

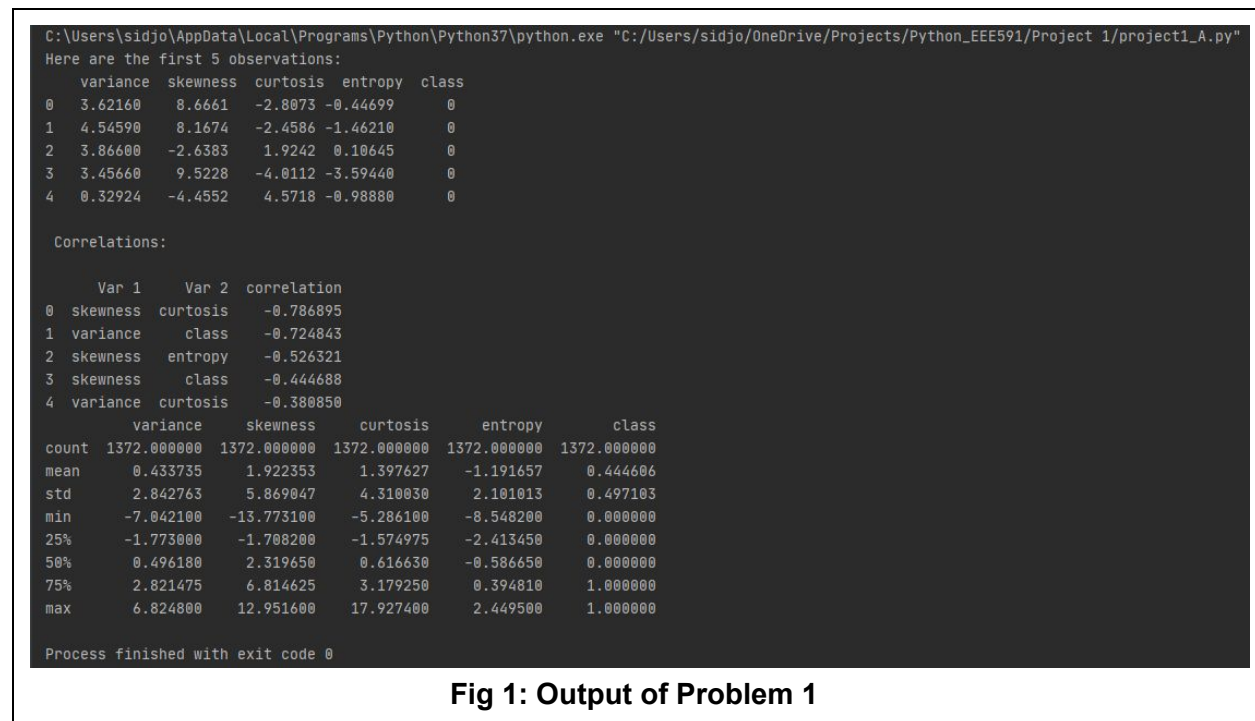
EEE 591 : Project 1
Name: Siddharth Nilesh Joshi
ASU Id: 1217923356

Project Report

Problem 1

Objective: The problem 1 of the project deals with a statistical analysis of the each variable in the provided csv file. A conclusive report was generated based on the correlations of each of the variables: variance, skewness, curtosis, entropy and class.

This analysis is presented in the correlation plot (figure 2) and variable pair-plot (figure 3). Substantial conclusions are derived from the Correlation Plot as discussed further.



The figure above displays the first 5 observations, correlations and statistical analysis of the data in three distinct tables.

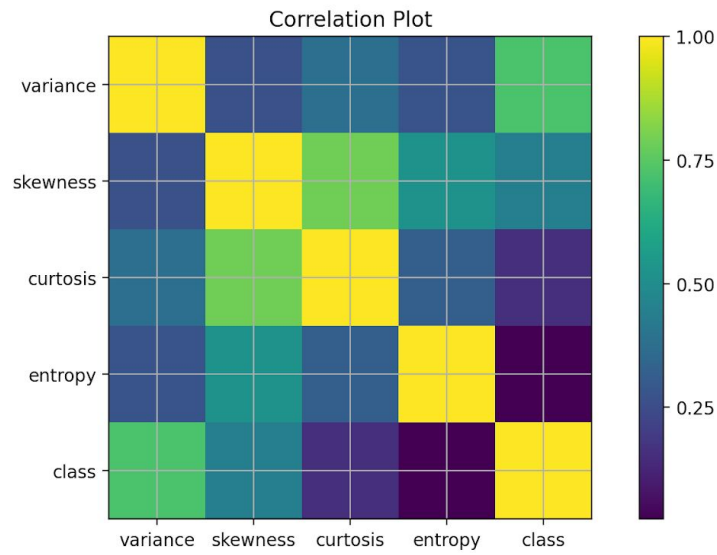


Fig 2: Correlation Plot between variables

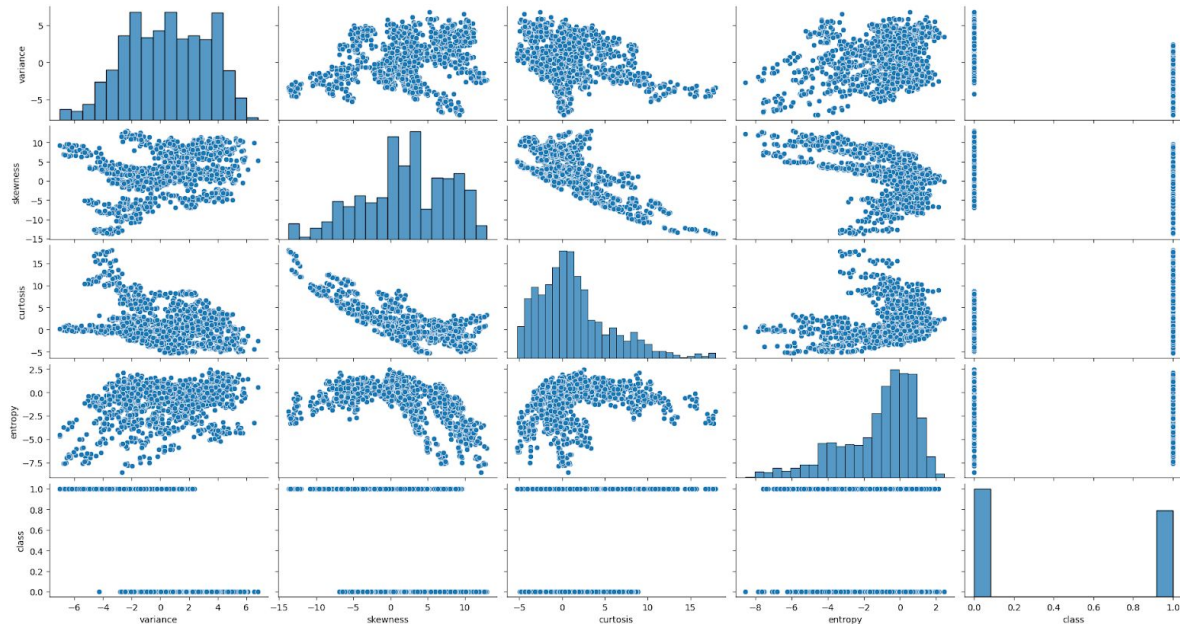


Fig 3: Variable Pair-Plot

As observed in the correlation plot, the colorbar on the right is the correlation measurement between variables, the highest being 1.00 along the diagonal (which is obvious). Observing in decreasing order, the optimal correlation-pairs are as follows: kurtosis-skewness > variance-class > entropy-skewness > class-skewness > kurtosis-variance > entropy-kurtosis > entropy-variance > skewness-variance > class-kurtosis > class-entropy. Hence the highest

correlation is observed in case of curtosis-skewness, and the lowest in class-entropy. This proves that using variable curtosis, skewness (with high correlation) would result in an accurately trained model which can predict accurate test results. The class-entropy variables (with low correlation) would lead to a decrease in the model accuracy, with far inaccurate results.

Problem 2:

Objective: The following problem deals with the training and testing of models based on numerous machine learning techniques namely: perceptron, logistic regression, support vector machines, decision tree learning, random forests and K-nearest neighbors.

Algorithm	Incorrectly Classified Samples	Accuracy (*100 %)	Combined incorrectly classified samples	Combined (train & test) accuracy (*100 %)
Perceptron	8	0.9806	28	0.9796
Logistic Regression	5	0.9879	12	0.9913
Support Vector machine	6	0.9854	17	0.9876
Decision Tree Learning	14	0.9660	20	0.9854
Random Forests	4	0.9903	20	0.9956
K-nearest neighbors	8	0.9806	23	0.9832

The table above elicits the incorrect classification during training of the model as well as along with test data, accuracy in testing and combined accuracy in testing and training. As is it clearly observed, a random forests algorithm proves to be most accurate in predicting if the bill is counterfeit (with 4 mis-classifications). For each of the algorithms, the data provided was standardized through the “standardize” function for testing and training. Random forest hereby yields great results while avoiding overfitting of data. It is hence, highly suitable for classification and regression problems.