

 [learn-co-curriculum](#) / [dsc-roc-curves-and-auc-lab](#) Public View license 0 stars  158 forks Star Watch ▼[Code](#) [Issues 1](#) [Pull requests](#) [Actions](#) [Projects](#) [Security](#) [Insights](#) solution ▼

...

This branch is [9 commits ahead](#), [9 commits behind](#) master.

lindseyberlin Fix typo ...

on Aug 24, 2020  10[View code](#) README.md

ROC Curves and AUC - Lab

Introduction

In this lab, you'll practice drawing ROC graphs, calculating AUC, and interpreting these results. In doing so, you will also further review logistic regression, by briefly fitting a model as in a standard data science pipeline.

Objectives

You will be able to:

- Create a visualization of ROC curves and use it to assess a model
- Evaluate classification models using the evaluation metrics appropriate for a specific problem

Train the model

Start by repeating the previous modeling steps we have discussed. For this problem, you are given a dataset `'mushrooms.csv'`. Your first job is to train a `LogisticRegression` classifier on the dataset to determine whether the mushroom is edible (e) or poisonous (p). The first column of the dataset `class` indicates whether or not the mushroom is poisonous or edible.

But first,

- Import the data
- Print the first five rows of the data
- Print DataFrame's `.info()`

```
import pandas as pd
```

```
# Load the data
```

```
df = pd.read_csv('mushrooms.csv')
```

```
# Data preview
```

```
print(df.head())
```

```
print('')
```

```
print(df.info())
```

```

class cap-shape cap-surface cap-color bruises odor gill-attachment \
0      p      x      s      n      t      p      f
1      e      x      s      y      t      a      f
2      e      b      s      w      t      l      f
3      p      x      y      w      t      p      f
4      e      x      s      g      f      n      f

```

```

gill-spacing gill-size gill-color ... stalk-surface-below-ring \
0           c      n      k ...                               s
1           c      b      k ...                               s
2           c      b      n ...                               s
3           c      n      n ...                               s
4           w      b      k ...                               s

```

```

stalk-color-above-ring stalk-color-below-ring veil-type veil-color \
0                       w                       w      p      w
1                       w                       w      p      w
2                       w                       w      p      w
3                       w                       w      p      w
4                       w                       w      p      w

```

	ring-number	ring-type	spore-print-color	population	habitat
0	o	p	k	s	u
1	o	p	n	n	g
2	o	p	n	n	m
3	o	p	k	s	u
4	o	e	n	a	g

[5 rows x 23 columns]

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8124 entries, 0 to 8123
Data columns (total 23 columns):
class                8124 non-null object
cap-shape             8124 non-null object
cap-surface          8124 non-null object
cap-color            8124 non-null object
bruises              8124 non-null object
odor                 8124 non-null object
gill-attachment      8124 non-null object
gill-spacing         8124 non-null object
gill-size            8124 non-null object
gill-color           8124 non-null object
stalk-shape          8124 non-null object
stalk-root           8124 non-null object
stalk-surface-above-ring 8124 non-null object
stalk-surface-below-ring 8124 non-null object
stalk-color-above-ring 8124 non-null object
stalk-color-below-ring 8124 non-null object
veil-type            8124 non-null object
veil-color           8124 non-null object
ring-number          8124 non-null object
ring-type            8124 non-null object
spore-print-color     8124 non-null object
population           8124 non-null object
habitat              8124 non-null object
dtypes: object(23)
memory usage: 1.4+ MB
None
```

The next step is to define the predictor and target variables. Did you notice all the columns are of type `object` ? So you will need to first create dummy variables.

- First, create a dummy variable for the `'class'` column. Make sure you drop the first level
- Drop the `'class'` column from `df` and then create dummy variables for all the remaining columns. Again, make sure you drop the first level

- Import `train_test_split`
- Split the data (`x` and `y`) into training and test sets with 25% in the test set. Set `random_state` to 42 to ensure reproducibility

```
# Define y
y = pd.get_dummies(df['class'], drop_first=True)
y = y['p']

# Define X
X = df.drop(columns=['class'], axis=1)
X = pd.get_dummies(X, drop_first=True)

# Import train_test_split
from sklearn.model_selection import train_test_split

# Split the data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=42)
```

- Fit the vanilla logistic regression model we defined for you to training data
- Make predictions using this model on test data

```
# Import LogisticRegression
from sklearn.linear_model import LogisticRegression

# Instantiate
logreg = LogisticRegression(fit_intercept=False, C=1e12, solver='liblinear')

# Fit the model to training data
model_log = logreg.fit(X_train, y_train)

# Predict on test set
y_hat_test = logreg.predict(X_test)
```

Calculate TPR and FPR

Next, calculate the false positive rate and true positive rate (you can use the built-in functions from `sklearn`):

```
# Import roc_curve, auc
from sklearn.metrics import roc_curve, auc

# Calculate the probability scores of each point in the training set
y_train_score = model_log.decision_function(X_train)
```

```
# Calculate the fpr, tpr, and thresholds for the training set
train_fpr, train_tpr, thresholds = roc_curve(y_train, y_train_score)

# Calculate the probability scores of each point in the test set
y_test_score = model_log.decision_function(X_test)

# Calculate the fpr, tpr, and thresholds for the test set
test_fpr, test_tpr, test_thresholds = roc_curve(y_test, y_test_score)
```

Draw the ROC curve

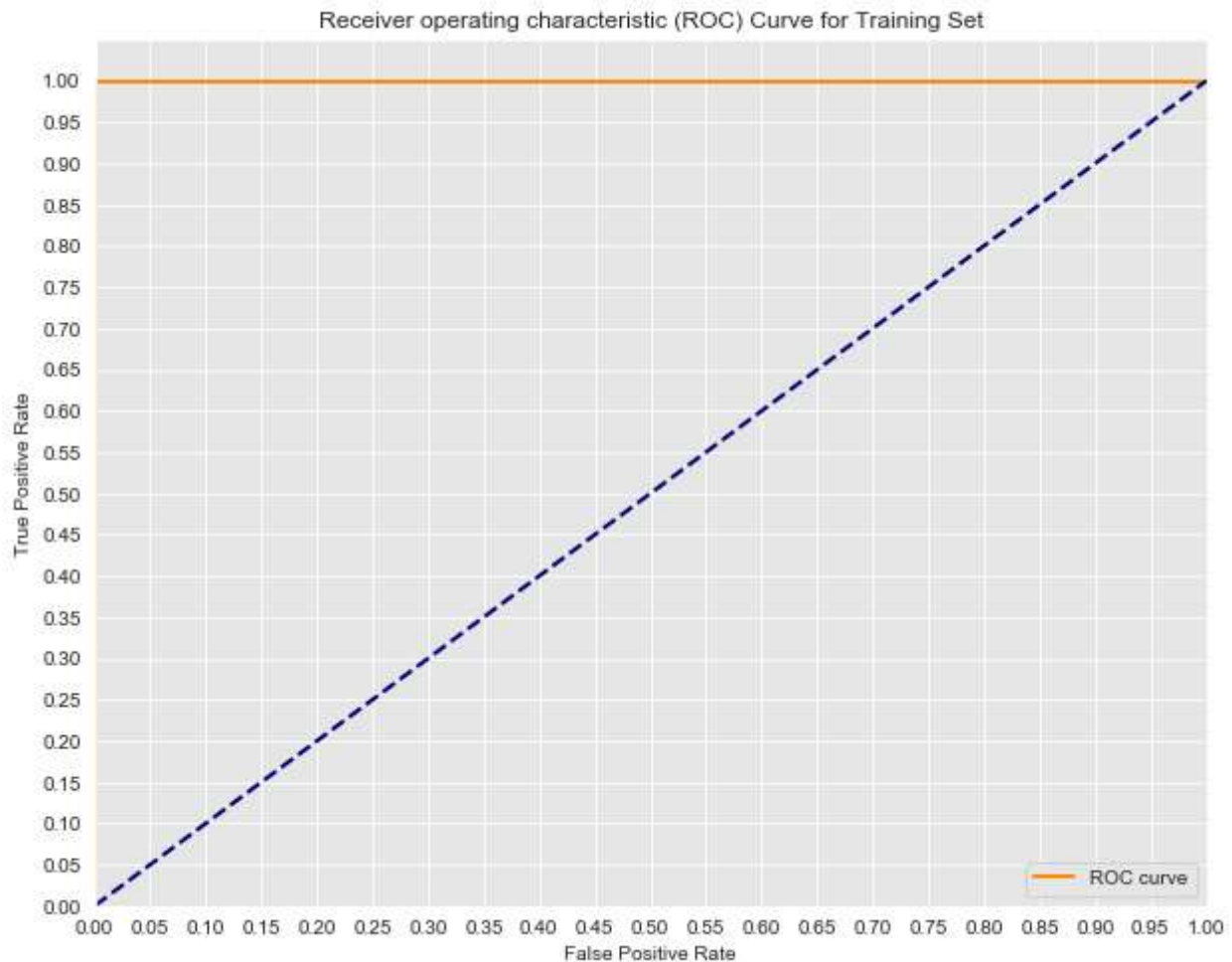
Next, use the false positive rate and true positive rate to plot the Receiver Operating Characteristic Curve for both the train and test sets.

```
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline

# Seaborn's beautiful styling
sns.set_style('darkgrid', {'axes.facecolor': '0.9'})

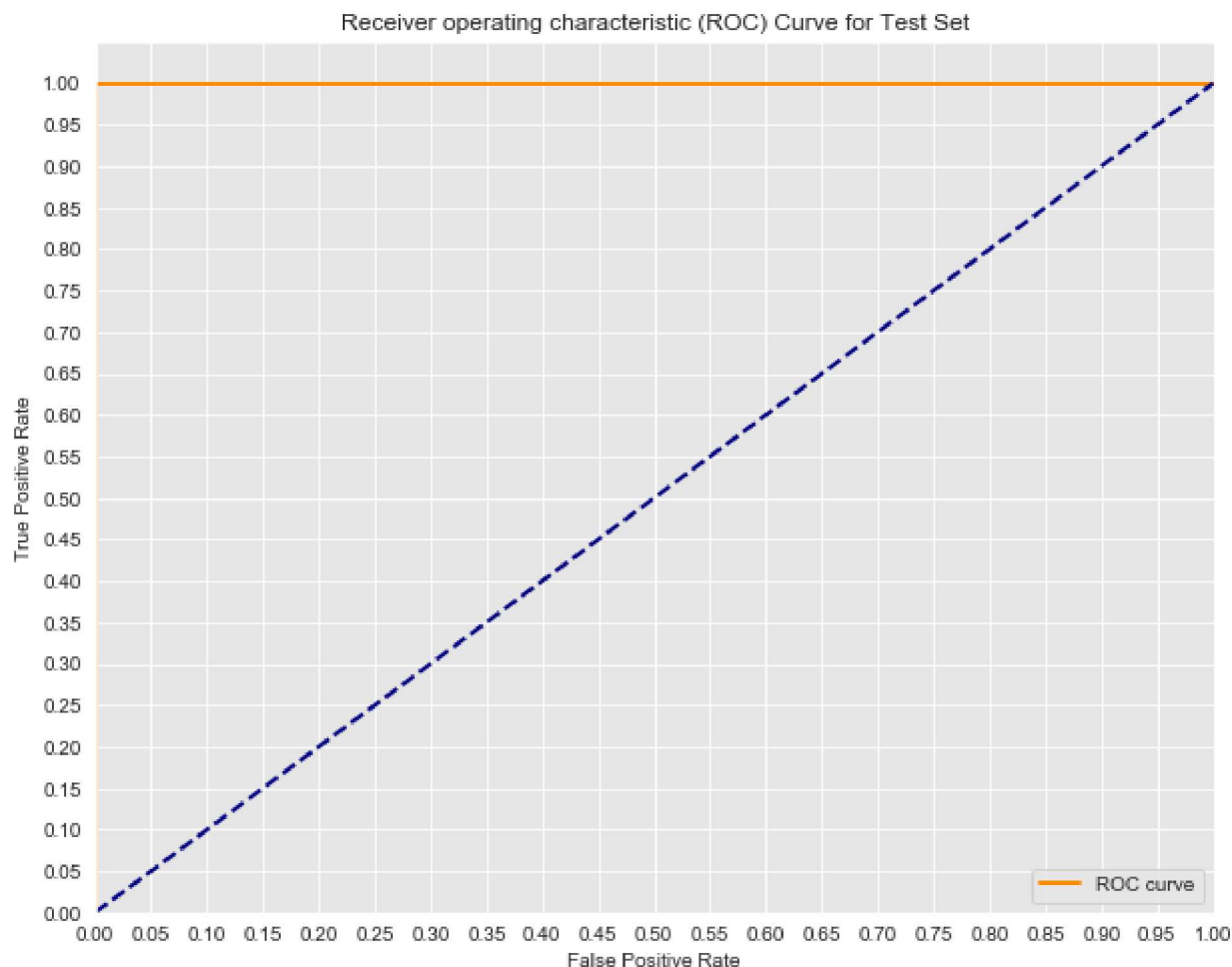
# ROC curve for training set
plt.figure(figsize=(10, 8))
lw = 2
plt.plot(train_fpr, train_tpr, color='darkorange',
         lw=lw, label='ROC curve')
plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.yticks([i/20.0 for i in range(21)])
plt.xticks([i/20.0 for i in range(21)])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic (ROC) Curve for Training Set')
plt.legend(loc='lower right')
print('Training AUC: {}'.format(auc(train_fpr, train_tpr)))
plt.show()
```

Training AUC: 1.0



```
# ROC curve for test set
plt.figure(figsize=(10, 8))
lw = 2
plt.plot(test_fpr, test_tpr, color='darkorange',
         lw=lw, label='ROC curve')
plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.yticks([i/20.0 for i in range(21)])
plt.xticks([i/20.0 for i in range(21)])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic (ROC) Curve for Test Set')
plt.legend(loc='lower right')
print('Test AUC: {}'.format(auc(test_fpr, test_tpr)))
print('')
plt.show()
```

Test AUC: 1.0

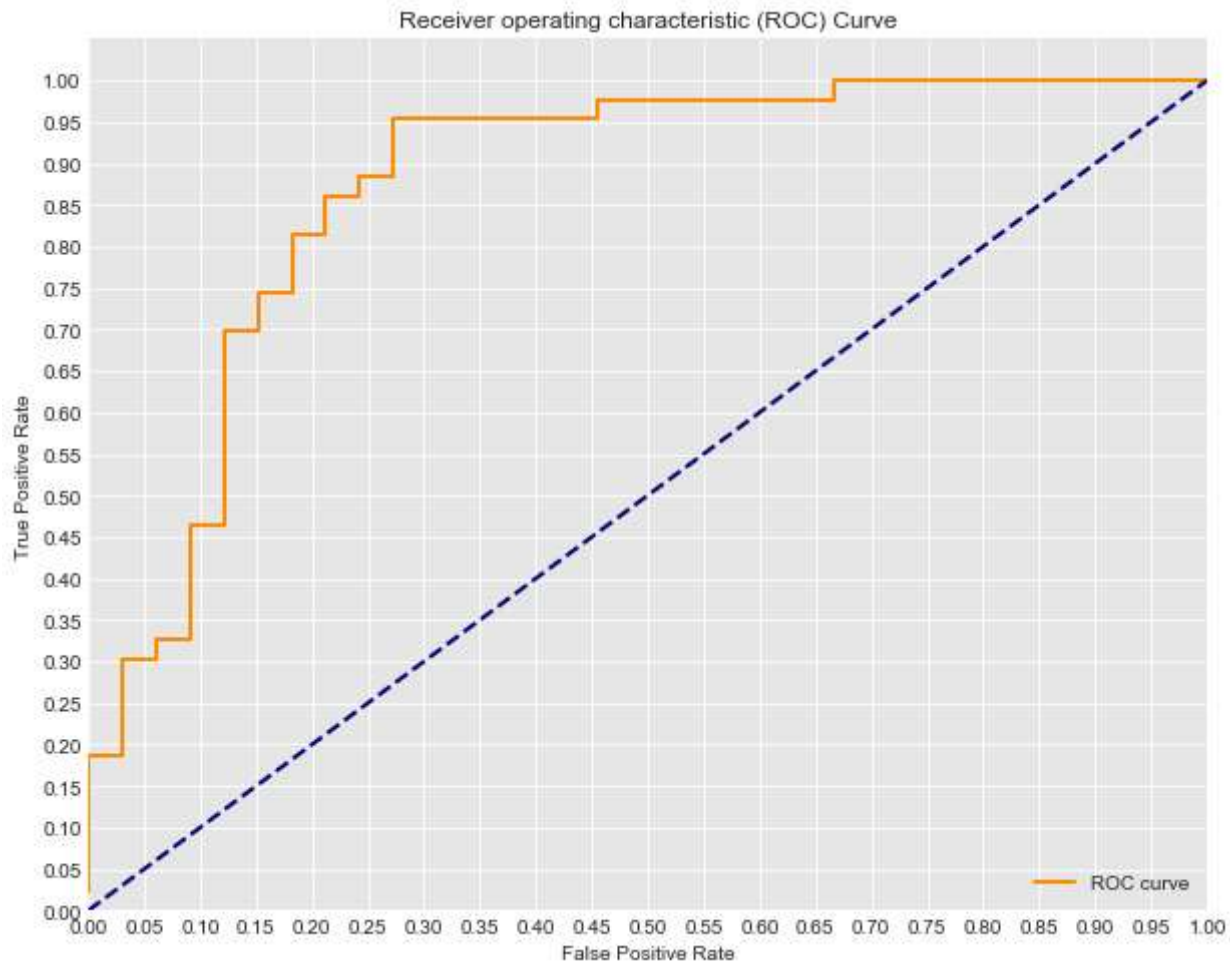


What do you notice about these ROC curves? Your answer here:

- # Both have an AUC of 1.0, indicating their performance is perfect.
- # Note that this is an extreme rarity!
- # Typically, if models perform this well it is too good to be true.

Interpret ROC curves

Look at the following ROC curve:



Think about the scenario of this model: predicting heart disease. If you tune the current model to have an 82% True Positive Rate, (you've still missed 20% of those with heart disease), what is the False positive rate?

```
# Write the approximate fpr when tpr = 0.8  
fpr = 0.17
```

If you instead tune the model to have a 95.2% True Postive Rate, what will the False Postive Rate be?

```
# Write the approximate fpr when tpr = 0.95  
fpr = 0.22
```

In the case of heart disease dataset, do you find any of the above cases acceptable? How would you tune the model? Describe what this would mean in terms of the number of patients falsely scared of having heart disease and the risk of missing the warning signs for those who do actually have heart disease.

Your answer here:


```
# With such an important decision, such as detecting heart disease, we would hope for
# The True positive weight is the more important of the two in this scenario.
# That is, the true positive rate determines the percentage of patients with heart disease
# The false positive rate is still very important, but it would be better to accidentally
# and warn them of potentially having heart disease than having missed warnings.
# That said, the false positive rate becomes rather unacceptably high once the true
# A .95 TPR indicates that out of 100 patients with heart disease we correctly warn
# At the same time, this has a FPR of nearly .25 meaning that roughly one in four times
# when they are actually healthy.
```

Summary

In this lab you further explored ROC curves and AUC, drawing graphs and then interpreting these results to lead to a more detailed and contextualized understanding of your model's accuracy.

Releases

No releases published

Packages

No packages published

Contributors 6



Languages

● Jupyter Notebook 100.0%