

Grayson Bass
ID Number: 801074960
Homework 2
Github: <https://github.com/gbass2/IntroToML>

Problem 1:

The data was split into 80% training and 20% validation. A random state of 42 was used to randomize the splitting of the data so that the data is split more evenly between the training and validation. This can be found in section In[2]. This helps prevent the training set or validation set from having more of one outcome. Standardization was used for the input scaling found in section In[3]. The LogisticRegression function apart of the sklearn library was used to create the logistic regression classifier. The metrics for the model are below and can be found in section In[4]:

Accuracy: 0.7532467532467533
Precision: 0.6491228070175439
Recall: 0.6727272727272727

Approximately 25% of the predicted data was incorrect predictions while approximately 75% were correct predictions. A precision of 0.649 means that when we predict that a patient has diabetes, we are correct 64.9% of the time. A recall of 0.672 means that we predicted 67.2% of the actual positive cases as true positives and 32.8% as false negatives. Figure 1 which is in section In[5] shows the confusion matrix using the linear regression model. Since the dataset has many more outcomes of 0 than 1, there are more true negative predictions on the validation set.

Problem 2:

Problem 2 was completed the same way as problem 1 except a Gaussian Naïve Bayes classifier from sklearn was used instead of the logistic regression classifier. The same split data was used and standardization was also used and the code can be found in section In[6]. The metrics for the model are below and can be found in section In[7]:

Accuracy: 0.7662337662337663
Precision: 0.6610169491525424
Recall: 0.7090909090909091

The accuracy was slightly improved when using the Naïve classifier. The recall and precision also increased compared to the logistic classification model. This means the model was able to predict slightly more true positives and slightly less false negatives. Based off of that, the Gaussian Naïve Bayes performed slightly better. This can be seen in the confusion matrix in figure 2 in section In[8].

Problem 3:

K-Fold cross validation was used with a logistic regression model and achieved by using the built-in libraries apart of sklearn. The data was not split into a training set and validation set in section In[9] since cross validation does this each iteration. The KFold function was used to set the number of folds used and cross_validate was used for the cross validation. The cross_validate function return the accuracy, precision, and recall for each iteration, so the mean was taken for each. The functions are

being used in section In[9]. The random state was also set to 42 like in problem 1. Both 5 and 10 folds were used. The metrics for both 5 and 10 folds are below and can be found in section In[10]:

```
Accuracy (K=5): 0.7682454800101859
Precision (K=5): 0.7155969191270859
Recall (K=5): 0.5726989972369619
Accuracy (K=10): 0.7707621326042379
Precision (K=10): 0.7173983781918565
Recall (K=10): 0.5800508774379743
```

Using 10 folds compared to 5 folds did improve the model, but the difference between 5 and 10 folds with the logistic regression model is less than 1%. When we have more folds like 10, we are increasing the number of iterations which means the validation set for the kth fold has less data, but this also means the training set for that fold has more data. Increasing the number of folds does not always improve the model but between 5 and 10 it slightly did. When comparing the results with problem 1 using the 10 folds, there was a little over a 2.15 % increase in accuracy, 7.8% increase in the precision, and a 9.36% decrease in the recall. Using K-Fold cross validation helped improve the metrics compared to just using logistic regression.

Problem 4:

The same approach was used as in problem 3 but the model was changed. K-Fold cross validation was used with a Gaussian Naïve Bayes model and achieved by using the built-in libraries apart of sklearn. Like problem 3, the data was not split into a training set and validation set in section In[11] since cross validation does this each iteration. The KFold function was used to set the number of folds used and cross_validate was used for the cross validation. The cross_validate function return the accuracy, precision, and recall for each iteration, so the mean was taken for each. The functions are being used in section In[11]. Both 5 and 10 folds were used. The metrics for both 5 and 10 folds are below and can be found in section In[12]:

```
Accuracy (K=5): 0.7539258127493421
Precision (K=5): 0.6645294767870302
Recall (K=5): 0.6011292482940126
Accuracy (K=10): 0.7512303485987697
Precision (K=10): 0.6537815101446814
Recall (K=10): 0.5938748312619281
```

There is not a significant improvement when using K-Fold cross validation when compared to just using the Naïve Bayes Classifier prediction. The recall is lower but this could be because a random state of 42 was used but the sklearn Gaussian Naïve Bayes classifier used in problem 2, section In[6], does not use random state. This is due to K-Fold cross validation being meaningless to Gaussian Naïve Bayes. Each iteration will be a separate model so if 10 folds are used then 10 independent models are created. This happens because the features for Naïve Bayes are independent of one another. So, it does not make sense to use K-Fold cross validation with Gaussian Naïve Bayes.