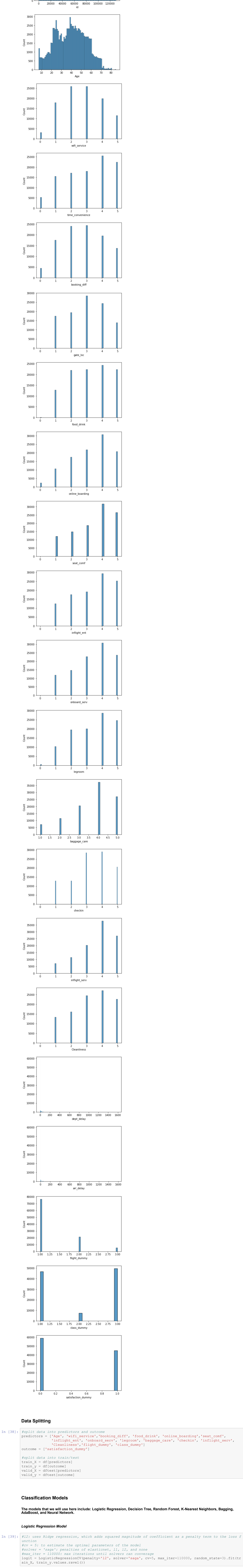
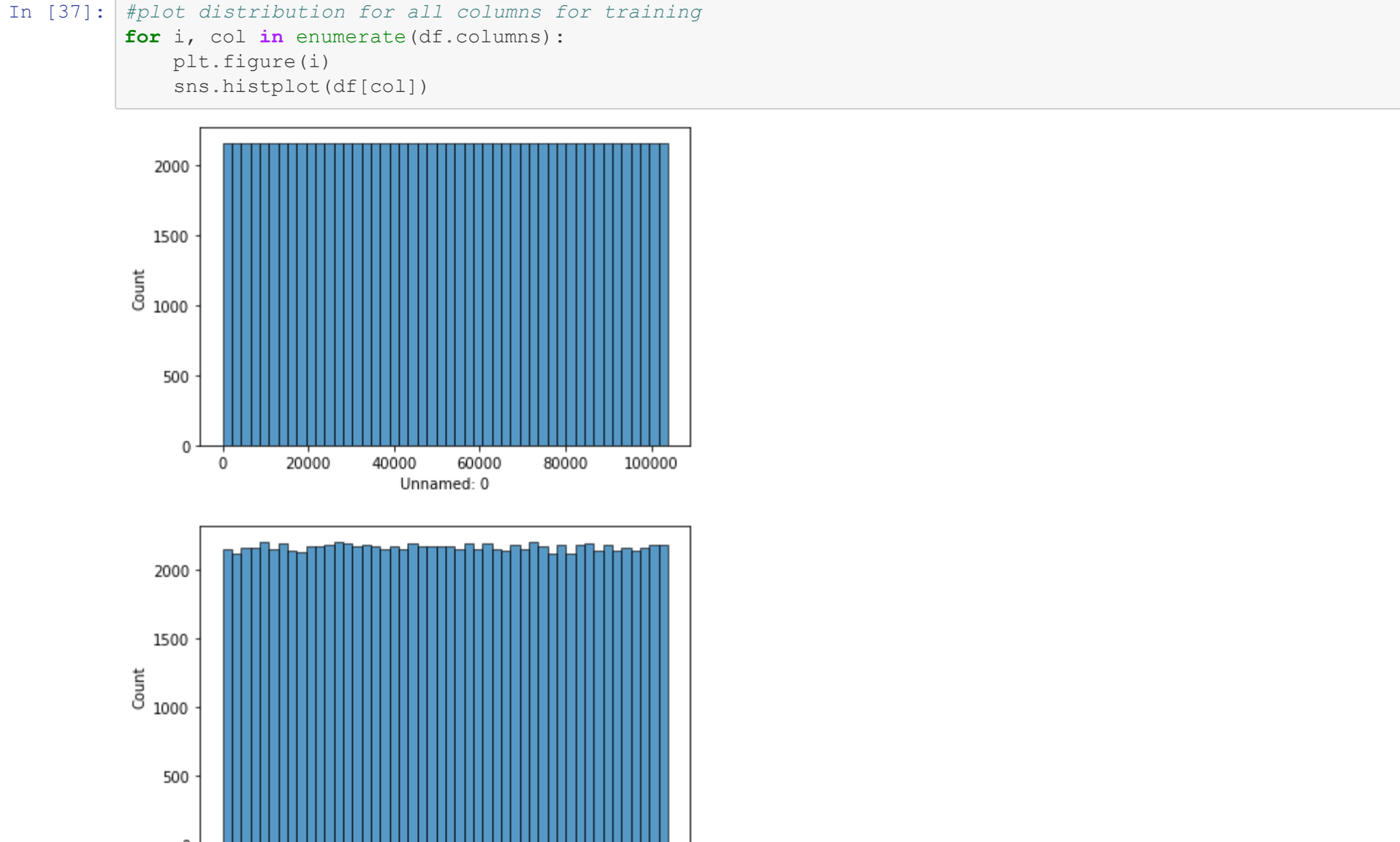


In [32]: #Viewing correlation between guest satisfaction and predictors.
#We see that the 5 predictors that are most correlated with customer satisfaction are online boarding, class, inflight_ent, seat_comf, and online_serv.
#The most correlated predictors with customer satisfaction are arr_delay, time_convenience, dept_delay, Unnamed: 0, and gate_loc. This is observation is valid because customers are less satisfied when they are inconvenienced with flight delays or flight times.
fig, ax = plt.subplots(figsize=(10,10))
fig.subplots('Correlation between Guest Satisfaction and Predictors',fontsize=20)
ax=sns.heatmap(df.corr()[['satisfaction_dummy']].sort_values("satisfaction_dummy"),vmax=1, vmin=-1, cma
p="rocket", annot=True, ax=ax;
ax.invert_yaxis()



Data Splitting

In [38]: #split data into predictors and outcome
predictors = ['Age', 'wifi_service', 'booking_diff', 'food_drink', 'online_boarding', 'seat_comf', 'inflight_ent', 'onboard_serv', 'legroom', 'baggage_care', 'checkin', 'inflight_serv', 'Cleanliness', 'flight_dummy', 'class_dummy']
outcome = ['satisfaction_dummy']
#split data into train/test
train_X = df[predictors]
train_y = df[outcome]
valid_X = df[predictors]
valid_y = df[outcome]

Classification Models

The models that we will use here include: Logistic Regression, Decision Tree, Random Forest, K-Nearest Neighbors, Bagging, AdaBoost, and Neural Network.

Logistic Regression Model

In [39]: #fit using Ridge regression, which adds squared magnitude of coefficient as a penalty term to the loss function
fcv = 5: to estimate the optimal parameters of the model
fsolver = 'saga': penalties of elasticnet, l1, l2, and none
#max_iter = 10000: max (iterations until) solvers can converge
Logit = LogisticRegressionCV(penalty='l2', solver='saga', cv=5, max_iter=10000, random_state=3).fit(train_X, train_y.values.ravel())

In [40]: logit_confusion = confusionMatrices(logit, 'Logistic regression')
Logistic regression - training results
Confusion Matrix (Accuracy 0.8434)

The training accuracy is 84.34% while the testing accuracy is 83.81%, which is not a huge difference showing that the model was not overtrained.

Decision Tree

In [41]: #pruned decision tree
dtree = DecisionTreeClassifier(random_state=3).fit(train_X, train_y)
tree_confusion = confusionMatrices(dtree, 'Decision Tree')

Decision Tree - training results
Confusion Matrix (Accuracy 0.9996)

The training accuracy is 99.96% and the test accuracy is 92.65%.

Bagging

In [42]: #using the decision classification tree as the base estimator, bagging is used to improve metrics
bagging = BaggingClassifier(dtree, random_state=3, max_estimator = 0.5, max_features = 0.5)
bagging.fit(train_X, train_y)
bag_confusion = confusionMatrices(bagging, 'Bagging')

Bagging - training results
Confusion Matrix (Accuracy 0.9552)

Training accuracy is 95.52% and test accuracy is 92.07%. Decision trees with bagging did not help with improvement in the accuracy.

AdaBoost

In [43]: #using the decision classification tree as the base estimator, adaboost is used to improve metrics
adaBoost = AdaBoostClassifier(n_estimators = 100, base_estimator = dtree, random_state=3)
adaBoost.fit(train_X, train_y)
ade_confusion = confusionMatrices(adaBoost, 'AdaBoost')

AdaBoost - training results
Confusion Matrix (Accuracy 0.9996)

Training accuracy is 99.96% and testing accuracy is 94.23%. AdaBoost improved accuracy in the decision tree.

Random Forest

In [44]: rf = RandomForestClassifier(n_estimators=100, random_state=3).fit(train_X, train_y.values.ravel())

In [45]: forest_confusion = confusionMatrices(rf, 'Random Forest')

Random Forest - training results
Confusion Matrix (Accuracy 0.9996)

The training accuracy is 99.96% and 94.83%. Random forest performed better than decision tree classifier by 2%.

k-NN

In [46]: #finding optimal k value
results = []
values = [1, 3, 5, 7, 9, 11, 13, 15, 17]
for k in values:
knn = KNeighborsClassifier(n_neighbors=k).fit(train_X, train_y)
results.append((k, knn.score(valid_X, knn.predict(valid_X))))
#accuracy: accuracy_score(valid_y, knn.predict(valid_X))
results = pd.DataFrame(results)
results

Out [46]:

0 k accuracy
1 3 0.898945
2 5 0.908846
3 7 0.910956
4 9 0.909147
5 11 0.898493
6 13 0.906300
7 15 0.906070
8 17 0.906375

In [47]: #creating knn algorithm with n_neighbors = 5
knn = KNeighborsClassifier(n_neighbors=5).fit(train_X, train_y)
knn_confusion = confusionMatrices(knn, 'k-NN Model')

k-NN Model - training results
Confusion Matrix (Accuracy 0.9360)

The training accuracy is 93.60% and test accuracy is 91.10%.

Linear Discriminant Analysis

In [48]: ldaModel = LinearDiscriminantAnalysis()
ldaModel.fit(train_X, train_y)
ldaConfusion = confusionMatrices(ldaModel, 'Linear Discriminant Analysis')

Linear Discriminant Analysis - training results
Confusion Matrix (Accuracy 0.8387)

The training accuracy is 83.87% and the testing accuracy is 83.34%.

Neural Network

In [49]: #scaling training features to 1.0
scaleInput = MinMaxScaler()
scaleInput.fit(train_X * 1.0)

#hidden_layer_sizes = 10: number of neurons in the ith hidden layer
#activation: logistic sigmoid function
#solver: weights optimization using stochastic gradient descent
neuralNet = MLPClassifier(hidden_layer_sizes = (10),
activation = 'logit',
solver = 'sgd',
max_iter = 3000,
random_state = 1)

neuralNet.fit(scaleInput.transform(train_X), train_y)
NNConfusion = confusionMatrices(neuralNet, 'Neural Network')

Neural Network - training results
Confusion Matrix (Accuracy 0.4926)

The neural network performs poorly on this data, with 49.26% accuracy on the training segment and 49.83% accuracy on the test segment.

Results

Models are evaluated using the test data for accuracy, precision, recall, f1 score, area under ROC curve, and cumulative gains.


```
In [50]: #Decision Tree
dtree_pred_test = dtree.predict(valid_X)
print('Accuracy score:')
print(accuracy_score(valid_y,dtree_pred_test))
dtree_roc = metrics.roc_curve(valid_y, dtree_pred_test)
dtree_auc = metrics.auc(dtree_roc[0], dtree_roc[1])
dtree_plot = metrics.RocCurveDisplay(dtree_roc[0], dtree_roc[1],
roc_auc=dtree_auc, estimator_name='Decision Tree')

Accuracy score:
0.9264705882352942
```

```
In [51]: print(classification_report(valid_y, dtree_pred_test))

              precision    recall  f1-score   support

    0.0         0.94         0.93         0.93         14573
    1.0         0.91         0.92         0.92         11403

 accuracy         0.93         0.93         0.93         25976
 macro avg        0.93         0.93         0.93         25976
weighted avg        0.93         0.93         0.93         25976
```

```
In [52]: #Bagging
bagging_pred_test = bagging.predict(valid_X)
print(accuracy_score(valid_y,bagging_pred_test))
bagging_roc = metrics.roc_curve(valid_y, bagging_pred_test)
bagging_auc = metrics.auc(bagging_roc[0], bagging_roc[1])
bagging_plot = metrics.RocCurveDisplay(bagging_roc[0], bagging_roc[1],
roc_auc=bagging_auc, estimator_name='Bagging')

0.920657530277179
```

```
In [53]: print(classification_report(valid_y, bagging_pred_test))

              precision    recall  f1-score   support

    0.0         0.92         0.94         0.93         14573
    1.0         0.92         0.89         0.91         11403

 accuracy         0.92         0.92         0.92         25976
 macro avg        0.92         0.92         0.92         25976
weighted avg        0.92         0.92         0.92         25976
```

```
In [54]: #Adaboost
adaboost_pred_test = adaboost.predict(valid_X)
print(accuracy_score(valid_y,adaboost_pred_test))
adaboost_roc = metrics.roc_curve(valid_y, adaboost_pred_test)
adaboost_auc = metrics.auc(adaboost_roc[0], adaboost_roc[1])
adaboost_plot = metrics.RocCurveDisplay(adaboost_roc[0], adaboost_roc[1],
roc_auc=adaboost_auc, estimator_name='Adaboost')

0.942928857406837
```

```
In [55]: print(classification_report(valid_y, adaboost_pred_test))

              precision    recall  f1-score   support

    0.0         0.94         0.96         0.95         14573
    1.0         0.94         0.92         0.93         11403

 accuracy         0.94         0.94         0.94         25976
 macro avg        0.94         0.94         0.94         25976
weighted avg        0.94         0.94         0.94         25976
```

```
In [56]: #Logistic Regression
logit_pred_test = logit.predict(valid_X)
print(accuracy_score(valid_y,logit_pred_test))
logit_roc = metrics.roc_curve(valid_y, logit_pred_test)
logit_auc = metrics.auc(logit_roc[0], logit_roc[1])
logit_plot = metrics.RocCurveDisplay(logit_roc[0], logit_roc[1],
roc_auc=logit_auc, estimator_name='Logistic Regression')

0.83811980289498
```

```
In [57]: print(classification_report(valid_y, logit_pred_test))

              precision    recall  f1-score   support

    0.0         0.84         0.87         0.86         14573
    1.0         0.83         0.79         0.81         11403

 accuracy         0.84         0.84         0.84         25976
 macro avg        0.84         0.83         0.83         25976
weighted avg        0.84         0.84         0.83         25976
```

```
In [58]: #Random Forest
rf_pred_test = rf.predict(valid_X)
print(accuracy_score(valid_y,rf_pred_test))
rf_roc = metrics.roc_curve(valid_y, rf_pred_test)
rf_auc = metrics.auc(rf_roc[0], rf_roc[1])
rf_plot = metrics.RocCurveDisplay(rf_roc[0], logit_roc[1],
roc_auc=rf_auc, estimator_name='Random Forest')

0.9482984293193717
```

```
In [59]: print(classification_report(valid_y, rf_pred_test))

              precision    recall  f1-score   support

    0.0         0.94         0.97         0.95         14573
    1.0         0.96         0.93         0.94         11403

 accuracy         0.95         0.95         0.95         25976
 macro avg        0.95         0.95         0.95         25976
weighted avg        0.95         0.95         0.95         25976
```

```
In [60]: #K-NN
knn_pred_test = knn.predict(valid_X)
print(accuracy_score(valid_y,knn_pred_test))
knn_roc = metrics.roc_curve(valid_y, knn_pred_test)
knn_auc = metrics.auc(knn_roc[0], knn_roc[1])
knn_plot = metrics.RocCurveDisplay(knn_roc[0], knn_roc[1],
roc_auc=knn_auc, estimator_name='k-Nearest Neighbors')

0.9109562673236834
```

```
In [61]: print(classification_report(valid_y, knn_pred_test))

              precision    recall  f1-score   support

    0.0         0.90         0.95         0.92         14573
    1.0         0.93         0.86         0.89         11403

 accuracy         0.91         0.91         0.91         25976
 macro avg        0.91         0.91         0.91         25976
weighted avg        0.91         0.91         0.91         25976
```

```
In [62]: #Linear Discriminant Analysis
ldaModel_pred_test = ldaModel.predict(valid_X)
print(accuracy_score(valid_y,ldaModel_pred_test))
ldaModel_roc = metrics.roc_curve(valid_y, ldaModel_pred_test)
ldaModel_auc = metrics.auc(ldaModel_roc[0], ldaModel_roc[1])
ldaModel_plot = metrics.RocCurveDisplay(ldaModel_roc[0], ldaModel_roc[1],
roc_auc=ldaModel_auc, estimator_name='Linear Discriminant Analysis')

0.834231593398522
```

```
In [63]: print(classification_report(valid_y, ldaModel_pred_test))

              precision    recall  f1-score   support

    0.0         0.84         0.86         0.85         14573
    1.0         0.82         0.79         0.81         11403

 accuracy         0.83         0.83         0.83         25976
 macro avg        0.83         0.83         0.83         25976
weighted avg        0.83         0.83         0.83         25976
```

```
In [64]: #Neural Network
neuralNet_pred_test = neuralNet.predict(valid_X)
print(accuracy_score(valid_y,neuralNet_pred_test))
neuralNet_roc = metrics.roc_curve(valid_y, neuralNet_pred_test)
neuralNet_auc = metrics.auc(neuralNet_roc[0], neuralNet_roc[1])
neuralNet_plot = metrics.RocCurveDisplay(neuralNet_roc[0], neuralNet_roc[1],
roc_auc=neuralNet_auc, estimator_name='Artificial Neural Network')

0.4983446258084386
```

```
In [65]: print(classification_report(valid_y, neuralNet_pred_test))

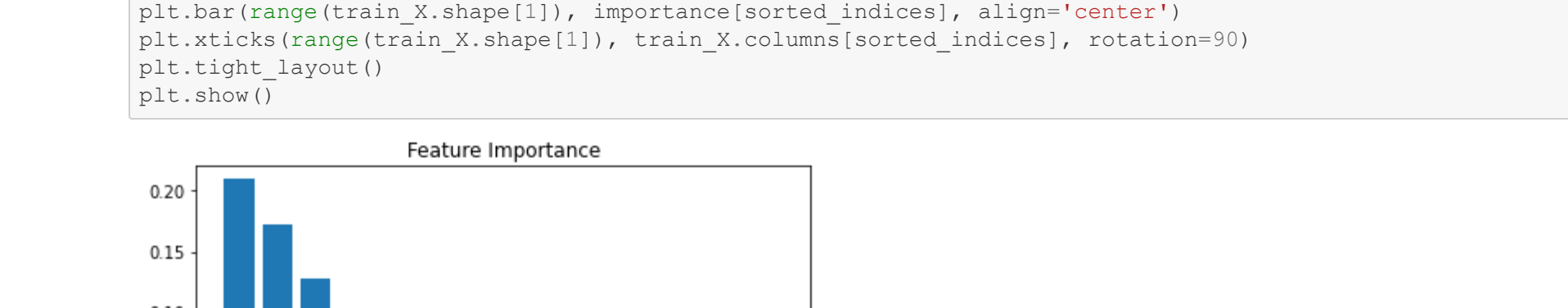
              precision    recall  f1-score   support

    0.0         0.82         0.14         0.23         14573
    1.0         0.47         0.96         0.63         11403

 accuracy         0.64         0.55         0.43         25976
 macro avg        0.66         0.50         0.41         25976
weighted avg        0.66         0.50         0.41         25976
```

```
In [66]: # Plotting ROC Curves for Models
fig, ax = plt.subplots(figsize=(12,8))
fig.suptitle('ROC Curves for Models', fontsize=12)
plt.plot([0, 1], [0, 1], linestyle = '--', color = '#174ab5')
plt.xlabel(' ', fontsize=12)
plt.ylabel(' ', fontsize=12)

logit_plot.plot(ax)
dtree_plot.plot(ax)
rf_plot.plot(ax)
knn_plot.plot(ax)
bagging_plot.plot(ax)
adaboost_plot.plot(ax)
ldaModel_plot.plot(ax)
neuralNet_plot.plot(ax)
plt.show()
```



The Random Forest model has the highest AUC, followed closely by Adaboost, which has strong performance when predicting true positives vs. false positives.

```
In [67]: from dmba import LiftChart
logisticchart = pd.Series(logit.predict_proba(valid_X)[: , 1])
logisticchart = logisticchart.sort_values(ascending=False)
dtreechart = pd.Series(dtree.predict_proba(valid_X)[: , 1])
dtreechart = dtreechart.sort_values(ascending=False)
rfchart = pd.Series(rf.predict_proba(valid_X)[: , 1])
rfchart = rfchart.sort_values(ascending=False)
knnchart = pd.Series(knn.predict_proba(valid_X)[: , 1])
knnchart = knnchart.sort_values(ascending=False)
baggingchart = pd.Series(bagging.predict_proba(valid_X)[: , 1])
baggingchart = baggingchart.sort_values(ascending=False)
adaboostchart = pd.Series(adaboost.predict_proba(valid_X)[: , 1])
adaboostchart = adaboostchart.sort_values(ascending=False)
ldaModelchart = pd.Series(ldaModel.predict_proba(valid_X)[: , 1])
ldaModelchart = ldaModelchart.sort_values(ascending=False)
neuralNetchart = pd.Series(neuralNet.predict_proba(scaleInput.transform(valid_X))[: , 1])
neuralNetchart = neuralNetchart.sort_values(ascending=False)

ax= gainsChart(logisticchart, figsize=(7,7))
gainsChart(dtreechart, color='palevioletred', ax=ax)
gainsChart(rfchart, color='green', ax=ax)
gainsChart(knnchart, color='purple', ax=ax)
gainsChart(baggingchart, color='orange', ax=ax)
gainsChart(adaboostchart, color='blue', ax=ax)
gainsChart(ldaModelchart, color='yellow', ax=ax)
gainsChart(neuralNetchart, color='red', ax=ax)

ax.set_title('Cumulative Gains Chart for Models')
colors = ['royalblue', 'palevioletred', 'green', 'purple', 'orange', 'blue', 'yellow', 'red']
lines = ['Logistic Regression', 'Decision Tree', 'Random Forest', 'k-NN', 'Bagging', 'Adaboost', 'Linear Discriminant Analysis', 'Artificial Neural Network']
labels=['Logistic Regression', 'Decision Tree', 'Random Forest', 'k-NN', 'Bagging', 'Adaboost', 'Linear Discriminant Analysis', 'Artificial Neural Network']

plt.legend(lines, labels)
```

```
Out [67]: <matplotlib.legend.Legend at 0x1f5cd34b5e0>
```



The Decision Tree, Random Forest, and Adaboost models have the strongest performance when examining cumulative gains.

Conclusion

```
In [68]: # Table of metrics
metrics = [{"Decision Tree", 0.926, 0.931, ['Bagging', 0.921, 0.92], ['Adaboost', 0.942, 0.94],
['Logistic Regression', 0.838, 0.83], ['Random Forest', 0.948, 0.95], ['k-NN', 0.911, 0.91],
['LDA', 0.833, 0.83], ['Neural Network', 0.498, 0.84]}]
metricdf = pd.DataFrame(metrics, columns = ['Model', 'Accuracy', 'AOC'])
metricdf
```

```
Out [68]:
```

	Model	Accuracy	AUC
0	Decision Tree	0.928	0.93
1	Bagging	0.921	0.92
2	Adaboost	0.942	0.94
3	Logistic Regression	0.838	0.83
4	Random Forest	0.948	0.95
5	k-NN	0.911	0.91
6	LDA	0.833	0.83
7	Neural Network	0.498	0.84

```
In [81]: # Plotting Importance by Feature
sorted_indices = np.argsort(importance)[::-1]
plt.title('Feature Importance')
plt.bar(range(train_X.shape[1]), importance[sorted_indices], align='center')
plt.xticks(range(train_X.shape[1]), train_X.columns[sorted_indices], rotation=90)
plt.tight_layout()
plt.show()
```



Using the random forest model, predicting whether a customer is satisfied or unsatisfied can be accomplished with high levels of accuracy, with the test data at 94.8% accuracy, the highest area under the ROC curve, and strong performance in cumulative gains.

Using the random forest model, the most influential features in determining if someone has a satisfactory experience is the quality of online boarding, the wifi service, what class they are in, their age, and the amount of legroom. Inflight entertainment and seat comfort are also important features, and are additional features airlines have a level of control over.

```
In [ ] :
```