# Put Kaggle data into colab

To open a Kaggle dataset in Google Colab, follow these step-by-step instructions. This process involves setting up Kaggle API credentials, configuring your Colab environment, and downloading the desired dataset.

## 1. Set Up Kaggle API Credentials

**a. Create a Kaggle Account:** If you haven't already, sign up for a <u>Kaggle</u> account.

**b. Generate Kaggle API Token:**

1. **Sign In:** Log in to your Kaggle account.

2. **Navigate to Account Settings:**

   - Click on your profile picture/avatar in the top-right corner.

   - Select **"My Account"** from the dropdown menu.

3. **Create API Token:**

   - Scroll down to the **"API"** section.

   - Click on **"Create New API Token"**.

   - This action will download a file named `kaggle.json` to your computer. This file contains your API credentials.

## 2. Upload `kaggle.json` to Colab

1. **Open Google Colab:**

   - Navigate to Google Colab and open a new or existing notebook.

2. **Upload the Credentials File:**

   - In the left sidebar, click on the **"Files"** tab.

   - Click the **"Upload"** button (represented by an upward-facing arrow icon).

   - Select and upload the downloaded `kaggle.json` file.

   *Alternatively,* you can upload the file programmatically using the following code snippet. This method prompts you to select the file from your local

machine:

```python
pythonCopy code
from google.colab import files
files.upload()  # Click to select and upload the kaggle.js
on file
```

## 3. Configure the Kaggle API in Colab

Run the following commands in a new code cell to set up the Kaggle API:

```python
pythonCopy code
# Install the Kaggle package if not already installed
!pip install -q kaggle

# Create a directory for Kaggle and move the kaggle.json file
there
import os
os.makedirs('/root/.kaggle/', exist_ok=True)

# Move the uploaded kaggle.json to the Kaggle directory
!cp kaggle.json /root/.kaggle/

# Set the permissions of the kaggle.json file to ensure it's
secure
!chmod 600 /root/.kaggle/kaggle.json
```

**Explanation of Commands:**

- `pip install -q kaggle` : Installs the Kaggle API package.

- `os.makedirs('/root/.kaggle/', exist_ok=True)` : Creates the `.kaggle` directory in the root folder if it doesn't exist.

- `cp kaggle.json /root/.kaggle/` : Copies the `kaggle.json` file to the `.kaggle` directory.

- `chmod 600 /root/.kaggle/kaggle.json` : Sets the file permissions to read/write for the owner only, ensuring security.

## 4. Verify the Kaggle API Setup

To ensure that the Kaggle API is set up correctly, you can list available datasets or competitions. For example:

```python
pythonCopy code
# List available Kaggle datasets to verify the setup
!kaggle datasets list
```

If the setup is correct, this command will display a list of datasets available on Kaggle without any errors.

## 5. Download the Desired Kaggle Dataset

Identify the dataset you wish to download. You can find datasets on Kaggle's Datasets page. Each dataset has a unique identifier in the format `username/dataset-name` .

**Example:** Suppose you want to download the "Titanic" dataset by paultimothymooney.

Use the following command to download it:

```python
pythonCopy code
# Replace 'username/dataset-name' with the actual dataset ide
ntifier
!kaggle datasets download -d paultimothymooney/titanic
```

**Notes:**

- **Dataset Identifier:** You can find the dataset identifier in the URL of the dataset page. For example, in `https://www.kaggle.com/paultimothymooney/titanic` , the identifier is `paultimothymooney/titanic` .

- **Downloading Specific Files:** If you only need specific files from the dataset, you can specify them using the `f` flag. For example:

```python
pythonCopy code
!kaggle datasets download -d paultimothymooney/titanic -f train.csv
```

## 6. Unzip the Downloaded Dataset

After downloading, the dataset is typically in a ZIP archive. Unzip it to access the files:

```python
pythonCopy code
# List the files in the current directory to identify the ZIP file
!ls

# Replace 'titanic.zip' with the actual ZIP filename if different
!unzip titanic.zip
```

**Notes:**

- The ZIP filename usually matches the dataset name, but you can verify it using the `!ls` command.
- If the dataset is large, you might want to unzip it to a specific directory.

**Handling Overwrites and Existing Files:** If you run the unzip command multiple times or if files already exist, you might encounter overwrite prompts. To suppress these prompts and overwrite existing files automatically, use:

```python
pythonCopy code
!unzip -o titanic.zip
```

## 7. Access and Use the Dataset in Your Colab Notebook

Once unzipped, you can access the dataset files using standard Python operations or libraries like pandas.

**Example: Loading a CSV File with Pandas**

```python
pythonCopy code
import pandas as pd

# Replace 'train.csv' with the actual filename
df = pd.read_csv('train.csv')

# Display the first few rows of the dataframe
df.head()
```

**Example: Listing Files in the Dataset Directory**

If the dataset was unzipped into a subdirectory, navigate to that directory first:

```python
pythonCopy code
# Change directory to the dataset folder if necessary
%cd titanic  # Replace 'titanic' with the actual folder name

# List files in the current directory
!ls
```

## 8. (Optional) Mount Google Drive for Persistent Storage

If you want to store the dataset in your Google Drive for persistent access across sessions, you can mount your Drive to Colab:

```python
pythonCopy code
from google.colab import drive
```

```
drive.mount('/content/drive')
```

**Steps After Mounting:**

1. **Authenticate:** Follow the prompted link to authorize Colab to access your Google Drive.

2. **Access Files:** Your Drive files are accessible under `/content/drive/My Drive/`.

3. **Copy Files to Drive:**

```python
pythonCopy code
# Copy all unzipped files to a folder in Google Drive
!cp -r path_to_unzipped_files /content/drive/My\ Drive/des
ired_folder_name/
```

*Replace `path_to_unzipped_files` with the actual path and `desired_folder_name` with your preferred folder name in Drive.*

## Important Security Considerations

- **Protect Your `kaggle.json` File:** The `kaggle.json` file contains sensitive API credentials. **Never share it** or expose it in shared or public notebooks.

- **Avoid Committing Credentials to Version Control:** If you're saving your notebook to a public repository (e.g., GitHub), ensure that the `kaggle.json` file is not included.

## Summary

By following these steps, you can seamlessly integrate Kaggle datasets into your Google Colab environment, enabling you to leverage Colab's computational resources for data analysis, machine learning, and other tasks. Here's a quick recap:

1. **Generate and upload** your `kaggle.json` API token.

2. **Configure** the Kaggle API in Colab.

3. **Download** and **unzip** the desired dataset.

4. **Access and use** the dataset within your notebook.

5. **(Optional)** Mount Google Drive for persistent storage.

Feel free to reach out if you encounter any issues or have further questions!