# WORKING WITH THE DATASET

| Date | 28 OCTOBER 2022 |
|---|---|
| Team ID | PNT2022TMID25946 |
| Project Name | ESTIMATION OF CROP YIELD PREDICTION USING DATA ANALYTICS |

## LOADING AND UNDERSTANDING THE DATASET:

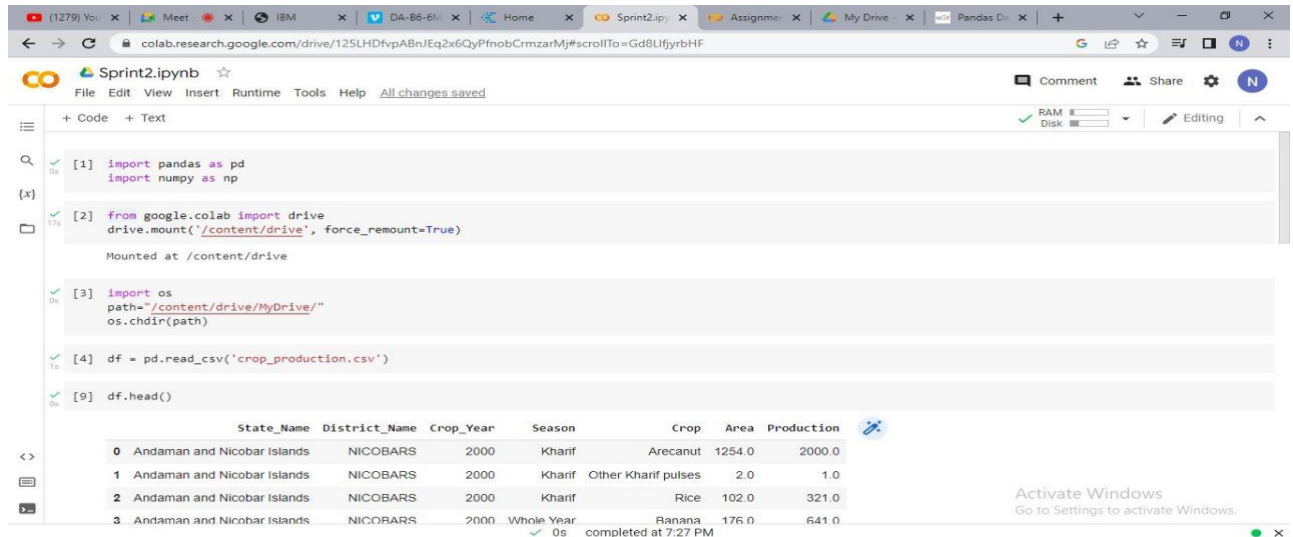The dataset is downloaded from the Kaggle website with reference to the hyperlink provided in the project flow.

## DESCRIPTION OF THE DATASET:

The given dataset was analyzed to get the knowledge about that dataset. It is done by writing a python code in Google Colab platform.

# CLEANING THE DATASET:

The dataset may contain null values, so it must cleaned before using and also cleaning process is important to perform accurate visualization.