

Introduction to PySpark

Para crear una conexión → Creamos una instancia de la clase `SparkContext`, requiere
unos argumentos que se pueden meter en un objeto `SparkConf()`
`SparkSession` → Interfaz con la conexión.

Imaginemos que el contexto se llama `SC`

Para ver la versión que está corriendo en el cluster → `SC.version`

La estructura básica de Spark es RDD (Resilient Distributed Dataset), usamos

`Spark DataFrame`

Optimizan directamente las operaciones.

Para crear una `SparkSession` →

```
from pyspark.sql import SparkSession
session = SparkSession.builder.getOrCreate()
```

Para ver las tablas creadas en el cluster → `spark.catalog.listTables()`

Se pueden realizar consultas SQL a estos tablos de la siguiente forma →

```
query = "SELECT ..."
(resultado = session.sql(query))
(resultado.show())
```

Nombre de la `SparkSession` creada

Este resultado se puede convertir a un dataframe de pandas →

```
df = resultado.toPandas()
```

```
spark_df.show()
```

Para crear un `Spark DataFrame` a partir de un `pd.DataFrame` → `spark_df = session.createDataFrame(pd.df)`

Para crear una tabla temporal en Spark a partir del dataframe → `spark_df.createOrReplaceTempView(nombre='...')`

Se pueden leer directamente archivos a `Spark DataFrame` → `spark_df = session.read.csv(ruta, header=True)`

Se le puede
fijar como nombre
de columna

Para crear una nueva columna → `spark_df = spark_df.withColumn('nueva columna', spark_df.veje_col * 5)`

El método `filter` funciona igual que un `WHERE` de SQL → `spark_df.filter(condicion)`

.. .. select SELECT → `spark_df.select('col.1', 'col.2', ...)`

Para hacer operaciones sobre columnas en `SELECT` y guardarlas como alias → `spark_df.select(spark_df.othme/60, alias='duracion')`

```
spark_df.selectExpr('ave.time/60 as duracion')
```

Para aplicar funciones de agregación como `min` o `max`, usamos el `groupBy` → `spark_df.groupBy().min(column)`

Después de `groupBy` podemos meter alias
`groupBy('clave').min(column)`

Para unir dataframes usamos el método `.join()` → `spark_df = spark_df_1.join(spark_df_2, on='key column', how='left')`

Introduction to PySpark

Para modificar el tipo de dato de una columna → `spark_df = spark_df.withColumn('nombre', spark_df.col.cast('integer'))`
El documento se llama double

Para convertir una columna categórica en un OneHotEncoder → `Paso 1 = indexer = StringIndexer(inputCol='nombre_col', outputCol='nombre_col')`
`Paso 2 = encoder = OneHotEncoder(" " " " " ")`

Para combinar varias columnas en 1 → `assembler = VectorAssembler(inputCols=['col.1', ..., 'col.3'], outputCol='nombre_col')`

Una vez tenemos la index, encoder y assembler los tenemos en 1 Pipeline →

```
from pyspark.ml import Pipeline
pipe = Pipeline(stages=[indexer, encoder, assembler])
```

Hay que realizar estas transformaciones antes de hacer el split train-test

Para aplicar estas transformaciones → `data_transform = pipe.fit(spark_df).transform(spark_df)`

Para realizar el split sobre los datos transformados → `train, test = data_transform.randomSplit([0.6, 0.4])`
80% train 20% test

Para crear una regresión logística →

```
from pyspark.ml.classification import LogisticRegression
lr = LogisticRegression()
```

Para evaluar modelos, cuenta con un evaluador →

```
from pyspark.ml.evaluation import BinaryClassificationEvaluator()
evaluador = Binary... (metricName='areaUnderROC')
```

Para buscar los parámetros óptimos, primero hay que crear un grid →

```
import pyspark.ml.tuning as tune
grid = tune.ParamGridBuilder()
grid = grid.addGrid(Parametro, valores)
grid = grid.build()
```

Para entrenar con validación cruzada →

```
cv = tune.CrossValidator(estimator = lr,
                          estimatorParamMaps = grid,
                          evaluator = evaluador)
```

PARA ENTRENAR

```
model = cv.fit(data.train)
Mejor = model.bestModel
```

Para realizar predicciones con el modelo →

```
resultados = modelo.transform(data.test)
```

" aplicar el evaluador sobre los resultados →

```
evaluador.evaluate(resultados)
```