

# Python Machine Learning [BASIC SHEET]

- Mean                                      Average value
- Median                                    Midpoint value
- Mode                                      Most common value

Standard deviation

**Sigma:  $\sigma$**

- Number that describes how spread out the values are.
- Low standard deviation means that most of the numbers are close to the mean (average) value.
- High standard deviation means that the values are spread out over a wider range.
- Standard deviation is the square root of the variance

Variance; Indicates how spread out the values are

**Sigma Square:  $\sigma^2$**

Percentiles; describes the value that a given percent of the values are lower than

## Data Distribution

```
import numpy
```

```
import matplotlib.pyplot as plt
```

```
x = numpy.random.uniform(0.0, 5.0, 250) # array containing 250 random floats between 0 and 5:
```

```
print(x)
```

## Histogram

```
x = numpy.random.uniform(0.0, 5.0, 250)
```

```
plt.hist(x, 5)
```

```
plt.show()
```

## Big Data Distributions

```
x = numpy.random.uniform(0.0, 5.0, 100000) # Create an array with 100000 random numbers
```

```
plt.hist(x, 100) # histogram with 100 bars:
```

```
plt.show()
```

## Normal Data Distribution

In probability theory this kind of data distribution is known as the *normal data distribution*, or the *Gaussian data distribution*, after the mathematician Carl Friedrich Gauss who came up with the formula of this data distribution.

```
import numpy
import matplotlib.pyplot as plt
x = numpy.random.normal(5.0, 1.0, 100000)
plt.hist(x, 100)
plt.show()
```

**Note:** A normal distribution graph is also known as the *bell curve* because of its characteristic shape of a bell.

## Histogram Explained

We use the array from the `numpy.random.normal()` method, with 100000 values, to draw a histogram with 100 bars. We specify that the mean value is 5.0, and the standard deviation is 1.0. Meaning that the values should be concentrated around 5.0, and rarely further away than 1.0 from the mean. And as you can see from the histogram, most values are between 4.0 and 6.0, with a top at approximately 5.0.

## Scatter Plot

A scatter plot is a diagram where each value in the data set is represented by a dot.

The Matplotlib module has a method for drawing scatter plots, it needs two arrays of the same length, one for the values of the x-axis, and one for the values of the y-axis:

```
import matplotlib.pyplot as plt
x = [5,7,8,7,2,17,2,9,4,11,12,9,6]          age of each car
y = [99,86,87,88,111,86,103,87,94,78,77,85,86] speed of each car
plt.scatter(x, y)
plt.show()
```

## Scatter Plot Explained

The x-axis represents ages, and the y-axis represents speeds. What we can read from the diagram is that the two fastest cars were both 2 years old, and the slowest car was 12 years old.

**Note:** It seems that the newer the car, the faster it drives, but that could be a coincidence, after all we only registered 13 cars.

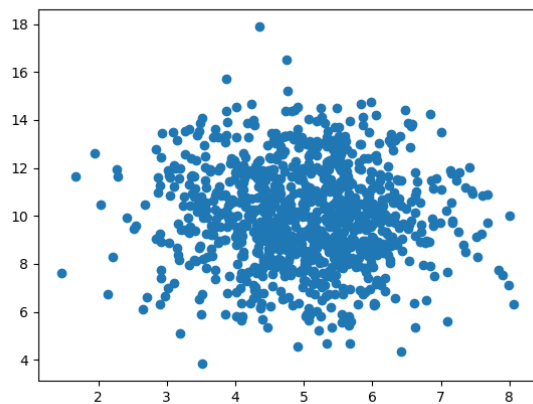
## Random Data Distributions

Create two arrays that are both filled with 1000 random numbers from a normal data distribution.

The first array will have the mean set to 5.0 with a standard deviation of 1.0.

The second array will have the mean set to 10.0 with a standard deviation of 2.0:

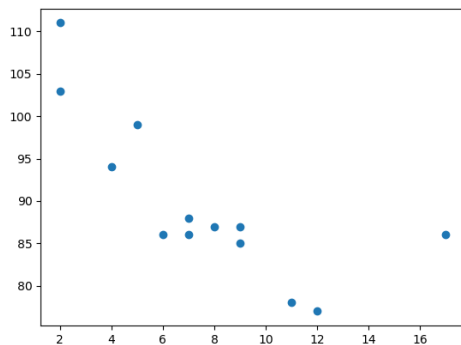
```
x = numpy.random.normal(5.0, 1.0, 1000)
y = numpy.random.normal(10.0, 2.0, 1000)
plt.scatter(x, y)
plt.show()
```



Regression; used when you try to find the relationship between variables. In Machine Learning, and in statistical modeling, that relationship is used to predict the outcome of future events.

Linear Regression; uses the relationship between the data-points to draw a straight line through all them. This line can be used to predict future values.

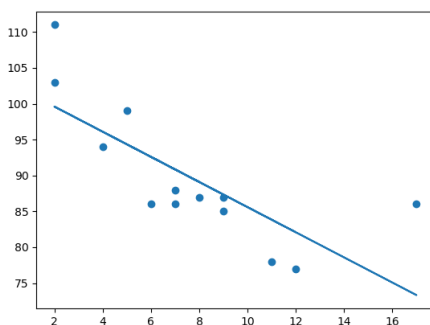
```
x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
plt.scatter(x, y)
plt.show()
```



```
x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
slope, intercept, r, p, std_err = stats.linregress(x, y)
```

```
def myfunc(x):
    return slope * x + intercept
```

```
mymodel = list(map(myfunc, x))
plt.scatter(x, y)
plt.plot(x, mymodel)
plt.show()
```



```
import matplotlib.pyplot as plt
from scipy import stats
x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
```

Execute a method that returns some important key values of Linear Regression:

```
slope, intercept, r, p, std_err = stats.linregress(x, y)
```

Create a function that uses the slope and intercept values to return a new value. This new value represents where on the y-axis the corresponding x value will be placed:

```
def myfunc(x):
    return slope * x + intercept
```

Run each value of the x array through the function

```
mymodel = list(map(myfunc, x))
```

Draw the original scatter plot:

```
plt.scatter(x, y)
```

Draw the line of linear regression:

```
plt.plot(x, mymodel)
```

Display the diagram:

```
plt.show()
```

## R for Relationship

It is important to know how the relationship between the values of the x-axis and the values of the y-axis is, if there are no relationship the linear regression can not be used to predict anything.

This relationship - the coefficient of correlation - is called r.

The r value ranges from -1 to 1, where 0 means no relationship, and 1 (and -1) means 100% related.

```
from scipy import stats
x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
slope, intercept, r, p, std_err = stats.linregress(x, y)
print(r)
```

**Note:** The result -0.76 shows that there is a relationship, not perfect, but it indicates that we could use linear regression in future predictions.

**# Predict the speed of a 10 years old car:**

```
from scipy import stats
```

```
x = [5,7,8,7,2,17,2,9,4,11,12,9,6]
```

```
y = [99,86,87,88,111,86,103,87,94,78,77,85,86]
```

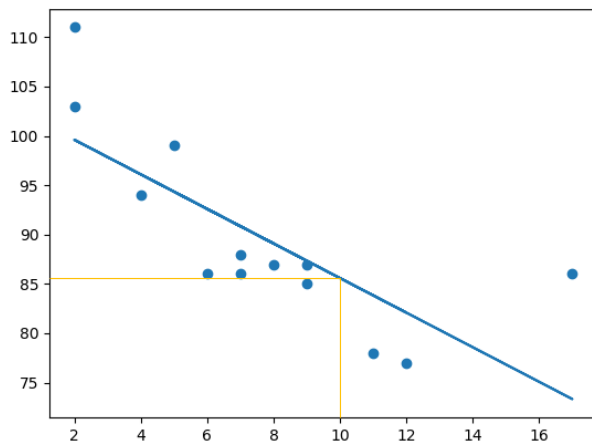
```
slope, intercept, r, p, std_err = stats.linregress(x, y)
```

```
def myfunc(x):
```

```
    return slope * x + intercept
```

```
speed = myfunc(10)
```

```
print(speed)
```



## Bad Fit?

Let us create an example where linear regression would not be the best method to predict future values. These values for the x- and y-axis should result in a very bad fit for linear regression:

```
x = [89,43,36,36,95,10,66,34,38,20,26,29,48,64,6,5,36,66,72,40]
```

```
y = [21,46,3,35,67,95,53,72,58,10,26,34,90,33,38,20,56,2,47,15]
```

```
slope, intercept, r, p, std_err = stats.linregress(x, y)
```

```
def myfunc(x):
```

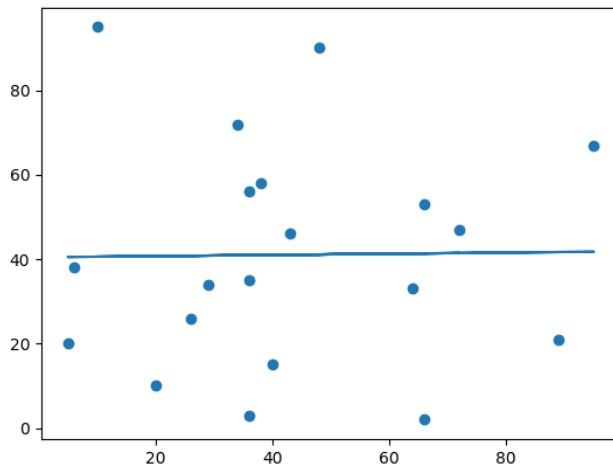
```
    return slope * x + intercept
```

```
mymodel = list(map(myfunc, x))
```

```
plt.scatter(x, y)
```

```
plt.plot(x, mymodel)
```

```
plt.show()
```



And the r for relationship; You should get a very low r value.

```
x = [89,43,36,36,95,10,66,34,38,20,26,29,48,64,6,5,36,66,72,40]
```

```
y = [21,46,3,35,67,95,53,72,58,10,26,34,90,33,38,20,56,2,47,15]
```

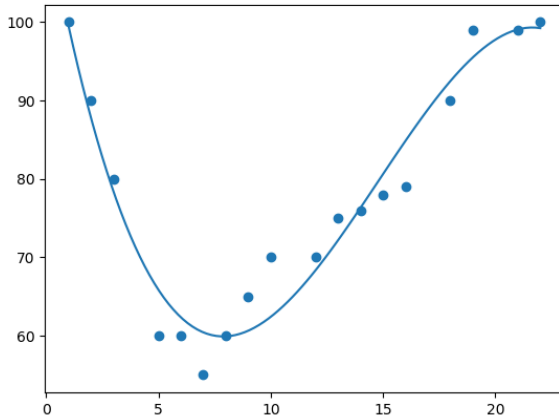
```
slope, intercept, r, p, std_err = stats.linregress(x, y)
```

```
print(r)
```

## Polynomial Regression

If your data points clearly will not fit a linear regression (a straight line through all data points), it might be ideal for polynomial regression.

Polynomial regression, like linear regression, uses the relationship between the variables  $x$  and  $y$  to find the best way to draw a line through the data points.



Python has methods for finding a relationship between data-points and to draw a line of polynomial regression. We will show you how to use these methods instead of going through the mathematic formula.

In the example below, we have registered 18 cars as they were passing a certain tollbooth.

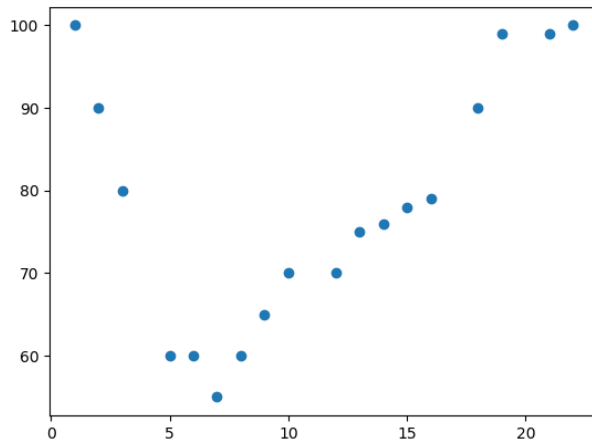
We have registered the car's speed, and the time of day (hour) the passing occurred.

The x-axis represents the hours of the day and the y-axis represents the speed:

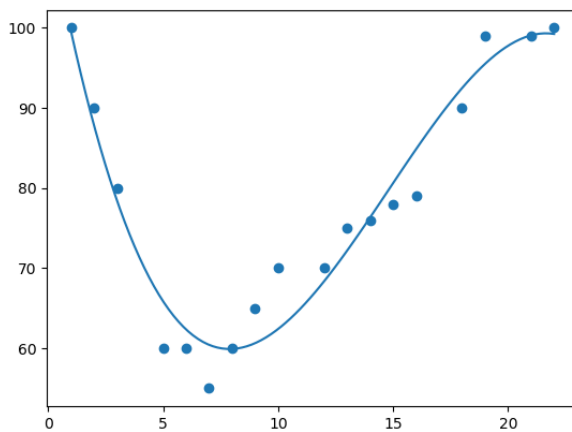


# Drawing a scatter plot:

```
import matplotlib.pyplot as plt
x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
plt.scatter(x, y)
plt.show()
```



```
import numpy
import matplotlib.pyplot as plt
x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
myline = numpy.linspace(1, 22, 100)
plt.scatter(x, y)
plt.plot(myline, mymodel(myline))
plt.show()
```



```
import numpy
import matplotlib.pyplot as plt

x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
```

Make a polynomial model	<code>mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))</code>
Then specify how the line will display	<code>myline = numpy.linspace(1, 22, 100)</code>
Draw the original scatter plot	<code>plt.scatter(x, y)</code>
Draw the line of polynomial regression	<code>plt.plot(myline, mymodel(myline))</code>
Display the diagram	<code>plt.show()</code>

## R-Squared

It is important to know how well the relationship between the values of the x- and y-axis is, if there are no relationship the polynomial regression can not be used to predict anything.

The relationship is measured with a value called the r-squared.

The r-squared value ranges from 0 to 1, where 0 means no relationship, and 1 means 100% related.

Python and the Sklearn module will compute this value for you, all you have to do is feed it with the x and y arrays:

How well does my data fit in a polynomial regression?

```
import numpy
from sklearn.metrics import r2_score
x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
print(r2_score(y, mymodel(x)))
```

**Note:** The result 0.94 shows that there is a very good relationship, and we can use polynomial regression in future predictions.

## Predict Future Values

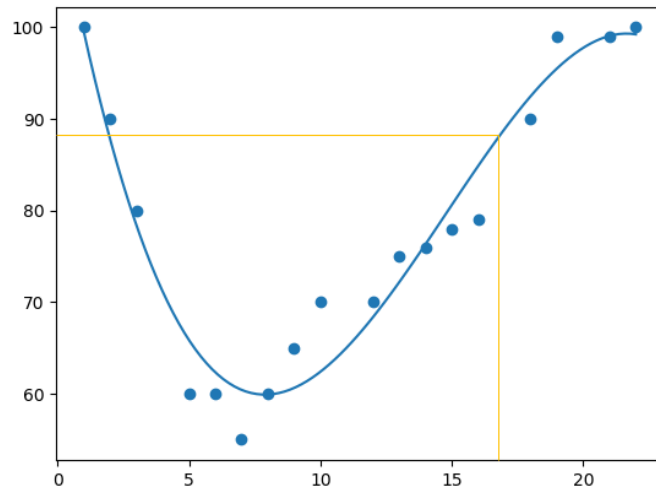
Example: Let us try to predict the speed of a car that passes the tollbooth at around 17 P.M:

```
mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
```

Predict the speed of a car passing at 17 P.M:

```
import numpy
from sklearn.metrics import r2_score
x = [1,2,3,5,6,7,8,9,10,12,13,14,15,16,18,19,21,22]
y = [100,90,80,60,60,55,60,65,70,70,75,76,78,79,90,99,99,100]
mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
speed = mymodel(17)
print(speed)
```

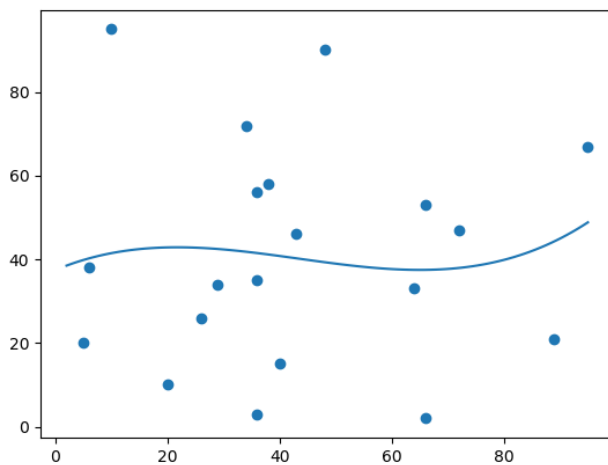
The example predicted a speed to be 88.87, which we also could read from the diagram:



## Bad Fit?

Let us create an example where polynomial regression would not be the best method to predict future values. These values for the x- and y-axis should result in a very bad fit for polynomial regression:

```
import numpy
import matplotlib.pyplot as plt
x = [89,43,36,36,95,10,66,34,38,20,26,29,48,64,6,5,36,66,72,40]
y = [21,46,3,35,67,95,53,72,58,10,26,34,90,33,38,20,56,2,47,15]
mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
myline = numpy.linspace(2, 95, 100)
plt.scatter(x, y)
plt.plot(myline, mymodel(myline))
plt.show()
```



And the r-squared value; you should get a very low r-squared value.

```
import numpy
from sklearn.metrics import r2_score
x = [89,43,36,36,95,10,66,34,38,20,26,29,48,64,6,5,36,66,72,40]
y = [21,46,3,35,67,95,53,72,58,10,26,34,90,33,38,20,56,2,47,15]
mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))
print(r2_score(y, mymodel(x)))
```

The result: 0.00995 indicates a very bad relationship, and tells us that this data set is not suitable for polynomial regression.

## Multiple Regression

Multiple regression is like linear regression, but with more than one independent value, meaning that we try to predict a value based on **two or more** variables.

We can predict the CO<sub>2</sub> emission of a car based on the size of the engine, but with multiple regression we can throw in more variables, like the weight of the car, to make the prediction more accurate.

```
import pandas # allows us to read csv files and return a DataFrame object.
```

```
df = pandas.read_csv("cars.csv")
```

```
X = df[['Weight', 'Volume']] # independent values
```

```
y = df['CO2'] # dependent values
```

```
from sklearn import linear_model # LinearRegression() method to create a linear regression object.
```

```
regr = linear_model.LinearRegression()
```

```
regr.fit(X, y)
```

This object has a method called fit() that takes the independent and dependent values as parameters and fills the regression object with data that describes the relationship

```
#predict the CO2 emission of a car where the weight is 2300kg, and the volume is 1300cm3:
```

```
predictedCO2 = regr.predict([[2300, 1300]])
```

```
import pandas
```

```
from sklearn import linear_model
```

```
df = pandas.read_csv("cars.csv")
```

```
X = df[['Weight', 'Volume']]
```

```
y = df['CO2']
```

```
regr = linear_model.LinearRegression()
```

```
regr.fit(X, y)
```

```
#predict the CO2 emission of a car where the weight is 2300kg, and the volume is 1300cm3:
```

```
predictedCO2 = regr.predict([[2300, 1300]])
```

```
print(predictedCO2) # [107.2087328]
```

We have predicted that a car with 1.3 liter engine, and a weight of 2300 kg, will release approximately 107 grams of CO<sub>2</sub> for every kilometer it drives.

## Coefficient

The coefficient is a factor that describes the relationship with an unknown variable.

Example: if  $x$  is a variable, then  $2x$  is  $x$  two times.  $x$  is the unknown variable, and the number 2 is the coefficient. In this case, we can ask for the coefficient value of weight against CO<sub>2</sub>, and for volume against CO<sub>2</sub>. The answer(s) we get tells us what would happen if we increase, or decrease, one of the independent values.

Print the coefficient values of the regression object:

```
import pandas

from sklearn import linear_model
df = pandas.read_csv("cars.csv")
X = df[['Weight', 'Volume']]
y = df['CO2']
regr = linear_model.LinearRegression()
regr.fit(X, y)
print(regr.coef_) # [0.00755095 0.00780526]
```

The result array represents the coefficient values of weight and volume.

- Weight: 0.00755095  
Volume: 0.00780526

These values tell us that if the weight increase by 1kg, the CO<sub>2</sub> emission increases by 0.00755095g. And if the engine size (Volume) increases by 1 cm<sup>3</sup>, the CO<sub>2</sub> emission increases by 0.00780526 g. We have already predicted that if a car with a 1300cm<sup>3</sup> engine weighs 2300kg, the CO<sub>2</sub> emission will be approximately 107g. What if we increase the weight with 1000kg?

```
import pandas
from sklearn import linear_model
df = pandas.read_csv("cars.csv")
X = df[['Weight', 'Volume']]
y = df['CO2']
regr = linear_model.LinearRegression()
regr.fit(X, y)
predictedCO2 = regr.predict([[3300, 1300]])
print(predictedCO2) # [114.75968007]
```

We have predicted that a car with 1.3 liter engine, and a weight of 3300 kg, will release approximately 115 grams of CO2 for every kilometer it drives.

Which shows that the coefficient of 0.00755095 is correct:

$$107.2087328 + (1000 * 0.00755095) = 114.75968$$

## Scale Features

When your data has different values, and even different measurement units, it can be difficult to compare them. What is kilograms compared to meters? Or altitude compared to time?

The answer to this problem is scaling. We can scale data into new values that are easier to compare.

Take a look at the table below, it is the same data set that we used in the multiple regression chapter, but this time the volume column contains values in *liters* instead of *cm<sup>3</sup>* (1.0 instead of 1000).

It can be difficult to compare the volume 1.0 with the weight 790, but if we scale them both into comparable values, we can easily see how much one value is compared to the other.

There are different methods for scaling data, in this tutorial we will use a method called standardization.

The standardization method uses this formula:

$$z = (x - u) / s$$

Where z is the new value, x is the original value, u is the mean and s is the standard deviation.

If you take the weight column from the data set above, the first value is 790, and the scaled value will be:

$$(790 - \text{1292.23}) / \text{238.74} = -2.1$$

If you take the volume column from the data set above, the first value is 1.0, and the scaled value will be:

$$(1.0 - \text{1.61}) / \text{0.38} = -1.59$$

Now you can compare -2.1 with -1.59 instead of comparing 790 with 1.0.

You do not have to do this manually, the Python sklearn module has a method called `StandardScaler()` which returns a `Scaler` object with methods for transforming data sets.



# Scale all values in the Weight and Volume columns:

```
import pandas
from sklearn import linear_model
from sklearn.preprocessing import StandardScaler
scale = StandardScaler()
df = pandas.read_csv("cars2.csv")
X = df[['Weight', 'Volume']]
scaledX = scale.fit_transform(X)
print(scaledX)
```

### **Predict CO2 Values**

The task in the Multiple Regression chapter was to predict the CO2 emission from a car when you only knew its weight and volume.

#Predict the CO2 emission from a 1.3 liter car that weighs 2300 kilograms

```
import pandas
from sklearn import linear_model
from sklearn.preprocessing import StandardScaler
scale = StandardScaler()
df = pandas.read_csv("cars2.csv")
X = df[['Weight', 'Volume']]
y = df['CO2']
scaledX = scale.fit_transform(X)
regr = linear_model.LinearRegression()
regr.fit(scaledX, y)
scaled = scale.transform([[2300, 1.3]])
predictedCO2 = regr.predict([scaled[0]])
print(predictedCO2) # [107.2087328]
```

## Evaluate Your Model

In Machine Learning we create models to predict the outcome of certain events, like in the previous chapter where we predicted the CO2 emission of a car when we knew the weight and engine size. To measure if the model is good enough, we can use a method called Train/Test.

### What is Train/Test

Train/Test is a method to measure the accuracy of your model.

It is called Train/Test because you split the the data set into two sets: a training set and a testing set.

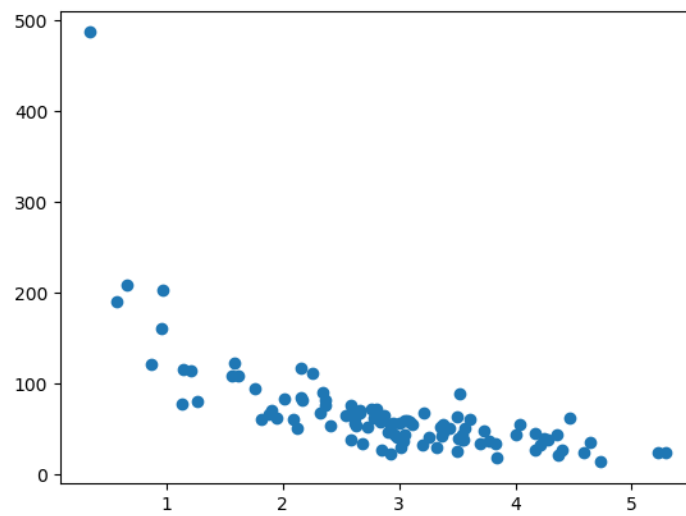
**80% for training, and 20% for testing.**

- You *train* the model using the training set.
- You *test* the model using the testing set.
  
- *Train* the model means *create* the model.
- *Test* the model means test the accuracy of the model.

### Start With a Data Set

```
import numpy
import matplotlib.pyplot as plt
numpy.random.seed(2)
x = numpy.random.normal(3, 1, 100)
y = numpy.random.normal(150, 40, 100) / x
plt.scatter(x, y)
plt.show()
```

- The x axis represents the number of minutes before making a purchase.
- The y axis represents the amount of money spent on the purchase.



### Split Into Train/Test

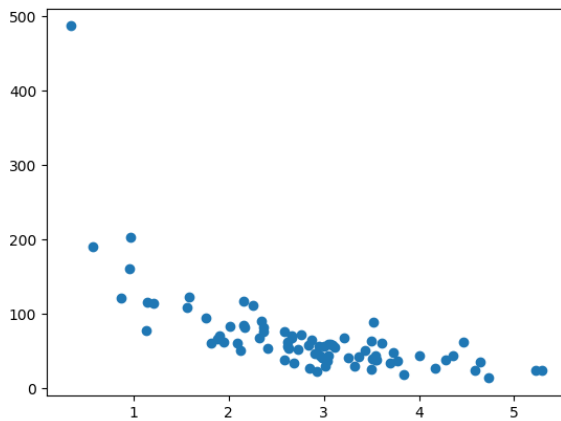
- The *training* set should be a random selection of 80% of the original data.
- The *testing* set should be the remaining 20%.

```
train_x = x[:80]  
train_y = y[:80]
```

```
test_x = x[80:]  
test_y = y[80:]
```

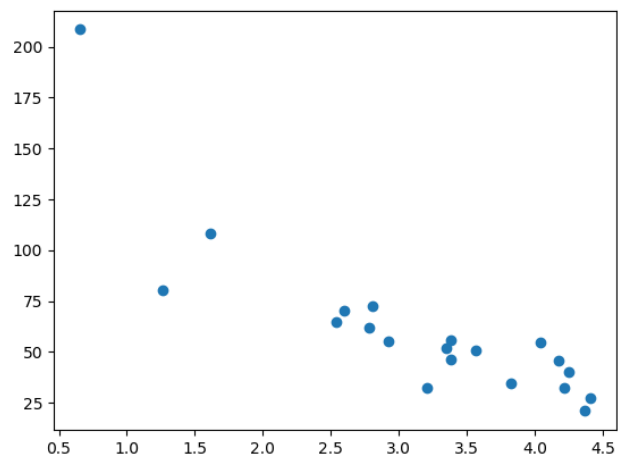
**Display the Training Set; Display the same scatter plot with the training set**

```
plt.scatter(train_x, train_y)  
plt.show()
```



Display the Testing Set; to make sure the testing set is not completely different, we will take a look at the testing set as well.

```
plt.scatter(test_x, test_y)  
plt.show()
```



## Fit the Data Set

What does the data set look like? In my opinion I think the best fit would be a [polynomial regression](#), so let us draw a line of polynomial regression.

To draw a line through the data points, we use the plot() method of the matplotlib module:

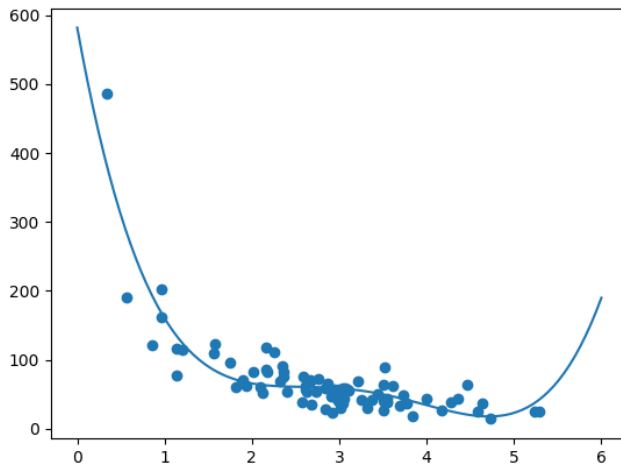
Draw a polynomial regression line through the data points:

```
import numpy
import matplotlib.pyplot as plt
numpy.random.seed(2)
x = numpy.random.normal(3, 1, 100)
y = numpy.random.normal(150, 40, 100) / x

train_x = x[:80]
train_y = y[:80]

test_x = x[80:]
test_y = y[80:]

mymodel = numpy.poly1d(numpy.polyfit(train_x, train_y, 4))
myline = numpy.linspace(0, 6, 100)
plt.scatter(train_x, train_y)
plt.plot(myline, mymodel(myline))
plt.show()
```



## R2

The R-squared score is a good indicator of how well my data set is fitting the model.

It measures the relationship between the x axis and the y axis, and the value ranges from 0 to 1, where 0 means no relationship, and 1 means totally related.

The sklearn module has a method called `r2_score()` that will help us find this relationship.

How well does my training data fit in a polynomial regression?

```
import numpy
from sklearn.metrics import r2_score
numpy.random.seed(2)

x = numpy.random.normal(3, 1, 100)
y = numpy.random.normal(150, 40, 100) / x

train_x = x[:80]
train_y = y[:80]

test_x = x[80:]
test_y = y[80:]

mymodel = numpy.poly1d(numpy.polyfit(train_x, train_y, 4))

r2 = r2_score(train_y, mymodel(train_x))

print(r2)
```

[Try it Yourself »](#)

**Note:** The result 0.799 shows that there is a OK relationship.

## Bring in the Testing Set

Let us find the R2 score when using testing data:

```
import numpy
from sklearn.metrics import r2_score
numpy.random.seed(2)
x = numpy.random.normal(3, 1, 100)
y = numpy.random.normal(150, 40, 100) / x
train_x = x[:80]
train_y = y[:80]
test_x = x[80:]
test_y = y[80:]
mymodel = numpy.poly1d(numpy.polyfit(train_x, train_y, 4))
r2 = r2_score(test_y, mymodel(test_x))
print(r2)
```

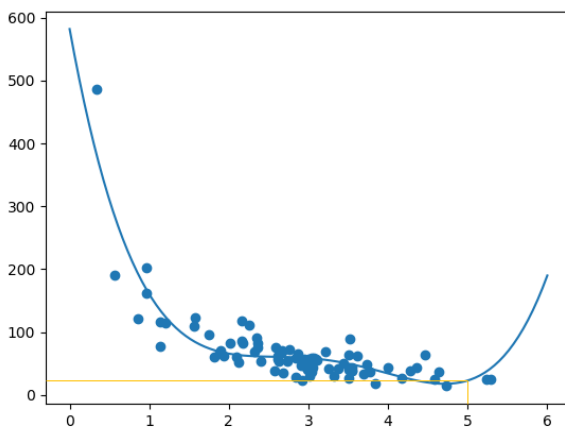
**Note:** The result 0.809 shows that the model fits the testing set as well, and we are confident that we can use the model to predict future values.

## Predict Values

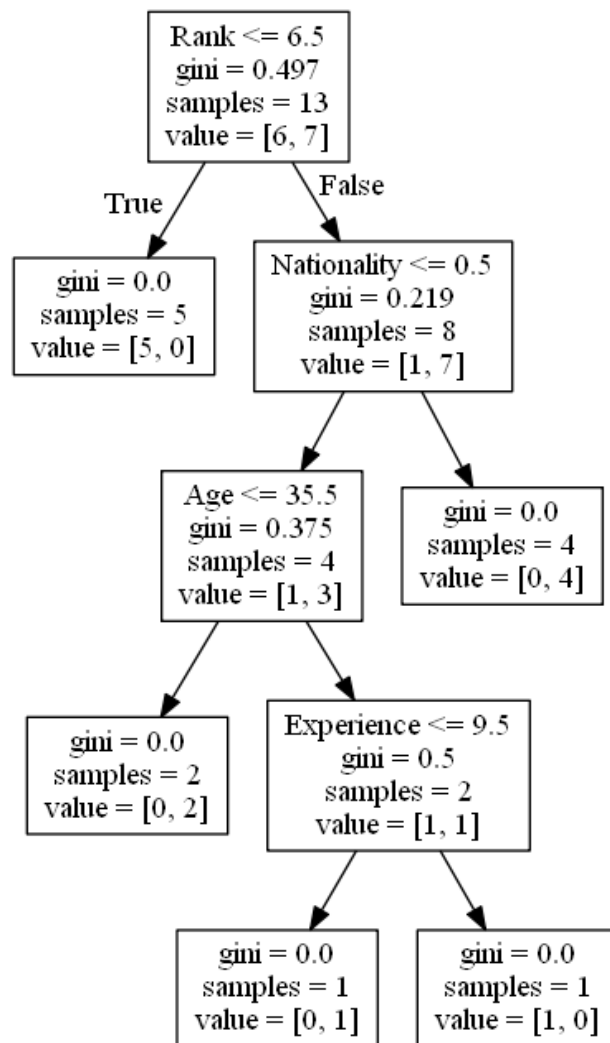
How much money will a buying customer spend, if she or he stays in the shop for 5 minutes?

```
print(mymodel(5))
```

The example predicted the customer to spend 22.88 dollars, as seems to correspond to the diagram:



### Decision Tree





```

import pandas
from sklearn import tree
import pydotplus
from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt
import matplotlib.image as pltimg
df = pandas.read_csv("shows.csv")
print(df)

```

To make a decision tree, all data has to be numerical.

We have to convert the non numerical columns 'Nationality' and 'Go' into numerical values.

Pandas has a map() method that takes a dictionary with information on how to convert the values.

```
{'UK': 0, 'USA': 1, 'N': 2}
```

Means convert the values 'UK' to 0, 'USA' to 1, and 'N' to 2.

Change string values into numerical values:

```

d = {'UK': 0, 'USA': 1, 'N': 2}
df['Nationality'] = df['Nationality'].map(d)
d = {'YES': 1, 'NO': 0}
df['Go'] = df['Go'].map(d)
print(df)

```

Then we have to separate the *feature* columns from the *target* column.

The feature columns are the columns that we try to predict *from*, and the target column is the column with the values we try to predict.

X is the feature columns, y is the target column:

```

features = ['Age', 'Experience', 'Rank', 'Nationality']
X = df[features]
y = df['Go']
print(X)
print(y)

```

Now we can create the actual decision tree, fit it with our details, and save a .png file on the computer:

Create a Decision Tree, save it as an image, and show the image:

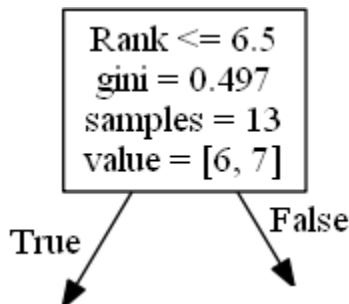
```
dtree = DecisionTreeClassifier()
dtree = dtree.fit(X, y)
data = tree.export_graphviz(dtree, out_file=None, feature_names=features)
graph = pydotplus.graph_from_dot_data(data)
graph.write_png('mydecisiontree.png')
```

```
img=pltimg.imread('mydecisiontree.png')
imgplot = plt.imshow(img)
plt.show()
```

## Result Explained

The decision tree uses your earlier decisions to calculate the odds for you to wanting to go see a comedian or not.

Let us read the different aspects of the decision tree:



Rank

Rank  $\leq 6.5$  means that every comedian with a rank of 6.5 or lower will follow the True arrow (to the left), and the rest will follow the False arrow (to the right).

gini = 0.497 refers to the quality of the split, and is always a number between 0.0 and 0.5, where 0.0 would mean all of the samples got the same result, and 0.5 would mean that the split is done exactly in the middle.

samples = 13 means that there are 13 comedians left at this point in the decision, which is all of them since this is the first step.

value = [6, 7] means that of these 13 comedians, 6 will get a "NO", and 7 will get a "GO".

Gini

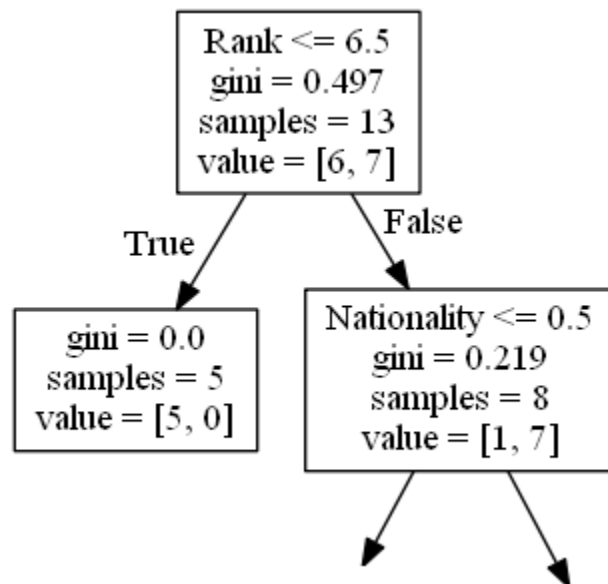
There are many ways to split the samples, we use the GINI method in this tutorial.

The Gini method uses this formula:

$$\text{Gini} = 1 - (x/n)^2 - (y/n)^2$$

Where x is the number of positive answers("GO"), n is the number of samples, and y is the number of negative answers ("NO"), which gives us this calculation:

$$1 - (7 / 13)^2 - (6 / 13)^2 = 0.497$$



The next step contains two boxes, one box for the comedians with a 'Rank' of 6.5 or lower, and one box with the rest.

**True - 5 Comedians End Here:**

gini = 0.0 means all of the samples got the same result.

samples = 5 means that there are 5 comedians left in this branch (5 comedian with a Rank of 6.5 or lower).

value = [5, 0] means that 5 will get a "NO" and 0 will get a "GO".

### False - 8 Comedians Continue:

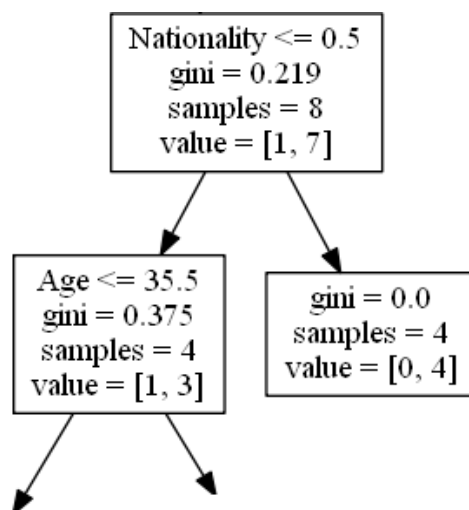
#### Nationality

Nationality  $\leq 0.5$  means that the comedians with a nationality value of less than 0.5 will follow the arrow to the left (which means everyone from the UK, ), and the rest will follow the arrow to the right.

gini = 0.219 means that about 22% of the samples would go in one direction.

samples = 8 means that there are 8 comedians left in this branch (8 comedian with a Rank higher than 6.5).

value = [1, 7] means that of these 8 comedians, 1 will get a "NO" and 7 will get a "GO".



### True - 4 Comedians Continue:

#### Age

Age  $\leq 35.5$  means that comedians at the age of 35.5 or younger will follow the arrow to the left, and the rest will follow the arrow to the right.

gini = 0.375 means that about 37,5% of the samples would go in one direction.

samples = 4 means that there are 4 comedians left in this branch (4 comedians from the UK).

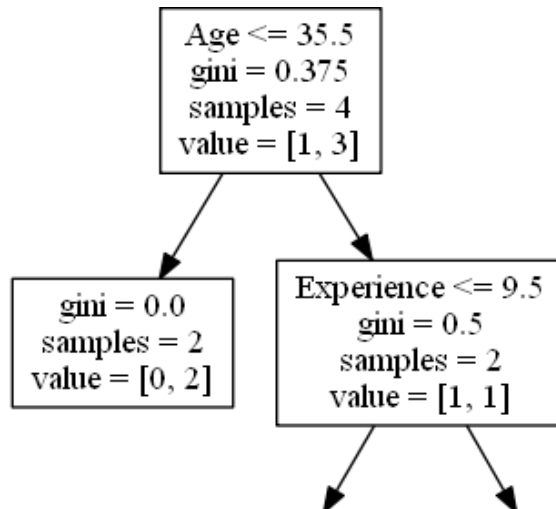
value = [1, 3] means that of these 4 comedians, 1 will get a "NO" and 3 will get a "GO".

### False - 4 Comedians End Here:

gini = 0.0 means all of the samples got the same result.

samples = 4 means that there are 4 comedians left in this branch (4 comedians not from the UK).

value = [0, 4] means that of these 4 comedians, 0 will get a "NO" and 4 will get a "GO".



### True - 2 Comedians End Here:

$\text{gini} = 0.0$  means all of the samples got the same result.

$\text{samples} = 2$  means that there are 2 comedians left in this branch (2 comedians at the age 35.5 or younger).

$\text{value} = [0, 2]$  means that of these 2 comedians, 0 will get a "NO" and 2 will get a "GO".

### False - 2 Comedians Continue:

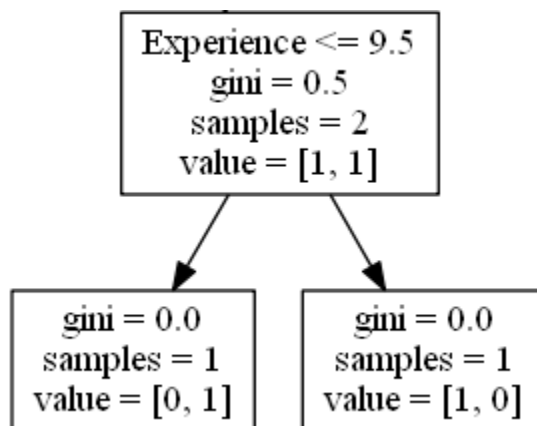
#### Experience

$\text{Experience} \leq 9.5$  means that comedians with 9.5 years of experience, or less, will follow the arrow to the left, and the rest will follow the arrow to the right.

$\text{gini} = 0.5$  means that 50% of the samples would go in one direction.

$\text{samples} = 2$  means that there are 2 comedians left in this branch (2 comedians older than 35.5).

$\text{value} = [1, 1]$  means that of these 2 comedians, 1 will get a "NO" and 1 will get a "GO".



### **True - 1 Comedian Ends Here:**

`gini = 0.0` means all of the samples got the same result.

`samples = 1` means that there is 1 comedian left in this branch (1 comedian with 9.5 years of experience or less).

`value = [0, 1]` means that 0 will get a "NO" and 1 will get a "GO".

### **False - 1 Comedian Ends Here:**

`gini = 0.0` means all of the samples got the same result.

`samples = 1` means that there is 1 comedians left in this branch (1 comedian with more than 9.5 years of experience).

`value = [1, 0]` means that 1 will get a "NO" and 0 will get a "GO".

### **Predict Values**

Example: Should I go see a show starring a 40 years old American comedian, with 10 years of experience, and a comedy ranking of 7?

Use `predict()` method to predict new values:

```
print(dtrees.predict([[40, 10, 7, 1]]))
```

What would the answer be if the comedy rank was 6?

```
print(dtrees.predict([[40, 10, 6, 1]]))
```

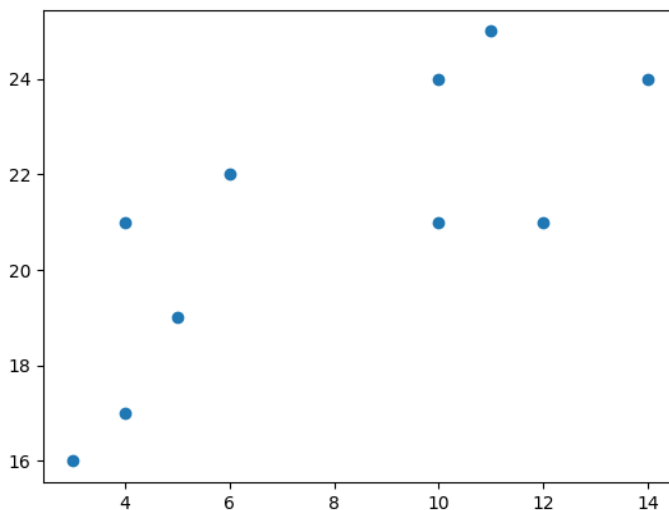
## K-means

K-means is an unsupervised learning method for clustering data points. The algorithm iteratively divides data points into K clusters by minimizing the variance in each cluster.

First, each data point is randomly assigned to one of the K clusters. Then, we compute the centroid (functionally the center) of each cluster, and reassign each data point to the cluster with the closest centroid. We repeat this process until the cluster assignments for each data point are no longer changing.

K-means clustering requires us to select K, the number of clusters we want to group the data into. The elbow method lets us graph the inertia (a distance-based metric) and visualize the point at which it starts decreasing linearly. This point is referred to as the "elbow" and is a good estimate for the best value for K based on our data.

```
import matplotlib.pyplot as plt
x = [4, 5, 10, 4, 3, 11, 14, 6, 10, 12]
y = [21, 19, 24, 17, 16, 25, 24, 22, 21, 21]
plt.scatter(x, y)
plt.show()
```



Now we utilize the elbow method to visualize the inertia for different values of K:

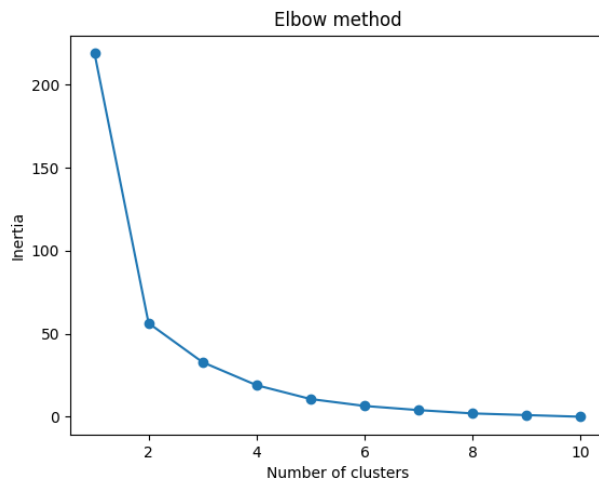
```
from sklearn.cluster import KMeans

data = list(zip(x, y))
inertias = []

for i in range(1,11):
    kmeans = KMeans(n_clusters=i)
    kmeans.fit(data)
    inertias.append(kmeans.inertia_)

plt.plot(range(1,11), inertias, marker='o')
plt.title('Elbow method')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.show()
```

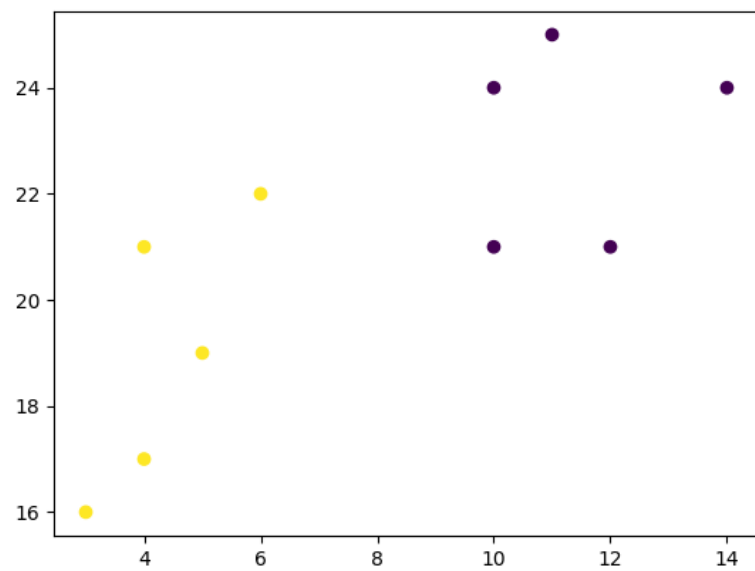
Result



The elbow method shows that 2 is a good value for K, so we retrain and visualize the result:

```
kmeans = KMeans(n_clusters=2)
kmeans.fit(data)
plt.scatter(x, y, c=kmeans.labels_)
plt.show()
```





```
import matplotlib.pyplot as plt

from sklearn.cluster import KMeans

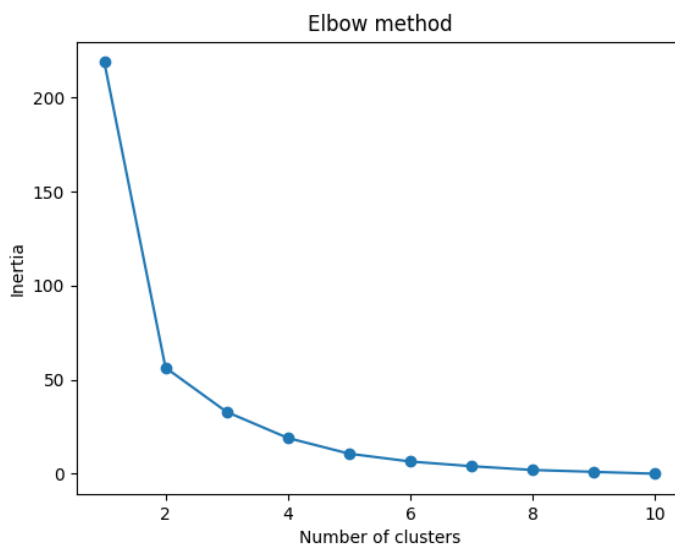
x = [4, 5, 10, 4, 3, 11, 14, 6, 10, 12]
y = [21, 19, 24, 17, 16, 25, 24, 22, 21, 21]

data = list(zip(x, y))
print(data)
```

```
[(4, 21), (5, 19), (10, 24), (4, 17), (3, 16), (11, 25), (14, 24), (6, 22), (10, 21), (12, 21)]
```

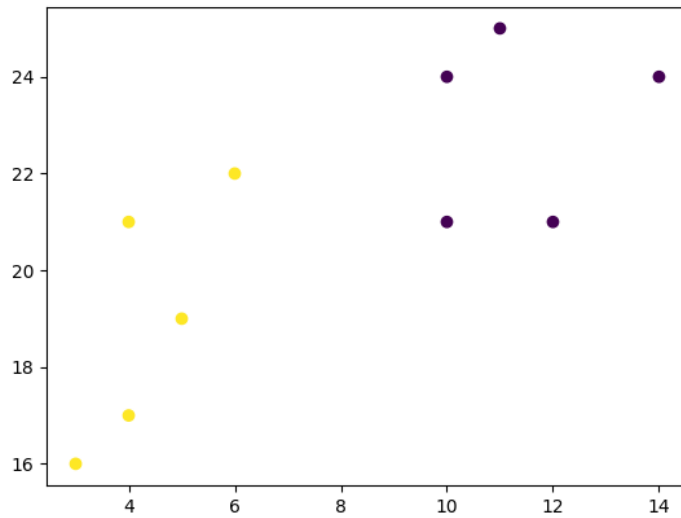
In order to find the best value for K, we need to run K-means across our data for a range of possible values. We only have 10 data points, so the maximum number of clusters is 10. So for each value K in range(1,11), we train a K-means model and plot the inertia at that number of clusters:

```
inertias = []
for i in range(1,11):
    kmeans = KMeans(n_clusters=i)
    kmeans.fit(data)
    inertias.append(kmeans.inertia_)
# for
plt.plot(range(1,11), inertias, marker='o')
plt.title('Elbow method')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.show()
```



We can see that the "elbow" on the graph above (where the inertia becomes more linear) is at K=2. We can then fit our K-means algorithm one more time and plot the different clusters assigned to the data:

```
kmeans = KMeans(n_clusters=2)
kmeans.fit(data)
plt.scatter(x, y, c=kmeans.labels_)
plt.show()
```



## AUC - ROC Curve

In classification, there are many different evaluation metrics. The most popular is accuracy, which measures how often the model is correct. This is a great metric because it is easy to understand and getting the most correct guesses is often desired. There are some cases where you might consider using another evaluation metric.

Another common metric is AUC, area under the receiver operating characteristic (ROC) curve. The Receiver operating characteristic curve plots the true positive (TP) rate versus the false positive (FP) rate at different classification thresholds. The thresholds are different probability cutoffs that separate the two classes in binary classification. It uses probability to tell us how well a model separates the classes.

## Imbalanced Data

Suppose we have an imbalanced data set where the majority of our data is of one value. We can obtain high accuracy for the model by predicting the majority class.

```
import numpy as np
from sklearn.metrics import accuracy_score, confusion_matrix, roc_auc_score, roc_curve
n = 10000
ratio = .95
n_0 = int((1-ratio) * n)
n_1 = int(ratio * n)
y = np.array([0] * n_0 + [1] * n_1)

# below are the probabilities obtained from a hypothetical model that always predicts the majority class
# probability of predicting class 1 is going to be 100%

y_proba = np.array([1]*n)
y_pred = y_proba > .5
print(f'accuracy score: {accuracy_score(y, y_pred)}')
cf_mat = confusion_matrix(y, y_pred)
print('Confusion matrix')
print(cf_mat)
print(f'class 0 accuracy: {cf_mat[0][0]/n_0}')
print(f'class 1 accuracy: {cf_mat[1][1]/n_1}')
```

Although we obtain a very high accuracy, the model provided no information about the data so it's not useful. We accurately predict class 1 100% of the time while inaccurately predict class 0 0% of the time. At the expense of accuracy, it might be better to have a model that can somewhat separate the two classes.

Example

```
# below are the probabilities obtained from a hypothetical model that doesn't always predict the mode
y_proba_2 = np.array(
    np.random.uniform(0, .7, n_0).tolist() +
    np.random.uniform(.3, 1, n_1).tolist()
)
y_pred_2 = y_proba_2 > .5

print(f'accuracy score: {accuracy_score(y, y_pred_2)}')
cf_mat = confusion_matrix(y, y_pred_2)
print('Confusion matrix')
print(cf_mat)
print(f'class 0 accuracy: {cf_mat[0][0]/n_0}')
print(f'class 1 accuracy: {cf_mat[1][1]/n_1}')
```

For the second set of predictions, we do not have as high of an accuracy score as the first but the accuracy for each class is more balanced. Using accuracy as an evaluation metric we would rate the first model higher than the second even though it doesn't tell us anything about the data.

In cases like this, using another evaluation metric like AUC would be preferred.

```
import matplotlib.pyplot as plt
```

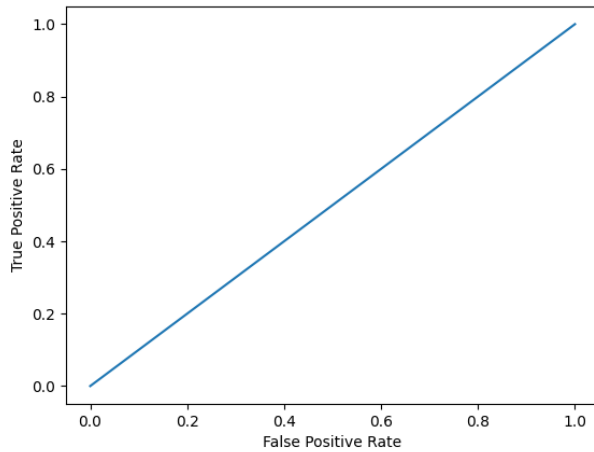
```
def plot_roc_curve(true_y, y_prob):
    """
    plots the roc curve based of the probabilities
    """

    fpr, tpr, thresholds = roc_curve(true_y, y_prob)
    plt.plot(fpr, tpr)
    plt.xlabel('False Positive Rate')
    plt.ylabel('True Positive Rate')
```

### Model 1:

```
plot_roc_curve(y, y_proba)
print(f'model 1 AUC score: {roc_auc_score(y, y_proba)}')
```

Result

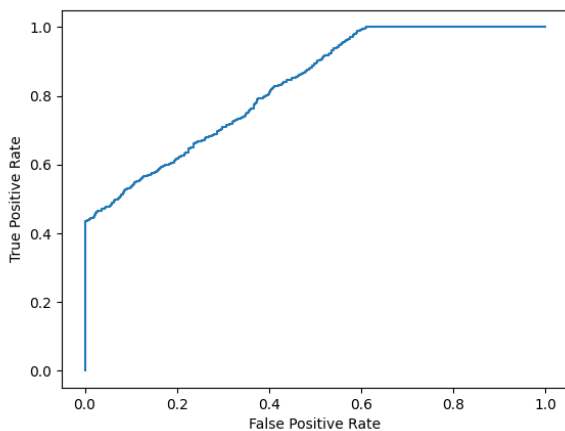


model 1 AUC score: 0.5

### Model 2:

```
plot_roc_curve(y, y_proba_2)
print(f'model 2 AUC score: {roc_auc_score(y, y_proba_2)}')
```

Result



model 2 AUC score: 0.8270551578947367

An AUC score of around .5 would mean that the model is unable to make a distinction between the two classes and the curve would look like a line with a slope of 1. An AUC score closer to 1 means that the model has the ability to separate the two classes and the curve would come closer to the top left corner of the graph.

## Probabilities

Because AUC is a metric that utilizes probabilities of the class predictions, we can be more confident in a model that has a higher AUC score than one with a lower score even if they have similar accuracies.

In the data below, we have two sets of probabilities from hypothetical models. The first has probabilities that are not as "confident" when predicting the two classes (the probabilities are close to .5). The second has probabilities that are more "confident" when predicting the two classes (the probabilities are close to the extremes of 0 or 1).

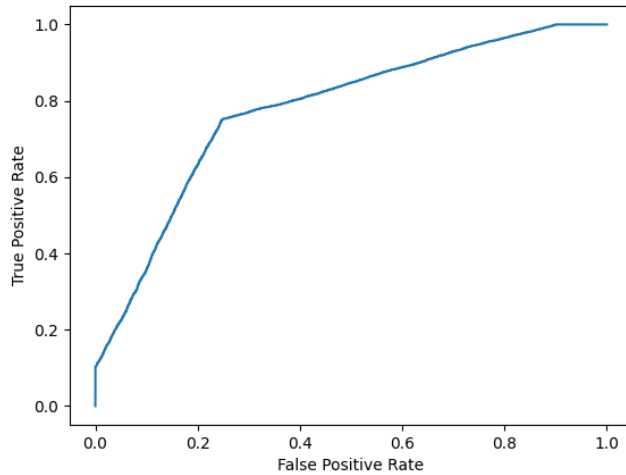
```
import numpy as np
n = 10000
y = np.array([0] * n + [1] * n)
#
y_prob_1 = np.array(
    np.random.uniform(.25, .5, n//2).tolist() +
    np.random.uniform(.3, .7, n).tolist() +
    np.random.uniform(.5, .75, n//2).tolist()
)
y_prob_2 = np.array(
    np.random.uniform(0, .4, n//2).tolist() +
    np.random.uniform(.3, .7, n).tolist() +
    np.random.uniform(.6, 1, n//2).tolist()
)

print(f'model 1 accuracy score: {accuracy_score(y, y_prob_1>.5)}')
print(f'model 2 accuracy score: {accuracy_score(y, y_prob_2>.5)}')

print(f'model 1 AUC score: {roc_auc_score(y, y_prob_1)}')
print(f'model 2 AUC score: {roc_auc_score(y, y_prob_2)}')
```

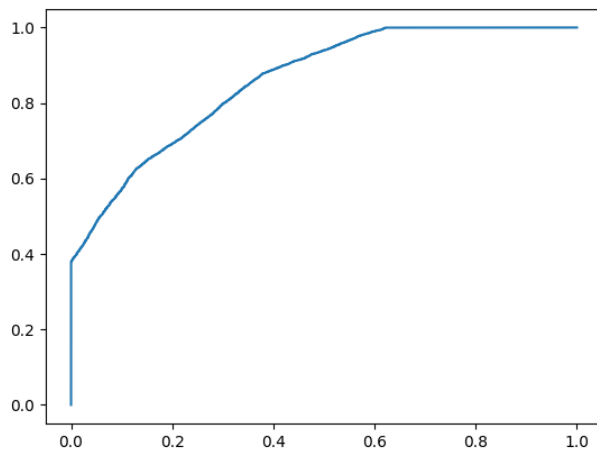
### Plot model 1:

```
plot_roc_curve(y, y_prob_1)
```



### Plot model 2:

```
fpr, tpr, thresholds = roc_curve(y, y_prob_2)  
plt.plot(fpr, tpr)
```



Even though the accuracies for the two models are similar, the model with the higher AUC score will be more reliable because it takes into account the predicted probability. It is more likely to give you higher accuracy when predicting future data.



## Statistics Methods

<code>statistics.harmonic_mean()</code>	Harmonic mean (central location) of the given data
<code>statistics.mean()</code>	Mean (average) of the given data
<code>statistics.median()</code>	Median (middle value) of the given data
<code>statistics.median_grouped()</code>	Median of grouped continuous data
<code>statistics.median_high()</code>	High median of the given data
<code>statistics.median_low()</code>	Low median of the given data
<code>statistics.mode()</code>	Mode (central tendency) of the given numeric or nominal data
<code>statistics.pstdev()</code>	Standard deviation from an entire population
<code>statistics.stdev()</code>	Standard deviation from a sample of data
<code>statistics.pvariance()</code>	Variance of an entire population
<code>statistics.variance()</code>	Variance from a sample of data

## Python Random Module

<code>seed()</code>	Initialize the random number generator
<code>getstate()</code>	Returns the current internal state of the random number generator
<code>setstate()</code>	Restores the internal state of the random number generator
<code>getrandbits()</code>	Returns a number representing the random bits
<code>randrange()</code>	Returns a random number between the given range
<code>randint()</code>	Returns a random number between the given range
<code>choice()</code>	Returns a random element from the given sequence
<code>choices()</code>	Returns a list with a random selection from the given sequence
<code>shuffle()</code>	Takes a sequence and returns the sequence in a random order
<code>sample()</code>	Returns a given sample of a sequence
<code>random()</code>	Returns a random float number between 0 and 1
<code>uniform()</code>	Returns a random float number between two given parameters

`triangular()`

Returns a random float number between two given parameters, you can also set a mode parameter to specify the midpoint between the two other parameters

`betavariate()`

Returns a random float number between 0 and 1 based on the Beta distribution (used in statistics)

`expovariate()`

Returns a random float number based on the Exponential distribution (used in statistics)

`gammavariate()`

Returns a random float number based on the Gamma distribution (used in statistics)

`gauss()`

Returns a random float number based on the Gaussian distribution (used in probability theories)

`lognormvariate()`

Returns a random float number based on a log-normal distribution (used in probability theories)

`normalvariate()`

Returns a random float number based on the normal distribution (used in probability theories)

`vonmisesvariate()`

Returns a random float number based on the von Mises distribution (used in directional statistics)

`paretovariate()`

Returns a random float number based on the Pareto distribution (used in probability theories)

`weibullvariate()`

Returns a random float number based on the Weibull distribution (used in statistics)