



Yotabites

YOTABITES

Big Data Summit KC 2017

SPARK WORKSHOP

Yotabites Consulting

bigdata@yotabites.com

Flat File Processing

1.	Getting Started with Apache Spark
1.1	Creating SparkContext <code>sc</code>
1.2	Import necessary Packages <pre>import sys import os import seaborn as sns import matplotlib.pyplot as plt import pandas as pd % matplotlib inline plt.rcParams["figure.figsize"] = [16,9] from pyspark import SparkContext from pyspark.streaming import StreamingContext import warnings warnings.filterwarnings("ignore") from operator import add import random from __future__ import print_function from IPython.display import Image,HTML,display from IPython.core.display import * display(HTML("<style>.container { width:75% !important; }</style>"))</pre>
2.	Creation of rdd and Loading a file from HDFS or S3

2.1	Create a RDD <code>rdd = sc.textFile("/tmp/meetup/meetup.csv")</code> <code>rdd.count()</code>
2.2	Another way to create an rdd <code>rdd1 = sc.parallelize([1,2,3,4,5])</code>
2.3	Display all the elements in the rdd <code>ondriver = rdd.collect()</code>
2.4	Display desired number of elements <code>prntSTR=rdd.take(5)</code> <code>print(pd.DataFrame(prntSTR))</code>
3.	Transformations on the RDD
3.1	Extract all the headers <code>pairRDD=rdd.map(lambda x: (x,x.split(",")[19]))</code> <code>pairRDD.first()</code>
3.2	Filter out all the rows whose state is kansas <code>filterRDD=pairRDD.filter(lambda pair: pair[1]=="KS")</code> <code>filterRDD.count()</code>
3.3	A flatMap on rdd <code>flatMapRDD=rdd.map(lambda x: x.split(",")).flatMap(lambda x:x)</code> <code>flatMapRDD.first()</code>

3.4	<p>Word count on the rdd</p> <pre>wordcountrdd = flatMapRDD\ .map(Lambda x: (x,1))\ .reduceByKey(Lambda x,y: x+y) wordcountrdd.take(5)</pre> <p>Note: To display a limited output use take(n), n is the number of desired outputs</p>
3.5	<p>RDD Lineage</p> <pre>wordcountrdd.toDebugString()</pre>
3.6	<p>Sort the output by Value</p> <pre>wordcountrdd.sortBy(Lambda value: value[1], ascending=False).take(5)</pre>
3.7	<p>Estimating Pi using the Monte Carlo Method</p> <pre>import random def inside(p): x, y = random.random(), random.random() return x*x + y*y < 1 NUM_SAMPLES=1000000 count = sc.parallelize(range(0, NUM_SAMPLES)) \ .filter(inside).count() print("Pi is roughly %f" % (4.0 * count / NUM_SAMPLES))</pre>
4.	<p>Let's perform some Actions on the rdd</p>
4.1	<p>CountbyKey action on rdd</p> <pre>wordcountrdd.map(Lambda line: (line,1)).countByKey()</pre>
4.2	<p>Return the first element in this RDD</p> <pre>wordcountrdd.first()</pre>

4.3	TakeOrdered(n) - Display number of desired elements in sequential order <code>wordcountrdd.takeOrdered(5)</code>
4.3	Output a Python RDD of key-value pairs (of form RDD[(K, V)]) to any Hadoop file system <code>wordcountrdd.saveAsSequenceFile("/tmp/bdkc_pyspark/bdkcSampleSequenceOutput")</code>
4.4	Save this RDD as a text file, using string representations of elements <code>wordcountrdd.saveAsTextFile("/tmp/bdkc_pyspark/bdkcSampleTextOutput")</code>
5.	DATAFRAMES A DataFrame is a Dataset organized into named columns. It is conceptually equivalent to a table in a relational database or a dataframe in R/Python, but with richer optimizations under the hood.
5.1	Load a CSV file and returns the result as a DataFrame <code>df = spark.read.csv("path")</code> Load a JSON file and returns the result as a DataFrame <code>jsondf=spark.read.format('json').Load('/tmp/bdkc_pyspark/satori_data/')</code>
5.2	Print the schema in a tree format <code>df.printSchema()</code>
5.3	Displays the content of the DataFrame <code>df.show()</code>
5.4	Filter and select operations on Dataframe <code>df.select(["group.group_state", "event.event_name"])\ .filter(df.group.group_state == "KS")\ .show()</code>

5.5	<p>Number of meetups in each city</p> <pre>df.groupby("group.group_city").count().orderBy("count",ascending=False).show()</pre>
5.6	<p>Number of members whose response was yes in each state</p> <pre>df.filter(df.response == 'yes')\ .groupby(["group.group_state", 'response'])\ .count()\ .show()</pre>
5.7	<p>Most Number of responses by each state</p> <pre>dataframe_df.filter(dataframe_df.group.group_state != 'null')\ .groupby(dataframe_df['group.group_state'],dataframe_df['response'])\ .count().orderBy("count",ascending=False)\ .show()</pre>
5.8	<p>Use Sql to find query the dataframe</p> <pre>df.createOrReplaceTempView("df") spark.sql("select event_time, count(*) as count \ from df \ group by event_time \ order by count desc")\ .show()</pre>

RDBMS

1.	Loading data from a JDBC source.
1.1	<p>Read data from Mysql.</p> <pre>jdbc_tblEmployees = spark.read \ .format("jdbc") \ .option("url", "jdbc:mysql://172.31.89.20:3306/employee") \ .option("dbtable", "tblEmployees") \ .option("user", "spark") \ .option("password", "spark") \ .option("driver", "com.mysql.jdbc.Driver") \ .load() jdbc_db.printSchema()</pre>

1.2	<p>Read another table tblpayemployeeparamdetails from dbo database</p> <pre> jdbc_employeepay = spark.read \ .format("jdbc") \ .option("url", "jdbc:mysql://172.31.89.20:3306/employee") \ .option("dbtable", "tblpayemployeeparamdetails") \ .option("user", "spark") \ .option("password", "spark") \ .option("driver", "com.mysql.jdbc.Driver") \ .load() jdbc_user.printSchema() </pre>
2.	<p>Joining the above dataframes with inner and leftsemi joins</p>
2.1	<p>Join the above two tables</p> <pre> joined_df = jdbc_tblemployees.join(jdbc_employeepay, jdbc_tblemployees['EmployeeNumber'] == jdbc_employeepay['EmployeeNumber'], "inner").drop(jdbc_tblemployees['EmployeeNum ber']) joined_df.printSchema() </pre>
2.2	<p>leftsemi join the above two tables</p> <pre> outer_join = jdbc_employeepay.join(jdbc_tblemployees, jdbc_tblemployees['EmployeeNumber'] == jdbc_employeepay['EmployeeNumber'], 'leftsemi')\ .drop(jdbc_employeepay['EmployeeNumber']) outer_join.printSchema() </pre>
3.	<p>Querying on Joined Dataframes</p>

3.1	Which department is earning the most <code>joined_df.groupBy("DepartmentCode").agg({'Amount': "mean"}).orderBy("avg(Amount)", ascending=False).show()</code>
3.2	Find Average Salaries of Employees <code>joined_df.groupby("EmployeeNumber").agg({'Amount': "mean"}).orderBy("avg(Amount)", ascending=False).show()</code>
3.3	Find Average Salaries of Employees in SQL <code>joined_df.createOrReplaceTempView("table")</code> <code>spark.sql("select EmployeeNumber, avg(Amount) as avg from table group by\ EmployeeNumber,Amount order by avg desc ").show()</code>
3.4	Another SQL query <code>spark.sql("select LocationType, Sex, count(*) as count from table group by \ LocationType,Sex order by count desc").show()</code>

SPARK STREAMING

1.	Spark Streaming
1.1	<p>Spark StreamingContext</p> <pre>ssc = StreamingContext(sc, 1)</pre>
1.2	<p>WordCount Spark Streaming</p> <pre>lines = ssc.textFileStream("/tmp/bdkc_pyspark/test.txt") counts = lines.flatMap(Lambda line: line.split(" "))\ .map(Lambda x: (x, 1))\ .reduceByKey(Lambda a, b: a+b) counts.pprint() counts.saveAsTextFiles("/tmp/bdkc_pyspark/output.txt") ssc.start() ssc.awaitTermination(10) ssc.stop()</pre>