


```
[20]: classifiers = [
# MultinomialNB(), # doesn't work
# KNeighborsClassifier(), # works
# SVC(kernel='rbf', C=0.001, probability=True), # took a long time... need to refresh memory
# SVC(kernel='linear'), # took a long time... need to refresh memory
# SVC(probability=True, nu=0.1), # took a long time... need to refresh memory
DecisionTreeClassifier(), # works
RandomForestClassifier(), # works
AdaBoostClassifier(), # works
GradientBoostingClassifier(), # works
GaussianNB(), # works
BernoulliNB(), # works
MLPClassifier(), # works
MLPClassifier(hidden_layer_sizes=[100, 100]), # works
LinearDiscriminantAnalysis(), # works
LogisticRegression(), # works
QuadraticDiscriminantAnalysis(), # works
]

log_cols=["Classifier", "Accuracy", "F1 Score", "ROC", "Precision", "Recall", "Log Loss", "Sensitivity", "Specificity", "Speci
log = pd.DataFrame(columns=log_cols)

for clf in classifiers:
    clf.fit(X_train, y_train)
    y_pred = clf.predict(X)
    name = clf.__class__.__name__

    print("="*60)
    print(name)
    print("****Results****")

    train_predictions = clf.predict(X_test)
    acc = accuracy_score(y_test, train_predictions)
    acc = acc_round(3)
    print("Accuracy: {:.4%}".format(acc))

    # coef_scores = X_scaled
    # coef_scores = clf.coef_
    # print(coef_scores)

    fbeta = fbeta_score(y_test, train_predictions, beta=1)
    fbeta = fbeta_round(3)
    print("F1 Score: {:.4%}".format(fbeta))

    roc = roc_auc_score(y_test, train_predictions)
    roc = roc_round(3)
    print("AUC (ROC) Score: {:.4%}".format(roc))

    precision = precision_score(y_test, train_predictions, average='binary')
    precision = precision_round(3)
    print("Precision Score: {:.4%}".format(precision))

    recall = recall_score(y_test, train_predictions)
    recall = recall_round(3)
    print("Recall Score: {:.4%}".format(recall))

    TP = confusion_matrix(y_test, train_predictions).ravel()[0]
    FN = confusion_matrix(y_test, train_predictions).ravel()[1]
    sensitivity = TP / (TP + FN)
    sensitivity = sensitivity_round(3)
    print("Sensitivity Score: {:.4%}".format(sensitivity))

    TN = confusion_matrix(y_test, train_predictions).ravel()[3]
    FP = confusion_matrix(y_test, train_predictions).ravel()[2]
    specificity = TN / (TN + FP)
    specificity = specificity_round(3)
    print("Specificity Score: {:.4%}".format(specificity))

    true_positives = confusion_matrix(y_test, train_predictions).ravel()[0]

    print(f"\nClassic train score: {np.round(clf.score(X_train, y_train),4)*100}")
    print(f"\nClassic test score: {np.round(clf.score(X_test, y_test),4)*100}")
    confusion_matrix = confusion_matrix(y_test, train_predictions)
    print(f"\nConfusion matrix: \n{confusion_matrix}\n")

    # train_predictions = clf.predict_proba(X_test)
    # ll = log_loss(y_test, train_predictions)
    # ll = ll_round(3)
    # print(f"Log Loss: {ll}")

    log_entry = pd.DataFrame([name, acc, f100, fbeta, roc, precision, recall, ll, sensitivity, specificity, true
log = pd.concat([log, log_entry])
# type(log)

# print("="*60)
# type(coef_scores)
# print(index)
# print(log_entry)
# type(log)

=====
KNeighborsClassifier
****Results****
Accuracy: 87.0000%
F1 Score: 0.929
AUC (ROC) Score: 0.568
Precision Score: 0.9
Recall Score: 0.96
Sensitivity Score: 0.177
Specificity Score: 0.946

Classic train score: 0.9107
Classic test score: 0.8699
Confusion matrix:
[[ 276 1287]
 [ 478 11523]]

=====
DecisionTreeClassifier
****Results****
Accuracy: 81.8000%
F1 Score: 0.897
AUC (ROC) Score: 0.591
Recall Score: 0.853
Sensitivity Score: 0.248
Specificity Score: 0.893

Classic train score: 0.9987
Classic test score: 0.8185
Confusion matrix:
[[ 1286 1176]
 [ 386 10715]]

=====
RandomForestClassifier
****Results****
Accuracy: 87.8000%
F1 Score: 0.933
AUC (ROC) Score: 0.576
Precision Score: 0.901
Recall Score: 0.968
Sensitivity Score: 0.184
Specificity Score: 0.968

Classic train score: 0.9986
Classic test score: 0.8775
Confusion matrix:
[[ 288 1279]
 [ 386 11615]]

=====
AdaBoostClassifier
****Results****
Accuracy: 89.3000%
F1 Score: 0.942
AUC (ROC) Score: 0.578
Precision Score: 0.901
Recall Score: 0.988
Sensitivity Score: 0.168
Specificity Score: 0.988

Classic train score: 0.8912
Classic test score: 0.893
Confusion matrix:
[[ 262 1301]
 [ 150 11831]]

=====
GradientBoostingClassifier
****Results****
Accuracy: 89.4000%
F1 Score: 0.943
AUC (ROC) Score: 0.585
Precision Score: 0.903
Recall Score: 0.986
Sensitivity Score: 0.184
Specificity Score: 0.986

Classic train score: 0.8945
Classic test score: 0.8937
Confusion matrix:
[[ 288 1275]
 [ 167 11834]]

=====
GaussianNB
****Results****
Accuracy: 77.0000%
F1 Score: 0.863
AUC (ROC) Score: 0.623
Precision Score: 0.917
Recall Score: 0.815
Sensitivity Score: 0.431
Specificity Score: 0.815

Classic train score: 0.7708
Classic test score: 0.7705
Confusion matrix:
[[ 373 890]
 [2223 9778]]

=====
BernoulliNB
****Results****
Accuracy: 89.4000%
F1 Score: 0.943
AUC (ROC) Score: 0.59
Precision Score: 0.904
Recall Score: 0.985
Sensitivity Score: 0.196
Specificity Score: 0.985

Classic train score: 0.8912
Classic test score: 0.8942
Confusion matrix:
[[ 306 1257]
 [ 178 11823]]

=====
C:\Users\veagina\AppData\Local\Packages\PythonSoftwareFoundation.Python.3.10_qbz5n2kfra8p0\LocalCache\local-pac
kages\Python310\site-packages\sklearn\neural_network\multilayer_perceptron.py:692: ConvergenceWarning: Stochastic
optimizers: Maximum iterations (200) reached and the optimization hasn't converged yet.
warnings.warn(

=====
MLPClassifier
****Results****
Accuracy: 88.9000%
F1 Score: 0.94
AUC (ROC) Score: 0.589
Precision Score: 0.904
Recall Score: 0.979
Sensitivity Score: 0.198
Specificity Score: 0.979

Classic train score: 0.8981
Classic test score: 0.8693
Confusion matrix:
[[ 310 1253]
 [ 249 11752]]

=====
C:\Users\veagina\AppData\Local\Packages\PythonSoftwareFoundation.Python.3.10_qbz5n2kfra8p0\LocalCache\local-pac
kages\Python310\site-packages\sklearn\neural_network\multilayer_perceptron.py:692: ConvergenceWarning: Stochastic
optimizers: Maximum iterations (200) reached and the optimization hasn't converged yet.
warnings.warn(

=====
MLPClassifier
****Results****
Accuracy: 86.9000%
F1 Score: 0.928
AUC (ROC) Score: 0.58
Precision Score: 0.902
Recall Score: 0.956
Sensitivity Score: 0.204
Specificity Score: 0.956

Classic train score: 0.9171
Classic test score: 0.8691
Confusion matrix:
[[ 319 1244]
 [ 532 11469]]

=====
LinearDiscriminantAnalysis
****Results****
Accuracy: 89.5000%
F1 Score: 0.943
AUC (ROC) Score: 0.589
Precision Score: 0.904
Recall Score: 0.986
Sensitivity Score: 0.192
Specificity Score: 0.986

Classic train score: 0.8919
Classic test score: 0.8946
Confusion matrix:
[[ 200 1263]
 [ 167 11834]]

=====
LogisticRegression
****Results****
Accuracy: 89.5000%
F1 Score: 0.943
AUC (ROC) Score: 0.587
Precision Score: 0.903
Recall Score: 0.987
Sensitivity Score: 0.186
Specificity Score: 0.987

Classic train score: 0.8917
Classic test score: 0.8946
Confusion matrix:
[[ 251 11844]]

=====
QuadraticDiscriminantAnalysis
****Results****
Accuracy: 71.5000%
F1 Score: 0.822
AUC (ROC) Score: 0.615
Precision Score: 0.917
Recall Score: 0.745
Sensitivity Score: 0.486
Specificity Score: 0.745

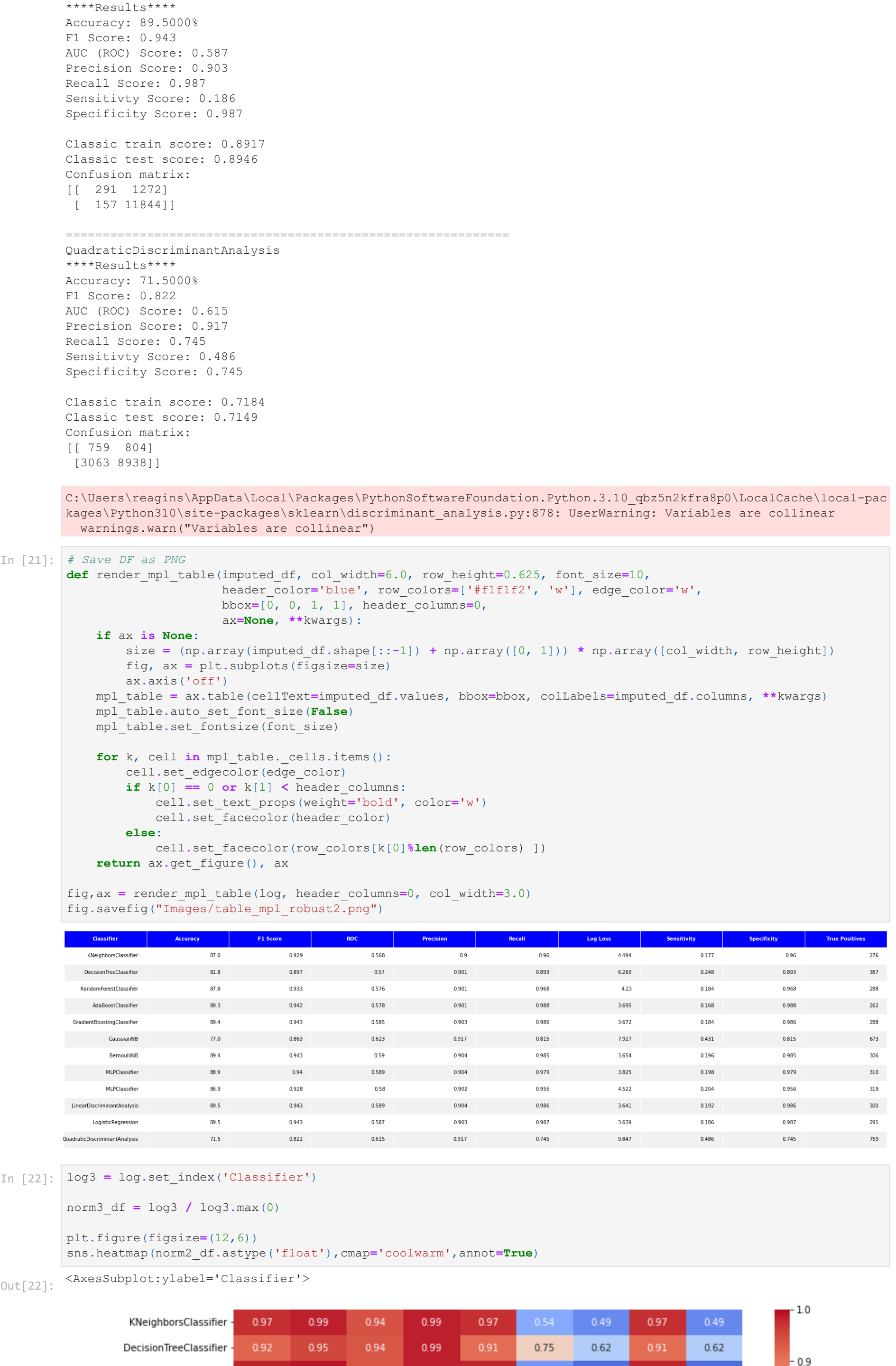
Classic train score: 0.7184
Classic test score: 0.7149
Confusion matrix:
[[ 759 804]
 [3063 8938]]

=====
C:\Users\veagina\AppData\Local\Packages\PythonSoftwareFoundation.Python.3.10_qbz5n2kfra8p0\LocalCache\local-pac
kages\Python310\site-packages\sklearn\discriminant_analysis.py:878: UserWarning: Variables are collinear
warnings.warn("Variables are collinear")
```

```
In [21]: # Save DF as PNG
def render_mpl_table(imputed_df, col_width=6.0, row_height=0.625, font_size=10,
fig, ax = plt.subplots(figsize=size)
bbox=[0, 0, 1, 1], header_columns=0,
ax=None, **kwargs):
    size = np.array(imputed_df.shape[:-1]) + np.array([0, 1]) * np.array([col_width, row_height])
    fig, ax = plt.subplots(figsize=size)
    ax.axis('tight')
    mpl_table = ax.table(cellText=imputed_df.values, bbox=bbox, colLabels=imputed_df.columns, **kwargs)
    mpl_table.auto_set_font_size(False)
    mpl_table.set_fontsize(font_size)

    for k, cell in mpl_table._cells.items():
        cell.set_edgewidth(edge_color)
        if k[0] == 0 or k[1] < header_columns:
            cell.set_text_props(weight='bold', color='w')
            cell.set_facecolor(row_colors[k[0]]*len(row_colors))
        else:
            cell.set_facecolor(row_colors[k[0]]*len(row_colors))
    return ax, ax.get_figure(), ax

fig, ax = render_mpl_table(log, header_columns=0, col_width=3.0)
fig.savefig(f'figures/table_mpl_{round2.png}')
```



```
In [22]: log3 = log.set_index('Classifier')
norm3_df = log3 / log3.max(0)
plt.figure(figsize=(12,6))
sns.heatmap(norm3_df,axtype='float',cmap='coolwarm',annot=True)

<AxesSubplot:ylabel='Classifier'>
```



```
In [24]: big_log_df = (log1 + log2 + log3) / 3
big_log_df
```

```
Out[24]:
```

```
Classifier
AdaBoostClassifier      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
BernoulliNB             NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
DecisionTreeClassifier  81.933333  0.897333  0.575333  0.901  0.893667  6.233333  0.257333  0.893667  406.666667
GaussianNB              NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
GradientBoostingClassifier  NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
KNeighborsClassifier     87.0  0.928667  0.576333  0.900333  0.959  4.491  0.194  0.959  306.333333
LinearDiscriminantAnalysis  NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
LogisticRegression      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
MLPClassifier            NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
MLPClassifier            NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
MLPClassifier            NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
QuadraticDiscriminantAnalysis  NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
RandomForestClassifier  NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN      NaN
```

```
In [ ]:
```

```
In [25]: plt.figure(figsize=(24,3))
#sns.barplot(X.columns,logreg_coef_[0])
#plt.xticks(rotation=60)
plt.title("Extracting the Feature Importance")
```

Further discussion for the group

- What further refinements to the dataset should we make as part of the EDA / cleanup?
 - Removing the pd.ys variable, for example
 - Dropping outliers
- How might the use of other classification algorithms and scalers affect the final predictions?
 - Algorithms like LogisticRegression, DecisionTree, RandomForest, KNeighbors, NaiveBayes, neural net, etc. -
 - Scalers like StandardScaler, MinMaxScaler, RobustScaler
- Playing with parameters, pipelines, gridsearches to maximize True Negatives and minimize False Negatives
 - That is, maximize deposit = 1 correct predictions and reducing deposit = 0 wrong predictions
 - Even if that means accidentally overpredicting the number of true deposits, better to try a bad path than miss a potential business opportunity
- Extending this to other predictions
 - e.g. predicting the 'default' variable, or some other categorical values
 - e.g. predicting a range for continuous values based on categorical values
- Best ways to impute missing data?

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```