

Data Collection and Preprocessing Phase

Date	15 July 2024
Team ID	740713
Project Title	Genetic Classification Of An Individual By Using Machine Learning
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description
Data Overview	Basic statistics, dimensions, and structure of the data.
Univariate Analysis	Exploration of individual variables (mean, median, mode, etc.).
Bivariate Analysis	Relationships between two variables (correlation, scatter plots).
Multivariate Analysis	Patterns and relationships involving multiple variables.
Outliers and Anomalies	Identification and treatment of outliers.
Data Preprocessing Code Screenshots	
Loading Data	<pre>dataset=pd.read_csv('/content/clinvar_conflicting.csv')</pre>

Handling Missing Data	<pre>] missing_values = dataset.isnull().sum() print(missing_values)</pre>
Data Transformation	<pre>IQR=[] IQR.append(dataset["X1"].quantile(0.75)-dataset["X1"].quantile(0.25)) IQR [13.0] UPPER=[] IQR.append(dataset["X1"].quantile(0.75)+(1.5)*IQR[0]) IQR [13.0, 25.906976744186046]</pre>
Feature Engineering	<pre>np.where(dataset["X1"]>155474899.66410196,155474899.66410196,np.where(dataset["X1"]<-165339573.33580733,-165339573.33580733,dataset["X1"])) array([-8.59302326, -8.59302326, -8.59302326, ..., -9.59302326, -9.59302326, -9.59302326])</pre>
Save Processed Data	<pre>from sklearn.ensemble import RandomForestClassifier import pickle # Save the model # Save the model with open('Genetic-classification.pkl', 'wb') as f: pickle.dump(rfc_model, f) # Change 'model' to 'rfc_model'</pre>