Raw Data Collected from Various Sources

Data Curation and Feature Engineering Tool

1 Data Imputation, Cleaning, & Exploration

2 Feature Engineering

3 Data Visualization

Preprocessed Data Ready for ML/DL Algorithms

# DCFE Web Application
# User Guide

## For Vadu HDSS Dataset

**Table of Contents**

## Introduction

In the process of data analysis, it is observed that 60% of the time is spent in cleaning and preprocessing the data. To overcome this drawback and to reduce the time spent in preprocessing we have developed a tool named Data Curation and Feature Engineering Tool (DCFE Tool). Epidemiologists, Researchers or any kind of Data analysts can use this web based application tool to upload, explore and analyze the epidemic/pandemic or health care datasets. The tool provides the provision for the user to upload datasets which are in .csv or .xlsx file format. The complete detail of the uploaded data like the size, No. of rows and columns, Percentage of missing value, Numerical and Categorical and date-time attributes can be viewed using 'Dataset Overview'. The 'data exploration' projects the statistical values such as mean, median, skewness, kurtosis, SD and No. of missing values of each attribute in the dataset.

Data imputation and Feature engineering are the two primary modules of this tool. The missing value present in the dataset is handled using the data imputation module where the missing values can be completely dropped or it can be filled using any suitable method like 'Attribute wise drop NaN', 'Attribute wise Fill NaN', Machine learning model based imputation like 'KNN' and 'Iterative imputation'. In Feature Engineering module, the user can use two ways of transformation. Both numerical and categorical data are transformed into new feature which are ready for any kind of algorithmic computation. The different encoding methods used here for transformation of both numerical and categorical features are Binning, Normalization, Log transform, Label, Ordinal, Binary, One-hot , Count frequency encoding.

The user can generate new features using Date-time features. If user wants to delete an unwanted feature, 'Remove unwanted feature' function can be used. If a user wants to roll back to the previous action performed 'Undo' function can be used.  The user can download his/her processed data using the function 'download processed data' or the user can also download the original uploaded dataset. For repeating the process with another dataset then the user can click 'Home' function to return to the front page of the tool.
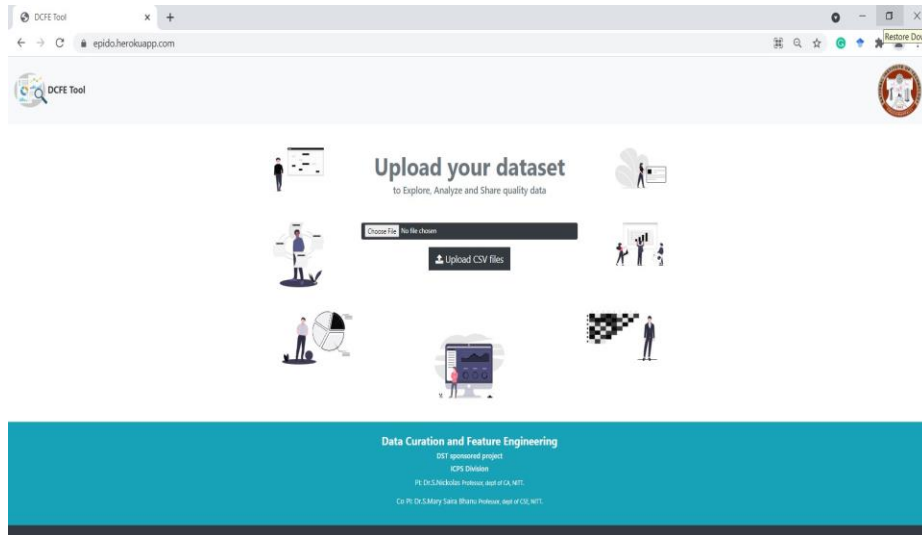
To visualize the uploaded data, the user can use the 'Visualization'. The function will give provision to use different charts such as bar, scatter, pie charts, which will take the features as inputs to display the counts of each category.

Further, the user can visualize the percentage of categorical and numerical features in the uploaded data using a pie chart. The user can also visualize the skewness, kurtosis and No. of missing value of each attribute in the dataset.

The correlation between the features in the dataset is shown using Heatmap. Finally, the outlier in each attribute of the data can be visualized using the box plot.


This document guides a user through the necessary steps to preprocess and download the curated data which can be used for any further computation with any Machine learning algorithm.

## Front page of the project



## DCFE Web Application Tool Link

Users can directly use our Web Application DCFE Tool by clicking on the following link: https://epido.herokuapp.com/

## Upload your dataset

The front page of our tool has a provision to upload datasets. The dataset must be in either .csv or .xlsx file format.

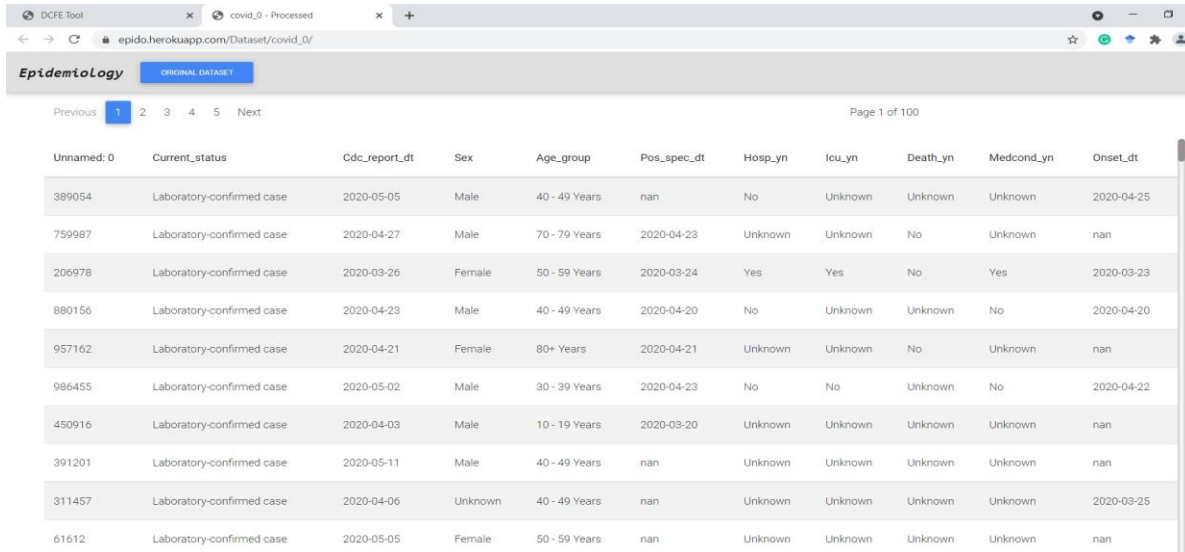## Choose File

If the user clicks 'Choose File' button, it will redirect the user to the file location to choose appropriate dataset file.

## Upload CSV files

Users can click this button to upload the dataset.

## Dataset view

**Epidemiology**   ORIGINAL DATASET

Previous 1 2 3 4 5 Next     Page 1 of 100

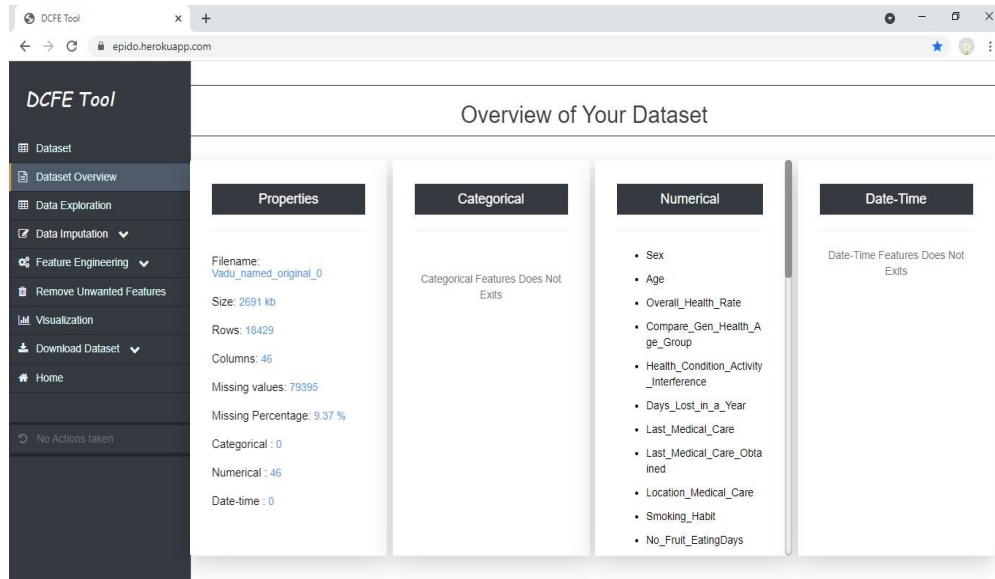| Unnamed: 0 | Current_status | Cdc_report_dt | Sex | Age_group | Pos_spec_dt | Hosp_yn | Icu_yn | Death_yn | Medcond_yn | Onset_dt |
|---|---|---|---|---|---|---|---|---|---|---|
| 389054 | Laboratory-confirmed case | 2020-05-05 | Male | 40 - 49 Years | nan | No | Unknown | Unknown | Unknown | 2020-04-25 |
| 759987 | Laboratory-confirmed case | 2020-04-27 | Male | 70 - 79 Years | 2020-04-23 | Unknown | Unknown | No | Unknown | nan |
| 206978 | Laboratory-confirmed case | 2020-03-26 | Female | 50 - 59 Years | 2020-03-24 | Yes | Yes | No | Yes | 2020-03-23 |
| 880156 | Laboratory-confirmed case | 2020-04-23 | Male | 40 - 49 Years | 2020-04-20 | No | Unknown | Unknown | No | 2020-04-20 |
| 957162 | Laboratory-confirmed case | 2020-04-21 | Female | 80+ Years | 2020-04-21 | Unknown | Unknown | No | Unknown | nan |
| 986455 | Laboratory-confirmed case | 2020-05-02 | Male | 30 - 39 Years | 2020-04-23 | No | No | Unknown | No | 2020-04-22 |
| 450916 | Laboratory-confirmed case | 2020-04-03 | Male | 10 - 19 Years | 2020-03-20 | Unknown | Unknown | Unknown | Unknown | nan |
| 391201 | Laboratory-confirmed case | 2020-05-11 | Male | 40 - 49 Years | nan | Unknown | Unknown | Unknown | Unknown | nan |
| 311457 | Laboratory-confirmed case | 2020-04-06 | Unknown | 40 - 49 Years | nan | Unknown | Unknown | Unknown | Unknown | 2020-03-25 |
| 61612 | Laboratory-confirmed case | 2020-05-05 | Female | 50 - 59 Years | nan | Unknown | Unknown | Unknown | Unknown | nan |

### Dataset
User can view the uploaded dataset by selecting the menu **'Dataset'**, where he/she can view the complete dataset page wise.

### Original Dataset
If the user wants to compare the processed file with the original file, then the button **'ORIGINAL DATASET'** can be clicked to download the original file.

## Dataset Overview



### Properties

If the user wants to know about the details of the uploaded .csv or .xlxs file, the menu **'Properties'** will provide all the necessary information like Filename, Size, Rows, Columns, Missing values, categorical, Numerical and Date-time features.

### Categorical

If the user wants to view how many categorical data is there in a dataset, the menu **'Categorical'** will expose those details.
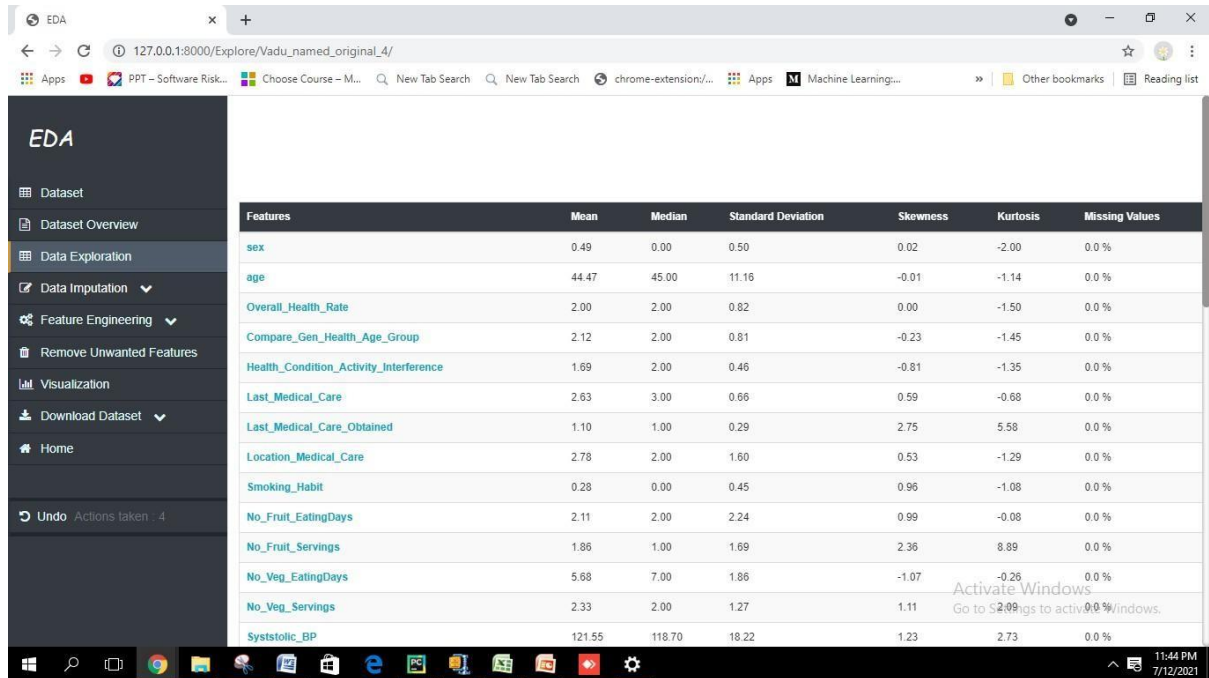
### Numerical

If the user wants to see how many numerical data present in a dataset, the menu **'Numerical'** will expose those details.

### Date-time

Once the file is uploaded the user is provided with the provision (**'Date-time')** to see whether the dataset contains Date-time features or not.

# Dataset Exploration



## Mean

If a user wants to know the mean/average values of each category, then the menu 'Mean' will generate the values for each feature.

## Median

If a user wants to know the median value of each feature, then the menu 'Median' will showcase the values.

## Standard deviation, Skewness, kurtosis

If a user wants to know the median values of each feature, then the menu '**Standard deviation, Skewness, kurtosis** will showcase the values.

## Missing values percentage

The user will come to know about the missing percentage for each attribute.

## Data Imputation-Complete Drop NaN



### Proceed

When the users click the link **'Complete Drop NaN'** it will open a Dialog Box showing the rules for dropping the NaN values. The users can click the button **'Proceed'** after reading the rules for dropping the NaN values.

### Close

The users can click **'close'** button after finish reading the mentioned rules.

## Data Imputation- Attribute Wise Drop NaN



## Select Features to Drop NaN values

If the user wants to drop a feature which has NaN values, then this menu should be clicked by the user. The check boxes are there for each feature. To drop the missing values in specific attributes/features when there are less than 15% of missing values in it. The features with less than 15% missing values are shown in bold embossed i.e enabled to check the check box.

### Check all
If a user wants to perform drop nan values on all attributes, then 'check all' option can be used.

### Drop
If user wants to do drop the feature after selecting the feature, the user has to click the 'Drop' button. Then a message will be displayed that the selected feature is deleted. (**Success! NaN values are dropped. Please refresh the page and see the changes**)

## Data Imputation- Attribute Wise Drop NaN



**Success! NaN values are dropped. Please refresh the page and see the changes**

The user can view the update after processing of 'Attribute wise drop NaN' module. The missing values in those attributes are deleted using 'Attribute wise drop' module. After dropping the missing values there are 0% of missing values in these attributes.
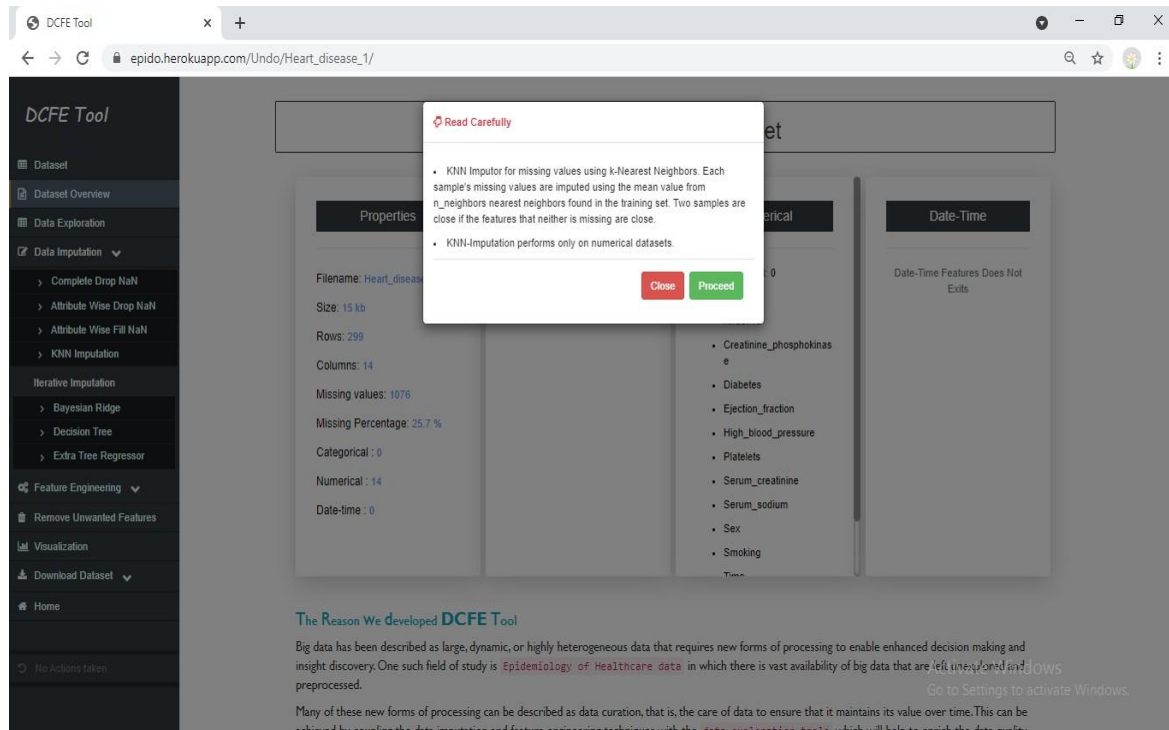
## Data Imputation- Attribute wise Fill NaN



## Fill

If the user wants to fill NaN, then they can click the value which is shown in the drop down list as 'Fill'. Another option shown in the drop down list is 'Replace'. So either Fill or Replace value can be chosen by the user to perform 'filling missing values' operation.

## Forward Fill, Backward Fill, Mean and Mode

The user is provided with four different methods (Forward Fill, Backward Fill, Mean and Mode) to fill NaN. The options are provided through drop down list.
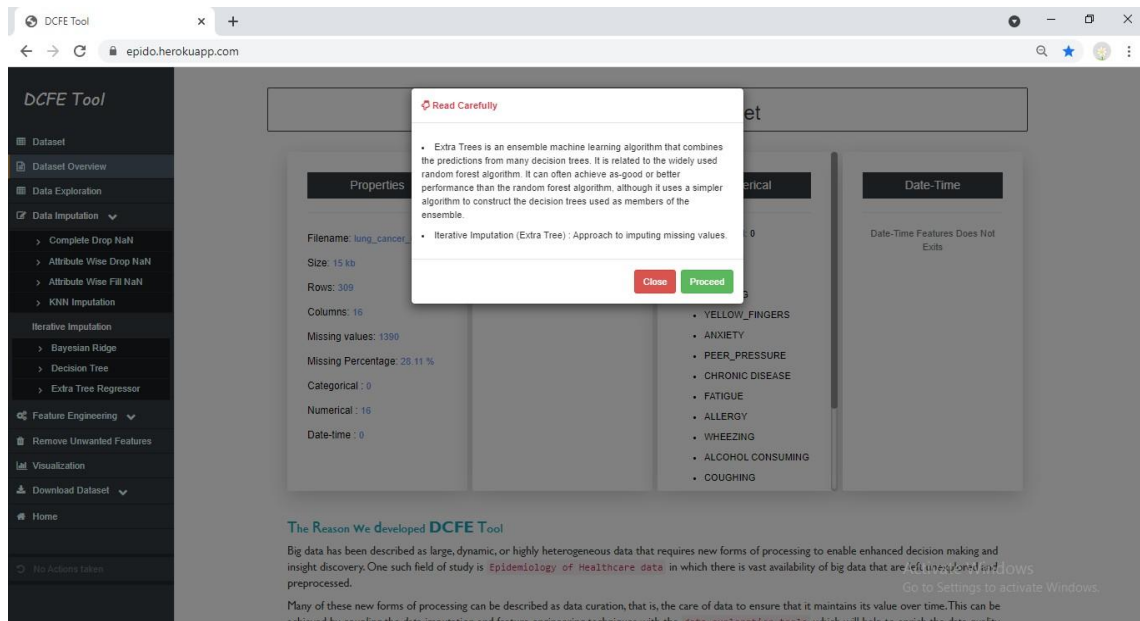
# Data Imputation- KNN Imputation



## Close

'Close' button is used to close the dialog box after the rules are read by any user.

## Proceed

If the user wants to fill missing values using KNN imputation, then click 'Proceed' button. The missing values are imputed using KNN imputation where the missing values are filled using average value of the 'n' nearest neighbor values.

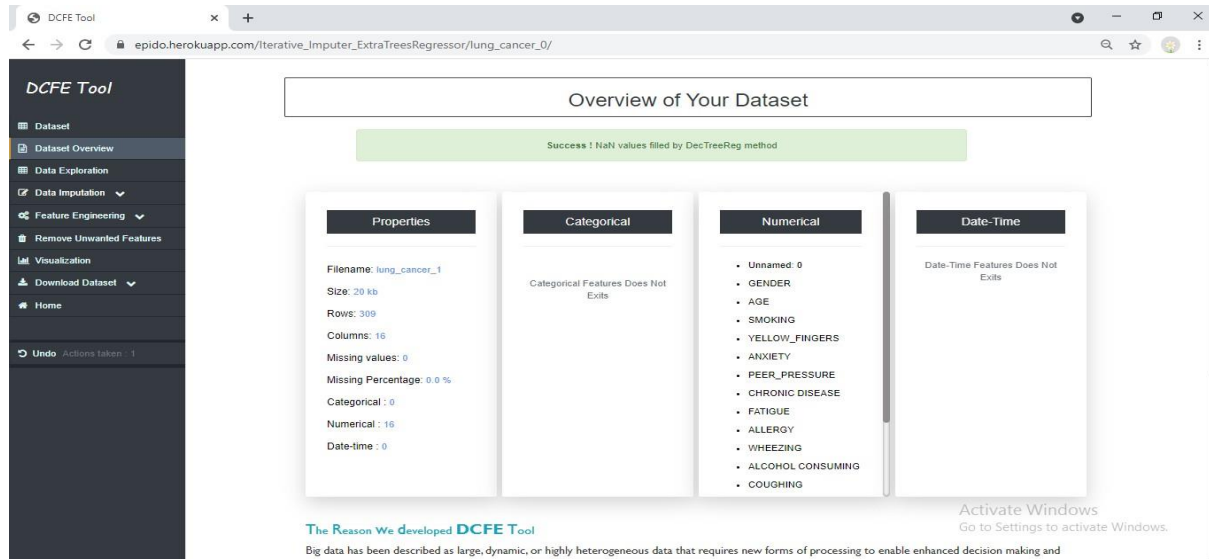## Iterative Imputation- Extra Tree Regressor Imputation



### Close
'Close' button is used to close the dialog box after the rules are read by any user.

### Proceed
If the user wants to fill missing values using Extra Tree Regressor imputation, then click 'Proceed' button. The missing values are imputed using Extra Tree Regressor.

# Iterative Imputation- Decision Tree Regressor Imputation



## Close
'Close' button is used to close the dialog box after the rules are read by any user.

## Proceed
If the user wants to fill missing values using Decision Tree Regressor imputation, then click 'Proceed' button. The missing values are imputed using Decision Tree Regressor imputation method. The changes can be viewed in the Dataset Overview in which Missing percentage will be shown as 'Zero'

## Feature Engineering- Binning



## Check All

The user can chose any attribute for doing Binning. Binning can be done in numerical data alone. Individual check option is also there for any user.

## Custom range Per Bin

The user can chose/customize the range of each bin for cut (Equal sized Discretization) method. The maximum and minimum values can be known through visualization of Box plot for every attribute.

Next method shown in the drop down list is qcut (Quantile based discretization) method. The user can also chose qcut method for Binning process.

## Proceed

If user wants to start any method for processing, then click 'Proceed' button for Binning any continuous values. After clicking 'proceed' button, it will redirect the user to Dataset view menu to see the changes in the dataset.

# Feature Engineering- Normalization



## Check All
The user can chose any attribute for doing Normalization. Normalization can be done in numerical data alone. Individual check option is also there for any user.

## Drop down list
The drop down list has three options for normalization of any numerical data. The user can click any method to perform normalization depending upon their uploaded data.

## Proceed
The user will click 'proceed' button for doing normalization. The updated values can be viewable from the Dataset menu option. Every change can also be reflected in visualization part.

## Feature Engineering- Min – Max normalization



## Check All

The user can chose any attribute for doing Min-Max Normalization. Normalization can be done in numerical data alone. Individual check option is also there for any user.

## Drop down list

The Min-Max option can be clicked by any user.

## Max and Min value

The user can identify the minimum and maximum value using Box –Plot which is shown in visualization part.

## Feature Engineering- Z Score normalization



## Check All

The user can choose any attribute for doing Z-score Normalization. Z-Score normalization can be done in numerical data alone. Individual check option is also there for any user.

## Proceed

The 'Proceed' button can provide the user to perform Z-score normalization on the checked attribute. The mean and mean absolute deviation is automatically assigned for chosen attribute.

# Feature Engineering-Decimal Scaling



## Check All

The user can chose any attribute for doing Decimal Scaling Normalization. Individual check option is also there for any user.

## Proceed

The 'Proceed' button can provide the user to perform Decimal Scaling normalization on the checked attribute. The value for Decimal Scaling is automatically generated.

# Feature Engineering-Log Transform



## Check All

The user can select any attribute for doing Log Transform. Individual check option is also there for any user.

## Proceed

The 'Proceed' button can provide the user to perform Log Transform on the selected attribute. The value is automatically generated which replaces each variable X with a log(X).

## Feature Engineering-Ordinal Encoding



## Check All

The user can choose any attribute for performing Ordinal Encoding. Individual check option is also there for any user. The categorical features alone displayed in the display list of features.

## Proceed

The **'Proceed'** button can provide the user to perform ordinal encoding on the checked attribute.

**Message: Success! Ordinal Encoding was done on selected features.** The user can see the mentioned message which is displayed after finishing the proceed button. The user can select the 'Dataset' menu option to know the changes made in the dataset.

## Ordinal Encoded Features

The user can view the newly generated feature with its new attribute name.

# Feature Engineering- Label Encoding



## Check All

The user can choose any attribute for performing Label Encoding. Individual check option is also there for any user. The categorical features alone displayed in the display list of features.

## Proceed

The **'Proceed'** button can provide the user to perform Label encoding on the checked attribute.

**Message: Success! Label Encoding was done on selected features.** The user can see the mentioned message which is displayed after finishing the proceed button. The user can select the 'Dataset' menu option to know the changes made in the dataset.

## Labelled Features

The user can view the newly generated feature with its new attribute name.

# Feature Engineering-One-Hot Encoding



## Check All

The user can choose any attribute for performing One-Hot Encoding. Individual check option is also there for any user. The categorical features alone displayed in the display list of features.
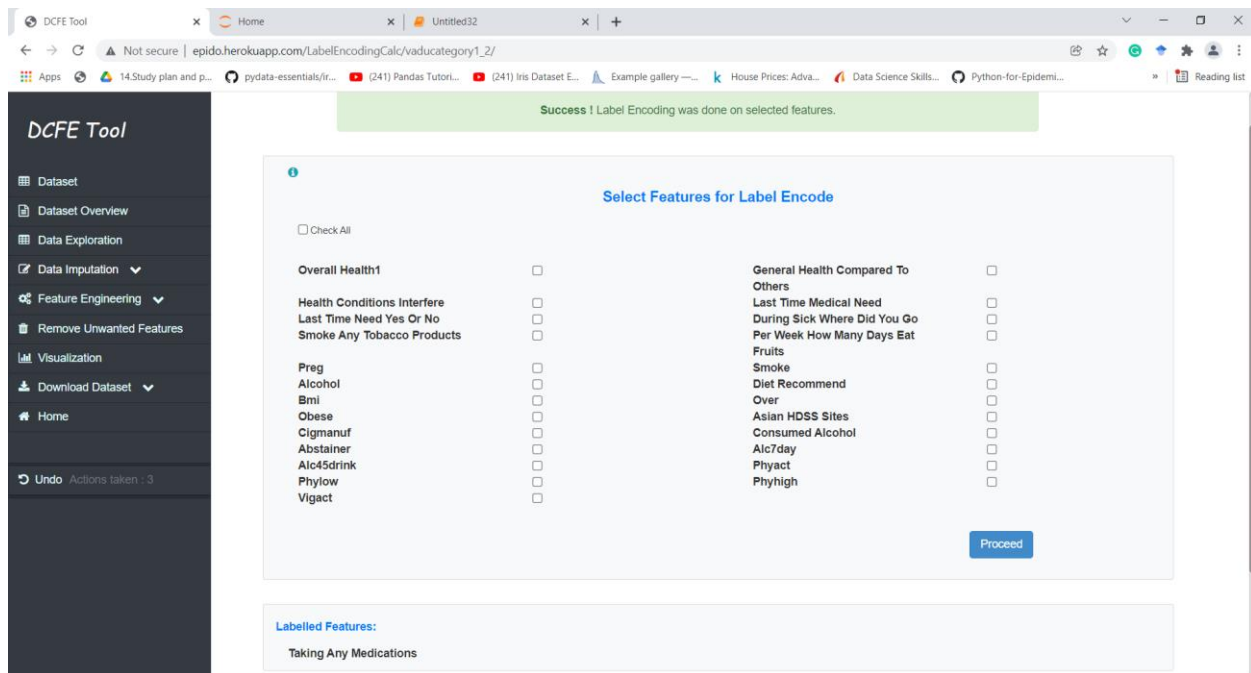
## Proceed

The **'Proceed'** button can provide the user to perform One-Hot encoding on the checked attribute.

**Message: Success! One-Hot Encoding was done on selected features.** The user can see the mentioned message which is displayed after finishing the proceed button. The user can select the 'Dataset' menu option to know the changes made in the dataset.

## One-Hot Encoded Features
**The user can view the newly generated feature with its new attribute name.**

## Feature Engineering- Count Frequency



### Check All

The user can choose any attribute for performing Count Frequency Encoding. Individual check option is also there for any user. The categorical features alone displayed in the display list of features.
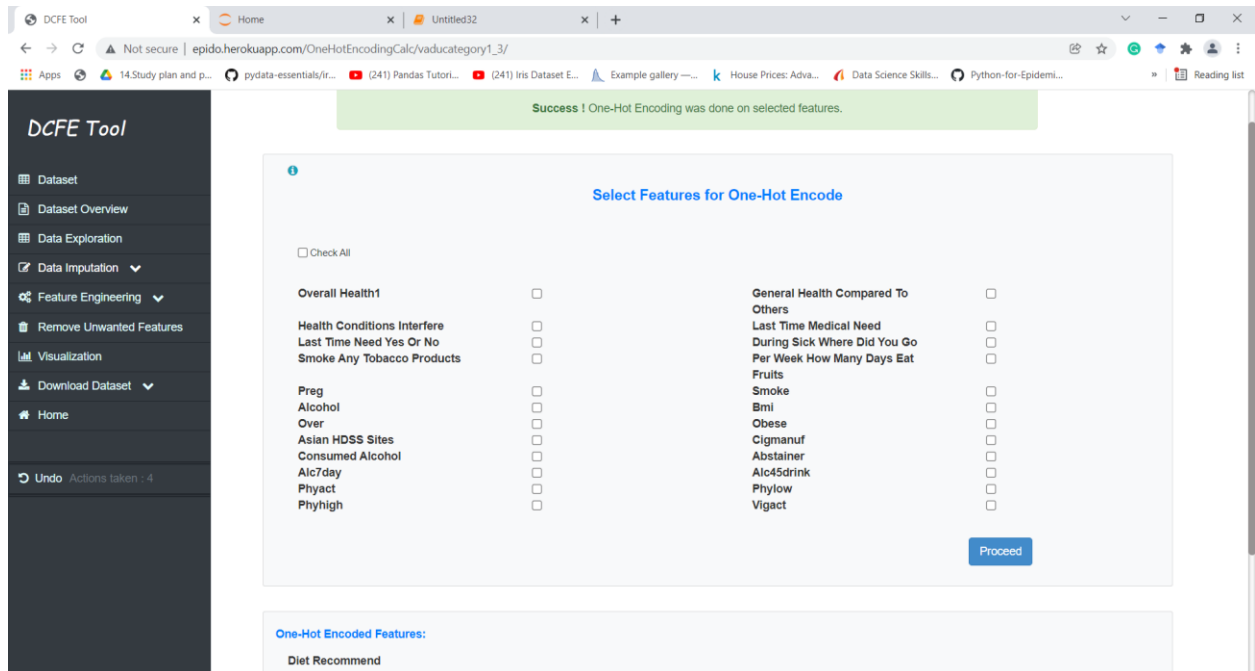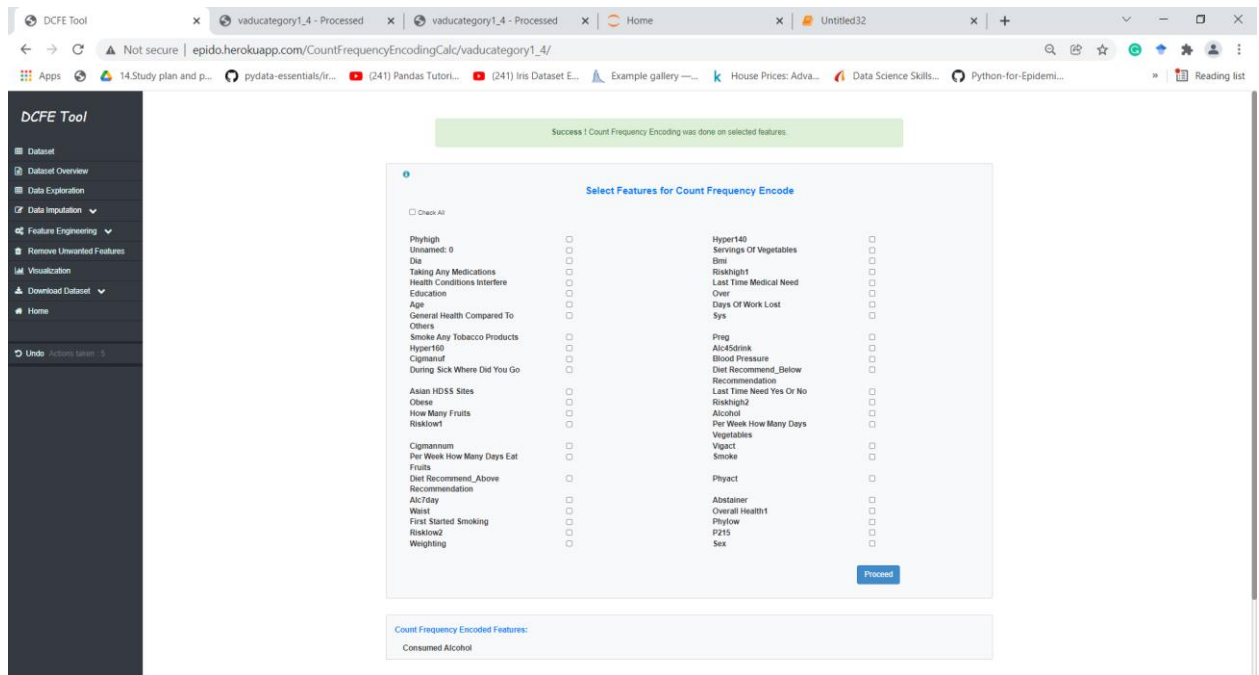
### Proceed

The **'Proceed'** button can provide the user to perform count frequency encoding on the checked attribute.

**Message: Success! Count frequency Encoding was done on selected features.** The user can see the mentioned message which is displayed after finishing the proceed button. The user can select the 'Dataset' menu option to know the changes made in the dataset.

### Count-frequency Encoded Features

The user can view the newly generated feature with its new attribute name.

# Feature Engineering-Binary Encoding



## Check All

The user can choose any attribute for performing Binary Encoding. Individual check option is also there for any user. The categorical features alone displayed in the display list of features.
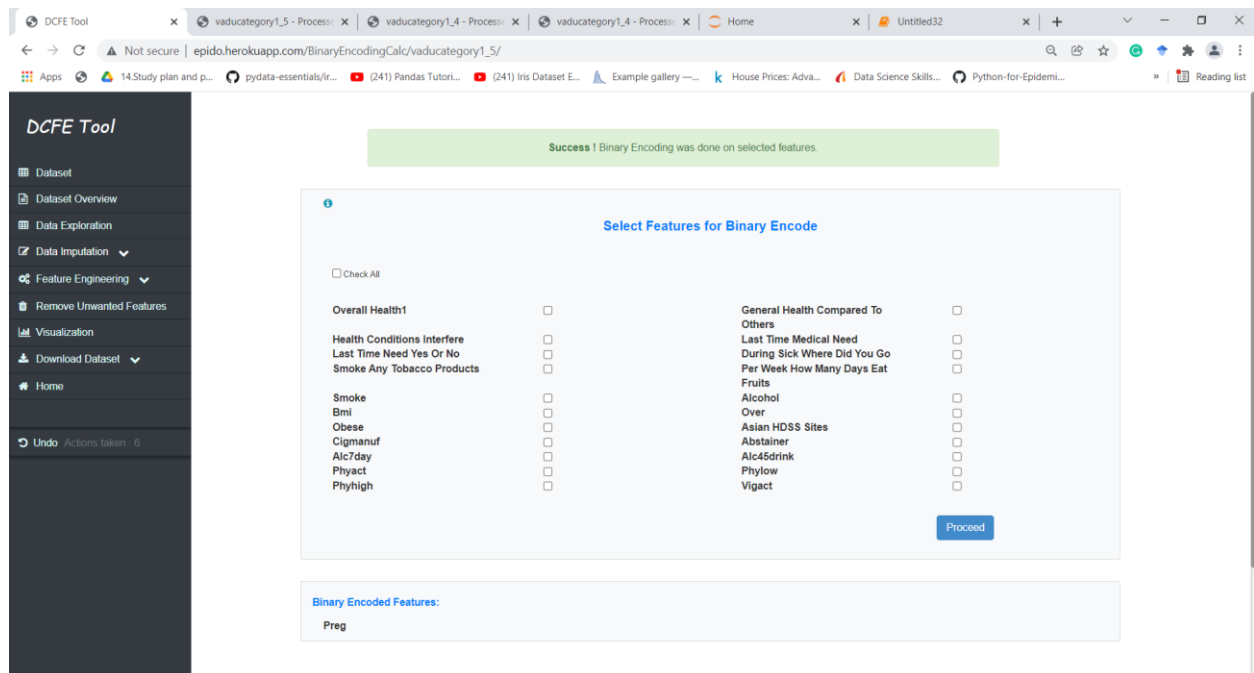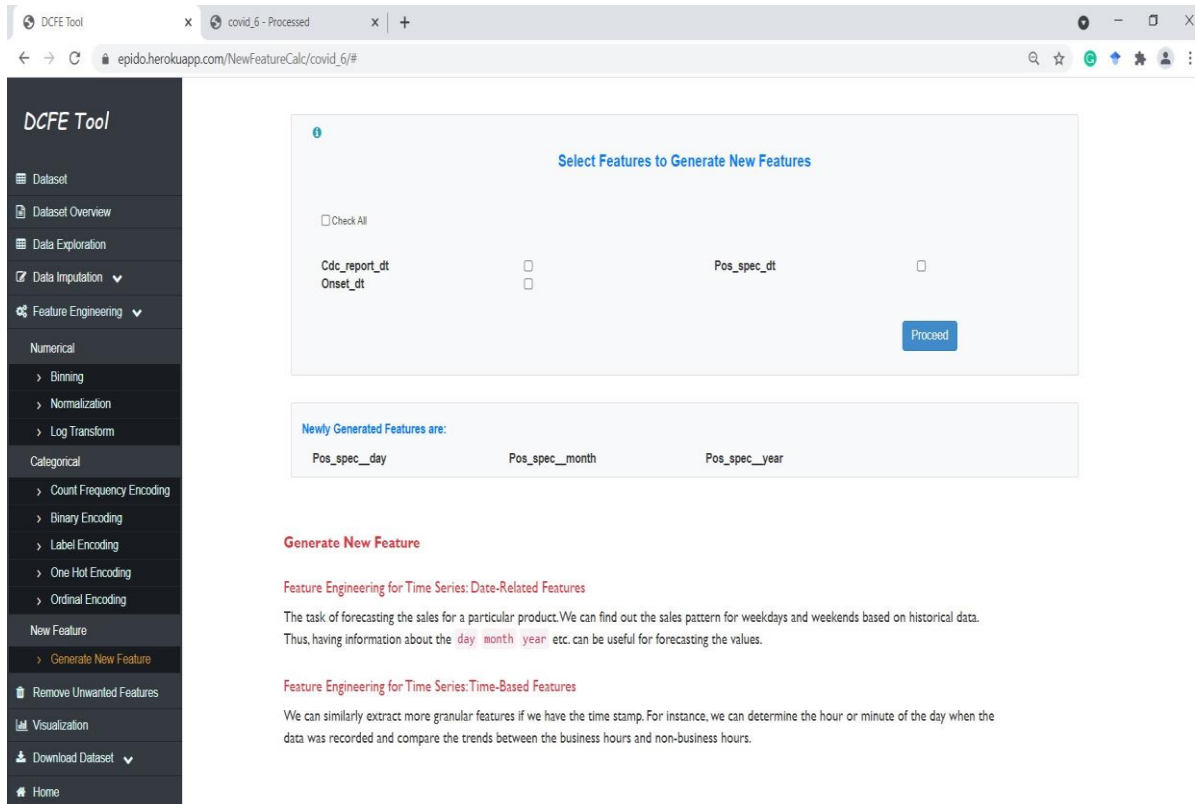
## Proceed

The **'Proceed'** button can provide the user to perform binary encoding on the checked attribute.

**Message: Success! Binary Encoding was done on selected features.** The user can see the mentioned message which is displayed after finishing the proceed button. The user can select the 'Dataset' menu option to know the changes made in the dataset.

## Binary Encoded Features

The user can view the newly generated feature with its new attribute name.

# Feature Engineering-New feature- Generate new feature



## Check All
The user can choose any attribute to perform new feature generation only on Date-Time features. Individual check option is also there for any user. The Date-Time features alone displayed in the display list of features.
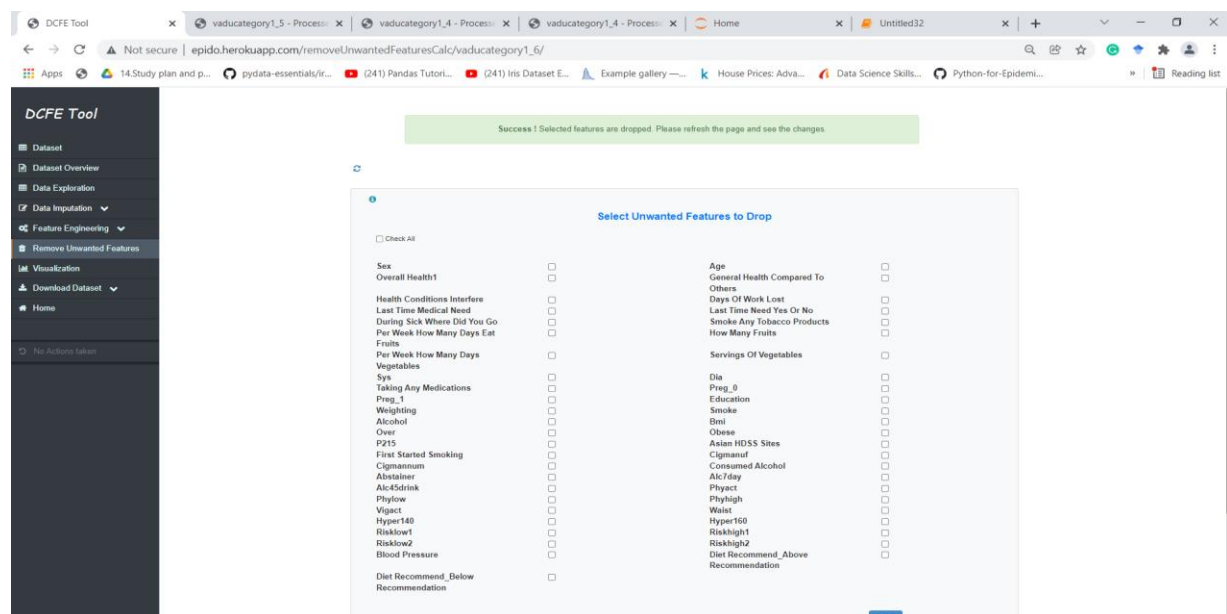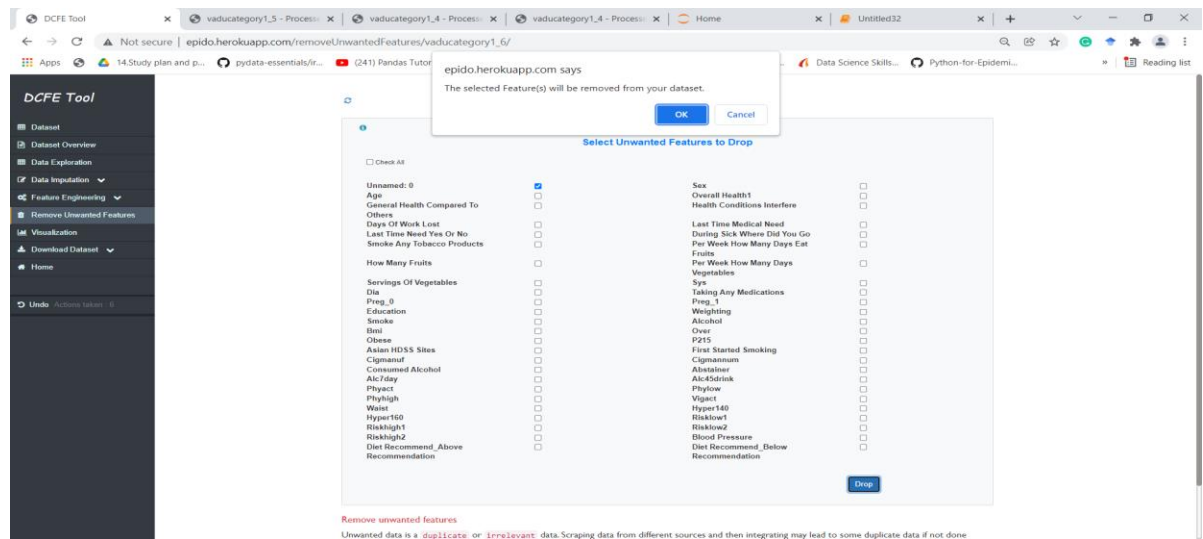
## Proceed

The **'Proceed'** button can provide the user to perform new feature generation on the checked attribute.

**Message: Success! New features are generated on selected features.** The user can see the mentioned message which is displayed after finishing the proceed button. The user can select the 'Dataset' menu option to know the changes made in the dataset.

**Newly generated Features are:**
**The user can view the newly generated feature with its new attribute name.**

## Remove unwanted feature





## Check All

The user can choose any attribute to perform 'remove unwanted feature' among all features. Individual check option is also there for any user. All features are displayed in the display list of features
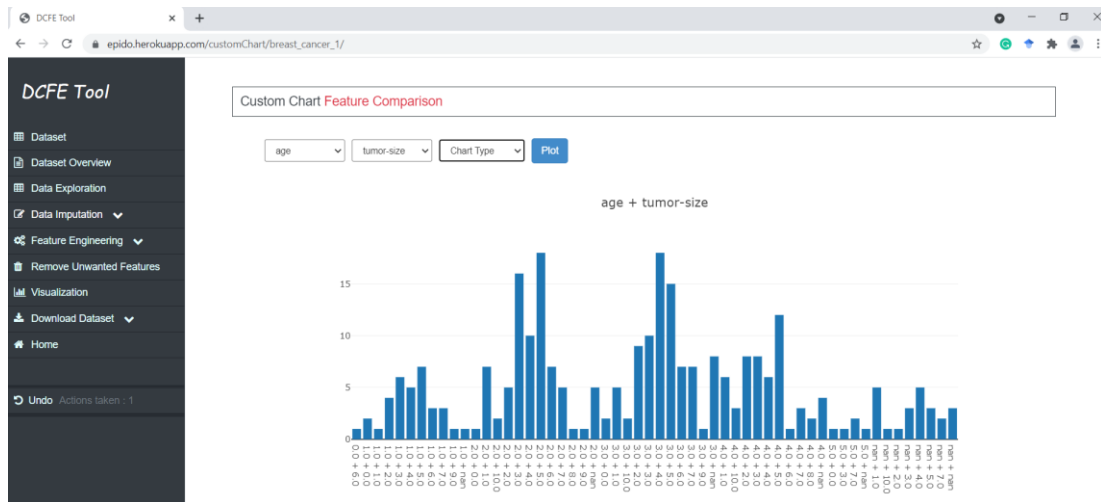
## Drop

**'Drop' button is provided to the user, to proceed with the operation of removing unwanted features.**

**Message: Success! Selected features are dropped. Please refresh the page and see the changes.**

**This message is shown to the user to know the status of the function.**

## Visualization-Feature Comparison



## Drop-Down list Feature Comparison

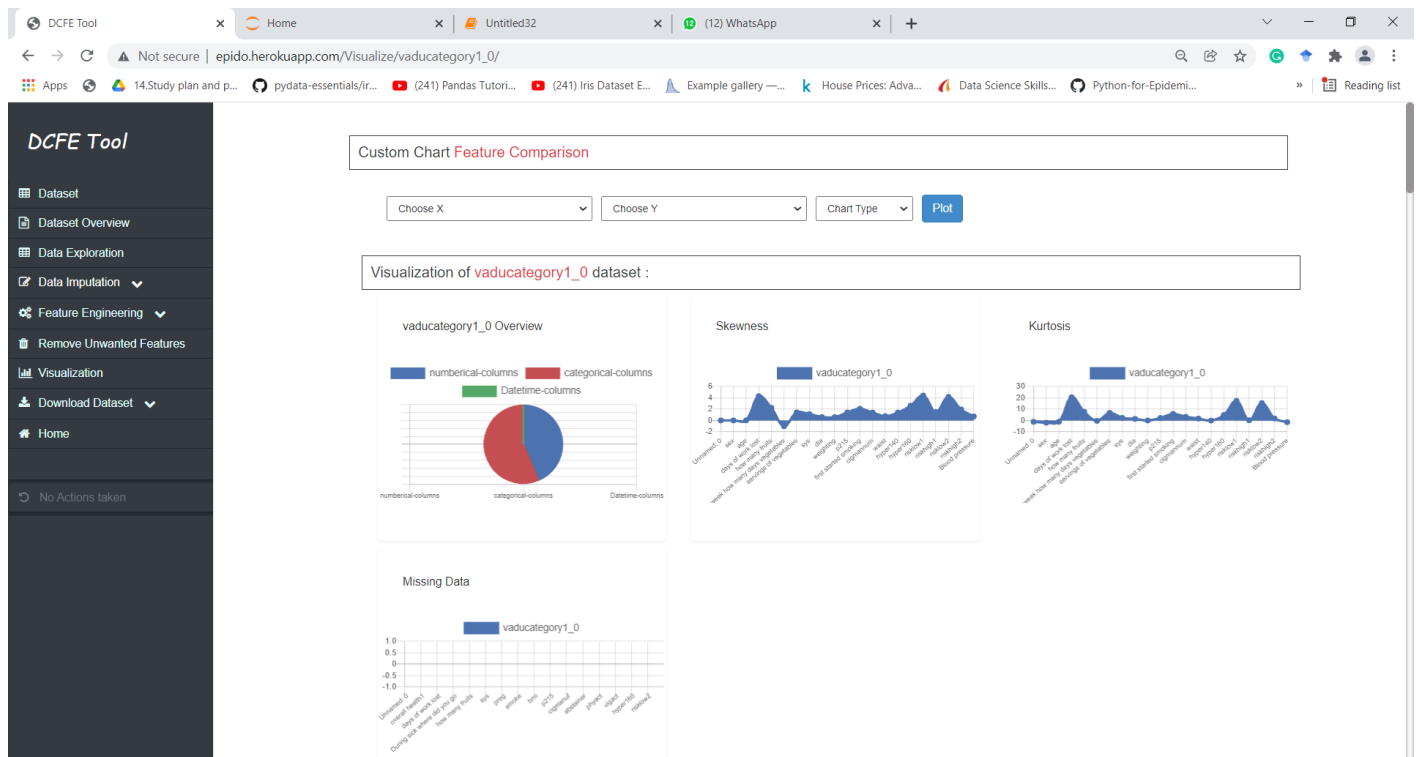The drop-down list is provided for the user to choose X axis and Y axis features.

## Chart Type
The chart type drop down list for the user is provided for selecting the various charts to plot such as Pie chart, Bar chart, Scatter chart and line chart.

Plot
**'Plot'** Button is assigned for triggering the function to perform with the selected chart type.

## Visualization- Data Exploration



### Dataset overview

If the user wants to visualize the full dataset i.e how many numerical and categorical features present in the dataset,that can be viewed here as scatter plot.

### Skewness
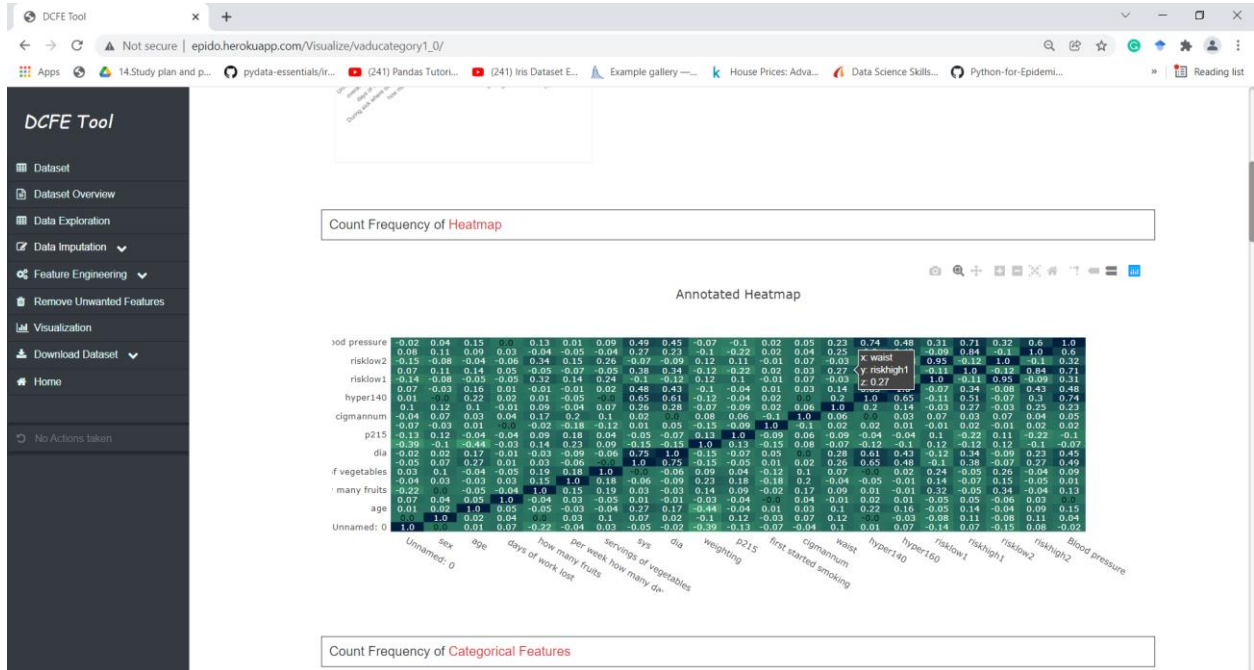The users are provided with the provision to view the skewness of each variable in the dataset.

### Kurtosis
The users are provided with the provision to view the kurtosis of each variable in the dataset.

### Missing Data

The users can view how many values are missed in each feature.

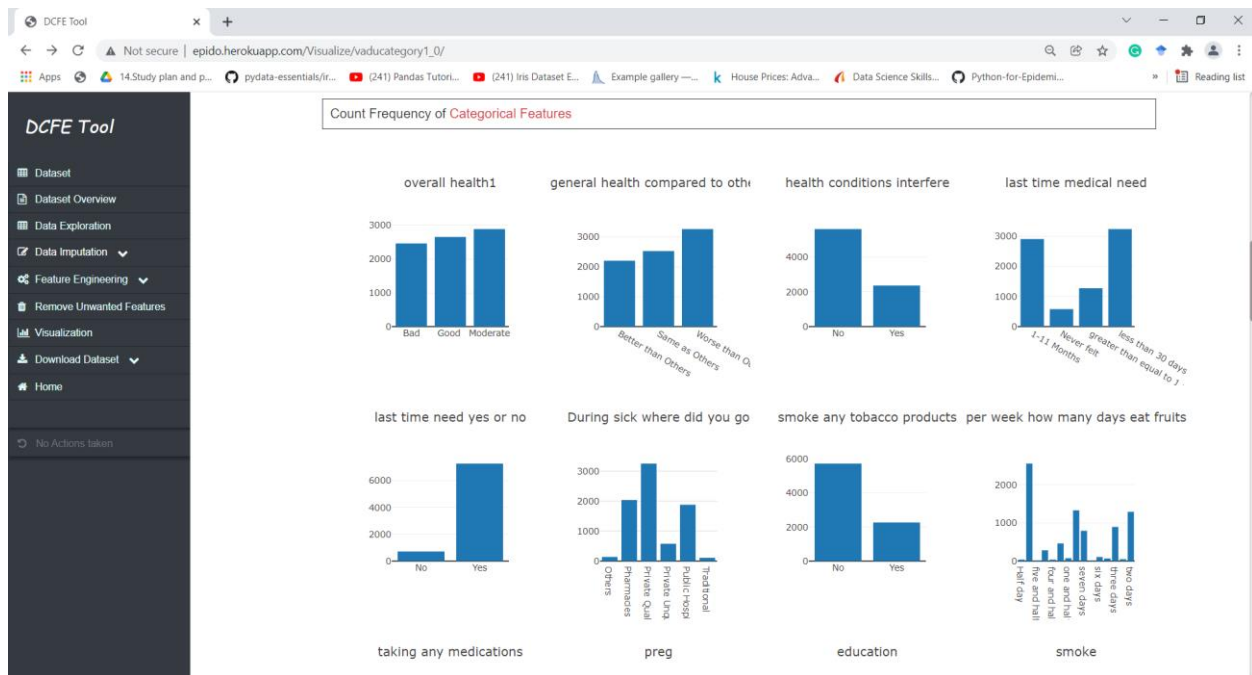## Visualization-Annotated Heat Map



If the user wants to know the correlation between every attribute, then the above heatmap is used to know the correlation between attributes.

Hovering of the cursor will display the exact percentage values.

The user can view X axis and Y axis selected features and Z axis will show the correlated value.
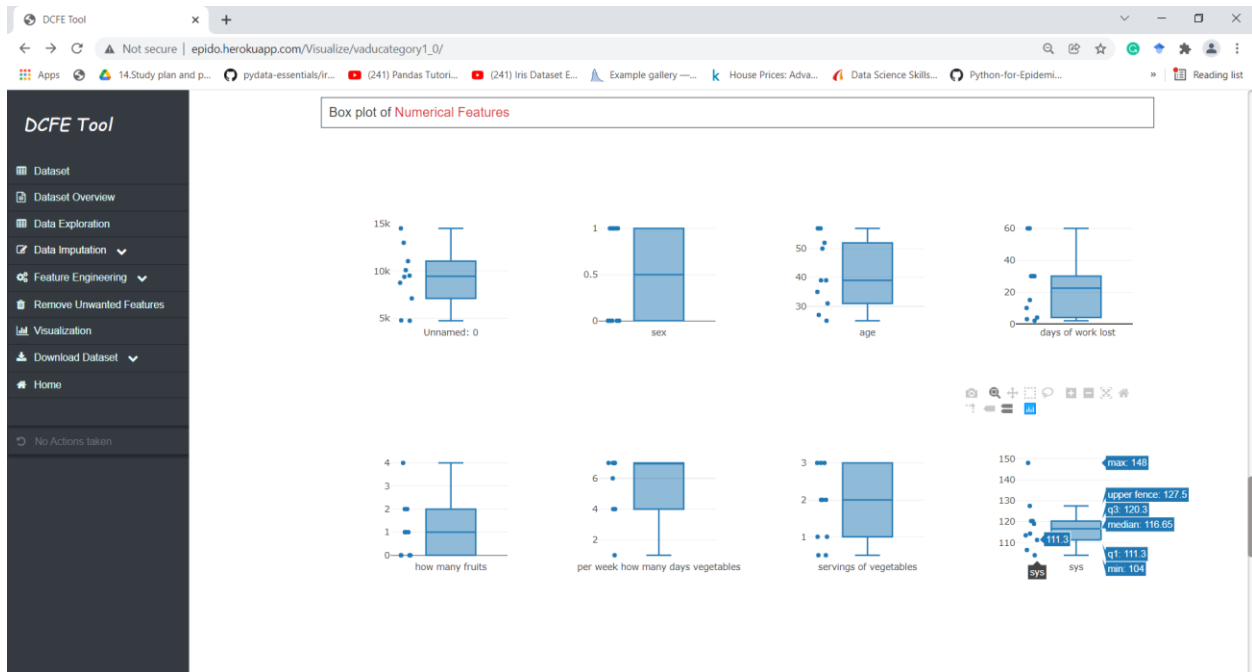
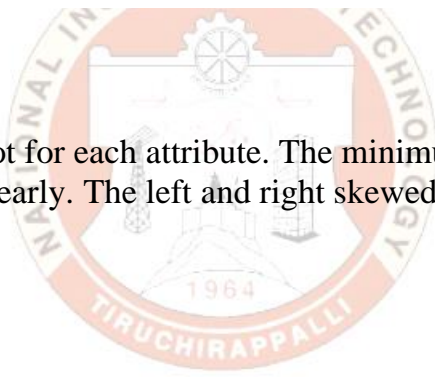# Visualization- Count Frequency of categorical features



The users can view the counts of each category of every attribute using Bar chart.

## Visualization-Box-Plot



The user can view the Box-Plot for each attribute. The minimum and maximum values of each attribute can be shown clearly. The left and right skewed values are also viewable to the user.
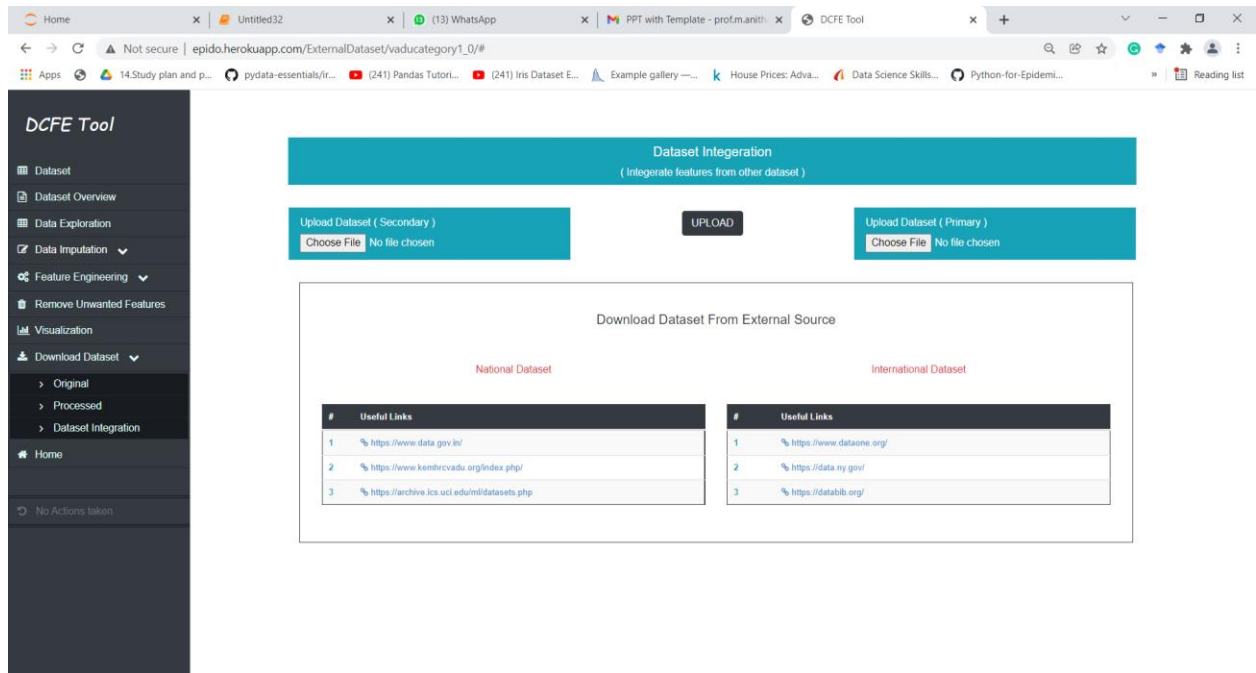
Download Dataset



## Original Dataset

If the user wants to download the original dataset for any comparison of values with the processed file, then they are provided with Download option of original dataset. The Downloaded dataset can be stored in the local file location by the user.

## Processed Dataset

If the user has finished every functions mentioned in the Tool and if it is in curated form, then it is ready to download.
The user is allowed to click the processed link and the file is downloaded and stored in the local file location itself for further usage.

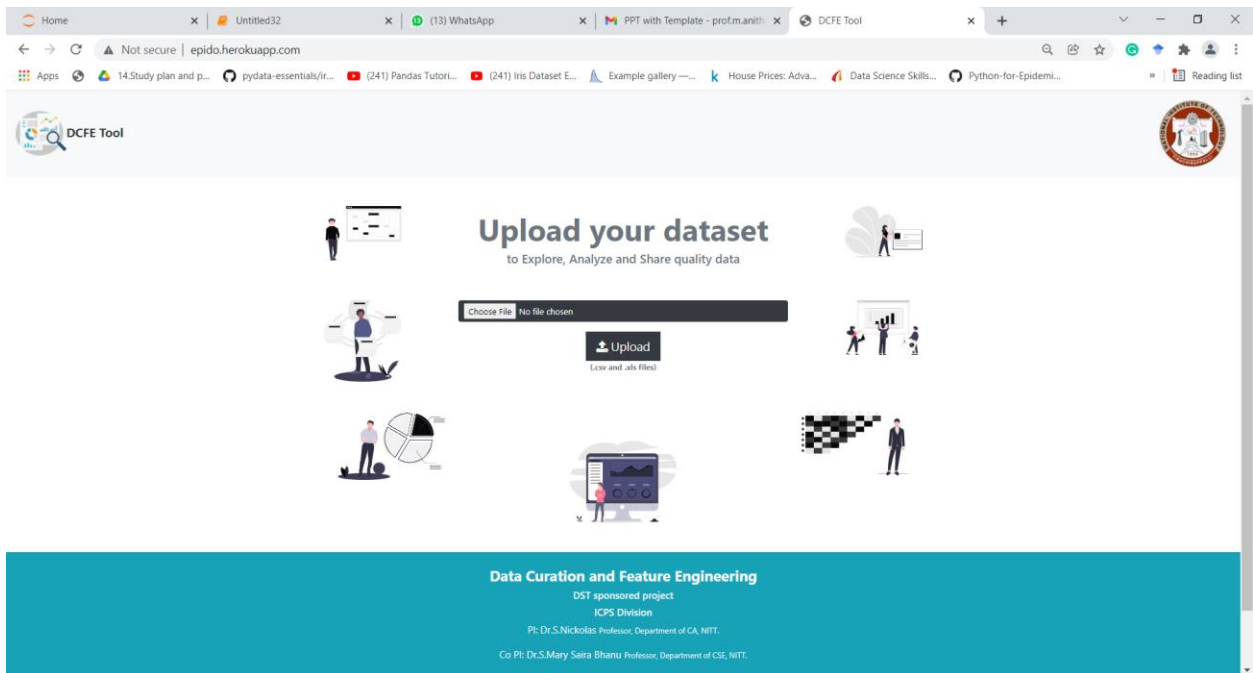Download Dataset-Dataset Integration



**Choose file**

The user can select .csv files for uploading as both secondary and primary. The file must have common attribute between them in-order to do integration of data.

**Download Dataset from External Source**

The dataset can be downloaded from external source and a provision is given in the Download Dataset menu itself for the user.

The user is provided with the provision of downloading National and International datasets very easily through Data collection interface

# Home



Home Link can be used by any user in-order to return to Home Page of the project. If the user is willing to move in case of uploading some other dataset, then this function link is useful for the user.

Undo



If the user has done different operations, but one point of time they may think to delete an action which may finished already. At that time, the user is provided with an option 'Undo' link. If the link is clicked one time the recent action will be deleted. At a time one action can be removed and a counter is also displayed along with the link.