# google-colab-setup-PySpark

November 14, 2025

## 1 Data Processing using Pyspark

```python
#configuración en google colab de spark y pyspark
from google.colab import drive
drive.mount('/content/gdrive')
```

```python
#instalar java y spark
!apt-get install openjdk-17-jdk-headless -qq > /dev/null
!wget -q https://downloads.apache.org/spark/spark-4.0.1/spark-4.0.1-bin-hadoop3.
 ↪tgz
!tar xf spark-4.0.1-bin-hadoop3.tgz
!pip install -q findspark
```

```python
import os
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-17-openjdk-amd64"
os.environ["SPARK_HOME"] = "/content/spark-4.0.1-bin-hadoop3"
```

```python
import findspark
findspark.init()
from pyspark.sql import SparkSession
spark = SparkSession.builder.master("local[*]").getOrCreate()
sc = spark.sparkContext
```

```python
# Load csv Dataset
df=spark.read.csv('sample_data/california_housing_test.
 ↪csv',inferSchema=True,header=True)
# desde S3
# df=spark.read.csv('s3a://bucke_name/datasets/sample_data.
 ↪csv',inferSchema=True,header=True)
```

```python
#columns of dataframe
df
```

```python
#check number of columns
len(df.columns)
```

```python
#number of records in dataframe
df.columns
```

```python
#shape of dataset
print((df.count(),len(df.columns)))
```

```python
#printSchema
df.printSchema()
```

```python
#fisrt few rows of dataframe
df.show(5)
```