



Date:

Introduction to PySpark.

SparkContext class.

SparkConf() constructor.

Spark - cluster computing.

- computation over clusters.

- data processing & computation \rightarrow parallel.

- master node $\swarrow \searrow$ worker nodes

Spark's core data structure \rightarrow RDD - Resilient Distributed Dataset

- Spark Dataframe abstraction built on RDD
easily optimized.

Create SparkSession Object from Spark Context

interface w/ connection.

conn'g to cluster

from pyspark.sql import SparkSession.

my_spark = SparkSession.builder.getOrCreate()

.catalog \Rightarrow data in cluster

.listTables() \Rightarrow names of tables.

x = spark.sql(query)

x.show()

x.toPandas() \Rightarrow Spark to pandas.

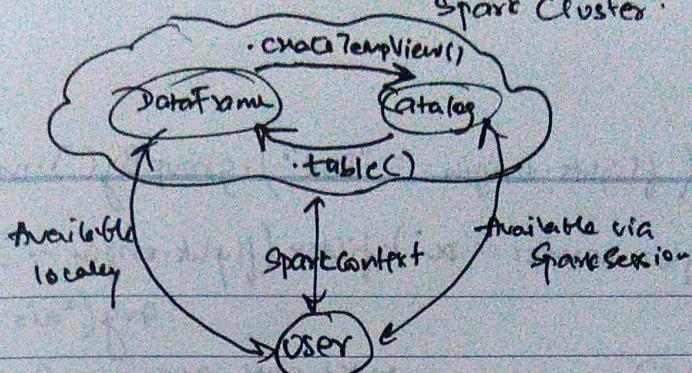
.createDataFrame() \Rightarrow pandas to spark. \rightarrow store locally,

.createTempView()

.createOrReplaceTempView()

Date: _____

Spark Cluster



pd = pd.DataFrame(np.random.randint(10))

spark_tmp = spark.createDataFrame(pd)

spark_tmp.createOrReplaceTempView("tmp")

file_path = " \u2191 \u2191 .csv"

airports = spark.read.csv(file_path, header=True)

airports.show()

Column Operations

• withColumn()

col-name [col-itself]

extract \Rightarrow df.colName

update \Rightarrow df = df.withColumn("newCol", df.oldCol + 1)

filtering

• filter()

flights.filter("air-time > 120").show()

Select -

• select()

df.colName • ("colName")

• alias \rightarrow rename column

• selectExpr

Grouped Data

• min
• max
• count

• groupBy()

df.groupBy().min("col").show()



Date :

```
flights.filter(flights.origin == "PDX").groupBy().min("distance").show()
flights.filter(flights.carrier == "DL").filter(flights.origin == "SEA").groupBy()
    .agg("air-time").show()
```

$\cdot \text{agg}() \rightarrow \text{column} \rightarrow \text{uses any agg-fun}$.

$\cdot \text{join}()$

$\begin{matrix} \text{2nd DF w/ 1st DF} \\ \text{on key how} \end{matrix}$

```
flights.join(airports, on = "dest", how = "leftouter")
```

Machine Learning Pipelines

PySpark.ml \rightarrow Transformer \rightarrow .transform() \rightarrow df.

Estimator \rightarrow .fit() \rightarrow model obj

Planes.withColumnRenamed(...)

Spark - Handles only numeric datatypes.

.cast() \leftarrow withColumn()

```
model_data = model_data.withColumn("arr-delay", model_data.arr-delay
    .cast("integer"))
```

PySpark.ml.features - one-hot vectors \Rightarrow categorical feature

① StringIndexer \rightarrow Estimator \rightarrow Transformer
str to obj $\quad \quad \quad$ df \rightarrow df.

② OneHotEncoder

StringIndexer(inputCol = "carrier", outputCol = "carrier_index")
OneHotEncoder(" " " carrier_fact")

Last step: combine all columns into one column.

VectorAssembler

Date: _____



VectorAssembler (inputCol: , outputCol:)

Pipeline → Combine all Estimators + Transformers.

Pipeline (stages = [- - - - .])

Piped data: flights Pipe . fit (model_data) . transform (model_data)

Logistic Regression

Hyperparameters → performance.

pyspark.ml.classification → LogisticRegression()

k-fold Cross Validation.

elastic NetParam regParam.

Evaluation

pyspark.ml.evaluation.

BinaryClassificationEvaluator (metricName = "areaUnderROC")

Grid

Look for optimal hyperparam.

pyspark.ml.tuning

ParamGridBuilder

• addGrid()

• build()

cv = tune.CrossValidator (estimator = lr,
EstimatorParamMaps = grid,
evaluator = evaluator)

grid = tune.ParamGridBuilder
grid = grid.addGrid (lr. regParam,

mp. array (0, -1, 0.01))

best_lr = lr.fit (training)

grid = grid.build()

Test

test_res = best_lr.transform (test)

print (evaluator.evaluate (test_res))