# NIT HACKATHON

Vertical :  Data Science

Topic : Website called 'Pipeline' for data preprocessing
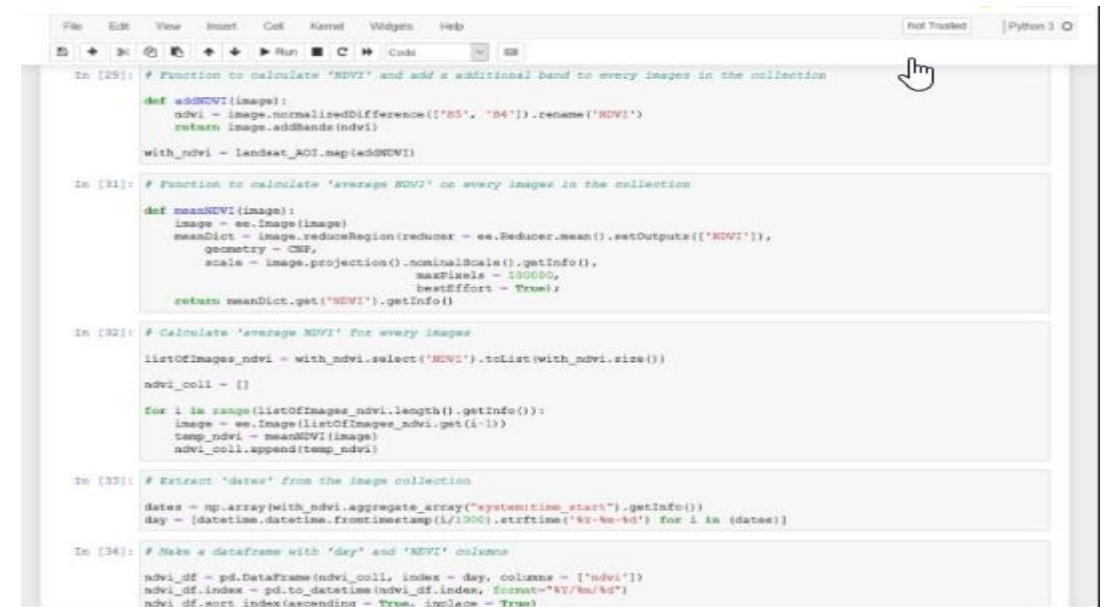
College : Shiv Nadar University Chennai

Team : Rhea Pandita (B.TECH 2nd Year)

Vedant Nair   (B.TECH 2nd Year)

# Aim of the project



➢ 'Pipeline' for Data Preprocessing is a website meant to aid in the data preprocessing of csv files.

➢ The usual process of preprocessing csv files includes writing 10-20 lines of code. Also, the methods applied differ from one file to another based on its contents.

➢ This website's aim is to reduce the time taken to do the same by usage of drop-down menus back-ended by common methods data analysts apply.

# Features available

- Uploading required file from the system or using drag and drop
- Dropdown menus to choose data preprocessing method from
- Following are the preprocessing dropdowns:
  - Null Value Treatment
  - Outlier Treatment
  - Dropping Data
  - Data Creation
  - Feature Scaling

# Scope of the project

- This project can be used to preprocess data that is in table form. Since a large number of industries are required to have their own repository of their business in the past and the present, this will facilitate a quick cleanup of the data

- Since most industries are rapidly moving forward in the direction of online enterprise, it is necessary for them to tabulate all changes that are encountered (even the minor ones). This increases the number of databases overall and its size, thus demanding a need for cleanup. This can greatly benefit the medical industry, the judicial system keep track of their cases etc and many more places.