**DSCI-549: Introduction to
Computational Thinking and
Data Science**

# USC Viterbi School
# of Engineering

## Homework 2

### Communication
Please contact TA **Yashika Goyal** (yashikag@usc.edu) for this homework and include "**DSCI549**" in the subject line of your email.

### Assignment
The purpose of this homework is to get familiar with Jupyter Notebooks and Google Colab, and run simple data analysis. You must complete the assignment using the Jupyter Notebook. Answer the questions in a separate word document, along with relevant screenshots. Then submit it on D2L.

For this assignment, you will be using four notebooks. You can access them here:

https://github.com/doctorningwangusc/DSCI549/tree/main/Homework%20Assignments/Assignment2_CaesarCypher_SimpleMath%26Statistics

1. Getting data in and out of Colab: Notebook1_Importing_Data_to_Colab (25 points)
   a. Download the dataset from the internet using the wget function and take a screenshot of your Colab directory with the appropriate file. [Include two screenshots: One displaying the wget cell output & other of the Colab directory with the downloaded file present in it.]
   b. Follow the instructions to upload the Forest Fires dataset from your local machine to Colab and take appropriate screenshots of the process.

2. Running functions inside Jupyter: Notebook2_Running_Functions (25 points)
   a. Enter a text of your choosing when prompted. You may copy and paste any text. Take a screenshot of the resulting text with Line Breaks.
   b. Run the Caesar Cypher. Take a screenshot of the ciphered text.

3. Summarizing data: Notebook3_Summarizing _Data (25 points)
   a. When prompted, enter a list of 20 numbers. [Provide a screenshot]
   b. Take a screenshot of the summary statistics for your dataset. What do Q1 & Q3 represent?
   c. Generate a boxplot for your data and include the plot here. How do you interpret boxplots in the context of the summary statistics from part b?

4. Visualizing data: Notebook4_Visualiazing_Data (25 points)
   a. When prompted enter a list containing at least 20 scores for the "fictitious" class and 5 bins. Include the resulting histogram here.
   b. Rerun the cell with a varying number of bins. What's the optimal number of bins considering the size of your dataset? Why? Include the histogram here.

---

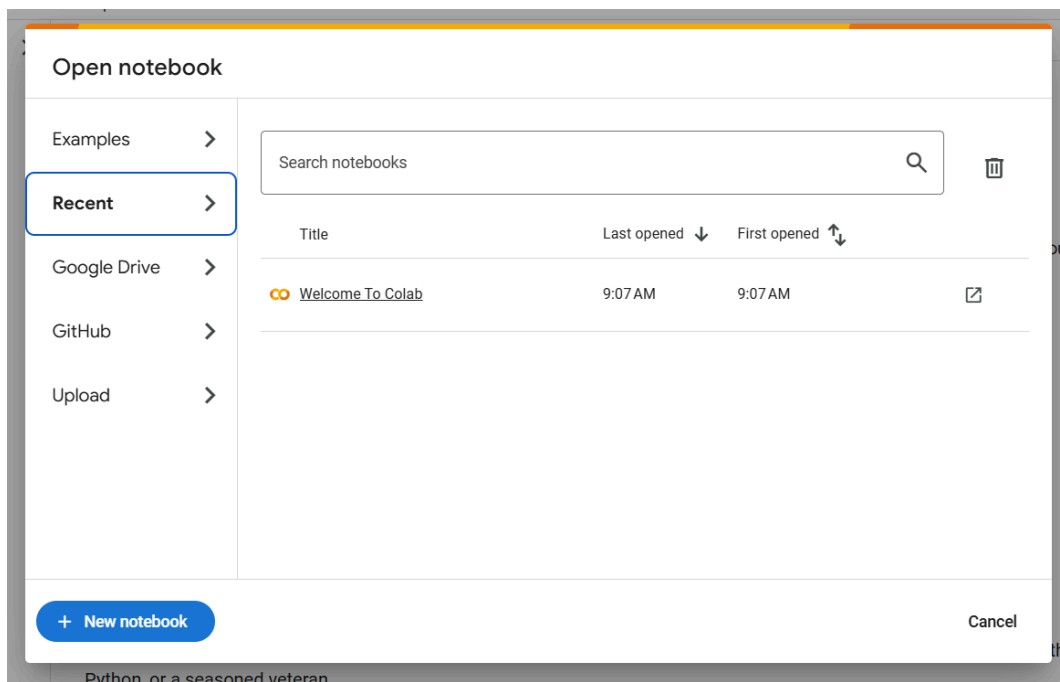# Getting Started with Google Colab

1.  To use Colab you should be logged in to your Google account. After this, click on the following link - https://colab.research.google.com

2.  You will see the following tab pop up when you open Colab.

    [Note: If this tab doesn't pop up, go to File → Open Notebook]
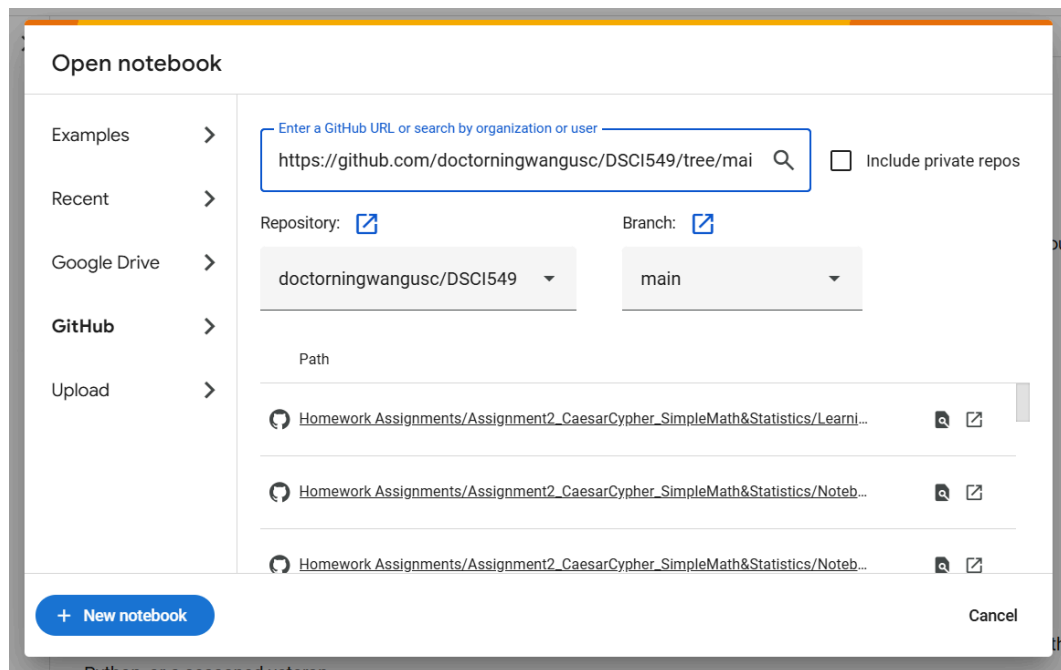


3.  You will be using Notebooks stored on a GitHub repository.
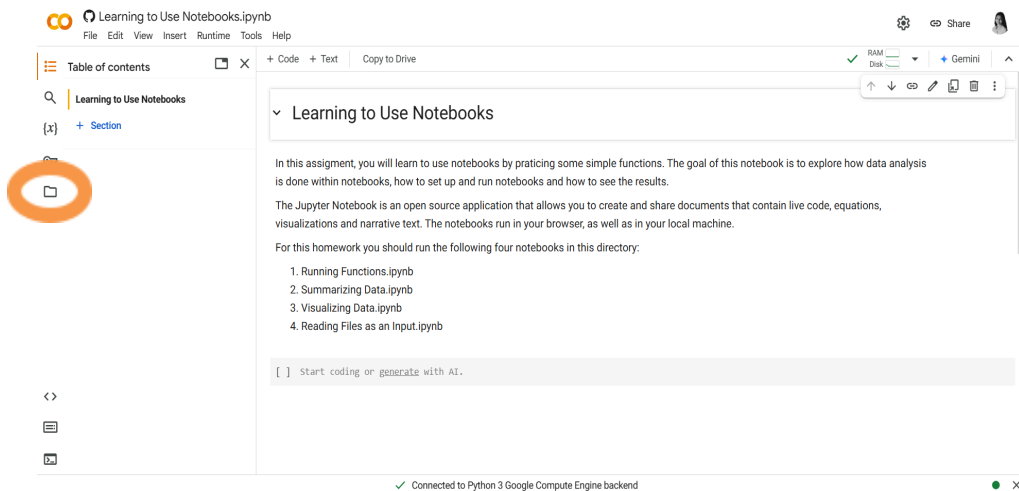    Click on the GitHub link.

4. Paste the following URL on the search bar:
   https://github.com/doctorningwangusc/DSCI549/tree/main/Homework%20Assignments/Assignment2_CaesarCypher_SimpleMath%26Statistics
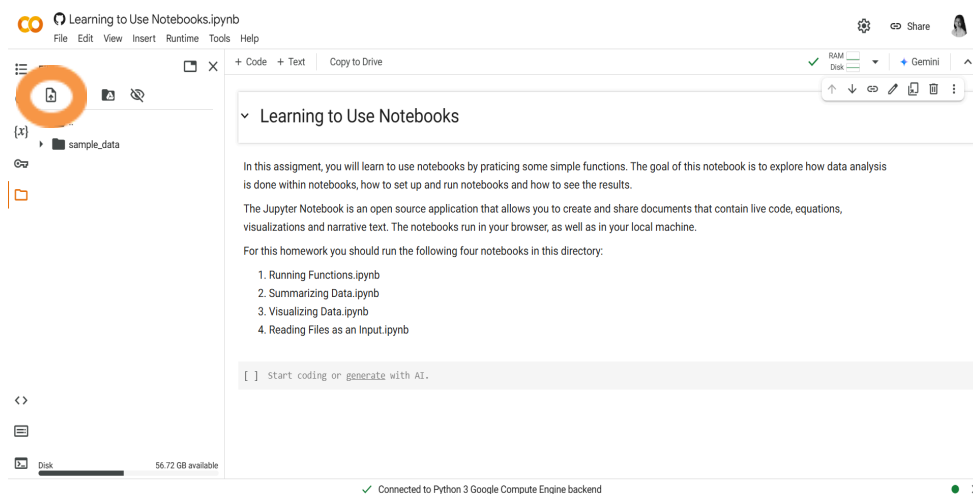   and click enter. This link should allow you access to all the notebooks on the repository. Scroll down and select the notebook you need to complete your homework. [You should be redirected to the Jupyter Notebook.]

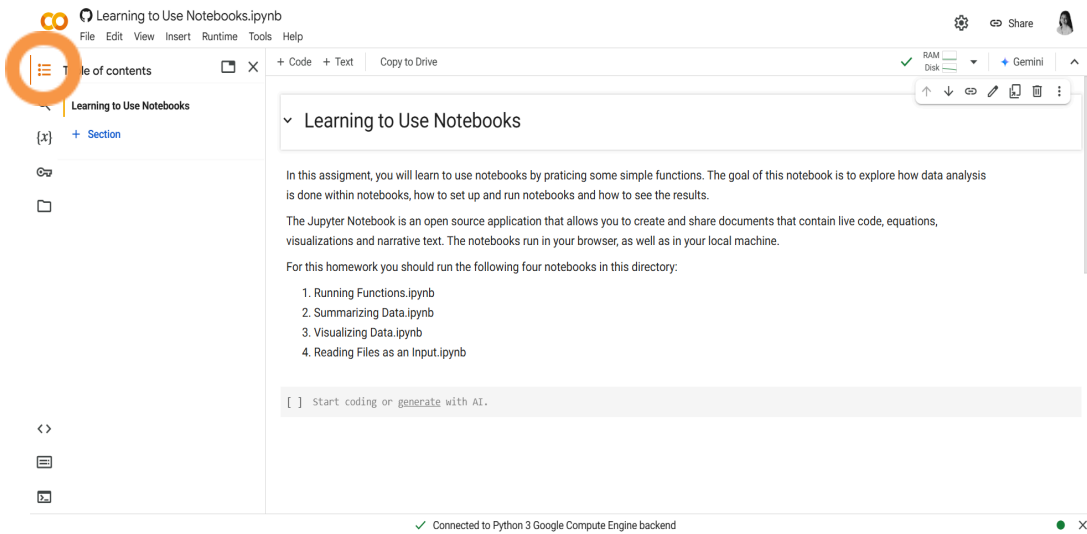5. Click on the folder icon and wait a few seconds for the runtime to connect.



6. The sample data folder contains datasets populated automatically by Google Colab. **DO NOT USE FOR YOUR HOMEWORK.**

   a. Some datasets are included as examples. To access them and download them into your Google Colab directory, run the corresponding cells in your notebook.
   b. If you want to upload a file from your local computer, click on 'Upload to session storage' in the menu bar.



**Note:** Uploaded files are deleted after 12 hours. You might need to upload it again to use the file after 12 hours.

7. You can click on the Contents button to view and access any part of the notebook.



8. Use the notebooks to answer the questions given in the homework assignment.