

```
import sys
import os
import pandas as pd
from pyspark.context import SparkContext
from pyspark.sql.session import SparkSession
```

In [3]:

```
!apt-get install openjdk-8-jdk-headless -qq > /dev/null
!wget -q https://archive.apache.org/dist/spark/spark-3.0.0/spark-3.0.0-bin-hadoop3.2.tgz
!tar xf spark-3.0.0-bin-hadoop3.2.tgz
```

```
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
os.environ["SPARK_HOME"] = "/content/spark-3.0.0-bin-hadoop3.2"
```

```
!pip install findspark
```

In [17]:

```
import findspark
findspark.init("/content/spark-3.0.0-bin-hadoop3.2")
```

In [7]:

```
data = pd.read_csv('mobile.csv')
print(data)
```

	id	battery_power	blue	...	three_g	touch_screen	wifi
0	1	1043	1	...	0	1	0
1	2	841	1	...	1	0	0
2	3	1807	1	...	0	1	1
3	4	1546	0	...	1	1	0
4	5	1434	0	...	1	0	1
...
995	996	1700	1	...	1	1	0
996	997	609	0	...	0	1	1
997	998	1185	0	...	1	0	0
998	999	1533	1	...	0	1	0
999	1000	1270	1	...	1	0	1

Dualsim data in the format of list

```
data = list(data['dual_sim'])
print(data)
```

[1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 1,
1, 1, 1, 1, 0, 0, 1, 1, 0, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0,
1, 1, 0, 1, 1, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1,
1, 1, 1, 0, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0,
0, 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 0,

3	6	235	1671	3911	15	8		8	1		1	1
	13	900	1		1.4		1	0	0	30	1	87
2	3	829	1893	439	6	2		20	1		0	0
	14	1190	1		2.2		1	5	0	19	0.9	158
5	15	227	1856	992	13	0		16	1		1	0
	15	630	0		1.8		0	8	1	51	0.9	193
8	9	1315	1323	2751	17	6		3	1		1	0
	16	1846	1		1		0	5	1	53	0.7	106
8	7	185	1832	563	9	5		10	1		0	1
	17	1985	0		0.5		1	14	1	26	1	163
2	17	613	1511	2083	13	3		14	1		1	0
	18	1042	0		2.9		0	5	1	48	0.2	186
4	15	335	532	2187	9	2		5	1		0	0
	19	1231	1		1.7		1	2	1	37	0.2	194
2	3	82	1771	3902	19	12		15	1		0	1
	20	1488	0		2.6		0	9	0	37	0.7	189
4	20	47	559	2524	5	0		6	0		0	0
+--++	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+
--+--	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+	+-----+
only showing top 20 rows												

Reading the dual_sim column in the pyspark

```
sc = SparkContext.getOrCreate()
spark = SparkSession(sc)

mobileConfig = spark.read.csv(
    "mobile.csv",
    header = "true"
)
mobileConfig.select("dual_sim").show()
```

only showing top 20 rows

In [21]:

```
[ '1', '1', '0', '1', '0', '1', '0', '1', '1', '0', '0', '0', '1', '1', '0', '0', '1', '0'
  '1', '0', '0', '1', '1', '1', '1', '1', '1', '0', '1', '1', '1', '1', '1', '0', '0' ]
```

