

---

## Nouveaux paradigmes de base de données

Wissem Inoubli – Nejat Arnik

IUT de Lens – BUT Info S5 – 2023–2024

---

### TP 5 : manipulation de données avec Apache Spark

---

## Exercice 1 Manipulation des RDDs (toutes les réponses doivent être avec les transformations ou actions sur les RDDs)

1. Créer un RDD (allData) à partir du fichier weblog.csv.

```
1 from pyspark.context import SparkContext
2 sc = SparkContext('local', 'TP Spark : correction')
3 alldata=sc.textFile("weblog.csv")
```

2. Déterminer le nombre de lignes de l’RDD créé.

```
1 print(alldata.count())
```

3. Créer un deuxième RDD (RDD200) à partir de l’RDD créé (allData) qui regroupe que les accès dont le code HTTP =200.

```
1 rdd200=alldata.map(lambda x:x.split(",")).filter(lambda x : x[3]=='200')
```

4. Calculer le pourcentage de succès (code200) des requêtes HTTP.

```
1 print(rdd200.count()/(alldata.coun()/100))
```

5. Déterminer pour chaque adresse IP le nombre d’accès.

```
1 pairIP=alldata.map(lambda x:(x.split(",")[0],1))
2 IpCount=pairIP.reduceByKey(lambda x,y:x+y)
3 IpCount.take(5)
```

6. Déterminez le nombre d’accès à la page login.php pour chaque adresse IP.

```
1 pairIP=alldata.filter(lambda x: "login.php" in x.split(",")[2]).map(lambda
    x:(x.split(",")[0],1)).reduceByKey(lambda x,y:x+y)
2 pairIP.take(5)
```

7. Trier et afficher les adresses IP selon leurs accès.

```
1 pairIP.sortBy(lambda x:x[1]).take(5)
```

8. Donnez les listes liste distinctes des codes HTTP disponibles dans le fichier.

```
1 alldata.map(lambda x: x.split(",")[3]).distinct().collect()
```

## Exercice 2 Manipulation des RDDs (toutes les réponses doivent être avec les transformations ou actions sur les RDDs)

1. Créer un RDD (vols) à partir du fichier vols1.csv.

```
1 vols=sc.textFile("vols1.csv")
```

2. Afficher les noms (distincts) de compagnies qui desservent les villes suivantes (mumbai, bangalore, delhi). La liste des villes doit être partagé avec toutes les partitions avec une variable de type broadcast.

```
1 broadcast_villes = sc.broadcast(set(["Mumbai", "Bangalore", "Delhi"]))
2 vols.map(lambda vol:vol.split(",")).filter(lambda vol: vol[3] in
    broadcast_villes.value).map(lambda x:x[1]).distinct().collect()
```

3. Calculer les revenus : par classe de vol, par compagnie.

```
1 VolswithoutHeader=vols.filter(lambda v:"airline" not in v )
2 projection = VolswithoutHeader.map(lambda line: line.split(',') ).map(lambda
    values: ((values[1], values[8]), float(values[11])))
3 prixComp = projection.reduceByKey(lambda x, y: x + y)
4 prixComp.take(15)
```

4. Calculer pour chaque compagnie la somme de toutes les durées de vols, et trier les compagnies selon cette somme.

```
1 projection = VolswithoutHeader.map(lambda line: line.split(',') ).map(lambda
    values: (values[1], float(values[9])))
2 DureeComp = projection.reduceByKey(lambda x, y: x + y).sortBy(lambda x:x[1])
3 DureeComp.take(15)
```

5. Créer les deux RDDs vols2 (vols2.csv) aéroport (aéroport.csv).

```
1 vols2=sc.textFile("vols2.csv")
2 aeroport=sc.textFile("aéroport.csv")
3 #les ent tes de chaque rdd
4 enteteVols = vols2.first()
5 enteteAeroport = aeroport.first()
```

6. A partir de l’RDD aéroport calculer le nombre d’aéroports pour chaque *state*.

```

1 aeroportByState = aeroport.filter(lambda a: a!=enteteAeroport).map(lambda line:
    line.split(',')).map(lambda values: (values[2], 1)).reduceByKey(lambda x,y:x+y)
2 aeroportByState.take(5)

```

7. Déterminez le nombre de vols partant de la ville de Nome après le 15 du mois (avec la transformation join()).

```

1 # Cr ation de paires (a roport_id , city) pour les a roports
2 part1 = aeroport.filter(lambda a: a!=enteteAeroport).map(lambda line:
    line.split(',')).filter(lambda a:a[1]=="Nome").map(lambda values: (values[0],
    values[1]))
3 # Filtre des vols et cr ation de paire (a roport_id , les info de vols)
4 part2 = vols2.filter(lambda v: v!=enteteVols).map(lambda line:
    line.split(',')).filter(lambda values: int(values[0]) > 15).map(lambda x:
    (x[3],x))
5 #jointure
6 join=part2.join(part1)
7 join.take(5)

```

8. Calculer le nombre de vols par jour de la semaine et par jour du mois.

```

1 volParSemaine = vols2.filter(lambda v: v!=enteteVols).map(lambda line:
    line.split(',')).map(lambda values: (values[1], 1))
2 volParMois = vols2.filter(lambda v: v!=enteteVols).map(lambda line:
    line.split(',')).map(lambda values: (values[0], 1))
3
4 NbVolParSemaine = volParSemaine.reduceByKey(lambda x, y: x + y)
5 NBVolParMois = volParMois.reduceByKey(lambda x, y: x + y)

```

## Exercice 3 Manipulation des DataFrames

1. Créer un dataframe df.vols à travers la création de un RDD (fichier vols1.csv).

```

1 from pyspark.sql import SparkSession
2 spark = SparkSession.builder.appName("FlightData").getOrCreate()
3 dataframe = spark.read.option("header", True).csv("vols1.csv")
4 dataframe.show(5)

```

2. Afficher le schéma du dataframe créé.

```

1 dataframe.printSchema()

```

3. Afficher les noms de compagnies.

```

1 dataframe.select("Airline").show()

```

4. Afficher les nombre de vols pour chaque compagnies.

```

1 dataframe.groupBy("Airline").count().show()

```

5. Afficher les compagnies (seulement les noms) pour lesquelles le nombre de vols assurés dépasse 10.

```
1 from pyspark.sql.functions import col
2 volCount=df.groupBy("Airline").count()
3 volCount.filter(col("count")>10).show()
```

6. Enregistrer le résultat de la question précédente dans un fichier csv.

```
1 volCount.write.csv("rst.csv", header=True)
```

7. Créer deux dataframes, df\_vols2 et df\_aeroport à travers la création de deux RDD (de deux fichiers vols2.csv et aeroport.csv).

```
1 vols2=sc.textFile("vols2.csv")
2 aeroport=sc.textFile("aeroport.csv")
3 #les entetes de chaque rdd
4 enteteVols = vols2.first()
5 enteteAeroport = aeroport.first()
```

8. Déterminez le nombre de vols partant de la ville de "Nome" après le 15 du mois.

```
1 enteteVols = vols2.first()
2 enteteAeroport = aeroport.first()
3 # Cr ation de paires (a aeroport_id , city) pour les a eroports
4 part1 = aeroport.filter(lambda a: a!=enteteAeroport).map(lambda line:
    line.split(',')).filter(lambda a:a[1]=="Nome").map(lambda values: (values[0],
    values[1]))
5 # Filtre des vols et cr ation de paire (a aeroport_id , les info de vols)
6 part2 = vols2.filter(lambda v: v!=enteteVols).map(lambda line:
    line.split(',')).filter(lambda values: int(values[0]) > 15).map(lambda x:
    (x[3],x))
7 #jointure
8 join=part2.join(part1)
9 join.take(5)
```

9. Calculer le nombre de vols par jour de la semaine et par jour du mois.

```
1 # Cr ation de paires RDD avec (jour de la semaine, 1) et (jour du mois, 1)
2 volParSemaine = vols2.filter(lambda v: v!=enteteVols).map(lambda line:
    line.split(',')).map(lambda values: (values[1], 1))
3 volParMois = vols2.filter(lambda v: v!=enteteVols).map(lambda line:
    line.split(',')).map(lambda values: (values[0], 1))
4
5 NbVolParSemaine = volParSemaine.reduceByKey(lambda x, y: x + y)
6 NBVolParMois = volParMois.reduceByKey(lambda x, y: x + y)
7 NbVolParSemaine.take(5)
```

10. Créer deux dataframes df\_vols3 et df\_aeroport2 directement à partir de deux fichiers vols2.csv et aeroport.csv.

```
1 from pyspark.sql import SparkSession
2 spark = SparkSession.builder.appName("FlightData").getOrCreate()
3
4 df_vols3 = spark.read.option("header", True).csv("vols2.csv")
5 df_aeroport2 = spark.read.option("header", True).csv("aeroport.csv")
```

11. Afficher les deux schémas de deux dataframes créés, puis les comparer avec les dataframes créés lors de la question 7.

```
1 df_vols3.printSchema()
2 df_aeroport2.printSchema()
3 df_vols3.show(5)
4 df_aeroport2.show(5)
```

## Exercice 4 Manipulation des DataFrames avec SparkSQL

1. Créer 3 tables temporaires vols1, vols2 et aeroport à partir de fichiers vols2.csv, vols1.csv et aeroport.csv.

```
1 from pyspark.sql import SparkSession
2 spark = SparkSession.builder.appName("FlightData").getOrCreate()
3 df_vols1 = spark.read.option("header", True).csv("vols1.csv")
4 df_vols2 = spark.read.option("header", True).csv("vols2.csv")
5 df_aeroport = spark.read.option("header", True).csv("aeroport.csv")
6 df_vols1.createOrReplaceTempView("vols1")
7 df_vols2.createOrReplaceTempView("vols2")
8 df_aeroport.createOrReplaceTempView("aeroport")
```

2. En utilisant la table vols1, afficher les nom de compagnies dont le nombre de stops est égala à "zero"

```
1 resultat = spark.sql("SELECT airline FROM vols1 WHERE stops = 'zero'")
2 resultat.show()
```

3. Avec la même table, afficher les noms de compagnies ainsi que la somme de toutes les durées de vols pour chaque compagnie.

```
1 resultat = spark.sql("SELECT airline , SUM(duration) AS somme_duree_vols FROM vols1
2 GROUP BY airline")
3 resultat.show()
```

4. Toujours Avec la même table, retourner pour chaque compagnie les revenus des vols économiques.

```
1 rev = spark.sql("SELECT airline , SUM(price) AS revenus_economiques FROM vols1
2 WHERE class = 'Economy' GROUP BY airline")
3 rev.show()
```

5. Avec les deux tables (vols2 et aeroports), déterminez le nombre de vols partant de la ville de "Nome" après le 15 du mois.

```
1 resultat = spark.sql("""
2     SELECT COUNT(*) AS nombre_de_vols
3     FROM vols2 v
4     JOIN aeroport a ON v.OriginAirportID = a.airport_id
5     WHERE a.city = 'Nome' AND v.DayOfMonth > 15
6 """)
7 resultat.show()
```

6. Avec les mêmes tables de la question précédente, calculer la somme de retard de départ pour tous aéroports de "Ak".

```
1 resultat = spark.sql("""
2     SELECT a.state, SUM(v.DepDelay) AS somme_retard_depart
3     FROM vols2 v
4     JOIN aeroport a ON v.OriginAirportID = a.airport_id
5     WHERE a.state = 'Ak'
6     GROUP BY a.state
7 """)
8 resultat.show()
```

7. Identifier le jour de la semaine ayant enregistré le plus grand retard au départ.

```
1 rst = spark.sql("""
2     SELECT DayOfWeek, SUM(DepDelay) AS somme_retard_depart
3     FROM vols2
4     GROUP BY DayOfWeek
5     ORDER BY somme_retard_depart DESC
6     LIMIT 1
7 """)
8 rst.show()
```

## Exercice 5 Manipulation des données dans cassandra avec Spark

1. Créer deux tables vols et aeroport dans une keyspace cassandra nommé "spark".
2. Charger les deux tables dans des Dataframes spark.
3. Afficher les schémas de deux tables.
4. A partir des dataframes créés et spark SQL répondez aux requêtes suivantes :
  1. Liste de aéroport de l'état "AL"
  2. Liste des noms de l'aéroport de départ et leurs nombres de vols.
  3. Afficher pour chaque aéroport le "name", "state", moyenne des retard de départ.