





```
Accuracy: 88.2000%
F1 Score: 0.935
AUC (ROC) Score: 0.614
Precision Score: 0.931
Recall Score: 0.961
Sensitivity Score: 0.267
Specificity Score: 0.961

Classic train score: 0.9193
Classic test score: 0.8817
Confusion matrix:
[[ 416 1140]
 [ 464 11544]]

=====
DecisionTreeClassifier
***Results***
Accuracy: 83.4000%
F1 Score: 0.906
AUC (ROC) Score: 0.602
Precision Score: 0.909
Recall Score: 0.903
Sensitivity Score: 0.3
Specificity Score: 0.903

Classic train score: 1.0
Classic test score: 0.8338
Confusion matrix:
[[ 457 1089]
 [ 1166 10842]]

=====
RandomForestClassifier
***Results***
Accuracy: 89.4000%
F1 Score: 0.943
AUC (ROC) Score: 0.602
Precision Score: 0.907
Recall Score: 0.981
Sensitivity Score: 0.222
Specificity Score: 0.981

Classic train score: 1.0
Classic test score: 0.8944
Confusion matrix:
[[ 346 1210]
 [ 1233 11793]]

=====
AdaBoostClassifier
***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: 0.8914
Classic test score: 0.8941
Confusion matrix:
[[ 286 1270]
 [ 166 11842]]

=====
GradientBoostingClassifier
***Results***
Accuracy: 89.7000%
F1 Score: 0.944
AUC (ROC) Score: 0.595
Precision Score: 0.905
Recall Score: 0.986
Sensitivity Score: 0.204
Specificity Score: 0.986

Classic train score: 0.8958
Classic test score: 0.8966
Confusion matrix:
[[ 317 1239]
 [ 164 11844]]

=====
GaussianNB
***Results***
Accuracy: 82.4000%
F1 Score: 0.897
AUC (ROC) Score: 0.685
Precision Score: 0.865
Recall Score: 0.905
Sensitivity Score: 0.505
Specificity Score: 0.865

Classic train score: 0.82
Classic test score: 0.8241
Confusion matrix:
[[ 786 1079]
 [ 1616 10932]]

=====
BernoulliNB
***Results***
Accuracy: 85.4000%
F1 Score: 0.917
AUC (ROC) Score: 0.666
Precision Score: 0.924
Recall Score: 0.931
Sensitivity Score: 0.422
Specificity Score: 0.91

Classic train score: 0.8543
Classic test score: 0.8544
Confusion matrix:
[[ 657 899]
 [ 1076 10932]]

C:\Users\resqina\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")

=====
MLPClassifier
***Results***
Accuracy: 87.9000%
F1 Score: 0.933
AUC (ROC) Score: 0.615
Precision Score: 0.931
Recall Score: 0.957
Sensitivity Score: 0.272
Specificity Score: 0.957

Classic train score: 0.9475
Classic test score: 0.8788
Confusion matrix:
[[ 424 1132]
 [ 512 11496]]

=====
MLPClassifier
***Results***
Accuracy: 85.5000%
F1 Score: 0.919
AUC (ROC) Score: 0.619
Precision Score: 0.912
Recall Score: 0.925
Sensitivity Score: 0.314
Specificity Score: 0.925

Classic train score: 0.9857
Classic test score: 0.8548
Confusion matrix:
[[ 488 1068]
 [ 901 11107]]

=====
LinearDiscriminantAnalysis
***Results***
Accuracy: 89.1000%
F1 Score: 0.94
AUC (ROC) Score: 0.623
Precision Score: 0.927
Sensitivity Score: 0.276
Specificity Score: 0.97

Classic train score: 0.8891
Classic test score: 0.8907
Confusion matrix:
[[ 430 1126]
 [ 356 11652]]

=====
LogisticRegression
***Results***
Accuracy: 89.5000%
F1 Score: 0.944
AUC (ROC) Score: 0.583
Precision Score: 0.903
Recall Score: 0.988
Sensitivity Score: 0.179
Specificity Score: 0.988

Classic train score: 0.8915
Classic test score: 0.8952
Confusion matrix:
[[ 278 1278]
 [ 143 11865]]

=====
QuadraticDiscriminantAnalysis
***Results***
Accuracy: 81.4000%
F1 Score: 0.89
AUC (ROC) Score: 0.681
Precision Score: 0.931
Recall Score: 0.853

C:\Users\resqina\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")

Classic train score: 0.833
Classic test score: 0.8124
Classic test score: 0.8136
Confusion matrix:
[[ 792 764]
 [ 1765 10243]]

In [18]:
# Save DF as PNG
def render_mpl_table(inputed_df, col_width=6.0, row_height=0.625, font_size=10,
                    header_color='#404066', row_color='#f1f1f2', 'w'), edge_color='w',
                    ax=None, **kwargs):
    if ax is None:
        size = np.array(inputed_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=inputed_df.values, bbox=bbox, colLabels=inputed_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_edgecolor(edge_color)
            if k[0] == 0 or k[1] < header_columns:
                cell.set_text_props(weight='bold', color='w')
                cell.set_facecolor(header_color)
            else:
                cell.set_facecolor(row_colors[k[0]]%len(row_colors))
        return ax.get_figure(), ax

fig,ax = render_mpl_table(log, header_columns=0, col_width=3.0)
fig.savefig("Images/table_mpl_standard1.png")

Classifier
```

Classifier	Accuracy	F1 Score	ROC	Precision	Recall	AUC	Log Loss	Sensitivity	Specificity	True Positive Count
DecisionTreeClassifier	83.4	0.906	0.602	0.903	0.903	0.602	4.084	0.301	0.903	457
RandomForestClassifier	89.4	0.943	0.602	0.907	0.981	0.602	3.702	0.222	0.981	346
AdaBoostClassifier	89.4	0.943	0.585	0.903	0.986	0.585	3.847	0.184	0.986	286
GradientBoostingClassifier	89.7	0.944	0.595	0.905	0.986	0.595	3.573	0.204	0.986	317
GaussianNB	82.4	0.897	0.685	0.865	0.905	0.685	4.008	0.505	0.865	786
BernoulliNB	85.4	0.917	0.666	0.924	0.931	0.666	3.923	0.422	0.931	657
MLPClassifier	87.9	0.933	0.615	0.931	0.957	0.615	4.126	0.272	0.957	424
MLPClassifier	85.5	0.919	0.619	0.912	0.925	0.619	3.924	0.314	0.925	488
LinearDiscriminantAnalysis	89.1	0.94	0.623	0.927	0.97	0.623	3.982	0.276	0.97	430
LogisticRegression	89.5	0.944	0.583	0.903	0.988	0.583	3.618	0.179	0.988	278
QuadraticDiscriminantAnalysis	81.4	0.89	0.681	0.931	0.853	0.681	4.44	0.853	0.931	792

log2 = log.set\_index('Classifier')

norm2\_df = log2 / log2.max(0)

sns.heatmap(norm2\_df.astype('float'), cmap='coolwarm')

<AxesSubplot:ylabel='Classifier'>

## RobustScaler

```
In [20]: # StandardScaler, MinMaxScaler, RobustScaler
Scaler = RobustScaler()
FX_scaled = scaler.fit_transform(inputed_df)

X = inputed_df.drop(columns='deposit')
y = inputed_df['deposit']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, shuffle=True)

scaler = RobustScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

```
In [21]: classifiers = [
    # LogisticRegression(), # doesn't work
    KNeighborsClassifier(3), # works
    # SVC(kernel="l2", C=0.001, probability=True), # took a long time... need to refresh memory
    # SVC(kernel="linear", # took a long time... need to refresh memory
    # NuSVC(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", "Speci
log = pd.DataFrame(columns=log_cols)

for clf in classifiers:
    clf.fit(X_train, y_train)
    # pred = clf.predict(X)
    name = clf.__class__.__name__
    print(name)
    print("=====")
    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))

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

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

    # train_predictions = clf.predict_proba(X_test)
    ll = log_loss(y_test, train_predictions)
    ll = ll_round(3)
    # print("Log Loss: {:.4%}".format(ll))

    log_entry = pd.DataFrame([name, acc,100, fbeta, roc, precision, recall, ll, sensitivity, specificity, true
log = pd.concat((log,log_entry))

# print("=="*30)
# type(coef_scores)
# print(index)
# print(log_entry)
# type(log)
```

```
=====
KNeighborsClassifier
***Results***
Accuracy: 88.0000%
F1 Score: 0.935
AUC (ROC) Score: 0.595
Precision Score: 0.906
Recall Score: 0.965
Sensitivity Score: 0.265
Specificity Score: 0.965

Classic train score: 0.9169
Classic test score: 0.8805
Confusion matrix:
[[ 351 1205]
 [ 416 11592]]

=====
DecisionTreeClassifier
***Results***
Accuracy: 83.5000%
F1 Score: 0.906
AUC (ROC) Score: 0.612
Precision Score: 0.911
Recall Score: 0.902
Sensitivity Score: 0.323
Specificity Score: 1.0

Classic train score: 1.0
Classic test score: 0.8353
Confusion matrix:
[[ 223 1233]
 [ 1180 10828]]

=====
RandomForestClassifier
***Results***
Accuracy: 89.3000%
F1 Score: 0.942
AUC (ROC) Score: 0.604
Precision Score: 0.907
Recall Score: 0.979
Sensitivity Score: 0.228
Specificity Score: 0.979

Classic train score: 0.9999
Classic test score: 0.893
Confusion matrix:
[[ 355 1201]
 [ 250 11758]]

=====
AdaBoostClassifier
***Results***
Accuracy: 89.3000%
F1 Score: 0.942
Precision Score: 0.904
Recall Score: 0.986
Sensitivity Score: 0.195
Specificity Score: 0.986

Classic train score: 0.8916
Classic test score: 0.8953
Confusion matrix:
[[ 304 1252]
 [ 168 11840]]

=====
GradientBoostingClassifier
***Results***
Accuracy: 89.6000%
F1 Score: 0.947
Precision Score: 0.906
Recall Score: 0.986
Sensitivity Score: 0.208
Specificity Score: 0.986

Classic train score: 0.8961
Classic test score: 0.8963
Confusion matrix:
[[ 223 1233]
 [ 174 11834]]

=====
GaussianNB
***Results***
Accuracy: 82.3000%
F1 Score: 0.896
AUC (ROC) Score: 0.684
Precision Score: 0.864
Recall Score: 0.903
Sensitivity Score: 0.503
Specificity Score: 0.865

Classic train score: 0.8187
Classic test score: 0.8232
Confusion matrix:
[[ 782 774]
 [ 1624 10384]]

=====
BernoulliNB
***Results***
Accuracy: 87.9000%
F1 Score: 0.933
AUC (ROC) Score: 0.649
Precision Score: 0.948
Recall Score: 0.35
Sensitivity Score: 0.948

Classic train score: 0.8764
Classic test score: 0.8794
Confusion matrix:
[[ 545 1011]
 [ 625 11383]]

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

=====
MLPClassifier
***Results***
Accuracy: 87.6000%
F1 Score: 0.931
AUC (ROC) Score: 0.623
Precision Score: 0.912
Recall Score: 0.951
Sensitivity Score: 0.294
Specificity Score: 0.951

Classic train score: 0.9438
Classic test score: 0.8759
Confusion matrix:
[[ 458 1098]
 [ 585 11423]]

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

=====
MLPClassifier
***Results***
Accuracy: 85.8000%
F1 Score: 0.921
AUC (ROC) Score: 0.613
Precision Score: 0.911
Recall Score: 0.931
Sensitivity Score: 0.294
Specificity Score: 0.931

Classic train score: 0.9933
Classic test score: 0.8584
Confusion matrix:
[[ 458 1098]
 [ 823 11185]]

=====
LinearDiscriminantAnalysis
***Results***
Accuracy: 89.1000%
F1 Score: 0.94
AUC (ROC) Score: 0.624
Precision Score: 0.912
Recall Score: 0.97
Sensitivity Score: 0.278
Specificity Score: 0.97

Classic train score: 0.8889
Classic test score: 0.8905
Confusion matrix:
[[ 433 1123]
 [ 362 11646]]

C:\Users\resqina\AppData\Local\Packages\PythonSoftwareFoundation.Python.3.10_qbz5n2kfra8p0\LocalCache\local-pac
kages\Python310\site-packages\sklearn\linear_model\logistic.py:814: ConvergenceWarning: lbfgs failed to converge
  warnings.warn(
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_1 = check_optimize_result(
C:\Users\resqina\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")

=====
LogisticRegression
***Results***
Accuracy: 89.3000%
F1 Score: 0.943
Precision Score: 0.984
Recall Score: 0.987
Sensitivity Score: 0.181
Specificity Score: 0.987

Classic train score: 0.8918
Classic test score: 0.8949
Confusion matrix:
[[ 282 1274]
 [ 152 11856]]

=====
QuadraticDiscriminantAnalysis
***Results***
Accuracy: 16.2000%
F1 Score: 0.108
AUC (ROC) Score: 0.513
Precision Score: 0.936
Recall Score: 0.057
Sensitivity Score: 0.97
Specificity Score: 0.057

Classic train score: 0.1692
Classic test score: 0.1618
Confusion matrix:
[[ 1508 47]
 [11322 686]]

In [22]: # Save DF as PNG
def render_mpl_table(inputed_df, col_width=6.0, row_height=0.625, font_size=10,
                    header_color='#404066', row_color='#f1f1f2', 'w'), edge_color='w',
                    ax=None, **kwargs):
    if ax is None:
        size = np.array(inputed_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=inputed_df.values, bbox=bbox, colLabels=inputed_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_edgecolor(edge_color)
            if k[0] == 0 or k[1] < header_columns:
                cell.set_text_props(weight='bold', color='w')
                cell.set_facecolor(header_color)
            else:
                cell.set_facecolor(row_colors[k[0]]%len(row_colors))
        return ax.get_figure(), ax

fig,ax = render_mpl_table(log, header_columns=0, col_width=3.0)
fig.savefig("Images/table_mpl_robust1.png")

Classifier
```

Classifier	Accuracy	F1 Score	ROC	Precision	Recall	AUC	Log Loss	Sensitivity	Specificity	True Positive Count
DecisionTreeClassifier	83.5	0.906	0.612	0.911	0.902	0.612	3.889	0.323	0.902	223
RandomForestClassifier	89.3	0.942	0.604	0.907	0.979	0.604	3.605	0.228	0.979	355
AdaBoostClassifier	89.3	0.942	0.595	0.903	0.986	0.595	3.624	0.195	0.986	304
GradientBoostingClassifier	89.6	0.947	0.606	0.906	0.986	0.606	3.498	0.208	0.986	355
GaussianNB	82.3	0.896	0.684	0.864	0.903	0.684	4.126	0.503	0.864	782
BernoulliNB	87.9	0.933	0.649	0.948	0.35	0.649	4.008	0.948	0.35	625
MLPClassifier	87.6	0.931	0.623	0.912	0.951	0.623	4.098	0.294	0.951	458
LinearDiscriminantAnalysis	89.1	0.94	0.624	0.912	0.97	0.624	3.982	0.278	0.97	433
LogisticRegression	89.3	0.943	0.984	0.987	0.987	0.984	3.618	0.181	0.987	282
QuadraticDiscriminantAnalysis	16.2	0.108	0.513	0.936	0.057	0.513	4.89	0.97	0.057	1508

log3 = log.set\_index('Classifier')

norm3\_df = log3 / log3.max(0)

sns.heatmap(norm3\_df.astype('float'), cmap='coolwarm')

<AxesSubplot:ylabel='Classifier'>

```
In [24]: #plt.Figure(figsize=(24,3))
#sns.heatmap(log.columns,logreg_coef_[0])
#plt.vlines(rotation=0)
#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?
  - Renaming the addvs 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, gridssearchs 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 classification
  - e.g. predicting a range for continuous values based on categorical values
- Best ways to impute missing data?

```
In [ ]:
```

```
In [ ]:
```

```
In [ ]:
```