# Implémentation parallèle de certains algorithmes de fouille de données avec le framework MapReduce

## Master 1 : Big Data Analytics

Ousmane DIA :
Abdoulaye Bara DIAW:

# PLAN

1-Classique du MapReduce

2-Parallélisation de l'algorithme K-means:
-Implémentation de K-means sur MapReduce
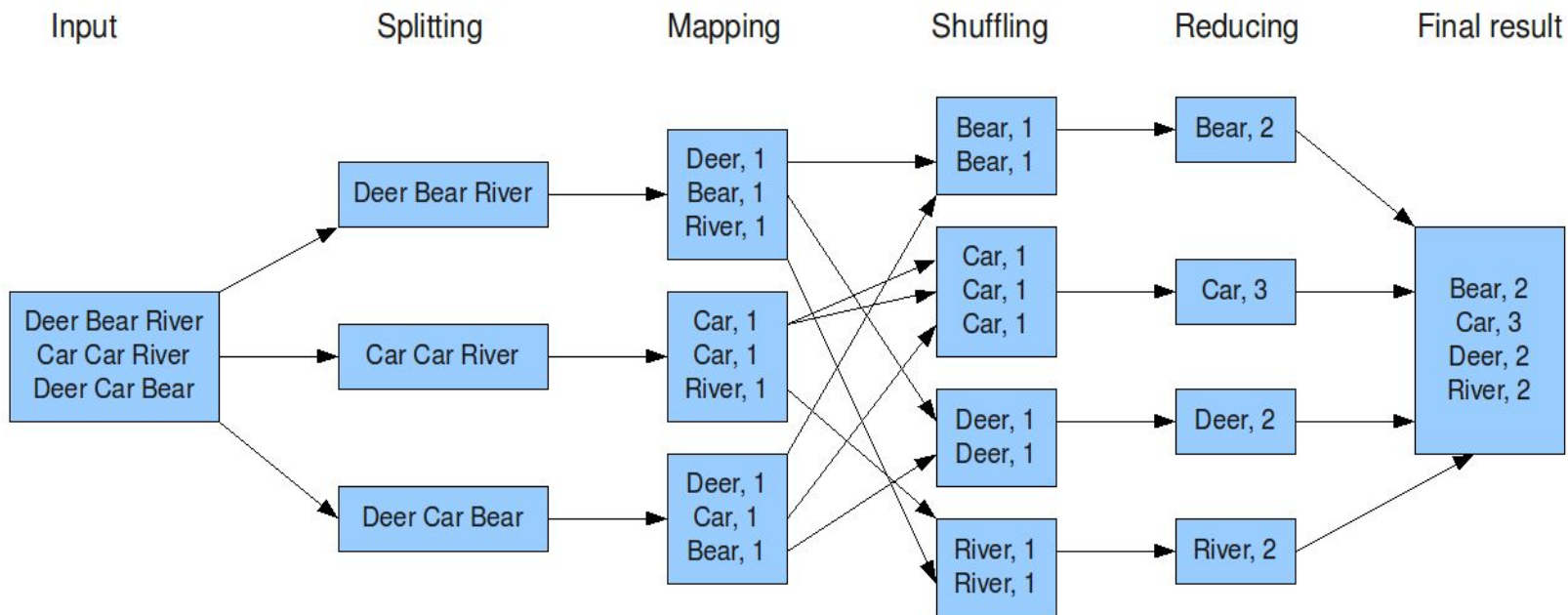
3-Parallélisation de l'algorithme Apriori
-Implémentation de Apriori sur MapReduce

# 1-Classique du Mapreduce

Description et méthode de Mapreduce

# Description et Méthode



The overall MapReduce word count process

| Input | Splitting | Mapping | Shuffling | Reducing | Final result |
|-------|-----------|---------|-----------|----------|--------------|

Deer Bear River
Car Car River
Deer Car Bear

Deer Bear River

Car Car River

Deer Car Bear

Deer, 1
Bear, 1
River, 1

Car, 1
Car, 1
River, 1

Deer, 1
Car, 1
Bear, 1

Bear, 1
Bear, 1

Car, 1
Car, 1
Car, 1

Deer, 1
Deer, 1

River, 1
River, 1

Bear, 2

Car, 3

Deer, 2

River, 2

Bear, 2
Car, 3
Deer, 2
River, 2

# 1-Classique du Mapreduce

Implémentation de Mapreduce

# Implémentation

Data: input.txt
id , age , sexe , adresse , salaire

```
 0,  25,homme,oran,28000
 1,  33,homme,oran,28000
 2,  46,homme,oran,54000
 3,  35,famme,oran,33000
 4,  23,famme,oran,25000
 5,  25,famme,mascara,25000
 6,  25,homme,oran,38000
 7,  33,homme,oran,38000
 8,  46,homme,oran,54000
 9,  35,famme,oran,33000
10,23,famme,oran,29000
11,25,famme,mascara,25000
12,25,homme,oran,28000
13,19,homme,oran,18000
14,46,homme,oran,45000
15,35,famme,oran,33000
16,23,famme,oran,23000
17,25,famme,mascara,21000
```

Calculer le maximum et le minimum du Salaire

## Implémentation

**Démarrage du CLuster Hadoop**
```
start-all.sh
```

**Création d'un répertoire HDFS**
```
hdfs dfs -mkdir /user
hdfs dfs -mkdir /user/ousmanealhayri/
hdfs  dfs  -mkdir  /user/ousmanealhayri/OusmaneAlhayri
alhayri1234
hdfs dfs -mkdir datainput
hdfs dfs -put data/input.txt datainput/
hdfs dfs -put data/input.txt datainput/
```

# Implémentation

Affichage des données : `hdfs dfs -head datainput/input.txt`

```
0,25,homme,oran,28000
1,33,homme,oran,28000
2,46,homme,oran,54000
3,35,famme,oran,33000
4,23,famme,oran,25000
5,25,famme,mascara,25000
6,25,homme,oran,38000
7,33,homme,oran,38000
8,46,homme,oran,54000
9,35,famme,oran,33000
10,23,famme,oran,29000
11,25,famme,mascara,25000
12,25,homme,oran,28000
13,19,homme,oran,18000
14,46,homme,oran,45000
15,35,famme,oran,33000
16,23,famme,oran,23000
17,25,famme,mascara,21000
```

## Implémentation

Mapper : `head ../data/input.txt | python mapper.py`

Reducer: `head ../data/input.txt |./mapper.py |sort |./reducer.py`



```
(base) ousmanealhayri@ousmanealhayri-ThinkPad-T450:
xt | python mapper.py
25      28000   1
33      28000   1
46      54000   1
35      33000   1
23      25000   1
25      25000   1
25      38000   1
33      38000   1
46      54000   1
35      33000   1
```



```
(base) ousmanealhayri@ousmanealhayri-ThinkPad-T450
xt |./mapper.py |sort |./reducer.py
23      25000   25000   1
25      38000   25000   3
33      38000   28000   2
35      33000   33000   2
46      54000   54000   2
```

Mapper Code Source:
https://github.com/Data-Mining-on-Hadoop-Mapreduce/Simple-Exemple-of-Mapreduce/blob/main/Code/mapper.py

Reducer Code Source:
https://github.com/Data-Mining-on-Hadoop-Mapreduce/Simple-Exemple-of-Mapreduce/blob/main/Code/reducer.py

# 2-Parallélisation de l'algorithme K-means

## Description et méthode de l'algo k-means

## Définition :

*Algorithme de clustering qui regroupe les éléments d'un dataset en groupes qui se ressemblent appelés **clusters.***

Pour ce faire, l'algorithme s'appuie sur deux éléments essentiels : **le centroïde et la distance**.

## Algorithme :

- Initialisation (choix au hasard de k centroïdes)
- Construction de k clusters
- Calcul des nouveaux centroïde
- On recommence jusqu'à ce qu'il y ait convergence

# 2-Parallélisation de l'algorithme K-means

## Implémentation de k-means sur Mapreduce

hdfs utilise le processeur en le subdivisant en noeuds.
Après notre première opération les nœuds ont été chargés, donc il faut les libérer pour pouvoir effectuer d'autres opérations.

Pour ce faire, nous procédons en trois étapes :

- Supprimer les fichiers temporaires:
```
sudo rm -R /tmp/*
```

- Supprimer les fichiers du datanote/
```
rm -rvf /hadoop/hdfs/datanode/*
```

- Formater et redémarrer
```
stop-all.sh
hdfs namenode -format
start-all.sh
jps
```

# Implémentation

Affichage des données : `hdfs dfs -head datainput/dataset.txt`

**Mapper :** `head ../data/dataset.txt | python mapper.py`



k-means Mapper Code Source :
https://github.com/Data-Mining-on-Hadoop-Mapreduce/K-means-Algorithm-on-Hadoop/blob/main/code/mapper.py

# Implémentation

**Reducer:** `head ../data/dataset.txt |./mapper.py |sort |./reducer.py`



```
Programmation paralléle/k-mean with Hadoop/code$ head -50 ../data/dataset.txt .
/mapper.py |sort |./reducer.py
10.1728565969406, 24.871711836734608
```

Reducer Code Source :
https://github.com/Data-Mining-on-Hadoop-Mapreduce/K-means-Algorithm-on-Hadoop/blob/main/code/reducer.py

# 2-Parallélisation de l'algorithme Apriori

## Description et méthode de Apriori

## Description et méthode

L'algorithme a priori fait référence à l'algorithme utilisé pour calculer les règles d'association entre les objets. Cela signifie savoir comment deux ou plusieurs objets sont liés les uns aux autres.

# Description et méthode

$$Support = \frac{frq(X,Y)}{N}$$

$$Rule: X \Rightarrow Y$$

$$Confidence = \frac{frq(X,Y)}{frq(X)}$$

$$Lift = \frac{Support}{Supp(X) \times Supp(Y)}$$

| Transactions | Itemset |
|---|---|
| T1 | {B, C, D, E} |
| T2 | {B, C, D} |
| T3 | {A, B, D} |
| T4 | {A, B, C, D, E} |
| T5 | {A, B, C} |
| T6 | {B, E} |

(a)

| (1-Itemsets) | Transactions |
|---|---|
| {A} | T3, T4, T5 |
| {B} | T1, T2, T3, T4, T5, T6 |
| {C} | T1, T2, T4, T5 |
| {D} | T1, T2, T3, T4 |
| {E} | T1, T4, T6 |

(b)

**Figure 1. An example to illustrate frequent itemsets mining**

# 2-Parallélisation de l'algorithme Apriori

## Implémentation de Apriori sur Mapreduce

# Implémentation

Affichage des données : `hdfs dfs -head datacsvinput/csv_dataset_sample.csv`

```
(base) ousmanealhayri@ousmanealhayri-ThinkPad-T450:~/Documents/Master Big Data/
Programmation paralléle/Apriori with Hadoop$ hdfs dfs -head datacsvinput/csv_da
taset_sample.csv
butter,butter
butter,ham
pepper
salt,bread,pepper,pepper,butter,rice,butter
ham,ham,corn,rice,bread,ham,butter
corn,ham,pepper,cheese,salt,rice,salt,ham
bread,corn,salt,bread,cheese,ham,bread,rice
corn,ham,salt,cheese,salt,salt,salt
salt,salt
rice
salt,corn,corn,rice,corn,salt,rice
ham,bread,ham,corn
pepper,butter,corn,salt,butter,pepper,bread,cheese
corn,salt,butter,cheese,ham
cheese,corn,pepper,ham,corn,salt
butter,bread,bread,ham,pepper,bread
bread,corn,pepper,bread,pepper,bread,butter,salt
corn,bread
corn,corn,salt,butter,rice,corn
```

## Implémentation

**Mapper** : `head -50 ../data/csv_dataset_sample.csv | python Apriori_mapper.py`



AprioriMapper Code source :
https://github.com/Data-Mining-on-Hadoop-Mapreduce/Apriori-Algorithm-on-Hadoop/blob/main/code/Apriori_mapper.py

## Implémentation

**Reducer :** `head ../data/csv_dataset_sample.csv |./Apriori_mapper.py |sort |./Apriori_reducer.py`

```
(base) ousmanealhayri@ousmanealhayri-ThinkPad-T450:~/Documents/Master Big Data/P
rogrammation paralléle/Apriori with Hadoop/code$ head ../data/csv_dataset_sample
.csv |./Apriori_mapper.py |sort |./Apriori_reducer.py
discarded        ('bread',),
discarded        ('butter',),
discarded        ('cheese',),
discarded        ('corn',),
discarded        ('ham',),
discarded        ('pepper',),
discarded        ('rice',),
discarded        ('salt',),
```

AprioriReducer Code Source:
https://github.com/Data-Mining-on-Hadoop-Mapreduce/Apriori-Algorithm-on-Hadoop/blob/main/code/Apriori_reducer.py

FIN