# PySpark Crash Course

January 27, 2021

```python
[5]: # Load Pyspark Pkgs
     import pyspark
```

```python
[7]: # Methods/Attrib
     dir(pyspark)
```

```
[7]: ['Accumulator',
      'AccumulatorParam',
      'BarrierTaskContext',
      'BarrierTaskInfo',
      'BasicProfiler',
      'Broadcast',
      'HiveContext',
      'MarshalSerializer',
      'PickleSerializer',
      'Profiler',
      'RDD',
      'RDDBarrier',
      'Row',
      'SQLContext',
      'SparkConf',
      'SparkContext',
      'SparkFiles',
      'SparkJobInfo',
      'SparkStageInfo',
      'StatusTracker',
      'StorageLevel',
      'TaskContext',
      '_NoValue',
      '__all__',
      '__builtins__',
      '__cached__',
      '__doc__',
      '__file__',
      '__loader__',
      '__name__',
      '__package__',
```

```
'__path__',
'__spec__',
'__version__',
'_globals',
'accumulators',
'broadcast',
'cloudpickle',
'conf',
'context',
'copy_func',
'files',
'find_spark_home',
'heapq3',
'java_gateway',
'join',
'keyword_only',
'profiler',
'rdd',
'rddsampler',
'resource',
'resultiterable',
'serializers',
'shuffle',
'since',
'sql',
'statcounter',
'status',
'storagelevel',
'taskcontext',
'traceback_utils',
'types',
'util',
'version',
'wraps']
```

**Working with DataFrames in PySpark**

- Read DataSet(CSV)
- Create DataFrame

Tips

- SparkSession
- SparkContext :sc
- SqlContext

```python
[8]:   # Create A SparkSession
       from pyspark.sql import SparkSession
```

```
spark = SparkSession.builder.appName("PySparkTut").getOrCreate()
```

[9]: 
```
!ls
```

'PySpark Crash Course.ipynb'   work

[10]: 
```
!wget https://raw.githubusercontent.com/Jcharis/common_ml_datasets_explorer_app/
→master/datasets/diamonds.csv
```

--2021-01-27 21:01:43--  https://raw.githubusercontent.com/Jcharis/common_ml_dat
asets_explorer_app/master/datasets/diamonds.csv
Resolving raw.githubusercontent.com (raw.githubusercontent.com)…
151.101.112.133
Connecting to raw.githubusercontent.com
(raw.githubusercontent.com)|151.101.112.133|:443… connected.
HTTP request sent, awaiting response… 200 OK
Length: 2772143 (2.6M) [text/plain]
Saving to: 'diamonds.csv'

diamonds.csv         100%[===================>]   2.64M  3.81MB/s    in 0.7s

2021-01-27 21:01:44 (3.81 MB/s) - 'diamonds.csv' saved [2772143/2772143]


[11]: 
```
!ls
```

 diamonds.csv  'PySpark Crash Course.ipynb'   work

[12]: 
```
# Read A DataSet without header
df = spark.read.csv('diamonds.csv')
```

[13]: 
```
# Preview dataset
df.show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|  _c0|      _c1|  _c2|    _c3|  _c4|  _c5|  _c6| _c7| _c8| _c9|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
```

```
| 0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93| 3.9|2.46|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
|  0.3|     Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
|  0.3|     Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

[62]:
```python
# Read A DataSet with header/column names
df = spark.read.csv('diamonds.csv',header=True)
```

[15]:
```python
df.show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
|  0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93| 3.9|2.46|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
|  0.3|     Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
|  0.3|     Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
|  0.3|Very Good|    J|    SI1| 62.7|   59|  351|4.21|4.27|2.66|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

[16]:
```python
# Columns
df.columns
```

```
[16]: ['carat', 'cut', 'color', 'clarity', 'depth', 'table', 'price', 'x', 'y', 'z']
```

```
[20]: # Shape (rows + columns)
      (df.count() ,len(df.columns))
```

```
[20]: (53940, 10)
```

```
[21]: # Number of columns
      len(df.columns)
```

```
[21]: 10
```

```
[22]: # Number of rows
      df.count()
```

```
[22]: 53940
```

```
[24]: # Descriptive Analysis
      df.describe().show()
```

```
+-------+----------------+---------+-----+-------+----------------+--------
---------+----------------+----------------+----------------+-------------
----+
|summary|           carat|      cut|color|clarity|           depth|
table|           price|               x|               y|
z|
+-------+----------------+---------+-----+-------+----------------+--------
---------+----------------+----------------+----------------+-------------
----+
|  count|           53940|    53940|53940|  53940|           53940|
53940|           53940|           53940|           53940|
53940|
|   mean|0.7979397478679852|    null| null|   null| 61.74940489432624|
57.45718390804603|3932.799721913237| 5.731157211716609|
5.734525954764462|3.5387337782723316|
| stddev|0.4740112444054196|    null| null|   null|1.4326213188336525|2.2344905
628213247|3989.439738146397|1.1217607467924915|1.1421346741235616|0.705698846949
9883|
|    min|             0.2|     Fair|    D|     I1|              43|
43|            1000|               0|               0|               0|
|    max|            5.01|Very Good|    J|   VVS2|              79|
95|            9999|            9.86|            9.94|            8.06|
+-------+----------------+---------+-----+-------+----------------+--------
---------+----------------+----------------+----------------+-------------
----+
```

```
[25]: # Pick a column & Get summary/describe a selected column
      df.describe('carat').show()
```

```
+-------+------------------+
|summary|             carat|
+-------+------------------+
|  count|             53940|
|   mean|0.7979397478679852|
| stddev|0.4740112444054196|
|    min|               0.2|
|    max|              5.01|
+-------+------------------+
```

```
[26]: # Preview the First Row
      df.first()
```

```
[26]: Row(carat='0.23', cut='Ideal', color='E', clarity='SI2', depth='61.5',
      table='55', price='326', x='3.95', y='3.98', z='2.43')
```

```
[31]: # Preview the first 10 rows
      # Like a list
      df.head(10)
```

```
[31]: [Row(carat='0.23', cut='Ideal', color='E', clarity='SI2', depth='61.5',
      table='55', price='326', x='3.95', y='3.98', z='2.43'),
       Row(carat='0.21', cut='Premium', color='E', clarity='SI1', depth='59.8',
      table='61', price='326', x='3.89', y='3.84', z='2.31'),
       Row(carat='0.23', cut='Good', color='E', clarity='VS1', depth='56.9',
      table='65', price='327', x='4.05', y='4.07', z='2.31'),
       Row(carat='0.29', cut='Premium', color='I', clarity='VS2', depth='62.4',
      table='58', price='334', x='4.2', y='4.23', z='2.63'),
       Row(carat='0.31', cut='Good', color='J', clarity='SI2', depth='63.3',
      table='58', price='335', x='4.34', y='4.35', z='2.75'),
       Row(carat='0.24', cut='Very Good', color='J', clarity='VVS2', depth='62.8',
      table='57', price='336', x='3.94', y='3.96', z='2.48'),
       Row(carat='0.24', cut='Very Good', color='I', clarity='VVS1', depth='62.3',
      table='57', price='336', x='3.95', y='3.98', z='2.47'),
       Row(carat='0.26', cut='Very Good', color='H', clarity='SI1', depth='61.9',
      table='55', price='337', x='4.07', y='4.11', z='2.53'),
       Row(carat='0.22', cut='Fair', color='E', clarity='VS2', depth='65.1',
      table='61', price='337', x='3.87', y='3.78', z='2.49'),
       Row(carat='0.23', cut='Very Good', color='H', clarity='VS1', depth='59.4',
      table='61', price='338', x='4', y='4.05', z='2.39')]
```

```
[32]: # Method 2: Useful Action with show()
      # Show first 10 datapoints
```

```
df.show(10)
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 10 rows
```

[33]:
```
# Get Last Rows
df.tail(5)
```

[33]: [Row(carat='0.72', cut='Ideal', color='D', clarity='SI1', depth='60.8',
    table='57', price='2757', x='5.75', y='5.76', z='3.5'),
     Row(carat='0.72', cut='Good', color='D', clarity='SI1', depth='63.1',
    table='55', price='2757', x='5.69', y='5.75', z='3.61'),
     Row(carat='0.7', cut='Very Good', color='D', clarity='SI1', depth='62.8',
    table='60', price='2757', x='5.66', y='5.68', z='3.56'),
     Row(carat='0.86', cut='Premium', color='H', clarity='SI2', depth='61',
    table='58', price='2757', x='6.15', y='6.12', z='3.74'),
     Row(carat='0.75', cut='Ideal', color='D', clarity='SI2', depth='62.2',
    table='55', price='2757', x='5.83', y='5.87', z='3.64')]

### 0.0.1 Selection of columns

- .select ###### Note
- Dot & Bracket Notation only gives the column name not the entire column
    - ['colA']*
    - .colA*

[35]:
```
# List all Columns
df.columns
```

[35]: ['carat', 'cut', 'color', 'clarity', 'depth', 'table', 'price', 'x', 'y', 'z']

[37]:
```
# Select A Column
df.select('carat').show()
```

```
+-----+
|carat|
+-----+
| 0.23|
| 0.21|
| 0.23|
| 0.29|
| 0.31|
| 0.24|
| 0.24|
| 0.26|
| 0.22|
| 0.23|
|  0.3|
| 0.23|
| 0.22|
| 0.31|
|  0.2|
| 0.32|
|  0.3|
|  0.3|
|  0.3|
|  0.3|
+-----+
only showing top 20 rows
```

[40]:
```python
# Select A Column irrespective of column word case
# will work irrespective of the case of the column once it is found within the
 ↪dataset
df.select('CARAT').show()
```

```
+-----+
|CARAT|
+-----+
| 0.23|
| 0.21|
| 0.23|
| 0.29|
| 0.31|
| 0.24|
| 0.24|
| 0.26|
| 0.22|
| 0.23|
|  0.3|
| 0.23|
| 0.22|
```

```
| 0.31|
|  0.2|
| 0.32|
|  0.3|
|  0.3|
|  0.3|
|  0.3|
+-----+
only showing top 20 rows
```

[41]: 
```python
# This is not as we would expect in pandas
# For Bracket Notation : pick column name not the entire columne
df['carat']
```

[41]: Column<b'carat'>

[44]: 
```python
# This is not as we would expect in pandas
# For Dot Notation : pick column name not the entire column
df.carat
```

[44]: Column<b'carat'>

[45]: 
```python
# Select Multiple Columns
df.select('carat','cut').show(5)
```

```
+-----+-------+
|carat|    cut|
+-----+-------+
| 0.23|  Ideal|
| 0.21|Premium|
| 0.23|   Good|
| 0.29|Premium|
| 0.31|   Good|
+-----+-------+
only showing top 5 rows
```

## 0.0.2 Column Filtering and Applying Conditions

- .filter
- .where

[46]: 
```python
# Filter of Columns
# Apply A Condition
df.show(10)
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
```

9

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 10 rows
```

[47]: 
```python
# Method 1:using filter
df.filter(df['cut'] == "Good").show()
```

```
+-----+----+-----+-------+-----+-----+-----+----+----+----+
|carat| cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+----+-----+-------+-----+-----+-----+----+----+----+
| 0.23|Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.31|Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
|  0.3|Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
|  0.3|Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
|  0.3|Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
|  0.3|Good|    I|    SI2| 63.3|   56|  351|4.26| 4.3|2.71|
| 0.23|Good|    F|    VS1| 58.2|   59|  402|4.06|4.08|2.37|
| 0.23|Good|    E|    VS1| 64.1|   59|  402|3.83|3.85|2.46|
| 0.31|Good|    H|    SI1|   64|   54|  402|4.29|4.31|2.75|
| 0.26|Good|    D|    VS2| 65.2|   56|  403|3.99|4.02|2.61|
| 0.26|Good|    D|    VS1| 58.4|   63|  403|4.19|4.24|2.46|
| 0.32|Good|    H|    SI2| 63.1|   56|  403|4.34|4.37|2.75|
| 0.32|Good|    H|    SI2| 63.8|   56|  403|4.36|4.38|2.79|
|  0.3|Good|    I|    SI1| 63.2|   55|  405|4.25|4.29| 2.7|
|  0.3|Good|    H|    SI1| 63.7|   57|  554|4.28|4.26|2.72|
| 0.26|Good|    E|   VVS1| 57.9|   60|  554|4.22|4.25|2.45|
|  0.7|Good|    E|    VS2| 57.5|   58| 2759|5.85| 5.9|3.38|
|  0.7|Good|    F|    VS1| 59.4|   62| 2759|5.71|5.76| 3.4|
|  0.7|Good|    H|   VVS2| 62.1|   64| 2767|5.62|5.65| 3.5|
| 0.71|Good|    E|    VS2| 59.2|   61| 2772| 5.8|5.88|3.46|
+-----+----+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

[48]: 
```python
# Method 1:using filter
df.filter(df.carat >= 0.7).show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|  0.7|    Ideal|    E|    SI1| 62.5|   57| 2757| 5.7|5.72|3.57|
| 0.86|     Fair|    E|    SI2| 55.1|   69| 2757|6.45|6.33|3.52|
|  0.7|    Ideal|    G|    VS2| 61.6|   56| 2757| 5.7|5.67| 3.5|
| 0.71|Very Good|    E|    VS2| 62.4|   57| 2759|5.68|5.73|3.56|
| 0.78|Very Good|    G|    SI2| 63.8|   56| 2759|5.81|5.85|3.72|
|  0.7|     Good|    E|    VS2| 57.5|   58| 2759|5.85| 5.9|3.38|
|  0.7|     Good|    F|    VS1| 59.4|   62| 2759|5.71|5.76| 3.4|
| 0.96|     Fair|    F|    SI2| 66.3|   62| 2759|6.27|5.95|4.07|
| 0.73|Very Good|    E|    SI1| 61.6|   59| 2760|5.77|5.78|3.56|
|  0.8|  Premium|    H|    SI1| 61.5|   58| 2760|5.97|5.93|3.66|
| 0.75|Very Good|    D|    SI1| 63.2|   56| 2760| 5.8|5.75|3.65|
| 0.75|  Premium|    E|    SI1| 59.9|   54| 2760|   6|5.96|3.58|
| 0.74|    Ideal|    G|    SI1| 61.6|   55| 2760| 5.8|5.85|3.59|
| 0.75|  Premium|    G|    VS2| 61.7|   58| 2760|5.85|5.79|3.59|
|  0.8|    Ideal|    I|    VS1| 62.9|   56| 2760|5.94|5.87|3.72|
| 0.75|    Ideal|    G|    SI1| 62.2|   55| 2760|5.87| 5.8|3.63|
|  0.8|  Premium|    G|    SI1|   63|   59| 2760| 5.9|5.81|3.69|
| 0.74|    Ideal|    I|   VVS2| 62.3|   55| 2761|5.77|5.81|3.61|
| 0.81|    Ideal|    F|    SI2| 58.8|   57| 2761|6.14|6.11| 3.6|
|  0.8|    Ideal|    F|    SI2| 61.4|   57| 2761|5.96|   6|3.67|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

[49]:
```
# Method 2: where
df.where(df['cut'] == 'Good').show()
```

```
+-----+----+-----+-------+-----+-----+-----+----+----+----+
|carat| cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+----+-----+-------+-----+-----+-----+----+----+----+
| 0.23|Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.31|Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
|  0.3|Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
|  0.3|Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
|  0.3|Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
|  0.3|Good|    I|    SI2| 63.3|   56|  351|4.26| 4.3|2.71|
| 0.23|Good|    F|    VS1| 58.2|   59|  402|4.06|4.08|2.37|
| 0.23|Good|    E|    VS1| 64.1|   59|  402|3.83|3.85|2.46|
| 0.31|Good|    H|    SI1|   64|   54|  402|4.29|4.31|2.75|
| 0.26|Good|    D|    VS2| 65.2|   56|  403|3.99|4.02|2.61|
| 0.26|Good|    D|    VS1| 58.4|   63|  403|4.19|4.24|2.46|
| 0.32|Good|    H|    SI2| 63.1|   56|  403|4.34|4.37|2.75|
| 0.32|Good|    H|    SI2| 63.8|   56|  403|4.36|4.38|2.79|
|  0.3|Good|    I|    SI1| 63.2|   55|  405|4.25|4.29| 2.7|
|  0.3|Good|    H|    SI1| 63.7|   57|  554|4.28|4.26|2.72|
```

```
| 0.26|Good|    E|   VVS1| 57.9|    60|   554|4.22|4.25|2.45|
|  0.7|Good|    E|    VS2| 57.5|    58|  2759|5.85| 5.9|3.38|
|  0.7|Good|    F|    VS1| 59.4|    62|  2759|5.71|5.76| 3.4|
|  0.7|Good|    H|   VVS2| 62.1|    64|  2767|5.62|5.65| 3.5|
| 0.71|Good|    E|    VS2| 59.2|    61|  2772| 5.8|5.88|3.46|
+-----+----+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

[51]:
```
# Method 2: where
# select certain columns
df.where(df['cut'] == 'Good').select('price','cut').show()
```

```
+-----+----+
|price| cut|
+-----+----+
|  327|Good|
|  335|Good|
|  339|Good|
|  351|Good|
|  351|Good|
|  351|Good|
|  402|Good|
|  402|Good|
|  402|Good|
|  403|Good|
|  403|Good|
|  403|Good|
|  403|Good|
|  405|Good|
|  554|Good|
|  554|Good|
| 2759|Good|
| 2759|Good|
| 2767|Good|
| 2772|Good|
+-----+----+
only showing top 20 rows
```

[52]:
```
# Unique Values
# df['cut'].unique()
df.select("cut").distinct().show()
```

```
+---------+
|      cut|
+---------+
|  Premium|
```

```
|    Ideal|
|     Good|
|     Fair|
|Very Good|
+---------+
```

### 0.0.3 How to Add Columns & Delete/Drop Columns

- .withColumn()
- .drop()

```python
[54]:  # Add Columns
       df.withColumn("carat10x",df['carat'] * 10).show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+---------------
--+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
carat10x|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+---------------
--+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|
326|3.95|3.98|2.43|2.3000000000000003|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
2.1|
| 0.23|     Good|    E|    VS1| 56.9|   65|
327|4.05|4.07|2.31|2.3000000000000003|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
2.9|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
3.1|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
2.4|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
2.4|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
2.6|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
2.2|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|
4|4.05|2.39|2.3000000000000003|
|  0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
3.0|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93|
3.9|2.46|2.3000000000000003|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
2.2|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
```

13

```
3.1|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
2.0|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
3.2|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
3.0|
|  0.3|     Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
3.0|
|  0.3|     Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
3.0|
|  0.3|Very Good|    J|    SI1| 62.7|   59|  351|4.21|4.27|2.66|
3.0|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+---------------
--+
only showing top 20 rows
```

[55]: `df.show()`

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
|  0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93| 3.9|2.46|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
|  0.3|     Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
|  0.3|     Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
|  0.3|Very Good|    J|    SI1| 62.7|   59|  351|4.21|4.27|2.66|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

```
[63]: df2 = df.withColumn("carat10x",df['carat'] * 10)
```

```
[65]: df2.show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+---------------
--+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
carat10x|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+---------------
--+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|
326|3.95|3.98|2.43|2.3000000000000003|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
2.1|
| 0.23|     Good|    E|    VS1| 56.9|   65|
327|4.05|4.07|2.31|2.3000000000000003|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
2.9|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
3.1|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
2.4|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
2.4|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
2.6|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
2.2|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|
4|4.05|2.39|2.3000000000000003|
|  0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
3.0|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93|
3.9|2.46|2.3000000000000003|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
2.2|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
3.1|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
2.0|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
3.2|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
3.0|
|  0.3|     Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
3.0|
|  0.3|     Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
```

```
  3.0|
|  0.3|Very Good|    J|    SI1| 62.7|   59|  351|4.21|4.27|2.66|
  3.0|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+---------------
--+
only showing top 20 rows
```

[67]: 
```python
# Delete/Drop A Column
df2.drop('carat10x').show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
|  0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93| 3.9|2.46|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
|  0.3|     Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
|  0.3|     Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
|  0.3|Very Good|    J|    SI1| 62.7|   59|  351|4.21|4.27|2.66|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

### 0.0.4 GroupBy

- value counts
- aggregate

[69]: 
```python
# Value Counts
# df['cut'].value_counts()
# df.groupby('cut')[].size()
df.groupBy('cut').count().show()
```

```
+---------+-----+
|      cut|count|
+---------+-----+
|  Premium|13791|
|    Ideal|21551|
|     Good| 4906|
|     Fair| 1610|
|Very Good|12082|
+---------+-----+
```

[70]:
```python
# More Groupby
df.groupBy('price').mean().show()
```

```
+-----+
|price|
+-----+
| 2904|
| 3210|
| 3414|
| 3606|
| 3959|
| 4032|
| 4821|
| 4937|
| 5325|
| 5645|
| 5925|
| 6194|
| 6240|
| 6613|
| 6731|
| 7273|
| 7711|
| 7762|
| 9009|
| 9030|
+-----+
only showing top 20 rows
```

[71]:
```python
#  Sum of A groupby
df.groupBy('price').sum().show()
```

```
+-----+
|price|
+-----+
| 2904|
```

```
| 3210|
| 3414|
| 3606|
| 3959|
| 4032|
| 4821|
| 4937|
| 5325|
| 5645|
| 5925|
| 6194|
| 6240|
| 6613|
| 6731|
| 7273|
| 7711|
| 7762|
| 9009|
| 9030|
+-----+
only showing top 20 rows
```

```python
# # Aggregation
# df.groupBy('carat').agg('col':'sum')
```

[73]:
```python
df.show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
|  0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93| 3.9|2.46|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
```

```
| 0.3|     Good|    J|    SI1| 63.4|    54|  351|4.23|4.29| 2.7|
| 0.3|     Good|    J|    SI1| 63.8|    56|  351|4.23|4.26|2.71|
| 0.3|Very Good|    J|    SI1| 62.7|    59|  351|4.21|4.27|2.66|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

[74]: `df.columns`

[74]: `['carat', 'cut', 'color', 'clarity', 'depth', 'table', 'price', 'x', 'y', 'z']`

[75]:
```python
# Rearrange Columns
df.select('carat', 'color', 'clarity', 'depth', 'table', 'price', 'x', 'y',␣
 ↪'z','cut').show()
```

```
+-----+-----+-------+-----+-----+-----+----+----+----+---------+
|carat|color|clarity|depth|table|price|   x|   y|   z|      cut|
+-----+-----+-------+-----+-----+-----+----+----+----+---------+
| 0.23|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|    Ideal|
| 0.21|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|  Premium|
| 0.23|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|     Good|
| 0.29|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|  Premium|
| 0.31|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|     Good|
| 0.24|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|Very Good|
| 0.24|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|Very Good|
| 0.26|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|Very Good|
| 0.22|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|     Fair|
| 0.23|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|Very Good|
|  0.3|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|     Good|
| 0.23|    J|    VS1| 62.8|   56|  340|3.93| 3.9|2.46|    Ideal|
| 0.22|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|  Premium|
| 0.31|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|    Ideal|
|  0.2|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|  Premium|
| 0.32|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|  Premium|
|  0.3|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|    Ideal|
|  0.3|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|     Good|
|  0.3|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|     Good|
|  0.3|    J|    SI1| 62.7|   59|  351|4.21|4.27|2.66|Very Good|
+-----+-----+-------+-----+-----+-----+----+----+----+---------+
only showing top 20 rows
```

[76]:
```python
# Assign DF to a New DF
new_df = df.select('carat', 'color', 'clarity', 'depth', 'table', 'price', 'x',␣
 ↪'y', 'z','cut')
```

[77]: `new_df`

```
[77]: DataFrame[carat: string, color: string, clarity: string, depth: string, table:
      string, price: string, x: string, y: string, z: string, cut: string]
```

```
[78]: # Check Datatype
      new_df.dtypes
```

```
[78]: [('carat', 'string'),
       ('color', 'string'),
       ('clarity', 'string'),
       ('depth', 'string'),
       ('table', 'string'),
       ('price', 'string'),
       ('x', 'string'),
       ('y', 'string'),
       ('z', 'string'),
       ('cut', 'string')]
```

```
[80]: # Check For the Schema
      df.printSchema()
```

```
root
 |-- carat: string (nullable = true)
 |-- cut: string (nullable = true)
 |-- color: string (nullable = true)
 |-- clarity: string (nullable = true)
 |-- depth: string (nullable = true)
 |-- table: string (nullable = true)
 |-- price: string (nullable = true)
 |-- x: string (nullable = true)
 |-- y: string (nullable = true)
 |-- z: string (nullable = true)
```

```
[81]: # Check type of DF
      type(df)
```

```
[81]: pyspark.sql.dataframe.DataFrame
```

### 0.0.5 Saving DataFrames as CSV,parquet etc

```
[83]: # Save
      new_df.write.format('csv').option('header','true').save("diamond_clean.csv")
```

```
[84]: !ls
```

```
diamond_clean.csv    diamonds.csv  'PySpark Crash Course.ipynb'    work
```

```
[ ]: # Save  as parquet
     new_df.write.format('parquet').save("diamond_clean.parquet")
```

### 0.0.6  Making SQL Queries

- parse in the spark.SparkContext
- sqlContext

```
[85]: from pyspark.sql import SQLContext
```

```
[92]: dir(spark)
```

```
[92]: ['Builder',
       '__class__',
       '__delattr__',
       '__dict__',
       '__dir__',
       '__doc__',
       '__enter__',
       '__eq__',
       '__exit__',
       '__format__',
       '__ge__',
       '__getattribute__',
       '__gt__',
       '__hash__',
       '__init__',
       '__init_subclass__',
       '__le__',
       '__lt__',
       '__module__',
       '__ne__',
       '__new__',
       '__reduce__',
       '__reduce_ex__',
       '__repr__',
       '__setattr__',
       '__sizeof__',
       '__str__',
       '__subclasshook__',
       '__weakref__',
       '_activeSession',
       '_convert_from_pandas',
       '_createFromLocal',
       '_createFromRDD',
       '_create_dataframe',
       '_create_from_pandas_with_arrow',
```

```
'_create_shell_session',
'_get_numpy_record_dtype',
'_inferSchema',
'_inferSchemaFromList',
'_instantiatedSession',
'_jsc',
'_jsparkSession',
'_jvm',
'_jwrapped',
'_repr_html_',
'_sc',
'_wrapped',
'builder',
'catalog',
'conf',
'createDataFrame',
'getActiveSession',
'newSession',
'range',
'read',
'readStream',
'sparkContext',
'sql',
'stop',
'streams',
'table',
'udf',
'version']
```

[93]:
```python
# Create A Spark Context From the Spark Session
sc = spark.sparkContext
```

[94]:
```python
# Parse into the SQLContext
sqlContext = SQLContext(sc)
```

[96]:
```python
# Register Current DataFrame As Temporal Table
df.registerTempTable("DiamondsTable")
```

[97]:
```python
# Making QUeries
sqlContext.sql('SELECT * FROM DiamondsTable').show()
```

```
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
|carat|      cut|color|clarity|depth|table|price|   x|   y|   z|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
| 0.23|    Ideal|    E|    SI2| 61.5|   55|  326|3.95|3.98|2.43|
| 0.21|  Premium|    E|    SI1| 59.8|   61|  326|3.89|3.84|2.31|
| 0.23|     Good|    E|    VS1| 56.9|   65|  327|4.05|4.07|2.31|
```

```
| 0.29|  Premium|    I|    VS2| 62.4|   58|  334| 4.2|4.23|2.63|
| 0.31|     Good|    J|    SI2| 63.3|   58|  335|4.34|4.35|2.75|
| 0.24|Very Good|    J|   VVS2| 62.8|   57|  336|3.94|3.96|2.48|
| 0.24|Very Good|    I|   VVS1| 62.3|   57|  336|3.95|3.98|2.47|
| 0.26|Very Good|    H|    SI1| 61.9|   55|  337|4.07|4.11|2.53|
| 0.22|     Fair|    E|    VS2| 65.1|   61|  337|3.87|3.78|2.49|
| 0.23|Very Good|    H|    VS1| 59.4|   61|  338|   4|4.05|2.39|
|  0.3|     Good|    J|    SI1|   64|   55|  339|4.25|4.28|2.73|
| 0.23|    Ideal|    J|    VS1| 62.8|   56|  340|3.93| 3.9|2.46|
| 0.22|  Premium|    F|    SI1| 60.4|   61|  342|3.88|3.84|2.33|
| 0.31|    Ideal|    J|    SI2| 62.2|   54|  344|4.35|4.37|2.71|
|  0.2|  Premium|    E|    SI2| 60.2|   62|  345|3.79|3.75|2.27|
| 0.32|  Premium|    E|     I1| 60.9|   58|  345|4.38|4.42|2.68|
|  0.3|    Ideal|    I|    SI2|   62|   54|  348|4.31|4.34|2.68|
|  0.3|     Good|    J|    SI1| 63.4|   54|  351|4.23|4.29| 2.7|
|  0.3|     Good|    J|    SI1| 63.8|   56|  351|4.23|4.26|2.71|
|  0.3|Very Good|    J|    SI1| 62.7|   59|  351|4.21|4.27|2.66|
+-----+---------+-----+-------+-----+-----+-----+----+----+----+
only showing top 20 rows
```

[98]:
```python
# Can also use it to work with DataFrames
dir(sqlContext)
```

[98]:
```
['__class__',
 '__delattr__',
 '__dict__',
 '__dir__',
 '__doc__',
 '__eq__',
 '__format__',
 '__ge__',
 '__getattribute__',
 '__gt__',
 '__hash__',
 '__init__',
 '__init_subclass__',
 '__le__',
 '__lt__',
 '__module__',
 '__ne__',
 '__new__',
 '__reduce__',
 '__reduce_ex__',
 '__repr__',
 '__setattr__',
 '__sizeof__',
```

```
      '__str__',
      '__subclasshook__',
      '__weakref__',
      '_conf',
      '_inferSchema',
      '_instantiatedContext',
      '_jsc',
      '_jsqlContext',
      '_jvm',
      '_sc',
      '_ssql_ctx',
      'cacheTable',
      'clearCache',
      'createDataFrame',
      'createExternalTable',
      'dropTempTable',
      'getConf',
      'getOrCreate',
      'newSession',
      'range',
      'read',
      'readStream',
      'registerDataFrameAsTable',
      'registerFunction',
      'registerJavaFunction',
      'setConf',
      'sparkSession',
      'sql',
      'streams',
      'table',
      'tableNames',
      'tables',
      'udf',
      'uncacheTable']
```

[99]:
```python
# Thanks For Watching
# Jesus Saves @JCharisTech
# Jesse E.Agbe(JCharis)
# 2021
```

[ ]: