


# Python For Data Science Cheat Sheet

## PySpark - SQL Basics

### PySpark & Spark SQL

Spark SQL is Apache Spark's module for working with structured data.



#### Initializing SparkSession

A SparkSession can be used create DataFrame, register DataFrame as tables, execute SQL over tables, cache tables, and read parquet files.

```
>>> from pyspark.sql import SparkSession
>>> spark = SparkSession \
    .builder \
    .appName("Python Spark SQL basic example") \
    .config("spark.some.config.option", "some-value") \
    .getOrCreate()
```

#### Creating DataFrames

##### From RDDs

```
>>> from pyspark.sql.types import *
Infer Schema
>>> sc = spark.sparkContext
>>> lines = sc.textFile("people.txt")
>>> parts = lines.map(lambda l: l.split(", "))
>>> people = parts.map(lambda p: Row(name=p[0],age=int(p[1])))
>>> peopledf = spark.createDataFrame(people)
Specify Schema
>>> people = parts.map(lambda p: Row(name=p[0],
                                     age=int(p[1].strip())))
>>> schemaString = "name age"
>>> fields = [StructField(field_name, StringType(), True) for
field_name in schemaString.split()]
>>> schema = StructType(fields)
>>> spark.createDataFrame(people, schema).show()
+----+----+
|  name  | age |
+----+----+
|   Mine |  28 |
|  Filip |  29 |
|Jonathan|  30 |
+----+----+
```

##### From Spark Data Sources

###### JSON

```
>>> df = spark.read.json("customer.json")
>>> df.show()
+-----+-----+-----+-----+-----+
| address | age | firstName | lastName | phoneNumbers |
+-----+-----+-----+-----+-----+
|[New York,10021,N.Y. ...]| 25 | John | Smith |[212 555-1234,ho ...]|
|[New York,10021,N.Y. ...]| 21 | Jane | Doe |[322 888-1234,ho ...]|
+-----+-----+-----+-----+-----+
```

>>> df2 = spark.read.load("people.json", format="json")

###### Parquet files

```
>>> df3 = spark.read.load("users.parquet")
```

###### TXT files

```
>>> df4 = spark.read.text("people.txt")
```

#### Inspect Data

|  |  |
|--|--|
| <pre>&gt;&gt;&gt; df.dtypes &gt;&gt;&gt; df.show() &gt;&gt;&gt; df.head() &gt;&gt;&gt; df.first() &gt;&gt;&gt; df.take(2) &gt;&gt;&gt; df.schema</pre> | <p>Return df column names and data types</p> <p>Display the content of df</p> <p>Return first n rows</p> <p>Return first row</p> <p>Return the first n rows</p> <p>Return the schema of df</p> |
|--|--|

#### Duplicate Values

```
>>> df = df.dropDuplicates()
```

#### Queries

|  |   |
|--|---|
| <pre>&gt;&gt;&gt; from pyspark.sql import functions as F <b>Select</b> &gt;&gt;&gt; df.select("firstName").show() &gt;&gt;&gt; df.select("firstName","lastName") \     .show() &gt;&gt;&gt; df.select("firstName",             "age",             explode("phoneNumber") \             .alias("contactInfo")) \     .select("contactInfo.type",             "firstName",             "age") \     .show() &gt;&gt;&gt; df.select(df["firstName"],df["age"]+ 1) \     .show() &gt;&gt;&gt; df.select(df['age'] &gt; 24).show() <b>When</b> &gt;&gt;&gt; df.select("firstName",             F.when(df.age &gt; 30, 1) \             .otherwise(0)) \     .show() &gt;&gt;&gt; df[df.firstName.isin("Jane", "Boris")] \     .collect() <b>Like</b> &gt;&gt;&gt; df.select("firstName",             df.lastName.like("Smith")) \     .show() <b>Startswith - Endswith</b> &gt;&gt;&gt; df.select("firstName",             df.lastName \             .startswith("Sm")) \     .show() &gt;&gt;&gt; df.select(df.lastName.endswith("th")) \     .show() <b>Substring</b> &gt;&gt;&gt; df.select(df.firstName.substr(1, 3) \             .alias("name")) \     .collect() <b>Between</b> &gt;&gt;&gt; df.select(df.age.between(22, 24)) \     .show()</pre> | <p>Show all entries in firstName column</p> <p>Show all entries in firstName, age and type</p> <p>Show all entries in firstName and age, add 1 to the entries of age</p> <p>Show all entries where age &gt;24</p> <p>Show firstName and 0 or 1 depending on age &gt;30</p> <p>Show firstName if in the given options</p> <p>Show firstName, and lastName is TRUE if lastName is like Smith</p> <p>Show firstName, and TRUE if lastName starts with Sm</p> <p>Show last names ending in th</p> <p>Return substrings of firstName</p> <p>Show age: values are TRUE if between 22 and 24</p> |
|--|---|

#### Add, Update & Remove Columns

##### Adding Columns

```
>>> df = df.withColumn('city',df.address.city) \
    .withColumn('postalCode',df.address.postalCode) \
    .withColumn('state',df.address.state) \
    .withColumn('streetAddress',df.address.streetAddress) \
    .withColumn('telePhoneNumber',
        explode(df.phoneNumber.number)) \
    .withColumn('telePhoneType',
        explode(df.phoneNumber.type))
```

##### Updating Columns

```
>>> df = df.withColumnRenamed('telePhoneNumber', 'phoneNumber')
```

##### Removing Columns

```
>>> df = df.drop("address", "phoneNumber")
>>> df = df.drop(df.address).drop(df.phoneNumber)
```

#### GroupBy

|  |  |
|--|--|
| <pre>&gt;&gt;&gt; df.groupBy("age") \     .count() \     .show()</pre> | <p>Group by age, count the members in the groups</p> |
|--|--|

#### Filter

|   |  |
|---|--|
| <pre>&gt;&gt;&gt; df.filter(df["age"]&gt;24).show()</pre> | <p>Filter entries of age, only keep those records of which the values are &gt;24</p> |
|---|--|

#### Sort

```
>>> peopledf.sort(peopledf.age.desc()).collect()
>>> df.sort("age", ascending=False).collect()
>>> df.orderBy(["age", "city"],ascending=[0,1]) \
    .collect()
```

#### Missing & Replacing Values

|  |  |
|--|--|
| <pre>&gt;&gt;&gt; df.na.fill(50).show() &gt;&gt;&gt; df.na.drop().show() &gt;&gt;&gt; df.na \     .replace(10, 20) \     .show()</pre> | <p>Replace null values</p> <p>Return new df omitting rows with null values</p> <p>Return new df replacing one value with another</p> |
|--|--|

#### Repartitioning

|  |   |
|--|---|
| <pre>&gt;&gt;&gt; df.repartition(10) \     .rdd \     .getNumPartitions() &gt;&gt;&gt; df.coalesce(1).rdd.getNumPartitions()</pre> | <p>df with 10 partitions</p> <p>df with 1 partition</p> |
|--|---|

#### Running SQL Queries Programmatically

##### Registering DataFrames as Views

```
>>> peopledf.createGlobalTempView("people")
>>> df.createTempView("customer")
>>> df.createOrReplaceTempView("customer")
```

##### Query Views

```
>>> df5 = spark.sql("SELECT * FROM customer").show()
>>> peopledf2 = spark.sql("SELECT * FROM global_temp.people") \
    .show()
```

#### Output

##### Data Structures

|   |   |
|---|---|
| <pre>&gt;&gt;&gt; rdd1 = df.rdd &gt;&gt;&gt; df.toJSON().first() &gt;&gt;&gt; df.toPandas()</pre> | <p>Convert df into an RDD</p> <p>Convert df into a RDD of string</p> <p>Return the contents of df as Pandas DataFrame</p> |
|---|---|

##### Write & Save to Files

```
>>> df.select("firstName", "city") \
    .write \
    .save("nameAndCity.parquet")
>>> df.select("firstName", "age") \
    .write \
    .save("namesAndAges.json",format="json")
```

#### Stopping SparkSession

```
>>> spark.stop()
```