

CLEANING DATA WITH SPARK

Spark Schema

- Es un esquema que indica el tipo de dato de cada columna del dataframe
- Ejemplo →


```
import pySpark.sql.types
schema = StructType([
    StructField('name', StringType(), True),
    ...
])
```

Annotations:
 - `import pySpark.sql.types`: ¿Puede ser así?
 - `name`: nombre
 - `StringType()`: tipo de dato
 - `True`: ¿Puede ser null?
- Analiza la lectura de datos
- Se aplica el esquema a la lectura de datos

`spark.df = spark.read.format('csv').load(nombre..., schema=...)`

Immutability and lazy processing

Spark DataFrames son inmutables.

Entendiendo parquet

CSV es lento de parsear

si no hay un esquema, debe leer toda la datos por crear uno

Robustos de CSV y Spark

Parquet

- Estos archivos incluyen automáticamente un schema.
- No necesita leer toda la información
- Como verbo → leer parquet →


```
df = spark.read.format('parquet').load('archivo.parquet')
df = spark.read.parquet('archivo.parquet')
```
- convertir archivo →


```
df.write.format('parquet').save('nombre.parquet')
df.write.parquet('nombre.parquet')
```
- Usos → Funciones bien con queries SQL (mejora el rendimiento).

DataFrame column operations I

Para crear condiciones

- Filter valores NO NULLS → `df[col1 != NULL]`
- " " NULLS → `df[col1 == NULL]`
- Filter strings que cambie un texto → `df[nombre != 'Juan']`

Annotation: con `is` funciona como con `not`

Funciones de transformación de strings

(import pySpark.sql.functions as F)

- En pandas es `reagrupados` → `F.group(column)`
- En dividir texto → `F.split(column, '!')`
- Para cambiar el tipo → `df[col1].cast(IntegerType())`
- Para interactuar con valores que son de tipo `ArrayType()`

Annotations:
 - `len` → `.size(col)`
 - Retorno un `df` de la lista → `.getItems(df)`

CLEANING DATA WITH SPARK

DataFrame column operations I

Operaciones condicionales

• Condicionales clauses →

`.when()` → Pone los valores cuando
`.otherwise()` → consecutivos.
`F.when(df.Age >= 18, 'Adult').otherwise('Minor')`

User defined functions

Ejemplo de como se define →

from pyspark.sql.functions import udf
df = df.withColumn('refin',

udf.func = udf(function, StringType())

Para usarlo

`df = df.withColumn('refin', udf.func(df.col))`

Partitioning and lazy processing

Para que los RDD en Spark vayan ordenados preparados para procesamiento es //

`pyspark.sql.functions.monotonically_increasing_id()`

Para ver el N° de Particiones que tiene un dataframe →

`df.rdd.getNumPartitions()`

Caching

Guardar los DataFrames en memoria o en disco

Reduce el uso de recursos

Mejora la velocidad.

Problemas → Si hay muchos datos igual no entra en memoria.

Para implementarlo se cachea antes de una accion →

`df = spark.read.csv(...)`
`df.cache().count()`

se pueden hacer tambien
sobre el.

Comprobar si esta cacheado → `df.isCached`

Operaciones

Para descacheado → `df.unpersist()`

CLEANING DATA WITH SPARK

data → DataFrame
Improving import performance

Parámetros importantes → N° de objetos (Archives, Network locations, etc)
→ Tamaño general de los objetos (Frecuencia mejor si son del mismo tamaño 128MB (optimo))

Frecuencias bien definidas mejor el funcionamiento

Para leer muchos archivos es más eficiente →

df = spark.read.csv('archivo_*.txt.gz')

Cluster configurations

Para leer la configuración del cluster

spark.conf.get(<conf name>)

" escribir "

spark.conf.set(<conf name>)

Opciones de deployment → Single node
→ Standalone
→ Managed (Yarn, Kubernetes...) → manejados por terceros

Responsabilidades del Driver del cluster → Asignación de tareas
→ Consolidación de resultados
→ Shared data access
→ Tips { - El driver node debe tener la memoria que valgan todos
- Fast local storage

" Worker del cluster → Ejecutar las tareas

Para leer algunas configuraciones de spark → Nombre → name = spark.conf.get('spark.app.name')
→ Puerto del driver (REP) → driver.port = spark.conf.get('spark.driver.port')
→ N° de particiones → partitions = spark.conf.get('spark.sql.shuffle.partitions')

Performance improvements

Para ver el plan que ha definido spark para realizar una tarea → df.explain()

User broadcasting (df.broadcast()) da una copia de todos los datos a cada worker, esto mejora el rendimiento de df.join()

Para usarlo → from pyspark.sql.functions import broadcast
df = df.b.join(broadcast(df.2))

CLEANING DATA WITH SPARK

Data Pipelines

- Como es un data Pipeline
1. Input
 2. Transformaciones
 3. Output
 4. Validation
 5. Analysis

Los Pipelines no estan definidos como un objeto en Spark, simplemente es el proceso

Data handling techniques

- Algunos posibles errores en los datos de entrada son
- Datos incorrectos
 - Nulos
 - Líneas corrotas
 - Headers
 - Estructuras unidades → Diferentes separadores
 - Datos no regulares

Parseador de CSV de Spark

- Puede remover líneas vacías
- " " comillas → `df=spark.read.csv('...', comment='#')`
- Usar los headers → `df=spark.read.csv('...', header=True)`
- Definimos el separador → `df=spark.read.csv('...', sep=',')`

→ Si el separador no es un String, `sep='#'` y guarda los datos en una columna `_c0`. Así podemos leer estructuras anidadas.

Data validation (comprobar que tiene un formato)

Validación usando join

- Mucho más rápido que recorrer todos los datos
- Ejemplo → `df1 = ...`
`df2 = ...`
`validado = df1.join(df2, df1.att == df2.att)`

El tercer argumento es el tipo de join (inner, outer...)

Analysis final y delivery

Usar calculos para dar sentido como medias, medianas etc.