

Name: Tan Nguyen

CS 422

INTRODUCTION TO MACHINE LEARNING

FALL 2023

ASSIGNMENT 3

Machine Learning Report: Credit Card Approval Prediction Using Logistic Regression and Naïve Bayes Models

Dataset Source:

The dataset from [Kaggle](#), titled "Credit Card Approval," consists of credit card application data submitted to a commercial bank. It's a classic dataset used in machine learning to build classification models that predict whether an application will be approved or denied based on various attributes.

| Field name | Description | Type |
|----------------|-----------------------------------|--------|
| Gender | Customer's gender | Cat |
| Age | Customer's age as of cut off date | Float |
| Debt | Amount of debt balance | Float |
| Married | Marital status | Cat |
| BankCustomer | Customer category | Cat |
| EducationLevel | Customer's education category | Cat |
| Ethnicity | Customer's ethnicity | Cat |
| YearsEmployed | Customer's years of employment | Float |
| PriorDefault | If default before | Bool |
| Employed | If currently employed | Bool |
| CreditScore | Customer's credit score | Int |
| DriverLicense | If customer has driving license | Bool |
| Citizen | Customer's citizenship | Cat |
| ZipCode | Primary's zip code | String |
| Income | Monthly income | Int |
| ApprovalStatus | If approved for credit card | Bool |

Data Processing:

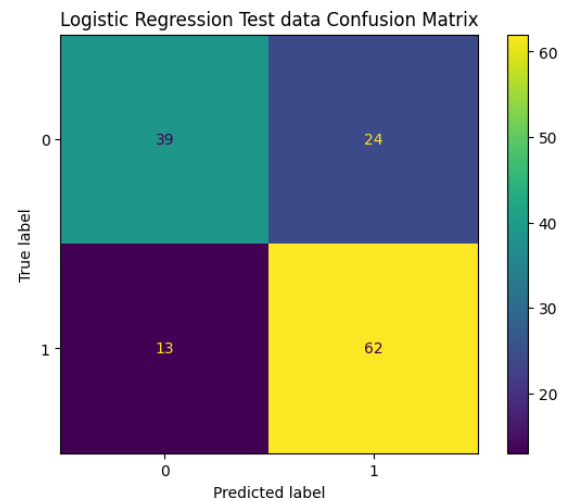
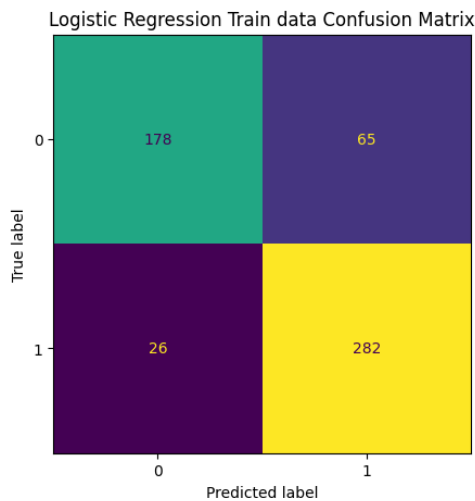
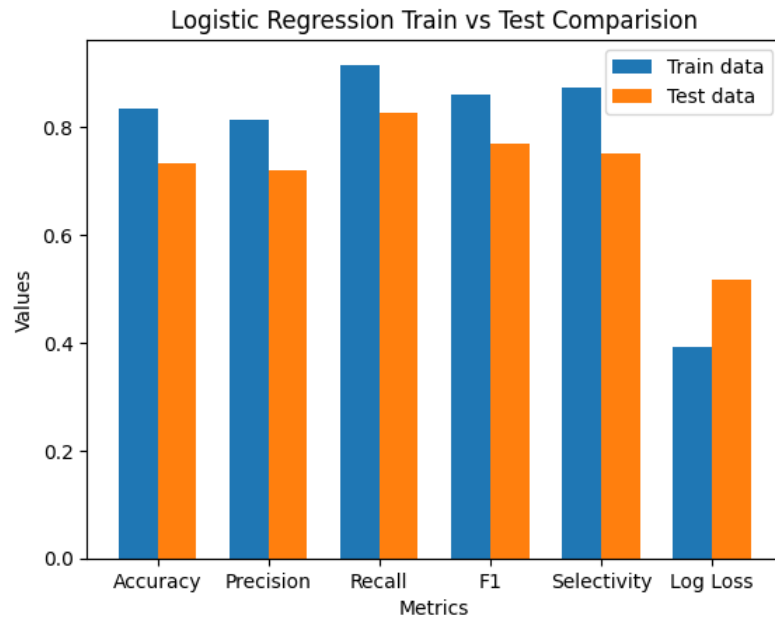
- Transform data set into pd.DataFrame
- As provided dataset is missing column headers, need to add in
- Encode all string-type features using `LabelEncoder()`
- The dataset was split into an 80/20 ratio for training and testing, respectively.

Parameter Vector:

W = [1.76183030e-01 2.19719163e-01 6.41410984e-04 -2.22811144e-02
6.36358240e-01 5.80457961e-01 -1.01392281e-01 -1.12581542e-02
-2.45389723e-01 -9.84176645e-01 -1.18861852e-01 -2.97046552e-01
6.33825132e-02 -7.63259961e-02 3.17690092e-03 -4.87591254e-04]

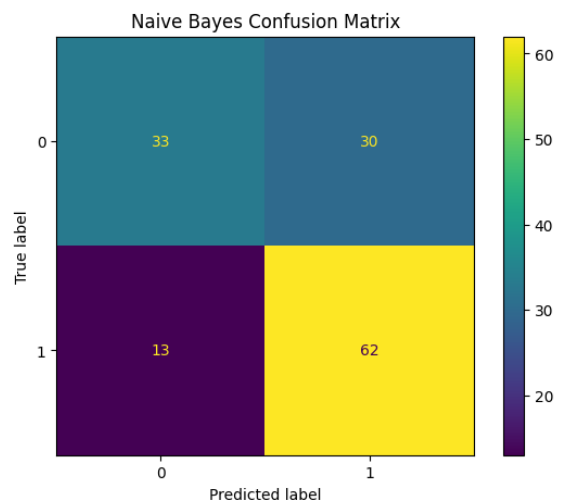
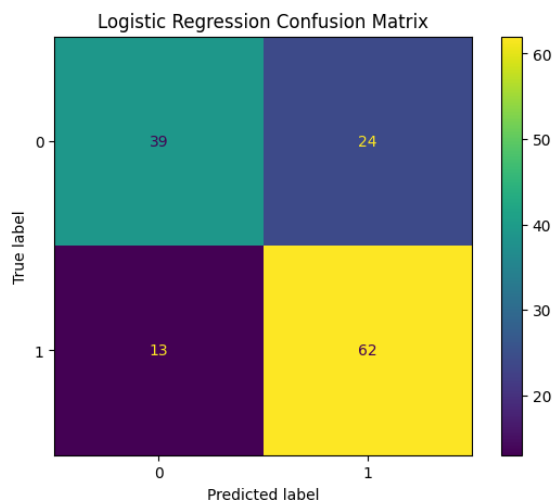
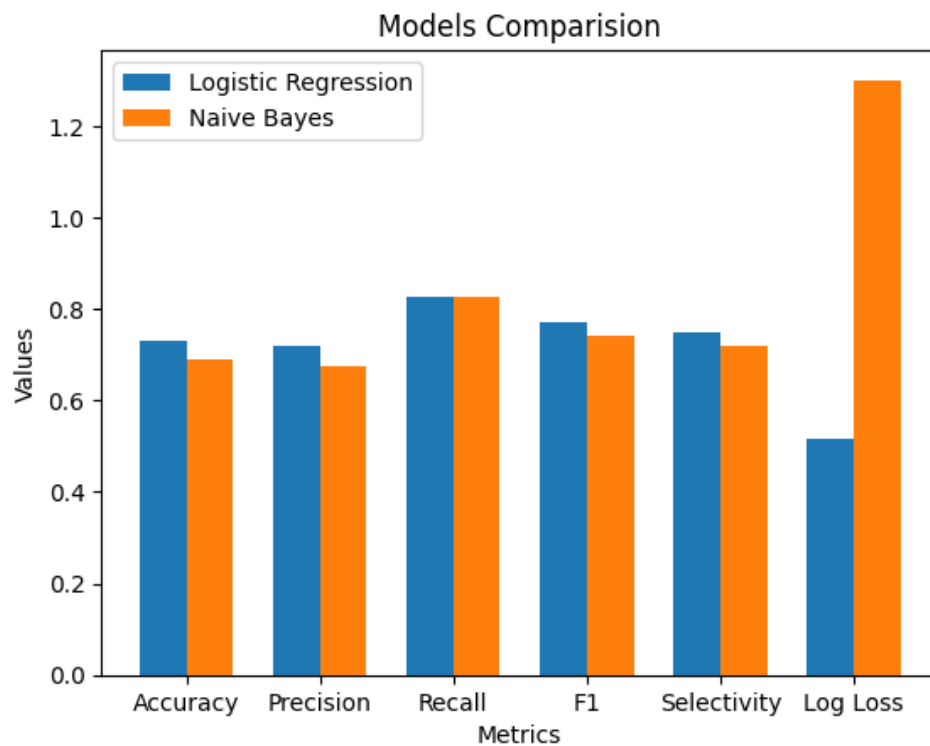
Intercept: 1.76183030e-01

Evaluation – Train vs Test data:



The Logistic Regression model demonstrates decent performance on both the training and testing datasets. However, it seems to be slightly overfitting the training data, as indicated by higher performance metrics for the training set compared to the test set. The false positive rate in the training data also reinforces this observation. Measures might be needed to mitigate this overfitting, such as regularization or further feature engineering.

Evaluation – Logistic Regression vs Naïve Bayes Gaussian:



Logistic regression slightly outperforms GaussianNB in the most of the metrics, with more consistent error rate, as it has lower log loss. The potential explanation would be the underlying assumptions of each model. Logistic regression might be better suited for this dataset as it does not assume features to be conditionally independent within one class as GaussianNB does.