

# TP1 :BIG DATA

Réalisé par : Mehdi KARECH

## Partie I : RDD sur databricks

1)2)3) création du vecteur aléatoire le charger dans un rdd et afficher le nombre de partition

```
1 import numpy as np
2 data = np.random.rand(20000)
3 rdd = sc.parallelize(data)
4 print(rdd.getNumPartitions())
5
```

8

Le nombre de partition égale 8 Cela veut dire qu'on peut lancer 8 taches sur 8 noeud en parallèles (une tache (task) par partition)

4)5) Chargement du fichier “text\_gradient.txt” et suppression des steps words de ce dernier

```
1 stop_words = ['the', 'a', 'of']
2 rdd2 = sc.textFile('FileStore/tables/text_gradient.txt')
3 rdd_filtred= rdd2.map(lambda x: x.replace(' the ', ' ')).map(lambda x: x.replace(' a ', ' ')).map(lambda x: x.replace(' of ', ' '))
4 rdd_filtred.take(5)
5
```

► (1) Spark Jobs

```
Out[4]: ['Gradient descent is first-order iterative optimization algorithm for ',
'finding local minimum differentiable function.',
'The idea is to take repeated steps in opposite direction of',
'the gradient (or approximate gradient) function at current point, ',
'because this is direction steepest descent. ']
```

Command took 0.64 seconds. Run on databricks-databricks-2.0.0-14 AM on TP

## Partie III : manipuler un dataset volumineux sur databricks

### 2) chargement dans un dataframe spark

```
1 sparkDf=spark.read.load('houses.csv')
2 sparkDf.printSchema()
```

► sparkDf: pyspark.sql.dataframe.DataFrame = [size: double, nb\_rooms: long ... 3 more fields]

root

```
-- size: double (nullable = true)
-- nb_rooms: long (nullable = true)
-- garden: long (nullable = true)
-- orientation: string (nullable = true)
-- price: double (nullable = true)
```

Mémoire consommée

```

1 from pyspark.serializers import PickleSerializer, AutoBatchedSerializer
2 def _to_java_object_rdd(rdd):
3
4     rdd = rdd._reserialize(AutoBatchedSerializer(PickleSerializer()))
5     return rdd.ctx._jvm.org.apache.spark.mllib.api.python.SerDe.pythonToJava(rdd._jrdd, True)
6
7 JavaObj = _to_java_object_rdd(sparkDf.rdd)
8
9 nbytes = sc._jvm.org.apache.spark.util.SizeEstimator.estimate(JavaObj)
10 print(nbytes)

```

4224424344

D'après le résultats la mémoire consommée est de 4224424344 bytes ce qui est équivalent à 4.2244 giga

### 3) Sélection des maisons qui ont plus de 2 chambres dans un nouveau dataframe

```

1 df_filtred = sparkDf.filter(sparkDf.nb_rooms > 2)
2 df_filtred.show(10)

```

► (1) Spark Jobs

► df\_filtred: pyspark.sql.dataframe.DataFrame = [size: double, nb\_rooms: long ... 3 more fields]

size	nb_rooms	garden	orientation	price
137.30027526131454	3	1	Ouest	312824.225304512
148.59258955996802	3	1	Est	312356.19881789363
192.67056190501847	3	1	Est	323977.13614897546
118.48163264872223	3	0	Nord	235242.31960432642
180.16142351158868	3	1	Est	320972.4814387185
167.61000386988792	3	1	Sud	340741.477043372
155.4553500321651	3	1	Nord	297374.6799684082
131.83054115223706	3	1	Est	307008.33851207944
200.06973888120532	3	1	Sud	348664.3313370723
174.8451659379434	3	1	Sud	342632.17409832845

only showing top 10 rows

### 4) Calcul de la surface moyenne et du prix moyen

Cmd 9

```

1 from pyspark.sql import functions as F
2 prix_et_surface_moyenne = sparkDf.agg(F.avg('price'),F.avg('size'))
3 print("le prix moyen du dataset original : ",prix_et_surface_moyenne.show())
4

```

► (2) Spark Jobs

► prix\_et\_surface\_moyenne: pyspark.sql.dataframe.DataFrame = [avg(price): double, avg(size): double]

avg(price)	avg(size)
267153.29091201985	150.04028343747882

### 5) méthodes envisagées pour réduire la taille du dataframe en mémoire

On peut appliquer une transformation sur la colonne orientation pour réduire la taille des entrées et ainsi réduire la taille en mémoire (on peut par exemple remplacer Ouest par O )

## Partie IV : Machine learning avec pyspark

### 2) Normalisation de la colonne size

```
1 from pyspark.ml.feature import StringIndexer, OneHotEncoder
2 from pyspark.sql.functions import stddev, mean, col
3 df_stats = sparkDf.select(
4     mean(col('size')).alias('mean'),
5     stddev(col('size')).alias('std')
6 ).collect()
7
8 mean = df_stats[0]['mean']
9 std = df_stats[0]['std']
10 sparkDf_normalized = sparkDf.withColumn('size', (F.col('size') - mean)/std)
11 sparkDf_normalized.show(5)
12 sparkDf.show(5)
13
```

▶ (4) Spark Jobs

▶ sparkDf\_normalized: pyspark.sql.dataframe.DataFrame = [size: double, nb\_rooms: long ... 3 more fields]

size	nb_rooms	garden	orientation	price
1.355121757813147	2	1	Est	309433.17891415354
-0.25540820679048665	3	1	Ouest	312824.225304512
-1.2914434589113652	1	0	Ouest	201769.8367561039
-0.02902297173783479	3	1	Est	312356.19881789363
0.8546401876526548	3	1	Est	323977.13614897546

pour normaliser la colonne ils suffit de faire (valeur - moyenne de la colonne) /écart type de la colonne

### 3) le one hot encoding de la colonne orientation

```

1 df_train, df_test = sparkDf.randomSplit([0.8, 0.2], seed=91)
2
3 string_indexer = StringIndexer(inputCol='orientation', outputCol='orientation'+Index')
4 one_hot_encoder = OneHotEncoder(inputCol=string_indexer.getOutputCol(),
5                                 outputCol='orientation'+OHE')
6

```

▶ df\_train: pyspark.sql.dataframe.DataFrame = [size: double, nb\_rooms: long ... 3 more fields]  
 ▶ df\_test: pyspark.sql.dataframe.DataFrame = [size: double, nb\_rooms: long ... 3 more fields]

Command took 0.47 seconds -- by em\_karech@esi.dz at 6/1/2021, 10:51:47 AM on TP

---

Cmd 13

```

1 label_to_index = StringIndexer(inputCol='price', outputCol='label')

```

Command took 0.05 seconds -- by em\_karech@esi.dz at 6/1/2021, 10:51:51 AM on TP

---

Cmd 14

```

1 string_model_indexer= string_indexer.fit(df_train)
2 display(string_model_indexer.transform(df_train))

```

▶ (3) Spark Jobs

	size	nb_rooms	garden	orientation	price	orientationIndex
1	0.00008298303231413229	1	1	Nord	30003.734081523573	1
2	0.00019212839873716805	3	1	Ouest	91008.64494749981	2
3	0.0012582259508064908	2	0	Est	17056.584577154976	0
4	0.0025701801760646963	2	1	Ouest	71115.50973101736	2
5	0.004157608722209716	1	0	Nord	-19813.295461324036	1
6	0.004203202703592979	2	0	Nord	188.7477264845802	1
7	0.004958708220783592	2	0	Est	17222.590444377183	0

Showing the first 1000 rows.

Command took 4.02 minutes -- by em\_karech@esi.dz at 6/1/2021, 10:52:01 AM on TP

#### 4) Groupage des colonnes dans un seul vecteur

```

1 from pyspark.ml.feature import VectorAssembler
2 numeric_cols = ["orientationOHE", "size", "nb_rooms", "garden"]
3 assembler_inputs = numeric_cols
4 vec_assembler = VectorAssembler(inputCols=assembler_inputs, outputCol='features')
5

```

Command took 0.08 seconds -- by em\_karech@esi.dz at 6/1/2021, 10:56:14 AM on TP

---

Cmd 16

```

1 from pyspark.ml.regression import LinearRegression
2 lr = LinearRegression(featuresCol = 'features', labelCol='label', regParam=0.3)
3
4
5

```

#### 5) Entraînement du modèle

```

1 from pyspark.ml import Pipeline
2
3 pipeline = Pipeline(stages=[string_indexer, one_hot_encoder, label_to_index, vec_assembler, lr])
4 pipeline_model = pipeline.fit(df_train)
5

```

On ne peut pas entraîner ce modèle avec la configuration actuelle de databricks on obtient une erreur de type Java heap space ce qui veut dire qu'on a pas assez de mémoire pour entraîner notre modèle

▶ (3) Spark Jobs

⊞ org.apache.spark.SparkException: Job aborted due to stage failure: Task 3 in stage 29.0 failed 1 times, most recent failure: Lost task 3.0 in stage 29.0 (TID 95) (ip-10-172-186-248.us-west-2.compute.internal executor driver): java.lang.OutOfMemoryError: Java heap space

Command took 8.16 minutes -- by em\_karech@esi.dz at 6/7/2021, 5:24:12 AM on TP

**On peut essayer avec un fichier d'une plus petite taille ou augmenter nos ressources**

## 6) Prédiction des valeurs

```
1 df_pred = pipeline_model.transform(df_test)
2 display(df_pred)|
```

## 7) Sauvegarder le résultat dans fichier

**df\_pred.write.csv(predcitions.csv)**