

Spark Streaming

Vélib ile de france

DOUNIA HARAG
CHAIMAE MOUH
AYOUB ERRADI
CHRISTIAN NZEGAING



Sommaire



1 Contexte

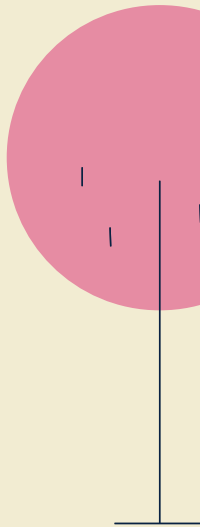
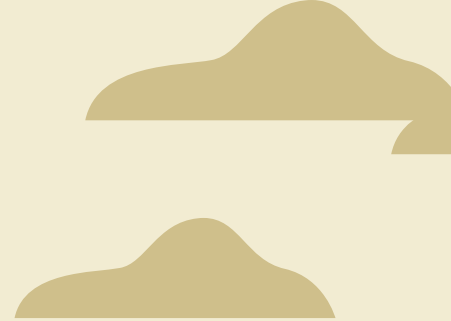
2 DATA



3 Application

3 Visualisation

4 Prédiction



Contexte :



- Automatisation du traitement des données de disponibilité des vélos Velib en île de France avec Apache Spark.
- Calcul des indicateurs de performance clés (KPI) pour chaque station.
- Consolidation des résultats pour une analyse approfondie.
- Visualisation des données sur Power BI

DATASET

Root

- stationcode : string
- name : string
- is_installed : string
- capacity : Int
- numdocksavailable : Int
- numbikesavailable : Int
- mechanical : Int
- ebike : Int
- is_renting : string
- is_returning : string
- duedate : string
- coordonnees_geo_long : string
- coordonnees_geo_lat : string
- nom_arrondissement_communes : string
- code_insee_commune : string



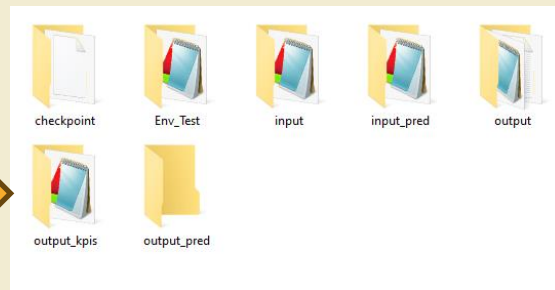
Traitement de données

1. Origine de dataset : data.gouv.fr

Ressources

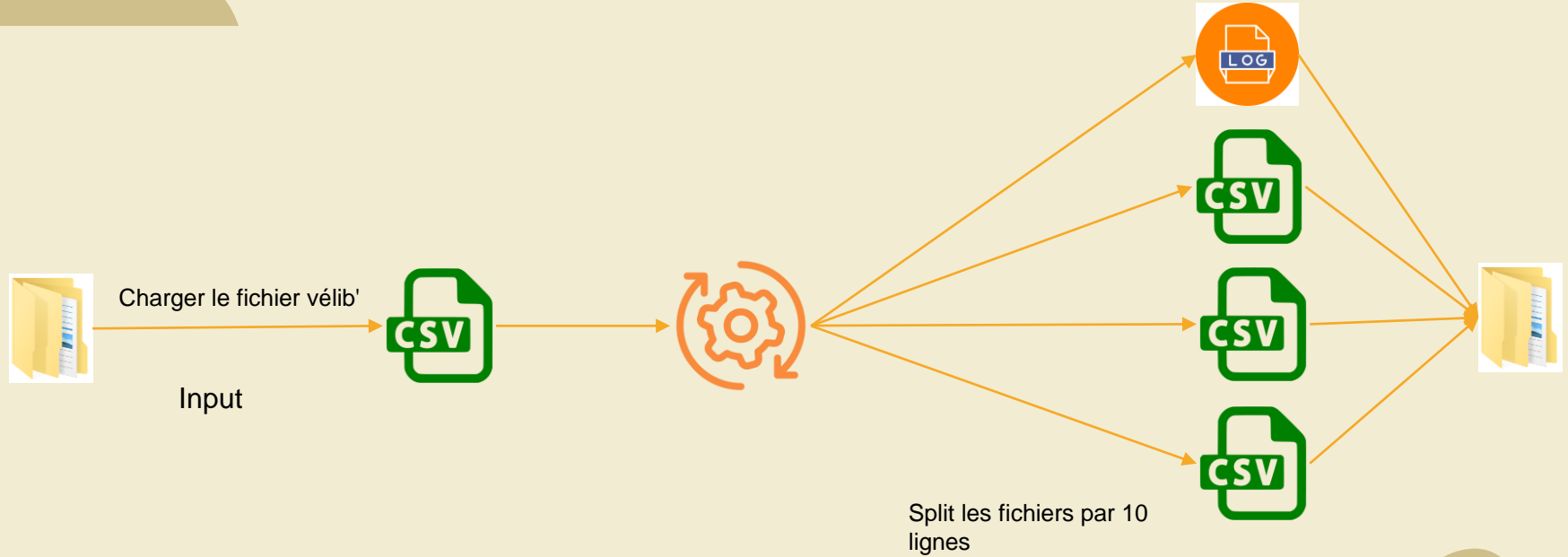
| | |
|---|--|
| <p>velib-disponibilite-en-temps-reel.zip</p> <p>🔄 inconnue Dernière modification</p> <p>🕒 100% Taux de disponibilité</p> <p>zip</p> <p>+ détails Télécharger</p> | <p>velib-disponibilite-en-temps-reel.json</p> <p>🔄 inconnue Dernière modification</p> <p>🕒 100% Taux de disponibilité</p> <p>json</p> <p>+ détails Télécharger</p> |
| <p>velib-disponibilite-en-temps-reel.geojson</p> <p>🔄 inconnue Dernière modification</p> <p>🕒 99.9% Taux de disponibilité</p> <p>geojson</p> <p>+ détails Télécharger</p> | <p>velib-disponibilite-en-temps-reel.csv</p> <p>🔄 inconnue Dernière modification</p> <p>🕒 100% Taux de disponibilité</p> <p>csv</p> <p>+ détails Télécharger</p> |

Stocker le fichier Velib-disponibilite-en-temps-reel.csv sur le répertoire input



2. Répertoire de données

Lire .Csv



Lire .Csv

```
scala, LireCSVVelib.scala
application.scala x main.scala x build.sbt x LireCSVVelib.scala x
1  import ...
5
6  object LireCSVVelib {
7    def main(args: Array[String]): Unit = {
8      // Créer une session Spark
9      val spark = SparkSession.builder
10        .appName("LireCSVVelib")
11        .master("local[*]") // Exécuter localement avec toutes les CPU disponibles
12        .getOrCreate()
13
14      // Lire le fichier CSV
15      val filePath = "C:\\Users\\DELL\\Documents\\input\\velib-disponibilite-en-temps-reel.csv"
16      val df = spark.read
17        .option("header", "true")
18        .option("inferSchema", "true")
19        .option("delimiter", ";") Choose schema // Spécifier le délimiteur correct
20        .csv(filePath)
21
22      // Définir le chemin du répertoire de sortie
23      val outputDir = "C:\\Users\\DELL\\Documents\\output"
24      val logFile = s"$outputDir\\log.txt"
25
26      val totalRows = df.count().toInt
27      val batchSize = 10
28
29      // Initialiser le fichier de log
```



Application

```
// Lire tous les fichiers du répertoire de sortie
val inputDir = "C:\\Users\\DELL\\Documents\\output"
val files = new File(inputDir).listFiles.filter(_.getName.endsWith(".csv")).sorted

val kpiDir = "C:\\Users\\DELL\\Documents\\output_kpis"
new File(kpiDir).mkdirs()

val outputFilePath = s"$kpiDir\\resultat_kpis.csv"

// Initialiser le fichier de sortie avec l'en-tête correct
val header = "stationcode;name;nom_arrondissement_communes;coordonnees_geo;duedate;total_bikes_available;avg_oc"
Files.write(Paths.get(outputFilePath), (header + "\\n").getBytes, StandardOpenOption.CREATE, StandardOpenOption.APPEND)

// Boucle pour traiter les fichiers périodiquement
var nextRunTime = LocalDateTime.now()
files.foreach { file =>
    var currentTime = LocalDateTime.now()
    while (currentTime.isBefore(nextRunTime)) {
        // Boucle d'attente active jusqu'à ce que l'heure prévue pour le prochain traitement soit atteinte
        currentTime = LocalDateTime.now()
    }

    val kpiDF = processFile(file.getAbsolutePath, spark)
```



KPIS

Disponibilité des Bornes et Vélos

- Taux d'occupation des stations : $(\text{numbikesavailable}/\text{capacity}) \times 100$
- Taux de disponibilité des bornes : $(\text{numdocksavailable}/\text{capacity}) \times 100$
- Nombre de vélos disponibles

Répartition des Types de Vélos

- Pourcentage de vélos mécaniques : $\text{mechanical}/\text{numbikesavailable} \times 100$
- Pourcentage de vélos électriques : $\text{ebike}/\text{numbikesavailable} \times 100$

Performance des Stations

- Stations Pleines
- Stations Vides

Application – Fichier Final CSV

```
// Lire tous les fichiers du répertoire de sortie
val inputDir = "C:\\Users\\chaimae\\Projet_Spark_streaming\\src\\output"
val files = new File(inputDir).listFiles().filter(_.getName.endsWith(".csv")).sorted

val kpiDir = "C:\\Users\\chaimae\\Projet_Spark_streaming\\src\\output_kpi"
new File(kpiDir).mkdirs()

val outputFilePath = s"$kpiDir/resultat_kpis.csv"

// Initialiser le fichier de sortie avec l'en-tête correct
val header = "stationcode;name;nom_arrondissement_communes;coordonnees_geo;duedate;total_bikes_available;avg_oc"
Files.write(Paths.get(outputFilePath), (header + "\n").getBytes, StandardOpenOption.CREATE, StandardOpenOption.

// Boucle pour traiter les fichiers périodiquement
var nextRunTime = LocalDateTime.now()
files.foreach { file =>
    var currentTime = LocalDateTime.now()
    while (currentTime.isBefore(nextRunTime)) {
        // Boucle d'attente active jusqu'à ce que l'heure prévue pour le prochain traitement soit atteinte
        currentTime = LocalDateTime.now()
    }

    val kpiDF = processFile(file.getAbsolutePath, spark)

    val outputDirPath = s"$kpiDir/resultat_kpis_temp"
    kpiDF.coalesce(1).write.mode("overwrite").option("header", "true").option("delimiter", ";").csv(outputDirPath)
```

- Lecture des fichiers
- Création du répertoire de sortie
- Initialisation du fichier de sortie
- Boucle de traitement périodique



visualisation



Prédiction

```
// Générer des lignes supplémentaires pour chaque station avec des disponibilités variables
val availabilityVariations = Seq(
  ("numdocksavailable", 5),
  ("numdocksavailable", -5),
  ("numbikesavailable", 5),
  ("numbikesavailable", -5),
  ("mechanical", 2),
  ("mechanical", -2),
  ("ebike", 2),
  ("ebike", -2)
)
```

Prédiction

```
// Appliquer les variations de disponibilité
val dfWithVariations = availabilityVariations.foldLeft(dfProcessed) { (accDf, variation) =>
  val (column, value) = variation
  accDf.union(
    dfProcessed.withColumn(column, col = $"$column" + lit(value))
  )
}

// Définir les variations temporelles
val timeVariations = Seq(
  1, 2, 3 // Variations temporelles en heures
)

// Appliquer les variations temporelles
val dfWithTimeVariations = timeVariations.foldLeft(dfWithVariations) { (accDf, hours) =>
  accDf.union(
    dfWithVariations.withColumn( colName = "duedate", col = $"duedate" + expr( expr = s"INTERVAL $hours HOUR"))
  )
}
```

Prédiction

```
// Appliquer les variations de disponibilité
val dfWithVariations = availabilityVariations.foldLeft(dfProcessed) { (accDf, variation) =>
  val (column, value) = variation
  accDf.union(
    dfProcessed.withColumn(column, col = $"$column" + lit(value))
  )
}

// Définir les variations temporelles
val timeVariations = Seq(
  1, 2, 3 // Variations temporelles en heures
)

// Appliquer les variations temporelles
val dfWithTimeVariations = timeVariations.foldLeft(dfWithVariations) { (accDf, hours) =>
  accDf.union(
    dfWithVariations.withColumn(colName = "duedate", col = $"$duedate" + expr(expr = s"INTERVAL $hours HOUR"))
  )
}
```

Prédiction

```
// Définir le modèle RandomForestRegressor pour la régression multi-cible
val rfMechanical = new RandomForestRegressor()
    .setLabelCol("mechanical")
    .setFeaturesCol("features")
    .setPredictionCol("pred_mechanical")

val rfEbike = new RandomForestRegressor()
    .setLabelCol("ebike")
    .setFeaturesCol("features")
    .setPredictionCol("pred_ebike")

// Créer des pipelines pour chaque modèle
val pipelineMechanical = new Pipeline().setStages(Array(assembler, rfMechanical))
val pipelineEbike = new Pipeline().setStages(Array(assembler, rfEbike))

// Diviser les données en ensembles d'entraînement et de test
val Array(trainingData, testData) = dfWithFeatures.randomSplit(Array(0.7, 0.3))
```

Prédiction

```
// Bornage des prédictions entre 0 et la capacité de la station
val predictionsBounded = predictions
  .withColumn( colName = "pred_mechanical", when( condition = $"pred_mechanical" < 0, value = 0)
    .when( condition = $"pred_mechanical" > $"capacity", value = $"capacity")
    .otherwise( value = $"pred_mechanical"))
  .withColumn( colName = "pred_ebike", when( condition = $"pred_ebike" < 0, value = 0)
    .when( condition = $"pred_ebike" > $"capacity", value = $"capacity")
    .otherwise( value = $"pred_ebike"))
```


Merci !

