





10.173.92.150 - PuTTY

```
spark-3.5.0-bin-hadoop3/conf/metrics.properties.template
spark-3.5.0-bin-hadoop3/conf/fairScheduler.xml.template
spark-3.5.0-bin-hadoop3/conf/log4j2.properties.template
spark-3.5.0-bin-hadoop3/LICENSE
spark-3.5.0-bin-hadoop3/bin/
spark-3.5.0-bin-hadoop3/bin/spark-sql
spark-3.5.0-bin-hadoop3/bin/spark-submit
spark-3.5.0-bin-hadoop3/bin/pyspark2.cmd
spark-3.5.0-bin-hadoop3/bin/beeline
spark-3.5.0-bin-hadoop3/bin/pyspark
spark-3.5.0-bin-hadoop3/bin/pyspark.cmd
spark-3.5.0-bin-hadoop3/bin/load-spark-env.sh
spark-3.5.0-bin-hadoop3/bin/sparkR.cmd
spark-3.5.0-bin-hadoop3/bin/spark-shell2.cmd
spark-3.5.0-bin-hadoop3/bin/load-spark-env.cmd
spark-3.5.0-bin-hadoop3/bin/run-example
spark-3.5.0-bin-hadoop3/bin/sparkR2.cmd
spark-3.5.0-bin-hadoop3/bin/beeline.cmd
spark-3.5.0-bin-hadoop3/bin/docker-image-tool.sh
spark-3.5.0-bin-hadoop3/bin/spark-sql.cmd
spark-3.5.0-bin-hadoop3/bin/sparkR
spark-3.5.0-bin-hadoop3/bin/spark-submit.cmd
spark-3.5.0-bin-hadoop3/bin/find-spark-home.cmd
spark-3.5.0-bin-hadoop3/bin/run-example.cmd
spark-3.5.0-bin-hadoop3/bin/spark-connect-shell
spark-3.5.0-bin-hadoop3/bin/spark-sql2.cmd
spark-3.5.0-bin-hadoop3/bin/spark-shell.cmd
spark-3.5.0-bin-hadoop3/bin/spark-class2.cmd
spark-3.5.0-bin-hadoop3/bin/spark-class
spark-3.5.0-bin-hadoop3/bin/spark-class.cmd
spark-3.5.0-bin-hadoop3/bin/spark-submit2.cmd
spark-3.5.0-bin-hadoop3/bin/spark-shell
spark-3.5.0-bin-hadoop3/bin/find-spark-home
EB3001 >sudo mv spark-3.5.0-bin-hadoop3 /opt/spark
EB3001 >vi ~/.bashrc
```

```
[1]+ Stopped vi ~/.bashrc
```

```
EB3001 >vi ~/.bashrc
```

```
[2]+ Stopped vi ~/.bashrc
```

```
EB3001 >vi ~/.bashrc
```

```
EB3001 >vi ~/.bashrc
```

```
[3]+ Stopped vi ~/.bashrc
```

```
EB3001 >source ~/.bashrc
```

```
EB3001 >pyspark
```

```
Python 3.10.12 (main, Nov 20 2023, 15:14:05) [GCC 11.4.0] on linux
```

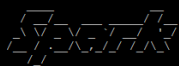
```
Type "help", "copyright", "credits" or "license" for more information.
```

```
Setting default log level to "WARN".
```

```
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
```

```
24/02/06 14:04:27 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
```

```
Welcome to
```



version 3.5.0

```
Using Python version 3.10.12 (main, Nov 20 2023 15:14:05)
```

```
Spark context Web UI available at http://assignment2:4040
```

```
Spark context available as 'sc' (master = local[*], app id = local-1707228269133).
```

```
SparkSession available as 'spark'.
```

```
>>>
```

## In this step-up sample data – Samsung models

```
10.173.92.150 - PuTTY
sample_data_df = spark.createDataFrame(
  sample_data,
  ['Id', 'Model', 'Year', 'ScreenSize', 'Storage', 'Weight']
)

# Show DataFrame
sample_data_df.show()
24/02/08 06:27:18 WARN SparkSession: Using an existing Spark session; only runtime SQL configurations will take effect.
>>>
>>> # Sample data
>>> sample_data = [
...   (1, 'Model_A', 2020, 15.6, 512, 3.5),
...   (2, 'Model_B', 2021, 13.3, 256, 2.8),
...   (3, 'Model_C', 2019, 14.0, 128, 3.0)
... ]
>>>
>>> # Create DataFrame
>>> sample_data_df = spark.createDataFrame(
...   sample_data,
...   ['Id', 'Model', 'Year', 'ScreenSize', 'Storage', 'Weight']
... )
>>>
>>> # Show DataFrame
>>> sample_data_df.show()
+-----+-----+-----+-----+-----+
| Id | Model | Year | ScreenSize | Storage | Weight |
+-----+-----+-----+-----+-----+
| 1 | Model_A | 2020 | 15.6 | 512 | 3.5 |
| 2 | Model_B | 2021 | 13.3 | 256 | 2.8 |
| 3 | Model_C | 2019 | 14.0 | 128 | 3.0 |
+-----+-----+-----+-----+-----+

>>> from pyspark.sql import SparkSession
>>>
>>> # Create a SparkSession
>>> spark = SparkSession.builder \
...   .appName("example_app") \
...   .getOrCreate()
>>>
>>> # Sample data
>>> sample_data = [
...   (1, 'Samsung_S21', 2021, '128GB', '5000mHA', 0.150),
...   (2, 'Samsung_S22', 2022, '256GB', '6000mHA', 0.208),
...   (3, 'Samsung_S23', 2023, '512GB', '7000mHA', 0.300)
... ]
>>>
>>> # Create DataFrame
>>> sample_data_df = spark.createDataFrame(
...   sample_data,
...   ['Id', 'Model', 'Year', 'Storage', 'Battery', 'Weight']
... )
>>>
>>> # Show DataFrame
>>> sample_data_df.show()
+-----+-----+-----+-----+-----+
| Id | Model | Year | Storage | Battery | Weight |
+-----+-----+-----+-----+-----+
| 1 | Samsung_S21 | 2021 | 128GB | 5000mHA | 0.15 |
| 2 | Samsung_S22 | 2022 | 256GB | 6000mHA | 0.208 |
| 3 | Samsung_S23 | 2023 | 512GB | 7000mHA | 0.3 |
+-----+-----+-----+-----+-----+

>>> █
```

## Creating DataFrames

```
10.173.92.150 - PuTTY
24/02/08 14:39:03 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Welcome to

      /_/_/  /_/_/  /_/_/  /_/_/  /_/_/
     /_/_/  /_/_/  /_/_/  /_/_/  /_/_/
    /_/_/  /_/_/  /_/_/  /_/_/  /_/_/
   /_/_/  /_/_/  /_/_/  /_/_/  /_/_/
  /_/_/  /_/_/  /_/_/  /_/_/  /_/_/
 /_/_/  /_/_/  /_/_/  /_/_/  /_/_/
/_/_/  /_/_/  /_/_/  /_/_/  /_/_/

version 3.5.0

Using Python version 3.10.12 (main, Nov 20 2023 15:14:05)
Spark context Web UI available at http://assignment2:4040
Spark context available as 'sc' (master = local[*], app id = local-1707403145266).
SparkSession available as 'spark'.
>>> sample_data.take(1)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'sample_data' is not defined
>>> sample_data.take(1)
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'sample_data' is not defined
>>> sample_data = sc.parallelize([
... (1, 'Samsung S21', 2020, '5.8"', '128GB', '3000mAh', 0.150)
... (2, 'Samsung S22', 2021, '6.2"', '256GB', '3200mAh', 0.126)
... (3, 'Samsung S23', 2022, '6.5"', '512GB', '3600mAh', 0.206)
... (4, 'Samsung S24', 2024, '7.1"', '1000GB', '5000mAh', 0.275)
... ])
>>> sample_data_df = spark.createDataFrame(
... sample_data
... [
... 'Id'
... 'Model'
... 'Year'
... 'ScreenSize'
... 'Storage'
... 'Baterly'
... 'Weight'
... ]
... )
>>> sample_data.take(1)
[(1, 'Samsung S21', 2020, '5.8"', '128GB', '3000mAh', 0.15)]
>>> sample_data_df.take(1)
[Row(Id=1, Model='Samsung S21', Year=2020, ScreenSize='5.8"', Storage='128GB', Baterly='3000mAh', Weight=0.15)]
>>> sample_data_df.show()
-----
| Id|      Model|Year|ScreenSize|Storage|  Baterly|Weight|
-----+-----+-----+-----+-----+-----+-----
|  1|Samsung S21|2020|   5.8"| 128GB|3000mAh| 0.15|
|  2|Samsung S22|2021|   6.2"| 256GB|3200mAh| 0.126|
|  3|Samsung S23|2022|   6.5"| 512GB|3600mAh| 0.206|
|  4|Samsung S24|2024|   7.1"|1000GB|5000mAh| 0.275|
-----

>>> sample_data_df.printSchema()
root
 |-- Id: long (nullable = true)
 |-- Model: string (nullable = true)
 |-- Year: long (nullable = true)
 |-- ScreenSize: string (nullable = true)
 |-- Storage: string (nullable = true)
 |-- Baterly: string (nullable = true)
 |-- Weight: double (nullable = true)
>>>
```

