# Using PySpark with Google Colab

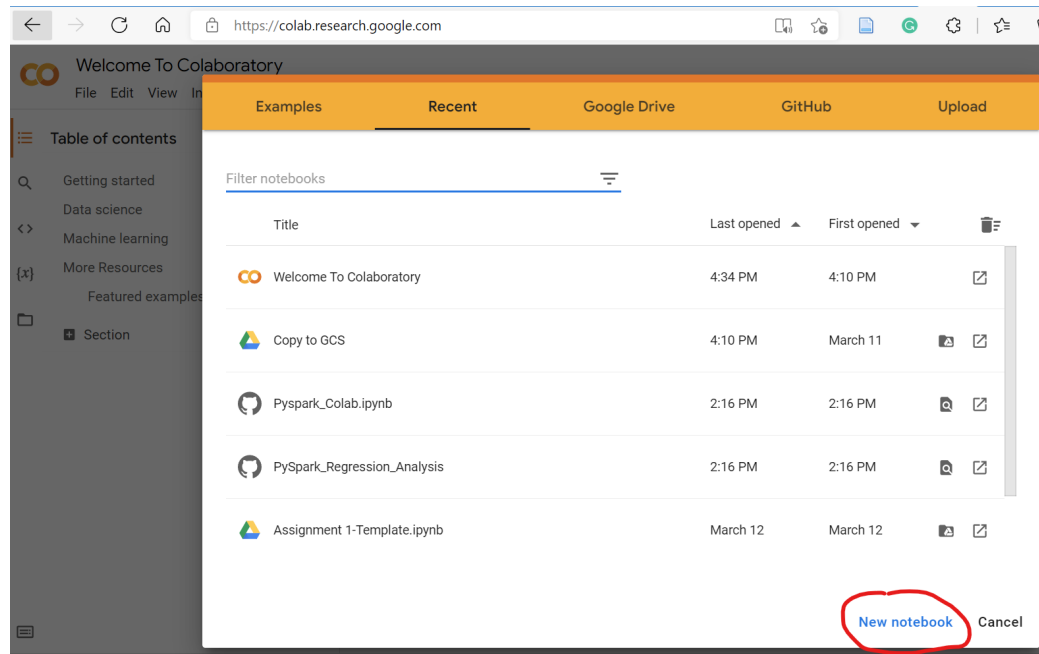Apache Spark (https://spark.apache.org/) is a multi-language engine for running data engineering, data science, and machine learning on a computer cluster or on a single machine.

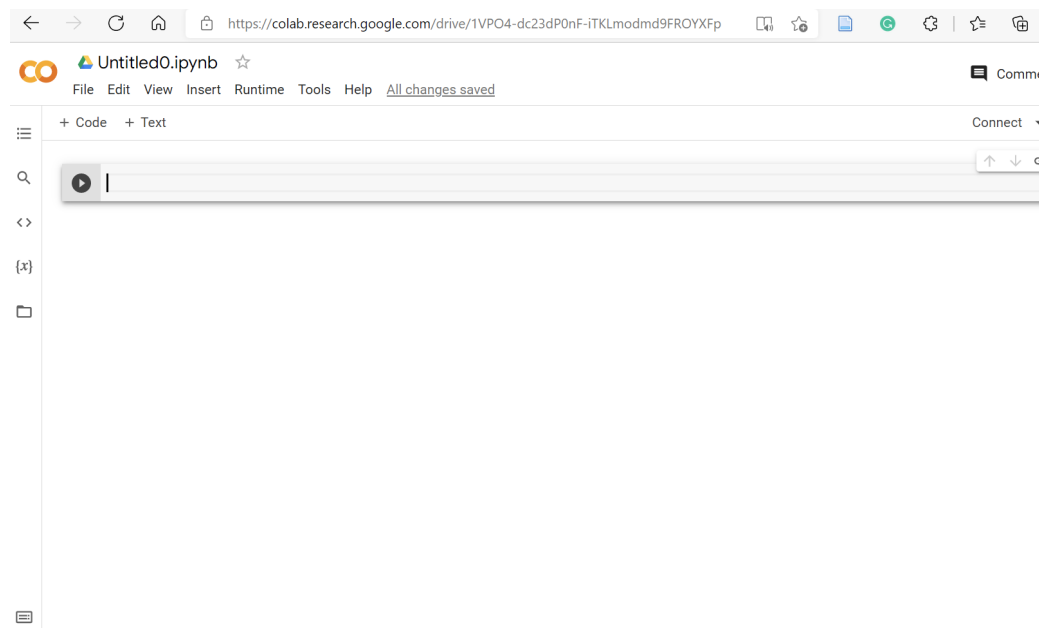PySpark (https://spark.apache.org/docs/latest/api/python/) is an interface for Apache Spark in Python. It allows writing Spark applications using Python APIs and provides the PySpark shell for interactively analyzing data in a distributed environment. PySpark supports most of Spark's features such as Spark SQL, DataFrame, Streaming, MLlib (Machine Learning), and Spark Core.

One of the easiest ways to use and try PySpark is through Google Colab. Colab is a product from Google Research and allows writing and executing Python code through the browser. Colab is a modified version of the Jupyter notebook hosted on Google servers and requires no setup while providing access free of charge to computing resources. It is particularly well suited to machine learning, data analysis, and education.

1. Go to https://colab.research.google.com/ and choose New notebook.



2. A new Jupiter notebook will be created, and you are ready to start writing your code.
3. You can click on the [+ Code] button to create new cells to enter your code.
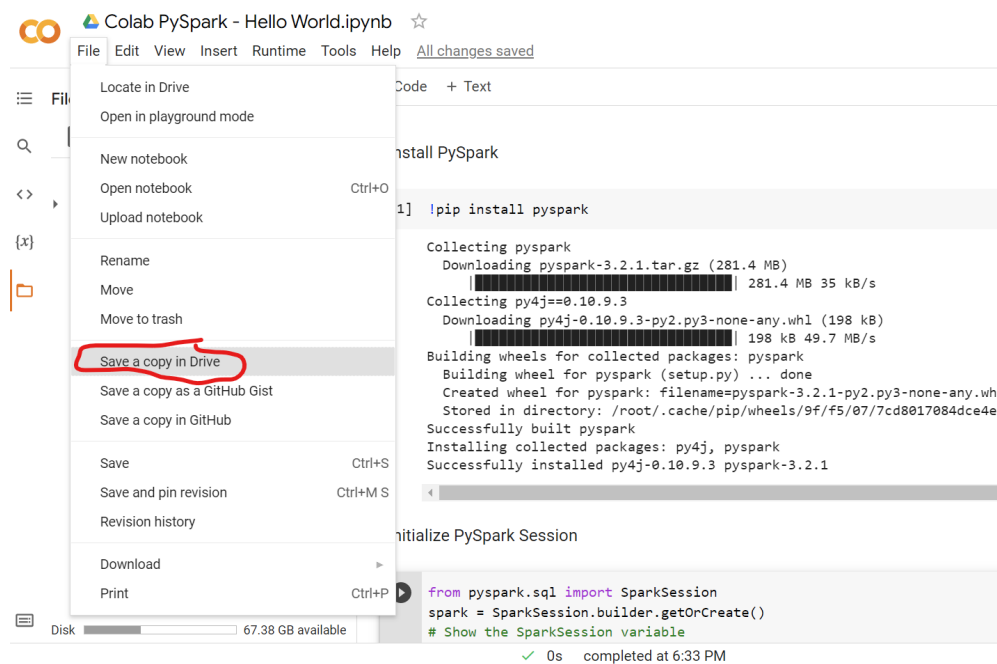
# Example notebook

You can use the following example notebook to start using the PySpark.
https://colab.research.google.com/drive/1VPO4-dc23dP0nF-iTKLmodmd9FROYXFp?usp=sharing

1. When you open the notebook, you can create your own copy using Save a copy to Drive.



2. Use the Run all command to run the notebook.