

## Index

Objectius	2
Descripció de la PAC	
Exercici 1 (2 punts)	
Exercici 2 (2 punts)	
Exercici 3 (2 punts)	
Exercici 4 (4 punts)	



#### **Objectius**

En aquest PAC es practicaran els conceptes del temari relacionat amb recomanació i clustering, en una vertent clarament pràctica enmarcada en l'ús de llibreries Python.

#### Descripció de la PAC

En aquesta PAC estudiarem en profunditat l'ús de la llibreria Surprise (http://www.surpriselib.com) com a eina en Python que facilita l'implementació de recomanadors basats en diferents tècniques d'agrupació (clustering).

Farem servir la coneguda base (<a href="https://grouplens.org/datasets/movielens/">https://grouplens.org/datasets/movielens/</a>) de dades Movielens-100k.

Important: Com farem un ús intensiu d'aquesta llibreria, aquesta PAC requereix implícitament un estudi i consulta actius de l'apartat "Documentation" de la pàgina web anteriorment esmentada.

Per a aquesta PAC, se us facilita un fitxer zip amb dos fitxers Python. Aquests dos fitxers hauran de ser completats per l'alumen omplint els diversos forats que s'hi indiquen amb el codi que correspongui.

NO MODIFIQUEU cap nom de funció, variable, ni tampoc dels propis fitxers. La PAC s'ha de completar tot fent servir la base que es proporciona, que en moltes ocasions pot guiar-vos en la manera com dur a terme els exercicis proposats.

Com podeu veure, al fitxer pac1.py tenim diverses definicions de funcions incompletes i un programa molt simple que serà executat (al final del fitxer). Al fitxer IAA\_predictors.py tenim, de manera anàloga, la definició d'unes classes sobre les que treballarem.

### Exercici 1 (2 punts)

a) (1 punt) Completeu la definició del mètode predict. Aquest mètode rep un dataset prèviament carregat, un identificador d'usuari i un identificador d'ítem, i la seva finalitat és preparar una sèrie d'algorismes de predicció (basats en diferents tècniques d'agrupació) per tal de dur a terme una predicció de la puntuació (rating) amb què l'usuari uid valorarà l'ítem iid. Com podeu veure al codi, aquest mètode realitzarà aquesta predicció partint dels algorismes SVD, KNNBasic, KNNWithMeans i NormalPredictor. Un cop definit el mètode, dueu a terme les prediccions que aquests algorismes faran de la valoració de l'usuari 777 de l'ítem 333.

Dins el bucle for del mètode predict(), per executar el mateix per a cadascun dels algoritmes de la variable algorithms inicialitzo l'algoritme, després entreno l'algoritme amb el trainset definit i finalment amb el mètode predict() genera la predicció, la qual es retorna s'imprimeix per pantalla.

Converteixo en string l'usuari 777 i l'ítem 333 abans d'utilitzar-los per cridar el mètode predict().



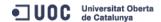
Codi al fitxer pac1.py, mètode predict():

```
def predict(data, uid, iid):
                                       trainset = data.build full trainset()
                                       # Set the algorithms to predict
algorithms = (SVD, KNNBasic, KNNWithMeans, NormalPredictor)
                                       for a in algorithms:
                                           *******************
                                           # Init the algorithm and train it with trainset (fit it)
# Store your prediction in a variable called "pred"
                                           algo.fit(trainset)
                                           pred= algo.predict(uid,iid, r_ui=None, verbose=True)
                                          ********************
                                           print (pred)
>>predict(data,str(uid),str(iid))
               item: 333
                                                  r_ui = None est = 3.98
                                                                              {'was_impossible': False}
                                item: 333
               user: 777
                                                  r_ui = None est = 3.98 {'was_impossible': False}
                Computing the msd similarity matrix...
               Done computing similarity matrix.
                                                  r_ui = None est = 3.83
                                                                             {'actual_k': 40, 'was_impossible': False}
               user: 777
                                item: 333
               user: 777
                                                  r_ui = None est = 3.83 {'actual_k': 40, 'was_impossible': False}
                                item: 333
               Computing the msd similarity matrix...
               Done computing similarity matrix.
                                                 r_ui = None est = 4.31 {'actual_k': 40, 'was_impossible': False}
r_ui = None est = 4.31 {'actual_k': 40, 'was_impossible': False}
r_ui = None est = 2.45 {'was_impossible': False}
               user: 777
                             item: 333
               user: 777
                                item: 333
                                item: 333
               user: 777
               user: 777
                                item: 333
                                               r_ui = None est = 2.45 {'was_impossible': False}
```

b) (1 punt) Consulteu la documentació i expliqueu com proporcionen la seva predicció els algorismes KNNWithMeans i NormalPredictor.

L'algorisme NormalPredictor calcula uns valor aleatoris basats en una distribució normal (gaussiana), la qual s'estima a partir d'un conjunt de dades d'entrenament utilitzant la estimació per màxima verosimilitat. Aquest mètode calcula el valor que maximitza la probabilitat d'obtenir el valor del conjunt d'entrenament.

L'algorisme KNNWithMeans per predir la puntuació que un usuari donara a un ítem, es cerquen els k usuaris que han puntuat ítems més similars a l'usuari del qual volem donar la predicció. A partir de les puntuacions dels k usuaris i similituds amb l'ítem a predir, es donara una puntuació a l'ítem. Els valors més propers s'obtenen a partir d'una funció de similitud com distància euclidiana, correlació de Pearson, MVDM o Hamming.



#### Exercici 2 (2 punts)

a) (1 punt) Completeu la definició del mètode multiple\_cv. Aquest mètode rep un dataset prèviament carregat, i el nombre de vies (o folds, per defecte 5) per tal de dur a terme una valoració creuada dels mateixos algorismes fets servir a l'Exercici 1. Expliqueu les característiques i els avantatges de la validació creuada i justifiqueu per què és important considerar diverses particions del trainset per validar la bondat del nostre sistema.

Característiques i avantatges de la validació creuada i perquè es importants considerar particions del trainset per validar.

La validació creuada es un sistema d'avaluació d'un sistema d'aprenentatge automàtic, que permet mesurar la qualitat dels resultats de l'algoritme utilitzat. Es divideix el dataset en k particions, per tal de calcular la mitjana i desviació de cadascuna de les k particions. Això permet comparar per tal d'estimar el grau de fiabilitat del resultat i millorar-los a partir de modificacions en els mètodes i els seus paràmetres. Els mètodes que hem utilitzat son RMSE (error mitjà quadrat) i MAE (error mitjà absolut), que calculen l'error i quan més petit el resultat es més fiable.

b) (1 punt) Identifiqueu el mètode que dona pitjors resultats i raoneu per què creieu que aquest mètode no és adient. En quines condicions seria més raonable utilitzar-ho com a predictor?

El mètode amb pitjors resultats s'identifica per el valor màxim de la mitjana del mètode d'avaluació, on NormalPredictor ha obtingut una mitjana d'error superior a la resta, amb valors RMSE=1.5203 i MAE=1.2211 respecte a la segona pitjor que es KNNBasic amb RMSE=0.9784 i MAE=0.7728.

Aquest mètode no es adient en casos on sigui requisit una alta fiabilitat com a l'àmbit de la medicina, on les prediccions respecte a malalties pot impactar en la salut de les persones.

Els casos on el mètode NormalPredictor, degut a que necessita menys poder computacional, es més raonable utilitzar-ho en aplicacions on es requereix prediccions en temps real i sigui suficient una predicció correcte.

Codi al fitxer pac1.py, mètode predict():



#### >>multiple\_cv(data,folds=5)

# 

	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	Std
RMSE (testset)	0.9528	0.9486	0.9484	0.9517	0.9506	0.9504	0.0017
MAE (testset)	0.7493	0.7494	0.7489	0.7485	0.7501	0.7493	0.0005
Fit time	0.47	0.50	0.49	0.50	0.49	0.49	0.01
Test time	3.03	2.95	3.03	3.00	2.98	3.00	0.03
Evaluating RMSE,	MAE of a	algorithm	NormalP	redictor	on 5 sp	lit(s).	

	Fold 1	Fold 2	Fe14 2	Fold 4	Fold F	Moon	Std
	roiu i	roiu z	roiu 3	roiu 4	roiu 3	mean	Stu
RMSE (testset)	1.5230	1.5176	1.5153	1.5180	1.5278	1.5203	0.0045
MAE (testset)	1.2227	1.2172	1.2172	1.2217	1.2267	1.2211	0.0036
Fit time	0.09	0.11	0.11	0.11	0.11	0.11	0.01
Tost time	0 13	0 16	0 16	0 16	0 11	0 15	0 02



#### Exercici 3 (2 punts)

Com ja heu observat, les valoracions registrades a MovieLens corresponen a valors discrets en l'interval [1,5], segurament obtinguts després que l'usuari valorés una determinada pel·lícula amb una quantitat d'estrelles.

Imaginem que treballem com a enginyers al portal web d'on s'han extret aquestes puntuacions, i des del departament d'user experience ens comuniquen que han realitzat un estudi que suggereix que hem de canviar aquesta manera de valorar les pel·lícules, de manera que ara, en comptes dels valors de l'1 al 5, seran valorades escollint una de les següents opcions:

["Molt dolenta", "Pobre", "No està malament", "Entretinguda", "Recomanable", "Molt bona", "Brillant"]

Quines consideracions hem de tenir en compte per a poder redissenyar el sistema de recomanació? Ens servirà el conjunt de dades d'entrenament que ja teníem? En cas d'haver d'introduir modificacions a la implementació del sistema, indiqueu-les, tot justificant cadascuna de les decisions que prendríeu.

Quines implicacions tindria el fet que al cap de 2 mesos ens indiquessin que hem d'afegir l'opció "Decepcionant" entre "Pobre" i "No està malament"?

Important: per realitzar aquest exercici no cal implementar cap programa, només un raonament justificat i complet.

Son necessàries les següents consideracions per poder redissenyat el sistema de recomanació:

- Les valoracions actuals son de tipus enter i en la nova son cadenes de caràcters (string).
- Augment de 5 a 6 possibles valors en les valoracions.

El conjunt de dades d'entrenament anterior ens servirà sempre que prèviament es normalitzin els valors a la nova escala utilitzada i s'ajuntin les dades.

Les modificacions a implementar al sistema serien les següents:

- Per a tot el conjunt de dades amb el nou tipus de valoracions, prèviament a entrenar-ho transformaria el tipus de les seves valoracions a un valor discret, per facilitar aprofitar l'actual programa que utilitza valors discrets.
- Un cop tenim el conjunt de dades antic i nou amb valors discrets, tant entrenament com test, normalitzaria les dades perquè tinguessin la mateixa escala de valors. D'aquesta forma l'antic conjunt es podrà incorporar al nou per tal de no perdre l'històric de valoracions i el sistema doni recomanacions més fiables que si ens basem només en el nou conjunt.

Si després de 2 mesos s'han d'afegir una opció entre Pobre i No està malament, significa que s'haurien tornar a normalitzar els conjunts de dades anteriors, per tal que s'adeqüessin al nou rang de valoracions i així utilitzar-los per entrenar el model.

["Molt dolenta", "Pobre", <u>«Decepcionant»,</u> "No està malament", "Entretinguda", "Recomanable", "Molt bona", "Brillant"]



#### Exercici 4 (4 punts)

Estudieu detingudament com crear un algorisme de predicció propi

(http://surprise.readthedocs.io/en/stable/building\_custom\_algo.html) i preneu com a exemple la definició dels algorismes que la llibreria Surprise porta implementats per defecte (per exemple,

https://github.com/NicolasHug/Surprise/blob/master/surprise/prediction\_algorithms/knns.py ).

En aquest exercici definirem un parell d'algorismes propis. Consulteu el fitxer IAA\_predictors.py, on hi ha definida una classe arrel SymmetricAlgo, de la qual heretaran els nostres algorismes, i a continuació tenim definicions incompletes de dos algorismes. En aquest context es demana:

a) (1.5 punts) Completar la definició de la classe Dummy\_IAA, de manera que el mètode estimate simplement retorni una valoració a l'atzar d'entre totes les valoracions fetes per l'usuari que es consulta.

Al mètode estimate amb la llista de tuples de self.yr amb els usuaris i valoracions, faig un recorregut per l'usuari u (user\_inner\_id) per tal de crear una llista 'Listux' amb les valoracions de l'usuari consultat. Només ens resta seleccionar una valoració de la llista 'ListAux' a l'atzar amb el mètode numpy.random.choice.

Codi al fitxer IAA\_predictors.py, clase Dummy\_IAA, mètode estimate(): >>run\_IAAPredictors(data,False, int(uid),int(iid))

b) (1.5 punts) Completar la definició de la classe KNN\_IAA, de manera que a la seva inicialització permeti rebre un paràmetre n, i el mètode estimate faci pràcticament el mateix que el que fa KNNBasic de Surprise, però amb la diferència que dels k\_neighbors triem n a l'atzar per calcular el promig ponderat (weighted average). En cas que k\_neighbors tingui una mida inferior a n, agafarem el valor de la mida de k\_neighbors com a n.

Codi al fitxer IAA\_predictors.py, clase KNN\_IAA, mètode estimate(): >>run\_IAAPredictors(data,False, int(uid),int(iid))



```
def estimate(self, u, i):
                                                                                                                          PS_D:\dropbox\uoc\iaa\pac1> python pac1.py
                                                                                                                         if not (self.trainset.knows_user(u) and self.trainset.knows_item(i)):
    raise PredictionImpossible('User and/or item is unkown.')
      x, v = self.switch(u, i)
                                                                                                                          \sharp Compute similarities between x and x2, where x2 describes all other
      # users that have also rated item Y.
neighbors = [(self.sim[x, x2], r) for (x2, r) in self.yr[y]]
k_neighbors = heapq.nlargest(self.k, neighbors, key=lambda t: t[0])
                                                                                                                         2.0
Computing the msd similarity matrix...
Done computing similarity matrix.
Evaluating KNN_IAA algorithm without cross-validation
(3.8014359950089598, {'actual_k': 34})
      #Compute weighted average from n random neighbors within k_neighbors
#First initialize variables
sum_sim = sum_ratings = actual_k = 0
kn_length = len(k_neighbors)
      #Obtenim _random_neighbors (n_random_neighbors)
#Si la mida de k_neighors es inferior a n, n=kn_length
if kn_length < self.n:
    n_random_neighbors = self.n</pre>
             ::
#$Sino seleccionem un número a l'atzar
#entre el mínim (self.n)i el màxim de neighbors (kn_length)
n_random_neighbors = np.random.randint(self.n,kn_length+1)
      #Selectionem a l'atzar n_random_neihbors de k_neighbors
k neighbors random = random.sample(k neighbors,n random neighbors)
      #Calculem la mitja de k_neighbors_random
for (sim, r) in k_neighbors_random:
    if sim > 0:
                   sum_sim += sim
                   sum_ratings += sim * r
actual_k += 1
      *****************
      if actual_k < self.min_k:
    raise PredictionImpossible('Not enough neighbors.')</pre>
      est = sum_ratings / sum_sim
      details = {'actual_k': actual_k}
return est, details
```

- c) (1 punt) Completeu el mètode run\_IAAPredictors de pac1.py, i executeu-ho. Aquest mètode es comportarà de dues maneres diferents, segons el paràmetre do\_cv:
- Si és True, realitzarà una validació creuada amb 5 vies (folds).
- Si és False, utilitzarà els paràmetres uid i iid per fer una predicció de valoració de l'ítem iid per part de l'usuari uid.

En ambdós casos, com podeu veure al codi, es faran servir els algorismes Dummy\_IAA i KNN\_IAA, un darrere l'altre.

Al executar el codi, amb do cv =False funciona, però amb do cv =True no.

Adjunto resultat de do\_cv=False