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
407
408 exit()
409 '''
410 #
411 # Test
412 #
413
414 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL
415
416 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":597, "seed":100}
417 trainNet(X, Y, DX, DY, TX, TY, config, user_size,train=False, save=True)

```



## 4.23. TEST 23

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:,len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
67     DEV = np.concatenate([DEV,USR_DEV])
68     TEST = np.concatenate([TEST,USR_TEST])
69
70 #Mezclar
71 if(shuffle):
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se ha aadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X,Y, DX,DY, TX, TY, config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112
113 num_sampled = 64 # Ejemplos negativos (Para el NSE)
114
115 model_name = 'model'
116 model_path = 'models/' + TEST_NAME + '/model'
117 complete_path = model_path + "/" + model_name
118
119 # Se almacena la loss de dev de las 5 ultimas epochs
120 slope_size = 100
121 train_hist = []
122 dev_hist = []
123
124 # Creacin del grafo de TF.
125 graph = tf.Graph()
126
127 with graph.as_default():
128
129     if(config["seed"]!=1):
130         tf.set_random_seed(config["seed"])
131
132     # nmero global de iteraciones
133     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
134
135     # Datos de entrada
136
137     # Array del tamaño del batch con las X
138     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
139
140     # Array del tamaño del batch con las Y asociadas
141     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
142
143     # Embeddings
144
145     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
146
147     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
148         / math.sqrt(song_emb_size)), name="O2")
149     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
150
151     # Operaciones
152
153     # Embedding de documento
154     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
155
156     # Embedding de cancion 1
157     ec1 = tf.nn.embedding_lookup(O2, train_dataset[:, 1])
158     ec1 = ec1 + B02
159
160     # Embedding de cancion 2
161     ec2 = tf.nn.embedding_lookup(O2, train_dataset[:, 2])
162     ec2 = ec2 + B02
163
164     # Cálculo de LOSS y optimizador—

```

```

164
165 embed = []
166 embed.append(ed)
167 embed.append(ec1)
168 embed.append(ec2)
169 embed = tf.concat(embed, 1)
170
171 # softmax weights, W and D vectors should be concatenated before applying softmax
172 weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
173     math.sqrt(song_size)), name="weights")
174 # softmax biases
175 biases = tf.Variable(tf.zeros([song_size]), name="biases")
176
177 # aClculo de LOSS y optimizador—
178
179 nce_loss = tf.nn.nce_loss(weights=weights,
180     biases=biases,
181     labels=train_labels,
182     inputs=embed,
183     num_sampled=num_sampled,
184     num_classes=song_size,
185     name="nce_loss")
186
187 batch_loss = nce_loss
188 nce_loss = tf.reduce_mean(nce_loss)
189
190 logits = tf.matmul(embed, tf.transpose(weights))
191 logits = tf.nn.bias_add(logits, biases)
192
193 top_v, top_i = tf.nn.top_k(logits, k=song_size)
194
195 # logits = tf.nn.softmax(logits, axis=0)
196
197 # optimizer =
198     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
199 train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
200     global_step=global_step)
201
202 # Inicializar variables
203 init = tf.global_variables_initializer()
204
205 # Crear objeto encargado de almacenar la red
206 saver = tf.train.Saver(max_to_keep=1)
207
208 config = tf.ConfigProto()
209 config.gpu_options.allow_growth = True
210 with tf.Session(graph=graph, config=config) as session:
211     #with tf.Session(graph=graph) as session:
212
213     # We must initialize all variables before we use them.
214     init.run()
215
216 if(save):
217     # Obtener el checkpoint
218     ckpt = tf.train.get_checkpoint_state(model_path)

```

```

216
217     # Si existe el checkpoint restaurar
218     if ckpt and ckpt.model_checkpoint_path:
219         saver.restore(session, ckpt.model_checkpoint_path)
220
221 #Si se desea entrenar
222 if(train):
223     #
224     # TRAIN
225     #
226
227     ITRS= int(len(X)/batch_size)
228     RES_ITMS = len(X) - (ITRS*batch_size)
229     if(RES_ITMS>0):ITRS+=1 #Si no es odivisin exacta, se naade una oiteracin
230
231     for e in range(epochs):
232         TRAIN_LOSS = 0.0
233
234         print("Epoch "+str(e))
235
236         for step in range(ITRS):
237
238             batch_inputs, batch_context =
239                 getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
240
241             feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
242
243             -, loss_val, batch_lss, gs = session.run([train_step,
244                 nce_loss,batch_loss,global_step], feed_dict=feed_dict)
245
246             TRAIN_LOSS += np.sum(batch_lss)
247
248             #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
249
250             if (save): saver.save(session, complete_path, global_step=global_step)
251
252             #
253             # DEV
254             #
255
256             if(len(DX)>0):
257
258                 ITRS_TEST = int(len(DX) / batch_size)
259                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
260                 TEST_LOSS = 0.0
261                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es odivisin exacta, se
262                     naade una oiteracin
263
264                 for step in range(ITRS_TEST):
265
266                     batch_inputs, batch_context =
267                         getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
268
269                     feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
270                     loss_val,batch_lss = session.run([nce_loss,batch_loss],
271                         feed_dict=feed_dict)

```

```

268
269         TEST_LOSS+=np.sum(batch_loss)
270
271         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
272         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
273
274     if (len(dev_hist) >= slope_size):
275         # Calcular la pendiente
276         print(getSlope(dev_hist[-slope_size:]))
277         print(dev_hist[-slope_size:])
278
279         if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
280
281             def printArray(data, data2):
282                 for i in range(len(data)):
283                     a=data[i]
284                     b=data2[i]
285                     print(str(a).replace(".", ",")+"\t"+str(b).replace(".", ","))
286
287             print("-" * 50)
288             print(e)
289             printArray(dev_hist, train_hist)
290             return
291 else:
292     #
293
294     # TEST FINAL
295     #
296
297     batch_size = 300
298
299     IN_TOP_5 = 0
300     IN_TOP_10 = 0
301     IN_TOP_20 = 0
302     IN_TOP_50 = 0
303     IN_TOP_100 = 0
304     IN_TOP_1000 = 0
305     real_positions = []
306
307     items = 0
308
309     ITRS_TEST = int(len(TX) / batch_size)
310     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
311     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
312     iteracin
313
314     for step in range(ITRS_TEST):
315         #tss = time.time()
316
317         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,
318             ITRS_TEST, RES_ITMS_TEST)
319         top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
320             batch_inputs, train_labels: batch_context})
321
322         for j in range(len(batch_context)):
323             real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
324             real_positions.append(real_position)

```

```

320
321     IN_TOP_5 += 1 if real_position <= 5 else 0
322     IN_TOP_10 += 1 if real_position <= 10 else 0
323     IN_TOP_20 += 1 if real_position <= 20 else 0
324     IN_TOP_50 += 1 if real_position <= 50 else 0
325     IN_TOP_100 += 1 if real_position <= 100 else 0
326     IN_TOP_1000 += 1 if real_position <= 1000 else 0
327
328     items += len(batch_inputs)*1.0
329
330     PTF = IN_TOP_5 / items
331     PTT = IN_TOP_10 / items
332     PTTW = IN_TOP_20 / items
333     PTFT = IN_TOP_50 / items
334     PTCIEN = IN_TOP_100 / items
335     PTMIL = IN_TOP_1000 / items
336
337     print(np.mean(real_positions), np.std(real_positions),
338           np.median(real_positions))
339     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
340           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
341           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
342
343     #print(time.time()-tss)
344
345     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
346
347     #-----
348
349     #####
350     # Llamadas
351     #####
352
353     TEST_NAME="test_073"
354
355     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
356
357     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
358     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
359     D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
360     W2V = np.load('embeddings/w2v/EMB_MATRIX')
361
362     #Se concatenan los embeddings de los usuarios para generar el perfil
363     D2V_EMB = np.concatenate((D2V_DM_PRESENT, D2V_DBOW_PRESENT), axis=1) # El reciente = Reciente
364     DM + Reciente DBOW
365
366     #Se cargan los datos de train de las canciones para conocer los IDs reales
367     DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
368
369     #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
370     ALL = set(list(range(W2V.shape[0])))
371     REAL = set(list(set(list(map(int, DATA)))))
372
373     #####
374     # TENSORFLOW
375     #####

```

```

375
376 #
377 # Grid Search + Early Stopping
378 #
379 '''
380 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL))
381 rates = [0.000001,1, 0.01, 0.0001]
382
383 for r in rates:
384     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
385
386     print("Learning Rate: "+str(r))
387     print("-" * 50)
388     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
389     print("-"*50)
390
391 exit()
392 '''
393 #
394 # Train
395 #
396 '''
397 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL))
398
399 #Para esta fase, se une TRAIN y DEV
400 X = np.concatenate((X,DX),axis=0)
401 Y = np.concatenate((Y,DY),axis=0)
402
403 DX=[]
404 DY=[]
405
406 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
407 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
408
409 exit()
410 '''
411 #
412 # Test
413 #
414
415 X, Y, DX, DY, TX, TY, user_size = getData-REAL_SONG_IDS=list-REAL))
416
417 config = {"batch_size": 1024, "learning_rate": 0.01, "epochs":1000, "seed":100}
418 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=False, save=True)

```



## 4.24. TEST 24

```

1  # -*- coding: utf-8 -*-
2
3  import os.path
4  import os
5  import tensorflow as tf
6  import numpy as np
7  import random
8  import time
9  import math
10 import pandas as pd
11 import pickle
12 from collections import Counter
13 from scipy.stats import linregress
14 from tensorflow.python.framework import ops as framework_ops
15
16 #####
17 # eMtodos
18 #####
19
20 def getData(REAL_SONG_IDS=[],data = "datos/Append/ARRAY_ARRAYS_REC.pkl",shuffle=True,
21            seed=100, train=.7):
22
23     DATA = np.load(data)
24
25     #DATA = DATA[:5];print("ELIMINAR FILTRO")
26
27     #User, cancion 1, cancion 2, cancion 3
28     TRAIN = np.empty([0, 4], dtype=np.int32)
29     DEV = np.empty([0, 4], dtype=np.int32)
30     TEST = np.empty([0, 4], dtype=np.int32)
31
32     for u in range(len(DATA)):
33
34         # Eliminar canciones escuchadas que no tienen embedding
35         USER_PLAYS = DATA[u]
36         USER_PLAYS = pd.DataFrame(list(map(int, USER_PLAYS)))
37         USER_PLAYS = USER_PLAYS.loc[USER_PLAYS[0].isin(REAL_SONG_IDS)][0].values
38
39         #Separar en train dev test
40         TRAIN_ITEMS = int(len(USER_PLAYS)*train)
41         DEV_ITEMS = int((len(USER_PLAYS)-TRAIN_ITEMS)/2)
42
43         TRAIN_DATA = USER_PLAYS[:TRAIN_ITEMS]
44         DEV_DATA = USER_PLAYS[TRAIN_ITEMS:TRAIN_ITEMS+DEV_ITEMS]
45         TEST_DATA = USER_PLAYS[TRAIN_ITEMS+DEV_ITEMS:]
46
47         # Crear ejemplos
48         USR_TRAIN = np.empty([len(TRAIN_DATA)-2, 4], dtype=np.int32)
49         USR_TRAIN[:,0] = u
50         USR_TRAIN[:,1] = TRAIN_DATA[:len(TRAIN_DATA)-2]
51         USR_TRAIN[:,2] = TRAIN_DATA[1:-1]
52         USR_TRAIN[:,3] = TRAIN_DATA[2:]
53
54         USR_DEV = np.empty([len(DEV_DATA) -2, 4], dtype=np.int32)
55         USR_DEV[:, 0] = u

```

```

55     USR_DEV[:, 1] = DEV_DATA[:,len(DEV_DATA)-2]
56     USR_DEV[:, 2] = DEV_DATA[1:-1]
57     USR_DEV[:, 3] = DEV_DATA[2:]
58
59     USR_TEST = np.empty([len(TEST_DATA)-2, 4], dtype=np.int32)
60     USR_TEST[:, 0] = u
61     USR_TEST[:, 1] = TEST_DATA[:,len(TEST_DATA) -2]
62     USR_TEST[:, 2] = TEST_DATA[1:-1]
63     USR_TEST[:, 3] = TEST_DATA[2:]
64
65     # Añadir al global
66     TRAIN = np.concatenate([TRAIN,USR_TRAIN])
67     DEV = np.concatenate([DEV,USR_DEV])
68     TEST = np.concatenate([TEST,USR_TEST])
69
70 #Mezclar
71 if(shuffle):
72     np.random.seed(seed)
73     np.random.shuffle(TRAIN)
74     np.random.shuffle(DEV)
75     np.random.shuffle(TEST)
76
77 return
78     TRAIN[:, :3], np.matrix(TRAIN[:, -1]).T, DEV[:, :3], np.matrix(DEV[:, -1]).T, TEST[:, :3], np.matrix(TEST
79
80 def getBatchData(X, Y, batch_size, step, ITRS, RES_ITMS):
81
82     # Si se ha aadi una iteracin, introducir los datos restantes
83     if (step == ITRS - 1 and RES_ITMS > 0):
84         batch_inputs = X[(ITRS - 1) * batch_size:]
85         batch_context = Y[(ITRS - 1) * batch_size:]
86     # Si no, es un batch normal
87     else:
88         st = step * batch_size
89         nd = (step + 1) * batch_size
90
91         batch_inputs = X[st:nd]
92         batch_context = Y[st:nd]
93     return(batch_inputs, batch_context)
94
95 def trainNet(X,Y, DX,DY, TX,TY,config, user_size, train=True, save=True):
96
97     def getSlope(data):
98         r = linregress(range(len(data)), data)
99         return r.slope
100
101     batch_size = config['batch_size']
102     learning_rate = config['learning_rate']
103     epochs = config['epochs']
104
105     user_size = user_size
106     song_size = W2V.shape[0]
107
108     user_emb_size = 128
109     song_emb_size = 64
110     hidden_size = 32

```

```

111 profile_size = user_emb_size + song_emb_size + song_emb_size
112 user_concat_size = user_emb_size * 2
113
114 num_sampled = 64 # Ejemplos negativos (Para el NSE)
115
116 model_name = 'model'
117 model_path = 'models/' + TEST_NAME + '/model'
118 complete_path = model_path + "/" + model_name
119
120 # Se almacena la loss de dev de las 5 ultimas epochs
121 slope_size = 100
122 train_hist = []
123 dev_hist = []
124
125 # Creacin del grafo de TF.
126 graph = tf.Graph()
127
128 with graph.as_default():
129
130     if(config["seed"]!=1):
131         tf.set_random_seed(config["seed"])
132
133     # Numero global de iteraciones
134     global_step = tf.Variable(0, dtype=tf.int32, trainable=False, name='global_step')
135
136     # Datos de entrada
137
138     # Array del tamaño del batch con las X
139     train_dataset = tf.placeholder(tf.int32, shape=[None, 3], name="train_inputs")
140
141     # Array del tamaño del batch con las Y asociadas
142     train_labels = tf.placeholder(tf.float32, shape=[None, 1], name="train_labels")
143
144     # Embeddings
145
146     D2V_EMB_VRB = tf.Variable(D2V_EMB, trainable=False, name="doc_embeddings")
147
148     T0 = tf.Variable(tf.truncated_normal([user_concat_size, user_emb_size], mean=0.0,
149         stddev=1.0 / math.sqrt(hidden_size)), name="T0")
150     B0 = tf.Variable(tf.zeros([user_emb_size]), name="B0")
151
152     O2 = tf.Variable(tf.truncated_normal([song_size, song_emb_size], mean=0.0, stddev=1.0
153         / math.sqrt(song_emb_size)), name="O2")
154     B02 = tf.Variable(tf.zeros([song_emb_size]), name="B02")
155
156     # Operaciones
157
158     # Embedding de documento
159     ed = tf.nn.embedding_lookup(D2V_EMB_VRB, train_dataset[:, 0])
160
161     #Transformar concat documento de 256 a 128
162     ed = tf.matmul(ed,T0)+B0
163
164     # Embedding de cancion 1
165     ec1 = tf.nn.embedding_lookup(O2, train_dataset[:, 1])

```

```

163     ec1 = ec1 + B02
164
165     # Embedding de cancion 2
166     ec2 = tf.nn.embedding_lookup(O2, train_dataset[:, 2])
167     ec2 = ec2 + B02
168
169     # aClculo de LOSS y optimizador—
170
171     embed = []
172     embed.append(ed)
173     embed.append(ec1)
174     embed.append(ec2)
175     embed = tf.concat(embed, 1)
176
177     # softmax weights, W and D vectors should be concatenated before applying softmax
178     weights = tf.Variable(tf.truncated_normal([song_size, profile_size], stddev=1.0 /
179         math.sqrt(song_size)), name="weights")
180     # softmax biases
181     biases = tf.Variable(tf.zeros([song_size]), name="biases")
182
183     # aClculo de LOSS y optimizador—
184
185     nce_loss = tf.nn.nce_loss(weights=weights,
186         biases=biases,
187         labels=train_labels,
188         inputs=embed,
189         num_sampled=num_sampled,
190         num_classes=song_size,
191         name="nce_loss")
192
193     batch_loss = nce_loss
194     nce_loss = tf.reduce_mean(nce_loss)
195
196     logits = tf.matmul(embed, tf.transpose(weights))
197     logits = tf.nn.bias_add(logits, biases)
198
199     top_v, top_i = tf.nn.top_k(logits, k=song_size)
200
201     # logits = tf.nn.softmax(logits, axis=0)
202
203     # optimizer =
204     tf.train.AdamOptimizer(learning_rate).minimize(nce_loss, global_step=global_step)
205     train_step = tf.train.GradientDescentOptimizer(learning_rate).minimize(nce_loss,
206         global_step=global_step)
207
208     # Inicializar variables
209     init = tf.global_variables_initializer()
210
211     # Crear objeto encargado de almacenar la red
212     saver = tf.train.Saver(max_to_keep=1)
213
214     config = tf.ConfigProto()
215     config.gpu_options.allow_growth = True
216     with tf.Session(graph=graph, config=config) as session:
217     #with tf.Session(graph=graph) as session:

```

```

215 # We must initialize all variables before we use them.
216 init.run()
217
218
219 if(save):
220     # Obtener el checkpoint
221     ckpt = tf.train.get_checkpoint_state(model_path)
222
223     # Si existe el checkpoint restaurar
224     if ckpt and ckpt.model_checkpoint_path:
225         saver.restore(session, ckpt.model_checkpoint_path)
226
227 #Si se desea entrenar
228 if(train):
229     # -----
230     # TRAIN
231     # -----
232
233     ITRS= int(len(X)/batch_size)
234     RES_ITMS = len(X) - (ITRS*batch_size)
235     if(RES_ITMS>0):ITRS+=1 #Si no es divisin exacta, se aade una iteracin
236
237     for e in range(epochs):
238         TRAIN_LOSS = 0.0
239
240         print("Epoch "+str(e))
241
242         for step in range(ITRS):
243
244             batch_inputs, batch_context =
245                 getBatchData(X,Y,batch_size,step,ITRS,RES_ITMS)
246
247             feed_dict = {train.dataset: batch_inputs, train.labels: batch_context}
248
249             ., loss_val, batch_loss, gs = session.run([train_step,
250                 nce_loss,batch_loss,global_step], feed_dict=feed_dict)
251
252             TRAIN_LOSS += np.sum(batch_loss)
253
254             #print('Epoch '+str(e)+' , last batch loss: '+str(loss_val))
255
256             if (save): saver.save(session, complete_path, global_step=global_step)
257
258             # -----
259             # DEV
260             # -----
261
262             if(len(DX)>0):
263
264                 ITRS_TEST = int(len(DX) / batch_size)
265                 RES_ITMS_TEST = len(DX) - (ITRS_TEST * batch_size)
266                 TEST_LOSS = 0.0
267                 if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se
268                     aade una iteracin

```

```

269
270         batch_inputs, batch_context =
                getBatchData(DX,DY,batch_size,step,ITRS_TEST,RES_ITMS_TEST)
271
272         feed_dict = {train_dataset: batch_inputs, train_labels: batch_context}
273         loss_val,batch_lss = session.run([nce_loss,batch_loss],
                feed_dict=feed_dict)
274
275         TEST_LOSS+=np.sum(batch_lss)
276
277         dev_hist.append(TEST_LOSS / (len(DX) * 1.0))
278         train_hist.append(TRAIN_LOSS / (len(X) * 1.0))
279
280         if (len(dev_hist) >= slope_size):
281             # Calcular la pendiente
282             print(getSlope(dev_hist[-slope_size:]))
283             print(dev_hist[-slope_size:])
284
285             if (getSlope(dev_hist[-slope_size:]) > -1e-5 or e==epochs-1):
286
287                 def printArray(data,data2):
288                     for i in range(len(data)):
289                         a=data[i]
290                         b=data2[i]
291                         print(str(a).replace(".",",")+ "\t"+str(b).replace(".",","))
292
293                 print("-" * 50)
294                 print(e)
295                 printArray(dev_hist,train_hist)
296                 return
297     else:
298         #
299
300     # TEST FINAL
301     #
302
303     batch_size = 300
304
305     IN_TOP_5 = 0
306     IN_TOP_10 = 0
307     IN_TOP_20 = 0
308     IN_TOP_50 = 0
309     IN_TOP_100 = 0
310     IN_TOP_1000 = 0
311     real_positions = []
312
313     items = 0
314
315     ITRS_TEST = int(len(TX) / batch_size)
316     RES_ITMS_TEST = len(TX) - (ITRS_TEST * batch_size)
317     if (RES_ITMS_TEST > 0): ITRS_TEST += 1 # Si no es divisin exacta, se aade una
318         iteracin
319
320     for step in range(ITRS_TEST):
321         #tss = time.time()
322
323         batch_inputs, batch_context = getBatchData(TX, TY, batch_size, step,

```

```

321         ITRS_TEST, RES_ITMS_TEST)
322     top_values, top_indx = session.run([top_v, top_i], feed_dict={train_dataset:
323         batch_inputs, train_labels: batch_context})
324
325     for j in range(len(batch_context)):
326         real_position = int(np.where(top_indx[j, :] == batch_context[j, 0])[0])
327         real_positions.append(real_position)
328
329         IN_TOP_5 += 1 if real_position <= 5 else 0
330         IN_TOP_10 += 1 if real_position <= 10 else 0
331         IN_TOP_20 += 1 if real_position <= 20 else 0
332         IN_TOP_50 += 1 if real_position <= 50 else 0
333         IN_TOP_100 += 1 if real_position <= 100 else 0
334         IN_TOP_1000 += 1 if real_position <= 1000 else 0
335
336     items += len(batch_inputs)*1.0
337
338     PTF = IN_TOP_5 / items
339     PTT = IN_TOP_10 / items
340     PTTW = IN_TOP_20 / items
341     PTFT = IN_TOP_50 / items
342     PTCIEN = IN_TOP_100 / items
343     PTMIL = IN_TOP_1000 / items
344
345     print(np.mean(real_positions), np.std(real_positions),
346           np.median(real_positions))
347     print(str(PTF).replace(".", ",") + "\t" + str(PTT).replace(".", ",") + "\t"
348           + str(PTTW).replace(".", ",") + "\t" + str(PTFT).replace(".", ",") + "\t"
349           + str(PTCIEN).replace(".", ",") + "\t" + str(PTMIL).replace(".", ","))
350
351     #print(time.time()-tss)
352
353     print(np.mean(real_positions), np.std(real_positions), np.median(real_positions))
354
355     #-----
356     #####
357     # Llamadas
358     #####
359     TEST_NAME="test_083"
360
361     os.environ['TF_CPP_MIN_LOG_LEVEL'] = '3'
362
363     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
364     D2V_DM_COMPLETE = np.load('embeddings/d2v_dm/complete/EMB_MATRIX')
365     D2V_DBOW_COMPLETE = np.load('embeddings/d2v_dbow/complete/EMB_MATRIX')
366
367     #Se obtienen cada uno de los embeddings aprendidos (incluidas las canciones)
368     D2V_DM_PRESENT = np.load('embeddings/d2v_dm/present/EMB_MATRIX')
369     D2V_DBOW_PRESENT = np.load('embeddings/d2v_dbow/present/EMB_MATRIX')
370
371     W2V = np.load('embeddings/w2v/EMB_MATRIX')
372
373     #Se concatenan los embeddings de los usuarios para generar el perfil
374     D2V_CONS = np.concatenate((D2V_DM_COMPLETE, D2V_DBOW_COMPLETE), axis=1) # El consolidado =

```

```

    Consolidado DM + Consolidado DBOW
375 D2V_PRES = np.concatenate((D2V_DM_PRESENT,D2V_DBOW_PRESENT),axis=1) # El presente = Presente
    DM + Presente DBOW
376
377 D2V_EMB = np.concatenate((D2V_CONS,D2V_PRES),axis=1) # El perfil completo CONS + PRES
378
379 #Se cargan los datos de train de las canciones para conocer los IDs reales
380 DATA = np.load("datos/Extend/ARRAY_ARRAYS_TRAIN.pkl")
381
382 #Obtener una lista con los ID de canciones que no aesten el la matriz de embeddings
383 ALL = set(list(range(W2V.shape[0])))
384 REAL = set(list(set(list(map(int, DATA)))))
385
386 #####
387 # TENSORFLOW
388 #####
389
390 #
391 # Grid Search + Early Stopping
392 #
393 '''
394 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
395 rates = [0.000001,1, 0.01, 0.0001]
396
397 for r in rates:
398     config = {"batch_size": 1024, "learning_rate": r, "epochs": 1000, "seed": 100}
399
400     print("Learning Rate: "+str(r))
401     print("-" * 50)
402     trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=False)
403     print("-"*50)
404
405 exit()
406 '''
407 #
408 # Train
409 #
410 '''
411 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))
412
413 #Para esta fase, se une TRAIN y DEV
414 X = np.concatenate((X,DX),axis=0)
415 Y = np.concatenate((Y,DY),axis=0)
416
417 DX=[]
418 DY=[]
419
420 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs":1000, "seed":100}
421 trainNet(X, Y, DX, DY, TX, TY, config, user_size, train=True, save=True)
422
423 exit()
424 '''
425 #
426 # Test
427 #
428
429 X, Y, DX, DY, TX, TY, user_size = getData(REAL_SONG_IDS=list(REAL))

```



```
430
431 config = {"batch_size": 1024, "learning_rate": 0.0001, "epochs":1000, "seed":100}
432 trainNet(X, Y, DX, DY, TX, TY, config, user_size,train=False, save=True)
```