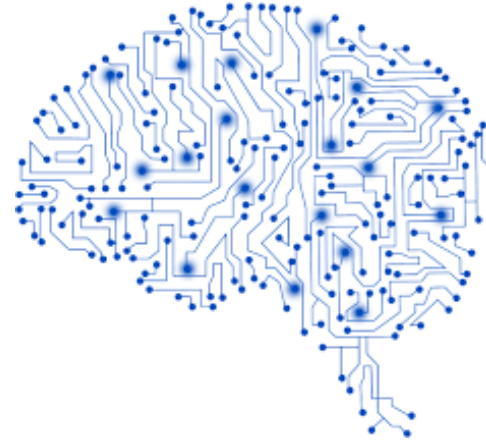




University of Minho
School of Engineering



Dados e Aprendizagem Automática

Linear & Logistic Regression

DAA @ MEI-1º/MiEI-4º – 1º Semestre

Bruno Fernandes, César Analide, Dalila Alves, Filipa Ferraz, Victor Alves

Contents

2

- Linear Regression
- Logistic Regression
- Hands On

3

Linear Regression

Linear Regression

4

□ Exercise:

- **Problem:** Development of a Machine Learning Model able to predict **house prices** for regions in the USA
- **Regression Approach:** Linear Regression approach to solve this problem
- **Dataset:** table with information regarding houses info. in regions of the United States, containing:
 - 'Avg. Area Income': Avg. Income of residents of the city house is located in.
 - 'Avg. Area House Age': Avg Age of Houses in same city
 - 'Avg. Area Number of Rooms': Avg Number of Rooms for Houses in same city
 - 'Avg. Area Number of Bedrooms': Avg Number of Bedrooms for Houses in same city
 - 'Area Population': Population of city house is located in
 - 'Price': Price that the house sold at
 - 'Address': Address for the house

	A	B	C	D	E	F	G
1	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
2	79545.45857	5.682861322	7.009188143	4.09	23086.8005	1059033.558	208
3	79248.64245	6.002899808	6.730821019	3.09	40173.07217	1505890.915	188
4	61287.06718	5.86588984	8.51272743	5.13	36882.1594	1058987.988	9127
5	63345.24005	7.188236095	5.586728665	3.26	34310.24283	1260616.807	USS
6	59982.19723	5.040554523	7.839387785	4.23	26354.10947	630943.4893	USNS

Linear Regression

5

Check out the data

We've been able to get some data from your neighbor for housing prices as a csv set, let's get our environment ready with the libraries we'll need and then import the data!

Import Libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

Check out the Data

```
USAhousing = pd.read_csv('USA_Housing.csv')
```

Linear Regression

6

```
USAhousing.head()
```

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price	Address
0	79545.458574	5.682861	7.009188	4.09	23086.800503	1.059034e+06	208 Michael Ferry Apt. 674\nLaurabury, NE 3701...
1	79248.642455	6.002900	6.730821	3.09	40173.072174	1.505891e+06	188 Johnson Views Suite 079\nLake Kathleen, CA...
2	61287.067179	5.865890	8.512727	5.13	36882.159400	1.058988e+06	9127 Elizabeth Stravenue\nDanielstown, WI 06482...
3	63345.240046	7.188236	5.586729	3.26	34310.242831	1.260617e+06	USS Barnett\nFPO AP 44820
4	59982.197226	5.040555	7.839388	4.23	26354.109472	6.309435e+05	USNS Raymond\nFPO AE 09386

```
USAhousing.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5000 entries, 0 to 4999
Data columns (total 7 columns):
Avg. Area Income      5000 non-null float64
Avg. Area House Age   5000 non-null float64
Avg. Area Number of Rooms  5000 non-null float64
Avg. Area Number of Bedrooms  5000 non-null float64
Area Population        5000 non-null float64
Price                  5000 non-null float64
Address                5000 non-null object
dtypes: float64(6), object(1)
memory usage: 273.5+ KB
```

```
USAhousing.describe()
```

	Avg. Area Income	Avg. Area House Age	Avg. Area Number of Rooms	Avg. Area Number of Bedrooms	Area Population	Price
count	5000.000000	5000.000000	5000.000000	5000.000000	5000.000000	5.000000e+03
mean	68583.108984	5.977222	6.987792	3.981330	36163.516039	1.232073e+06
std	10657.991214	0.991456	1.005833	1.234137	9925.650114	3.531176e+05
min	17796.631190	2.644304	3.236194	2.000000	172.610686	1.593866e+04
25%	61480.562388	5.322283	6.299250	3.140000	29403.928702	9.975771e+05
50%	68804.286404	5.970429	7.002902	4.050000	36199.406689	1.232669e+06
75%	75783.338666	6.650808	7.665871	4.490000	42861.290769	1.471210e+06
max	107701.748378	9.519088	10.759588	6.500000	69621.713378	2.469066e+06

```
USAhousing.columns
```

```
Index(['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',  
      'Avg. Area Number of Bedrooms', 'Area Population', 'Price', 'Address'],  
      dtype='object')
```

Linear Regression

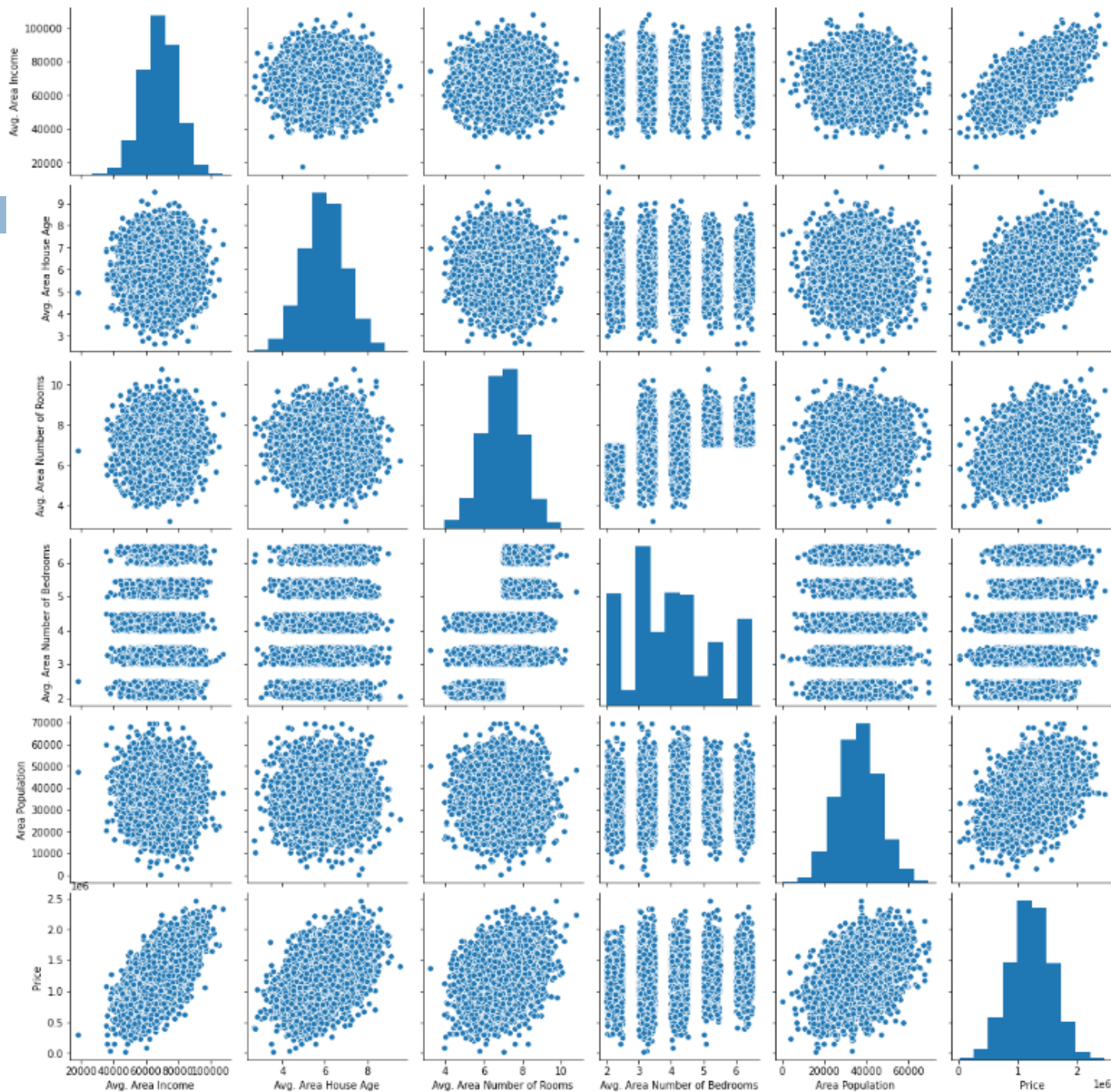
7

EDA

Let's create some simple plots to check out the data!

```
sns.pairplot(USAhousing)
```

```
<seaborn.axisgrid.PairGrid at 0x7f521d1c3ad0>
```

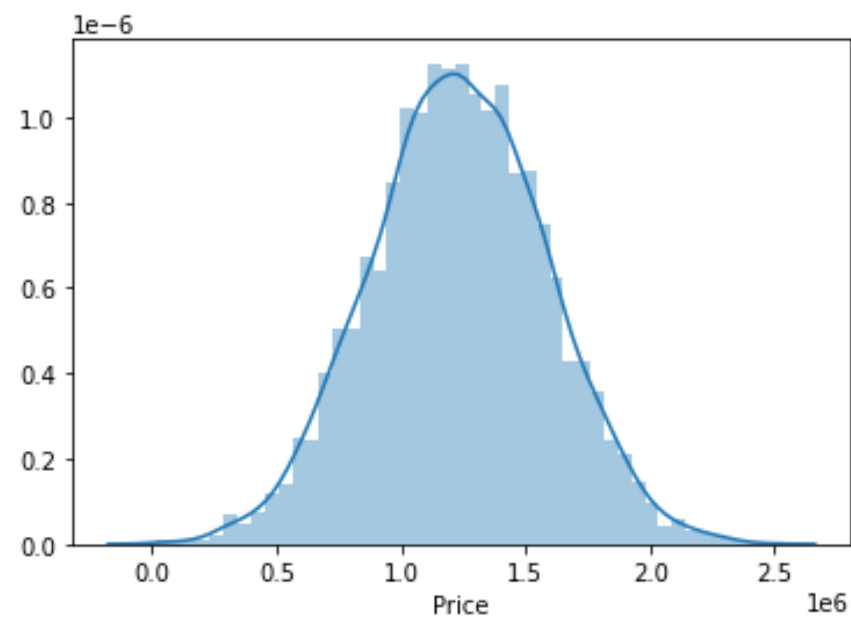


Linear Regression

8

```
sns.histplot(USAhousing['Price'])
```

<AxesSubplot:xlabel='Price'>



```
sns.heatmap(USAhousing.corr(numeric_only=True))
```

<AxesSubplot:>



Linear Regression

9

Training a Linear Regression Model

Let's now begin to train our regression model! We will need to first split up our data into an X array that contains the features to train on, and a y array with the target variable, in this case the Price column. We will toss out the Address column because it only has text info that the linear regression model can't use.

X and y arrays

```
X = USAhousing[['Avg. Area Income', 'Avg. Area House Age', 'Avg. Area Number of Rooms',  
               'Avg. Area Number of Bedrooms', 'Area Population']]  
y = USAhousing['Price']
```

Train Test Split

Now let's split the data into a training set and a testing set. We will train our model on the training set and then use the test set to evaluate the model.

```
from sklearn.model_selection import train_test_split
```

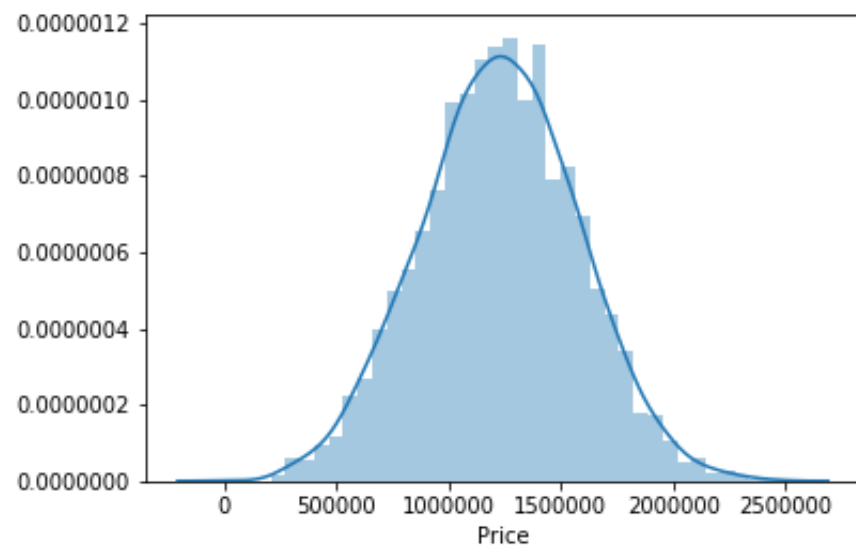
```
x_train, x_test, y_train, y_test = train_test_split(X, y, test_size=0.4, random_state=101)
```

Linear Regression

10

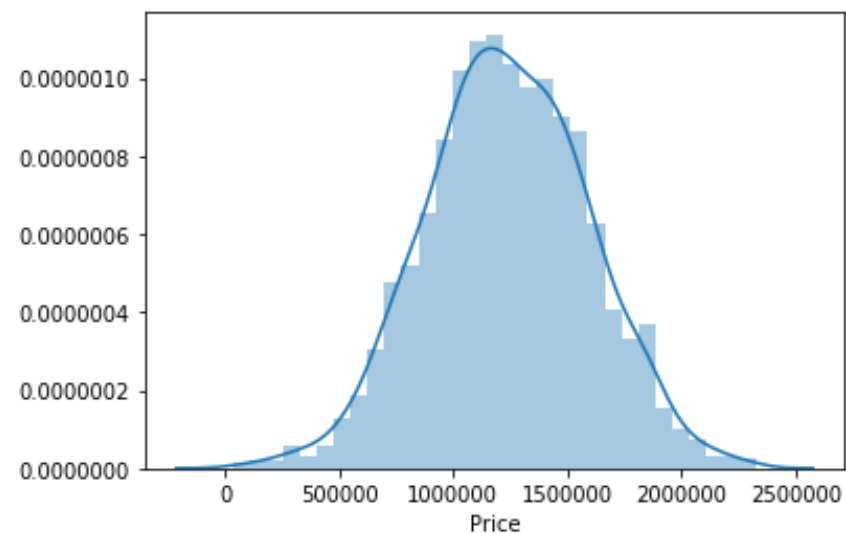
```
sns.histplot(y_train)
```

<matplotlib.axes._subplots.AxesSubplot at 0x1c6b394ed48>



```
sns.histplot(y_test)
```

<matplotlib.axes._subplots.AxesSubplot at 0x1c6b3ed28c8>



Linear Regression

11

Creating and Training the Model

```
from sklearn.linear_model import LinearRegression
```

```
lm = LinearRegression()
```

```
lm.fit(X_train,y_train)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)
```

Model Evaluation

Let's evaluate the model by checking out it's coefficients and how we can interpret them.

```
# print the intercept  
print(lm.intercept_)
```

```
-2640159.79685
```

```
coeff_df = pd.DataFrame(lm.coef_,X.columns,columns=['Coefficient'])  
coeff_df
```

	Coefficient
Avg. Area Income	21.528276
Avg. Area House Age	164883.282027
Avg. Area Number of Rooms	122368.678027
Avg. Area Number of Bedrooms	2233.801864
Area Population	15.150420

Interpreting the coefficients:

- Holding all other features fixed, a 1 unit increase in Avg. Area Income is associated with an *increase of \$21.52* *.
- Holding all other features fixed, a 1 unit increase in Avg. Area House Age is associated with an *increase of \$164883.28* *.
- Holding all other features fixed, a 1 unit increase in Avg. Area Number of Rooms is associated with an *increase of \$122368.67* *.
- Holding all other features fixed, a 1 unit increase in Avg. Area Number of Bedrooms is associated with an *increase of \$2233.80* *.
- Holding all other features fixed, a 1 unit increase in Area Population is associated with an *increase of \$15.15* *.

Linear Regression

12

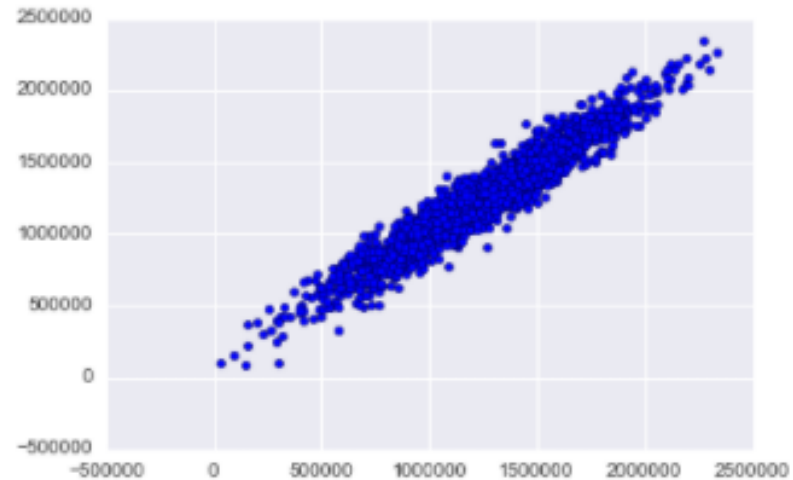
Predictions from our Model

Let's grab predictions off our test set and see how well it did!

```
predictions = lm.predict(X_test)
```

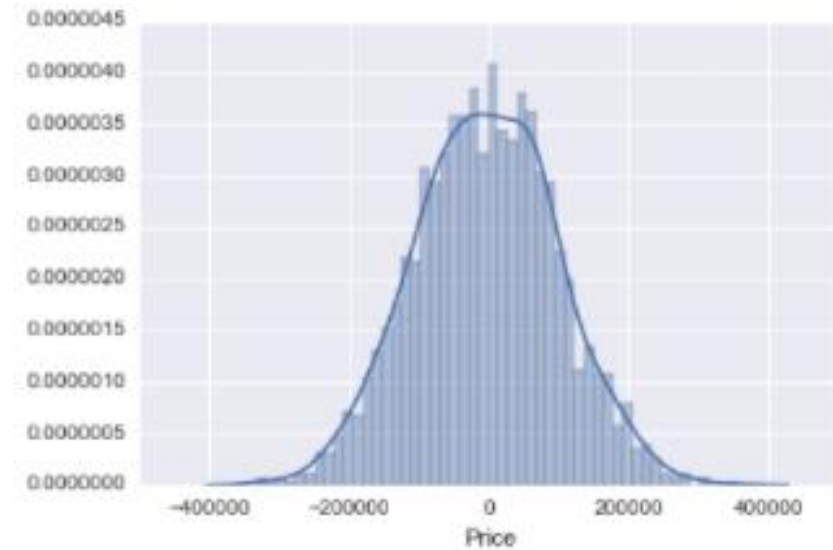
```
plt.scatter(y_test, predictions)
```

```
<matplotlib.collections.PathCollection at 0x142622c88>
```



Residual Histogram

```
sns.histplot((y_test-predictions),bins=50);
```



Linear Regression

13

Regression Evaluation Metrics

Here are three common evaluation metrics for regression problems:

Mean Absolute Error (MAE) is the mean of the absolute value of the errors:

$$\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Mean Squared Error (MSE) is the mean of the squared errors:

$$\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Root Mean Squared Error (RMSE) is the square root of the mean of the squared errors:

$$\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Comparing these metrics:

- **MAE** is the easiest to understand, because it's the average error.
- **MSE** is more popular than MAE, because MSE "punishes" larger errors, which tends to be useful in the real world.
- **RMSE** is even more popular than MSE, because RMSE is interpretable in the "y" units.

All of these are **loss functions**, because we want to minimize them.

```
from sklearn import metrics
```

```
print('MAE:', metrics.mean_absolute_error(y_test, predictions))  
print('MSE:', metrics.mean_squared_error(y_test, predictions))  
print('RMSE:', np.sqrt(metrics.mean_squared_error(y_test, predictions)))
```

MAE: 82288.2225191

MSE: 10460958907.2

RMSE: 102278.829223

14

Logistic Regression

Logistic Regression

15

□ Exercise:

- **Problem:** use machine learning to create a model that predicts which passengers survived the **Titanic** shipwreck
- **Classification Approach:** Logistic Regression approach to solve this problem
- **Dataset:** table with information regarding passengers' information, including:

Variable	Definition	Key
survival	Survival	0 = No, 1 = Yes
pclass	Ticket class	1 = 1st, 2 = 2nd, 3 = 3rd
sex	Sex	
Age	Age in years	
sibsp	# of siblings / spouses aboard the Titanic	
parch	# of parents / children aboard the Titanic	
ticket	Ticket number	
fare	Passenger fare	
cabin	Cabin number	
embarked	Port of Embarkation	C = Cherbourg, Q = Queenstown, S = Southampton

Logistic Regression

16

Import Libraries

Let's import some libraries to get started!

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import time
%matplotlib inline
```

The Data

Let's start by reading in the titanic_train.csv file into a pandas dataframe.

```
train = pd.read_csv('titanic_train.csv')
```

```
train.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

Logistic Regression

17

Exploratory Data Analysis

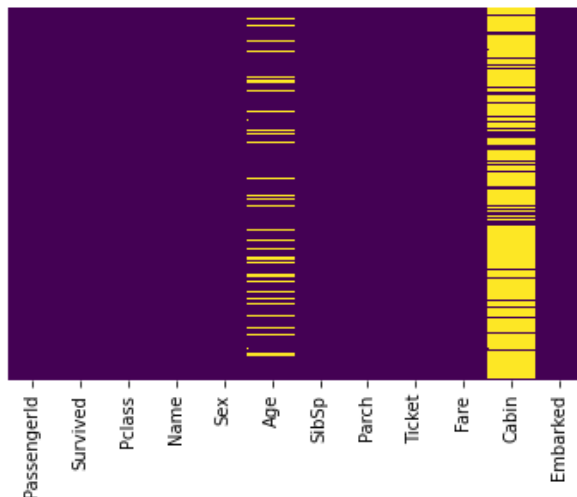
Let's begin some exploratory data analysis! We'll start by checking out missing data!

Missing Data

We can use seaborn to create a simple heatmap to see where we are missing data!

```
sns.heatmap(train.isnull(),yticklabels=False,cbar=False,cmap='viridis')
```

<AxesSubplot:>

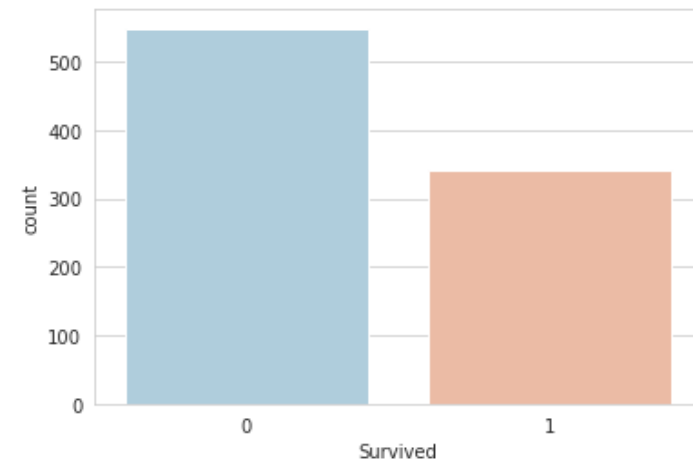


Roughly 20 percent of the Age data is missing. The proportion of Age missing is likely small enough for reasonable replacement with some form of imputation. Looking at the Cabin column, it looks like we are just missing too much of that data to do something useful with at a basic level. We'll probably drop this later, or change it to another feature like "Cabin Known: 1 or 0"

Let's continue on by visualizing some more of the data! Check out the video for full explanations over these plots, this code is just to serve as reference.

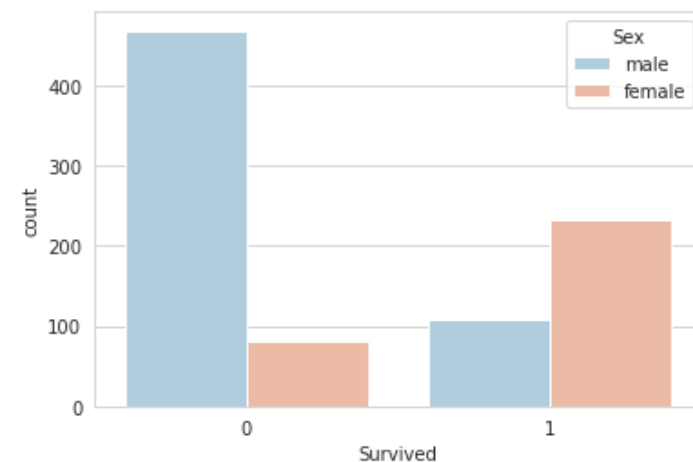
```
sns.set_style('whitegrid')
sns.countplot(x='Survived',data=train,palette='RdBu_r')
```

<AxesSubplot:xlabel='Survived', ylabel='count'>



```
sns.set_style('whitegrid')
sns.countplot(x='Survived',hue='Sex',data=train,palette='RdBu_r')
```

<AxesSubplot:xlabel='Survived', ylabel='count'>

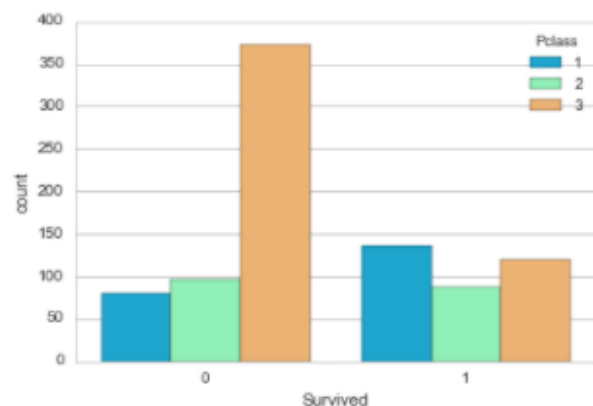


Logistic Regression

18

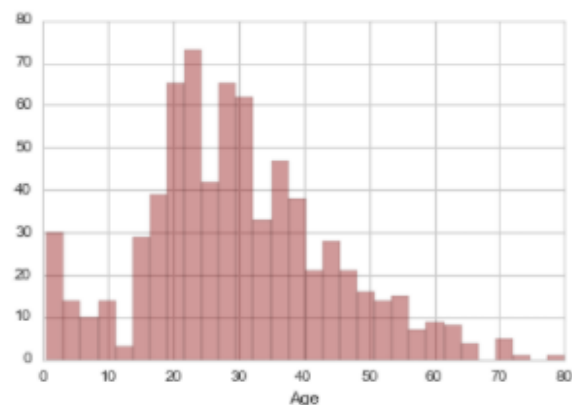
```
sns.set_style('whitegrid')  
sns.countplot(x='Survived', hue='Pclass', data=train, palette='rainbow')
```

<matplotlib.axes._subplots.AxesSubplot at 0x11b130f28>



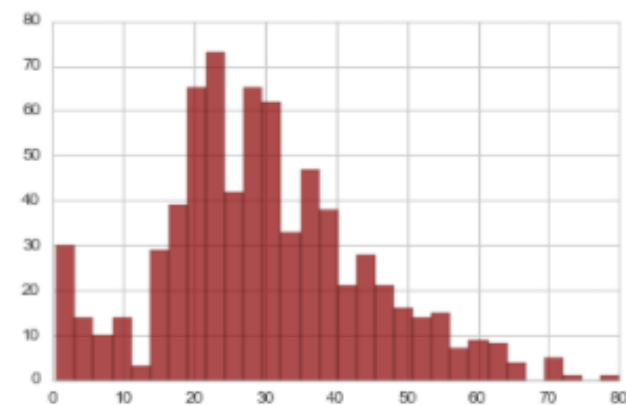
```
sns.histplot(train['Age'].dropna(), kde=False, color='darkred', bins=30)
```

<matplotlib.axes._subplots.AxesSubplot at 0x11c16f710>



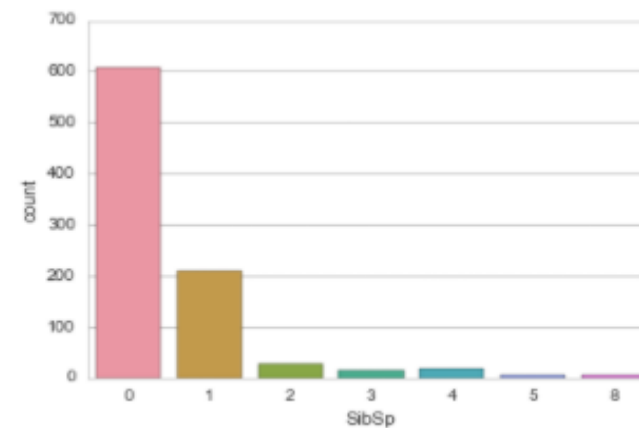
```
train['Age'].hist(bins=30, color='darkred', alpha=0.7)
```

<matplotlib.axes._subplots.AxesSubplot at 0x11b127ef0>



```
sns.countplot(x='SibSp', data=train)
```

<matplotlib.axes._subplots.AxesSubplot at 0x11c4139e8>

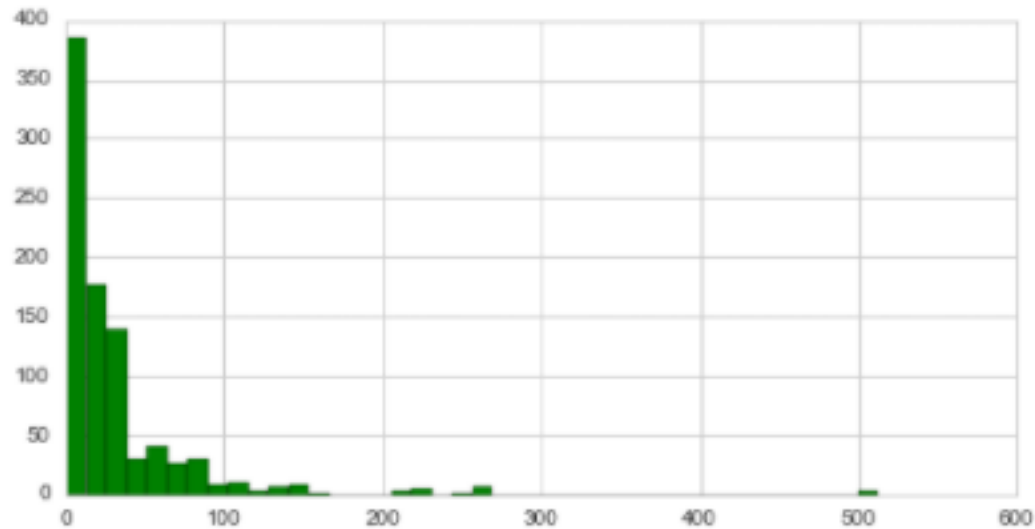


Logistic Regression

19

```
train['Fare'].hist(color='green',bins=40,figsize=(8,4))
```

<matplotlib.axes._subplots.AxesSubplot at 0x113893048>



Logistic Regression

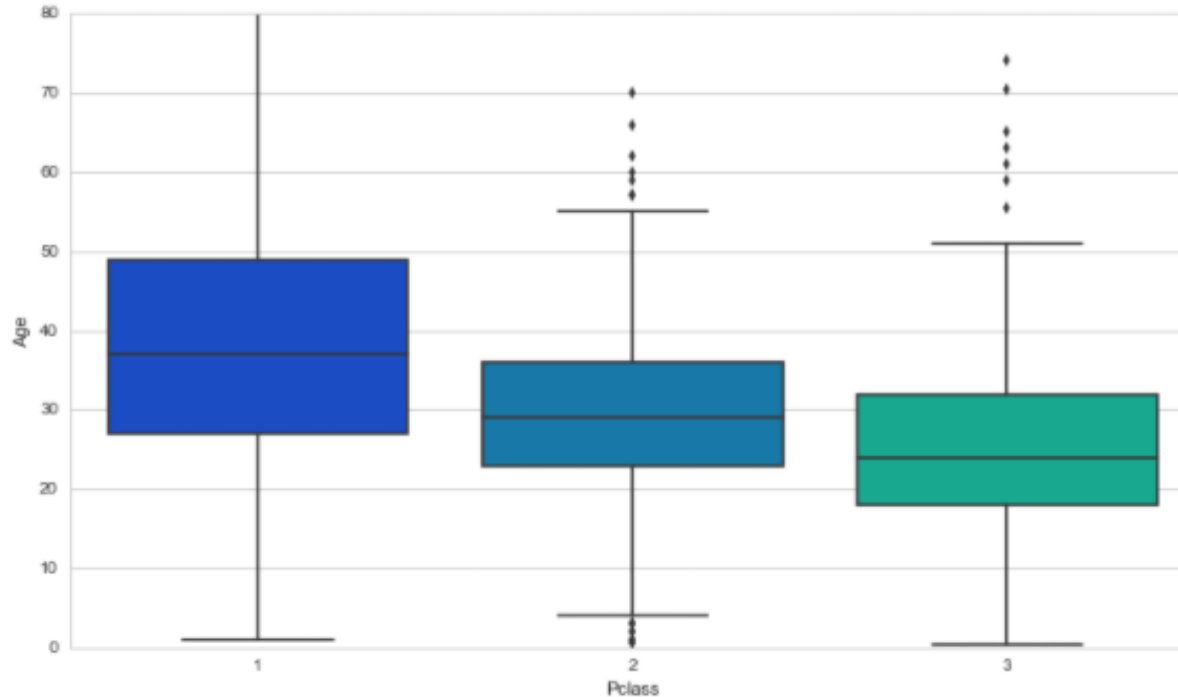
20

Data Cleaning

We want to fill in missing age data instead of just dropping the missing age data rows. One way to do this is by filling in the mean age of all the passengers (imputation). However we can be smarter about this and check the average age by passenger class. For example:

```
plt.figure(figsize=(12, 7))  
sns.boxplot(x='Pclass', y='Age', data=train, palette='winter')
```

<matplotlib.axes._subplots.AxesSubplot at 0x11c901cc0>



We can see the wealthier passengers in the higher classes tend to be older, which makes sense. We'll use these average age values to impute based on Pclass for Age.

Logistic Regression

21

```
def impute_age(cols):
    Age = cols[0]
    Pclass = cols[1]

    if pd.isnull(Age):

        if Pclass == 1:
            return 37

        elif Pclass == 2:
            return 29

        else:
            return 24

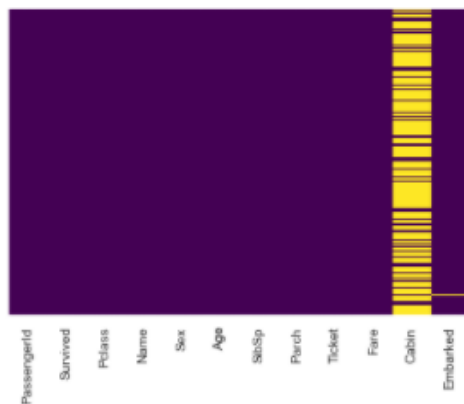
    else:
        return Age
```

Now apply that function!

```
train['Age'] = train[['Age', 'Pclass']].apply(impute_age,axis=1)
```

Now let's check that heat map again!

```
sns.heatmap(train.isnull(),yticklabels=False,cbar=False,cmap='viridis')
<matplotlib.axes._subplots.AxesSubplot at 0x11c4dae10>
```



Great! Let's go ahead and drop the Cabin column and the row in Embarked that is NaN.

```
train.drop('Cabin',axis=1,inplace=True)
```

```
train.head()
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	S

```
train.dropna(inplace=True)
```

Logistic Regression

22

Converting Categorical Features

We'll need to convert categorical features to dummy variables using pandas! Otherwise our machine learning algorithm won't be able to directly take in those features as inputs.

```
train.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 889 entries, 0 to 890
Data columns (total 11 columns):
PassengerId    889 non-null int64
Survived       889 non-null int64
Pclass         889 non-null int64
Name           889 non-null object
Sex            889 non-null object
Age            889 non-null float64
SibSp          889 non-null int64
Parch          889 non-null int64
Ticket         889 non-null object
Fare           889 non-null float64
Embarked       889 non-null object
dtypes: float64(2), int64(5), object(4)
memory usage: 83.3+ KB
```

```
sex = pd.get_dummies(train['Sex'],drop_first=True)
embark = pd.get_dummies(train['Embarked'],drop_first=True)
```

```
train.drop(['Sex','Embarked','Name','Ticket'],axis=1,inplace=True)
```

```
train = pd.concat([train,sex,embark],axis=1)
```

```
train.head()
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare	male	Q	S
0	1	0	3	22.0	1	0	7.2500	1.0	0.0	1.0
1	2	1	1	38.0	1	0	71.2833	0.0	0.0	0.0
2	3	1	3	28.0	0	0	7.9250	0.0	0.0	1.0
3	4	1	1	35.0	1	0	53.1000	0.0	0.0	1.0
4	5	0	3	35.0	0	0	8.0500	1.0	0.0	1.0

Great! Our data is ready for our model!

Logistic Regression

23

Building a Logistic Regression model

Let's start by splitting our data into a training set and test set.

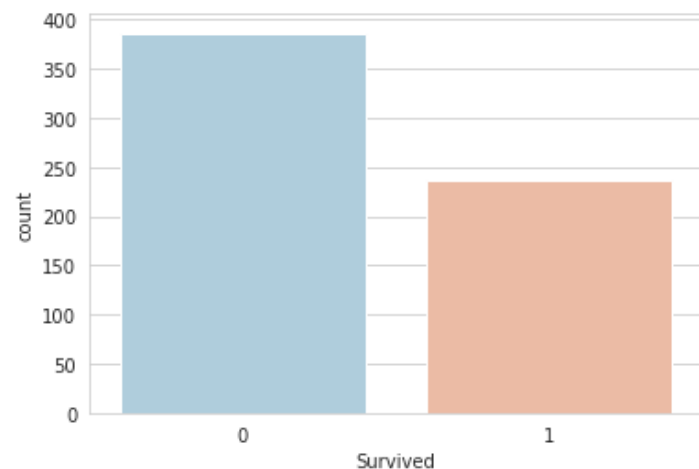
Train Test Split

```
from sklearn.model_selection import train_test_split
X = train.drop('Survived',axis=1)
y = train['Survived']
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    test_size=0.30,
                                                    random_state=101)
```

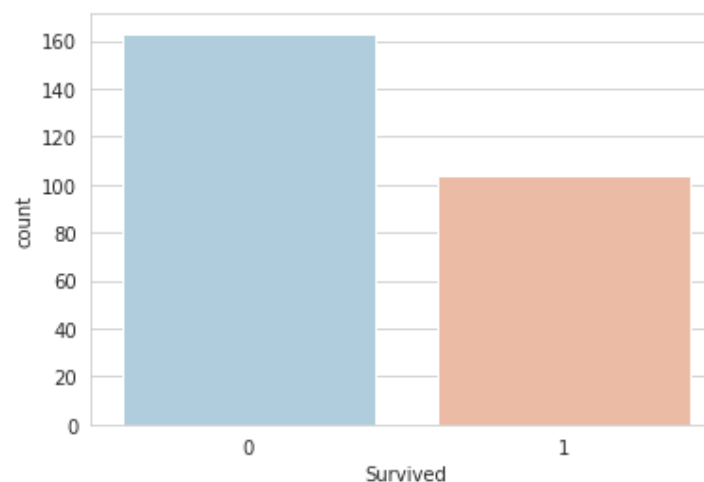
```
sns.set_style('whitegrid')
sns.countplot(x='Survived', data = pd.DataFrame(y_train,columns=['Survived'])) ,palette='RdBu_r')
```

<AxesSubplot:xlabel='Survived', ylabel='count'>



```
sns.set_style('whitegrid')
sns.countplot(x='Survived', data = pd.DataFrame(y_test,columns=['Survived'])) ,palette='RdBu_r')
```

<AxesSubplot:xlabel='Survived', ylabel='count'>



Logistic Regression

24

Training and Predicting

`sklearn.linear_model.LogisticRegression(penalty='l2', *, dual=False, tol=0.0001, C=1.0, fit_intercept=True, intercept_scaling=1, class_weight=None, random_state=None, solver='lbfgs', max_iter=100, multi_class='auto', verbose=0, warm_start=False, n_jobs=None, l1_ratio=None)`

LogisticRegression' **solvers**:

- Small datasets, 'liblinear' is a good choice, whereas 'sag' and 'saga' are faster for large ones
- Multiclass problems, only 'newton-cg', 'sag', 'saga' and 'lbfgs' handle multinomial loss
- 'liblinear' is limited to one-versus-rest schemes.

Supported **penalties** by solver:

- 'newton-cg' - ['l2', 'none']
- 'lbfgs' - ['l2', 'none']
- 'liblinear' - ['l1', 'l2']
- 'sag' - ['l2', 'none']
- 'saga' - ['elasticnet', 'l1', 'l2', 'none']

```
from sklearn.linear_model import LogisticRegression
```


Logistic Regression

25

logmodel1 - LogisticRegression(random_state=2022, solver='newton-cg')

```
starttime = time.process_time()

logmodel1 = LogisticRegression(random_state=2022, solver='newton-cg')
print(logmodel1)
logmodel1.fit(X_train,y_train)

endtime = time.process_time()
print(f"Time spent: {endtime - starttime} seconds")
```

LogisticRegression(random_state=2022, solver='newton-cg')
Time spent: 0.09568306500000023 seconds

```
predictions1 = logmodel1.predict(X_test)
```

logmodel2 - LogisticRegression(random_state=2022, solver='lbfgs')

```
starttime = time.process_time()

logmodel2 = LogisticRegression(random_state=2022, solver='lbfgs')
print(logmodel2)
logmodel2.fit(X_train,y_train)

endtime = time.process_time()
print(f"Time spent: {endtime - starttime} seconds")
```

LogisticRegression(random_state=2022)
Time spent: 0.06592618100000003 seconds

```
predictions2 = logmodel2.predict(X_test)
```

logmodel3 - LogisticRegression(random_state=2022, solver='liblinear')

```
starttime = time.process_time()

logmodel3 = LogisticRegression(random_state=2022, solver='liblinear')
print(logmodel3)
logmodel3.fit(X_train,y_train)

endtime = time.process_time()
print(f"Time spent: {endtime - starttime} seconds")
```

LogisticRegression(random_state=2022, solver='liblinear')
Time spent: 0.011570596999999516 seconds

```
predictions3 = logmodel3.predict(X_test)
```

Logistic Regression

26

Evaluation

We can check precision, recall, f1-score using classification report!

```
from sklearn.metrics import classification_report, ConfusionMatrixDisplay

print("With 'newton-cg': \n", classification_report(y_test, predictions1))
print("With 'lbfgs': \n", classification_report(y_test, predictions2))
print("With 'liblinear': \n", classification_report(y_test, predictions3))
```

With 'newton-cg':

	precision	recall	f1-score	support
0	0.82	0.91	0.86	163
1	0.84	0.68	0.75	104
accuracy			0.82	267
macro avg	0.83	0.80	0.81	267
weighted avg	0.83	0.82	0.82	267

With 'lbfgs':

	precision	recall	f1-score	support
0	0.79	0.91	0.85	163
1	0.81	0.62	0.71	104
accuracy			0.80	267
macro avg	0.80	0.77	0.78	267
weighted avg	0.80	0.80	0.79	267

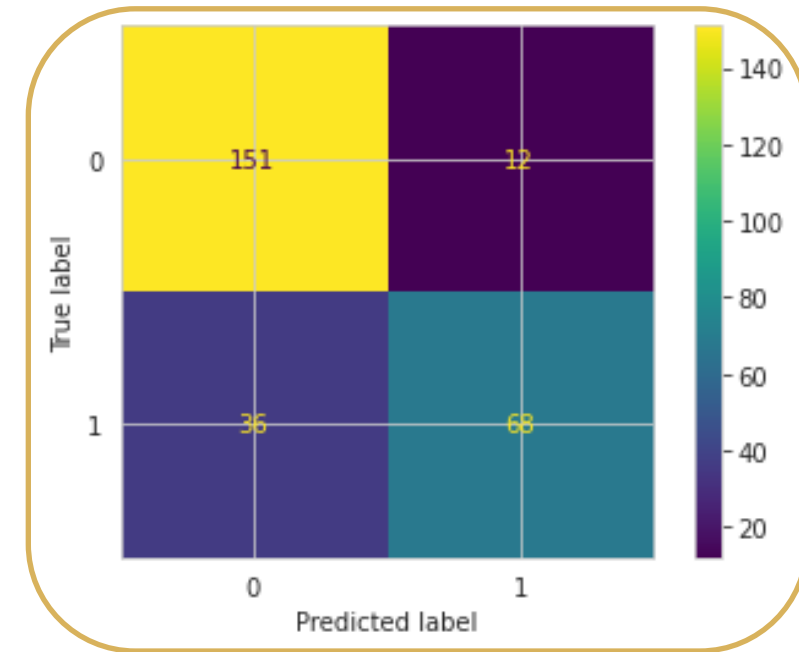
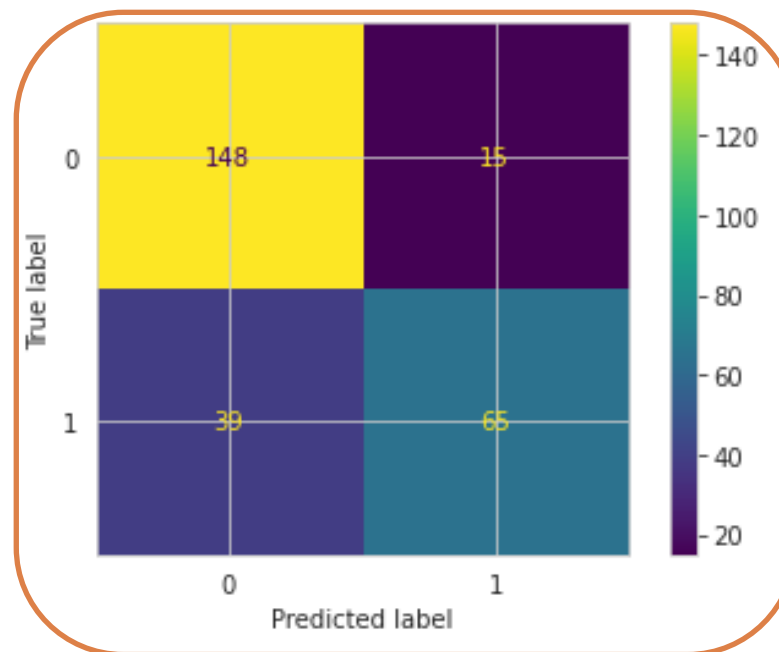
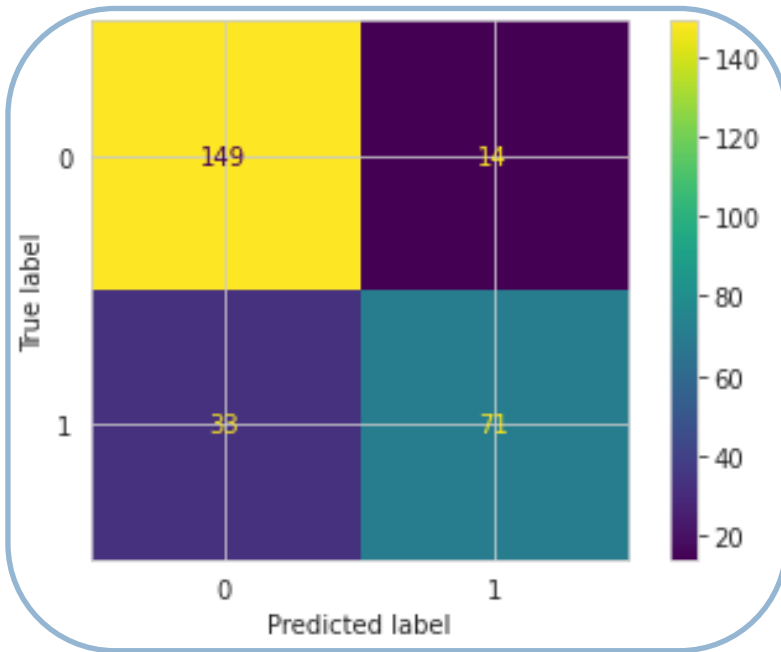
With 'liblinear':

	precision	recall	f1-score	support
0	0.81	0.93	0.86	163
1	0.85	0.65	0.74	104
accuracy			0.82	267
macro avg	0.83	0.79	0.80	267
weighted avg	0.82	0.82	0.81	267

Logistic Regression

27

```
#Get the confusion matrix  
ConfusionMatrixDisplay.from_predictions(y_test, predictions1)  
ConfusionMatrixDisplay.from_predictions(y_test, predictions2)  
ConfusionMatrixDisplay.from_predictions(y_test, predictions3)  
plt.show()
```



Hands On

28

Hands On

The screenshot displays the Spyder Python IDE interface. The main editor window shows a Python script for a Mean-Shift algorithm. The script defines a `Mean_Shift` class with an `__init__` method and a `fit` method. The `fit` method initializes centroids and iteratively updates them based on the data points, using a radius-based approach to find clusters. The script is located at `C:\data\PythonWorkspace\dev\meanshift_algorithm.py`.

The Variable explorer on the right shows the following variables:

Name	Type	Size	Value
batch_size	int	1	100
mnist	contrib.learn.python.learn.datasets.base.Datasets	3	Datasets object of...
n_classes	int	1	10
n_nodes_h1	int	1	500
n_nodes_h2	int	1	500
n_nodes_h3	int	1	500

The IPython console at the bottom shows the output of the script, including the loss and accuracy over 10 epochs:

```
See 'tf.nn.softmax_cross_entropy_with_logits_v2'.  
Epoch 0 completed out of 10 loss: 1666037.4677734375  
Epoch 1 completed out of 10 loss: 377809.3128890991  
Epoch 2 completed out of 10 loss: 201302.4857263565  
Epoch 3 completed out of 10 loss: 119427.91378033161  
Epoch 4 completed out of 10 loss: 72651.25679710507  
Epoch 5 completed out of 10 loss: 45327.621502393486  
Epoch 6 completed out of 10 loss: 31955.17812934518  
Epoch 7 completed out of 10 loss: 23664.35610633137  
Epoch 8 completed out of 10 loss: 18248.740643078025  
Epoch 9 completed out of 10 loss: 19962.00065876091  
Accuracy: 0.9511
```