

Pandas to PySpark Basic Command Conversion Cheatsheet

Action	Pandas <i>Assuming import of numpy as np</i>	PySpark <i>Assuming import of pyspark.sql.functions as fn</i>
Load a CSV	<code>df = pd.read_csv('file.csv')</code>	<code>df = spark.read.options(header = True, inferSchema = True).csv('file.csv')</code>
view entire dataframe	<code>df</code>	<code>display(df)</code>
View the head of a dataframe	<code>df.head(5)</code>	<code>df.show(5)</code> or <code>display(df.head(5))</code>
Graphing	<code>df.hist()</code> <i>(uses matplotlib)</i>	<code>display(df)</code> # then use graph options of display
View column names	<code>df.columns</code>	<code>df.columns</code>
View column types	<code>df.dtypes</code>	<code>df.dtypes</code> (may also want <code>df.schema</code>)
Get dataframe shape	<code>df.shape</code>	<code>print((df.count(), len(df.columns)))</code>
Get summary stats	<code>df.describe()</code>	<code>df.describe().show()</code>
Change columns names	<code>df.columns = ['a', 'b', 'c']</code>	<code>df.toDF('a', 'b', 'c')</code>
Rename a column	<code>df.rename(columns = {'old', 'new'})</code>	<code>df.withColumnRenamed('old', 'new')</code>
Drop a column	<code>df.drop('col1', axis = 1)</code>	<code>df.drop('col1')</code>
Add a column	<code>df['col_3'] = df['col_1'] + ['df.col_2']</code>	<code>df.withColumn('col_3', (df.col_1 + df.col_2))</code>
Fill Nulls with zero	<code>df.fillna(0)</code>	<code>df.fillna(0)</code>
Log transform	<code>df['log_a'] = np.log(df.a)</code>	<code>df.withColumns('log_a', fn.log(df.a))</code>
Extract a subset of columns of a dataframe	<code>df[['col1', 'col3']]</code>	<code>df.select('col1', 'col3')</code>
Extract a subset of rows by condition	<code>df.loc[df['col_1'] == some_value]</code>	<code>df.filter(fn.col('col_1') == some_value)</code>
Concatenate two dataframes vertically	<code>df = pd.concat([df1, df2, df3], axis = 0)</code>	<code>df = df1.unionAll(df2)</code>
Aggregation	<code>df.groupby(['a', 'b']).agg({'a': 'mean', 'b' = 'min'})</code>	<code>df.groupBy(['a', 'b']).agg({'a': 'mean', 'b' = 'min'})</code>
Pivot data and sum (could also count etc)	<code>pivot_table(df, values='col_3', index=['col_1'], columns=['col_2'], aggfunc=np.sum)</code>	<code>df.groupBy('col_1').pivot('col_2').sum('col_3')</code>
Row conditional	<code>df['cond'] = df.apply(lambda x: 1 if x.a > 10 else 2 if x.b < 5 else 3, axis = 1)</code>	<code>df.withColumn('cond', fn.when(df.a > 10, 1).when(df.b < 5, 2).otherwise(3))</code>
Apply a function	<code>df['squared'] = df.a.apply(lambda x: x**2)</code>	<code>func = fn.udf(lambda x: x**2, IntegerType())</code> <code>df.withColumn('squared', func(df.a))</code>
Join dataframes (pandas join command joins on keys which PySpark doesn't have)	<code>left.merge(right, on='id')</code> <code>left.merge(right, left_on='a', right_on='b')</code>	<code>left.join(right, on='id')</code> <code>left.join(right, left.a == right.b)</code>
Write a CSV to file	<code>df.to_csv('filename.csv')</code>	<code>df.write.csv('filename.csv')</code>
Convert between pyspark and pandas	<code>sqlContext.createDataFrame(pandas_df)</code>	<code>df.toPandas()</code>