

# Base de données NOSQL et Big Data

## Projet semestrielle: MapReduce

### Introduction :

L'objectif de cette analyse est d'explorer les données historiques sur les prix des actions de 500 grandes entreprises au cours des cinq dernières années. Les données sont extraites de Yahoo Finance et incluent des métriques quotidiennes telles que l'ouverture, la fermeture, le volume, etc. L'analyse vise à fournir des insights pour les investisseurs, les chercheurs et les passionnés d'apprentissage automatique.

### Étape 1: Exploration des Données

Dans cette étape, nous avons chargé le jeu de données à l'aide de la bibliothèque Pandas en Python. Nous avons affiché les premières lignes du jeu de données pour comprendre sa structure et vérifié la présence de valeurs manquantes et effectué des statistiques de base.

```
1 import pandas as pd
2 import numpy as np
3
4 # Charger Le jeu de données
5 df = pd.read_csv('stock_details_5_years.csv')
6
7 # Afficher Les premières lignes du jeu de données
8 print(df.head())
9
10 # Vérifier Les valeurs manquantes
11 print(df.isnull().sum())
12
13 # Statistiques de base
14 print(df.describe())
```

### Étape 2: Nettoyage des Données

Dans cette étape, nous avons supprimé les lignes avec des valeurs manquantes, éliminé les doublons et converti la colonne de date au format datetime. Nous avons également assuré que les types de données étaient appropriés.

```
1 # Supprimer Les lignes avec des valeurs manquantes
2 df = df.dropna()
3
4 # Vérifier et supprimer Les doublons
5 df = df.drop_duplicates()
6
7 # Convertir La colonne de date au format datetime
8 df['Date'] = pd.to_datetime(df['Date'])
9
10 # Assurer que Les types de données sont appropriés
11 df['Volume'] = df['Volume'].astype(int)
12
13 # Sauvegarder Le jeu de données nettoyé
14 df.to_csv('cleaned_dataset.csv', index=False)
```

### Étape 3: Visualisation des Données

Nous avons utilisé la bibliothèque Matplotlib pour visualiser les prix des actions d'une entreprise spécifique au fil du temps.

```
1 import matplotlib.pyplot as plt
2
3 #
4     Visualiser les prix des actions pour une entreprise spécifique
5 company_data = df[df['Company'] == 'AAPL'] #
6     Remplacer 'AAPL' par le symbole de l'entreprise souhaitée
7 plt.plot(company_data['Date'], company_data['Close'])
8 plt.title('Prix des actions au fil du temps')
9 plt.xlabel('Date')
10 plt.ylabel('Prix de clôture')
11 plt.show()
```

### ■ MapReduce avec Hadoop

Nous avons conçu une tâche MapReduce pour calculer le prix de clôture moyen par entreprise sur l'ensemble du jeu de données. Le script MapReduce est écrit en Java et comprend deux parties principales : le mappage et la réduction. Le mappage extrait les paires clé-valeur à partir des données, tandis que la réduction agrège ces paires pour calculer la moyenne.

#### ● Mapper.py:

```
1 #!/usr/bin/env python
2
3 import sys
4
5 # Entrée : Date,Open,High,Low,Close,Volume,Dividends,Stock Splits,Company
6 for line in sys.stdin:
7     # Divise la ligne d'entrée en jetons
8     tokens = line.strip().split(",")
9
10    # Vérifie si la ligne a suffisamment de jetons
11    if len(tokens) >= 9:
12        company = tokens[8]
13        closing_price = float(tokens[4])
14
15        # Sortie paire clé-valeur : (company, closing_price)
16        print("{}\t{}".format(company, closing_price))
17
```

## ● Reducer.py:

```
1  #!/usr/bin/env python
2
3  import sys
4
5  current_company = None
6  total_closing_price = 0.0
7  count = 0
8
9  # Lit l'entrée à partir de la saisie standard
10 for line in sys.stdin:
11     # Divise la ligne d'entrée en clé et valeur
12     company, closing_price = line.strip().split("\t")
13     closing_price = float(closing_price)
14
15     # Vérifie si l'entreprise a changé
16     if current_company != company:
17         # Résultat de sortie pour l'entreprise précédente
18         if current_company:
19             average_closing_price = total_closing_price / count
20             print("{}\t{}".format(current_company, average_closing_price))
21
22         # Réinitialise les variables pour la nouvelle entreprise
23         current_company = company
24         total_closing_price = closing_price
25         count = 1
26     else:
27         # Met à jour les variables pour l'entreprise en cours
28         total_closing_price += closing_price
29         count += 1
30
31 # Résultat de sortie pour la dernière entreprise
32 if current_company:
33     average_closing_price = total_closing_price / count
34     print("{}\t{}".format(current_company, average_closing_price))
35
```

## ■ MapReduce avec multiprocessing (Python)

Nous avons également implémenté une version MapReduce en utilisant le module multiprocessing en Python. Cela permet de tirer parti du traitement parallèle sur une machine locale.

### Étape 4: MapReduce Design

Dans cette étape, nous concevons le processus MapReduce en définissant les fonctions de mappage et de réduction ainsi que la stratégie de division des données en morceaux pour le traitement parallèle.

```
1  import matplotlib.pyplot as plt
2
3  #
4  # Visualiser les prix des actions pour une entreprise spécifique
5  company_data = df[df['Company'] == 'AAPL'] #
6  # Remplacer 'AAPL' par le symbole de l'entreprise souhaitée
7  plt.plot(company_data['Date'], company_data['Close'])
8  plt.title('Prix des actions au fil du temps')
9  plt.xlabel('Date')
10 plt.ylabel('Prix de clôture')
11 plt.show()
```

### Étape 5: MapReduce Implementation

Dans cette étape, nous mettons en œuvre le processus MapReduce en utilisant le module multiprocessing en Python. Cela permet de tirer parti du traitement parallèle sur une machine locale.

```
1 # Créer un pool de travailleurs
2 with Pool(num_chunks) as pool:
3     # Phase de map :
4     #
5     Appliquer la fonction de map à chaque morceau en parallèle
6     mapped_results = pool.map(map_function, data_chunks)
7
8     # Phase de réduction :
9     # Combinaison des résultats de la phase de map
10    final_result = reduce_function(mapped_results)
```

### Étape 6: MapReduce Execution

Dans cette étape, nous exécutons le processus MapReduce en affichant le résultat final, qui peut être utilisé pour des analyses plus approfondies.

```
1 # Display the final result
2 print(final_result)
```

### Conclusion

Ce rapport a couvert toutes les étapes du processus d'analyse des données sur les prix des actions, de l'exploration initiale à l'application de MapReduce avec la version parallèle en utilisant le module multiprocessing en Python. Ces méthodes offrent une approche scalable pour traiter des volumes importants de données. Les résultats peuvent être utilisés pour prendre des décisions éclairées en matière d'investissement, d'analyse de marché et de modélisation prédictive.