

## ▼ BDLE : Graphes en DataFrame

### ▼ Construction du graphe

```
1 #fonction utilisée par la suite pour calculer les poids des arcs
2 #df: dataframe, source: nom de la colonne contenant les noeuds source des arcs
3 # poids: poids initial avant normalisation, n: nombre maximum d'arcs à garder p
4 from pyspark.sql.functions import row_number, sum
5 from pyspark.sql import Window
6
7 def calcul_poids(df, source, poids, n):
8
9     window = Window.partitionBy(source).orderBy(col(poids).desc())
10
11     filterDF = df.withColumn("row_number", row_number().over(window)) \
12         .filter(col("row_number") <= n) \
13         .drop(col("row_number"))
14
15     tmpDF = filterDF.groupBy(col(source)).agg(sum(col(poids)).alias("sum_" + poi
16
17     finalDF = filterDF.join(tmpDF, source, "inner") \
18         .withColumn("norm_" + poids, col(poids) / col("sum_" + poids)) \
19         .cache()
20     return finalDF
```

### ▼ Construction des liens pondérés entre utilisateurs et chansons

- Construire un DataFrame userTrack à partir de data pour stocker les arcs entre utilisateurs et chansons. Pour chaque utilisateur (userId) on ajoute un arc vers une chanson (trackId) avec un poids égal au nombre total de fois que l'utilisateur à écouté la chanson. Utiliser la fonction calcul\_poids pour garder pour chaque utilisateur les 100 chansons avec le poids le plus élevé et normaliser les poids des arcs gardés.
- Affichage du résultat: garder uniquement les arcs qui ont les 5 plus grandes valeurs possibles des poids (utilisez la fonction rank() et la fenêtre over(window)). Afficher 20 lignes du résultat, en triant le résultat par ordre décroissant des poids, ensuite par ordre croissant des userId et des artistId.



```

8 #calculer le poids final en utilisant calcul_poids
9 userTrack = calcul_poids(userTrack,"userId","count",100).select("userId","trackId")
10
11 window = Window.orderBy(col("norm_count").desc())
12
13 userTrackList = userTrack.withColumn("position", rank().over(window)).where("position < 100")
14
15 for val in userTrackList:
16     print("%s %s %s" % val)
17
18 userTrack.count()#résultat: 210675

```

## Construction des liens pondérés entre utilisateurs et artistes

- Construire un DataFrame userArtist à partir de data pour stocker les arcs entre utilisateurs et artistes. Pour chaque utilisateur (userId) on ajoute un arc vers un artiste (artistId) avec un poids égal au nombre total de fois que l'utilisateur à écouté des chansons de cet artiste. Utiliser la fonction calcul\_poids pour garder pour chaque utilisateur au plus 100 artistes avec le poids le plus élevé et normaliser les poids des arcs gardés.
- Affichage du résultat: garder uniquement les arcs qui ont les 5 plus grandes valeurs possibles des poids (utilisez la fonction rank() et la fenêtre over(window)). Afficher 20 lignes du résultat, en triant le résultat par ordre décroissant des poids, ensuite par ordre croissant des userId et des artistId.

```

1 #poids=nombre de fois qu'un utilisateurs a écouté des morceaux de cet artiste
2 #regrouper les données par userId et artistId
3 userArtist = data.groupBy("userId","artistId").count()
4
5 #calculer le poids final en utilisant calcul_poids
6 userArtist = calcul_poids(userArtist,"userId","count",100).select("userId","artistId","weight")
7
8 window = Window.orderBy(col("norm_count").desc())
9
10 userArtistList = userArtist.withColumn("position", rank().over(window)).where("position < 20")
11 # à compléter ....
12
13 for val in userArtistList:
14     print("%s %s %s" % val)
15

```

- Affichage du résultat: garder uniquement les arcs qui ont les 5 plus grandes valeurs possibles des poids (utilisez la fonction `rank()` et la fenêtre `over(window)`). Afficher 20 lignes du résultat, en triant le résultat par ordre décroissant des poids, ensuite par ordre croissant des `artistId` et des `trackId`.\*

```

1 # poids de l'arc: nombre de fois qu'un track de l'artiste a été écouté par tous
2 artistTrack = data.groupBy("artistId","trackId").count()
3
4 #calculer le poids final en utilisant calcul_poids
5 artistTrack = calcul_poids(artistTrack,"artistId","count",100).select("artistId",
6                               "trackId","norm_count")
7
8 window = Window.orderBy(col("norm_count").desc())
9
10 artistTrackList = artistTrack.withColumn("position", rank().over(window)).\
11 where("position < 5").select("artistId","trackId","norm_count").orderBy("artistId",
12                                   "trackId","norm_count").show(20)
13
14 for val in artistTrackList:
15     print("%s %s %s" % val)
16
17 artistTrack.count() #résultat: 35408

```

## Construction des liens pondérés entre chansons

- Construire un DataFrame `trackTrack` à partir de `data` pour stocker les arcs entre les chansons. Le poids total d'un arc entre `trackId1` et `trackId2` est le nombre total d'utilisateurs qui ont écouté à la fois `trackId1` et `trackId2` dans un laps de temps de 10 minutes (à noter que le graphe est non orienté, `trackTrack` contient à la fois une entrée pour (`trackId1`, `trackId2`) et une entrée pour (`trackId2`, `trackId1`)). Utiliser la fonction `calcul_poids` pour garder pour chaque chanson au plus 100 chansons avec le poids le plus élevé et normaliser les poids gardés.

```

6
7 #pour chaque couple de trackId le nombre d'utilisateurs qui les ont écoutés ense
8
9
10 #calculer le poids final en utilisant calcul_poids
11 track1=data.withColumnRenamed('trackId',"trackId1").\
12 withColumnRenamed("timestamp","timestamp1").withColumnRenamed("userId","userId1'
13
14 trackTrack = data.join(track1,( (data.userId == track1.userId1) &(data.trackId
15 trackTrack=calcul_poids(trackTrack,"trackId","count",100).select('trackId',"trac
16
17 ""window = Window.orderBy(col("norm_count").desc())
18
19 trackTrackList = trackTrack.withColumn("position", rank().over(window))\
20     # à compléter ....
21
22 for val in trackTrackList:
23     print("%s %s %s" % val)""
24
25 trackTrack.count() #résultat: 136257

136257

```

## Construction du graphe final

Construire un DataFrame graphe pour stocker tous les noeuds et tous les liens calculés précédemment. Le dataframe contiendra une colonne 'source' (identifiant du noeud source), une colonne 'destination' et une colonne 'poids'. Les colonnes 'source' et 'dest' contiennent à la fois des identifiants utilisateur, chanson et artiste. La colonne 'poids' contient les poids des arcs

recommandation à chaque itération du calcul.

$$x[i] = \alpha * u[i] + (1-\alpha) * \sum(xant[j] * poids[j][i])$$

- poids[j][i] : poids de l'arc entre j à i
- u[i]: valeur de personnalisation, u[10]=1 et u[i]=0 si i !=10
- xant[j] : valeur du score du noeud j à l'itération précédente (x0[10]=1 et x0[i]=0 si i

On considère  $\alpha=0.15$  et on effectue le calcul pour 5 itérations (maxiter=5)

## Calcul du vecteur de recommandation x

```

1 import pandas as pd
2 from pyspark.sql.functions import when
3
4 user = 10
5 d=0.85
6 alpha=(1-d)
7 maxiter = 5
8 # Construire le vecteur d'importance initial
9 x0 = spark.createDataFrame(pd.DataFrame([(user,1)], columns=["id","rank"]))
10

```

```
1 # Définir une fonction qui prend comme argument une liste d'utilisateurs triés
2 # par ordre croissant (users) et retourne une liste de couples ordonnés d'utilis
3
4 from pyspark.sql.functions import udf
5 from pyspark.sql.types import *
6
7 def parse_string(users):
8     results=[]
9     for u in users[:-1]:
10         for v in users[1:]:
11             if ( u < v ):
12                 results.append(str(u)+"_"+str(v))
13     return results
14
15 parse_string_udf = udf(parse_string, ArrayType(StringType()))
```

```
1 # Implémentation de l'algorithme de calcul de triangles
2
```

