

Unit 3 – Correlation and Regression

Change variables as needed to observe how the change in data points impacts correlation and regression.

Covariance correlation

The initial code shows a strong positive correlation between the two variables.

```
# calculate the Pearson's correlation between two variables
from numpy import mean
from numpy import std
from numpy import cov
from numpy.random import randn
from numpy.random import seed
from matplotlib import pyplot as plt
import seaborn as sns

from scipy.stats import pearsonr
# seed random number generator
seed(1)

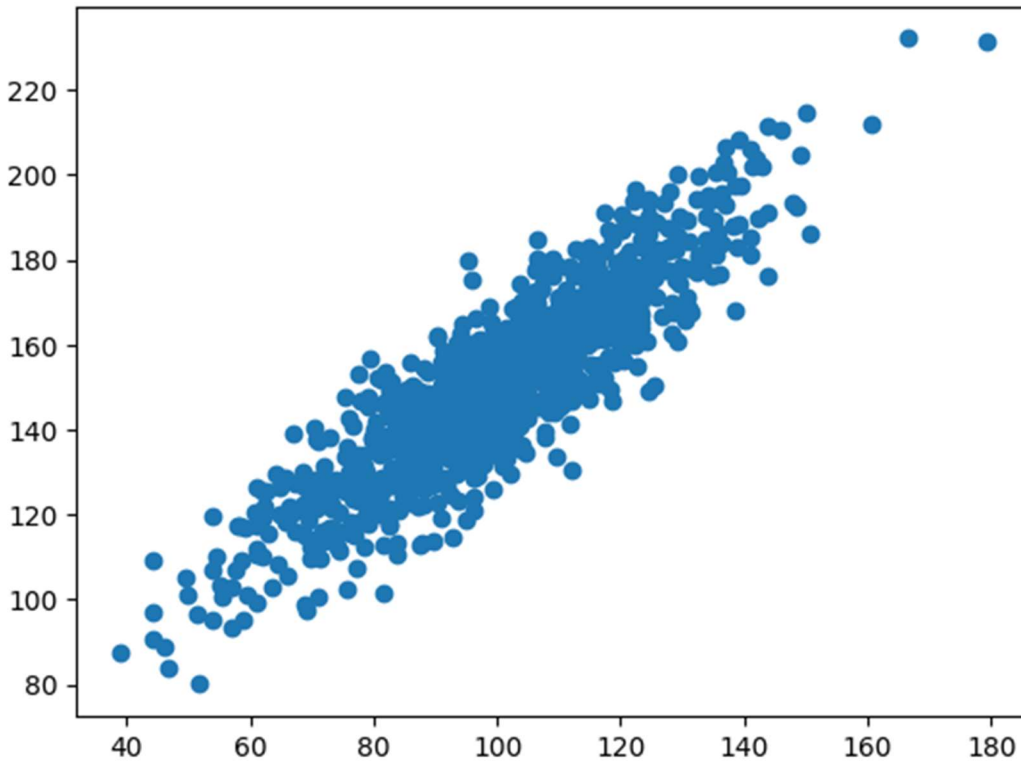
# prepare data
data1 = 20 * randn(1000) + 100
data2 = data1 + (10 * randn(1000) + 50)

# calculate covariance matrix
covariance = cov(data1, data2)

# calculate Pearson's correlation
corr, _ = pearsonr(data1, data2)

# plot
plt.scatter(data1, data2)
plt.show()

# summarize
print('data1: mean=%.3f stdv=%.3f' % (mean(data1), std(data1)))
print('data2: mean=%.3f stdv=%.3f' % (mean(data2), std(data2)))
print('Covariance: %.3f' % covariance[0][1])
print('Pearsons correlation: %.3f' % corr)
```



```
data1: mean=100.776 stdv=19.620  
data2: mean=151.050 stdv=22.358  
Covariance: 389.755  
Pearsons correlation: 0.888
```

When changing data and multiplying the data1 figure by 2 instead of 20 we can see less of a strong positive correlation, which now shows a correlation figure of 0.208 instead of 0.888 as seen with the prior result

```
# calculate the Pearson's correlation between two variables  
from numpy import mean
```

```
from numpy import std
from numpy import cov
from numpy.random import randn
from numpy.random import seed
from matplotlib import pyplot as plt
import seaborn as sns

from scipy.stats import pearsonr
# seed random number generator
seed(1)

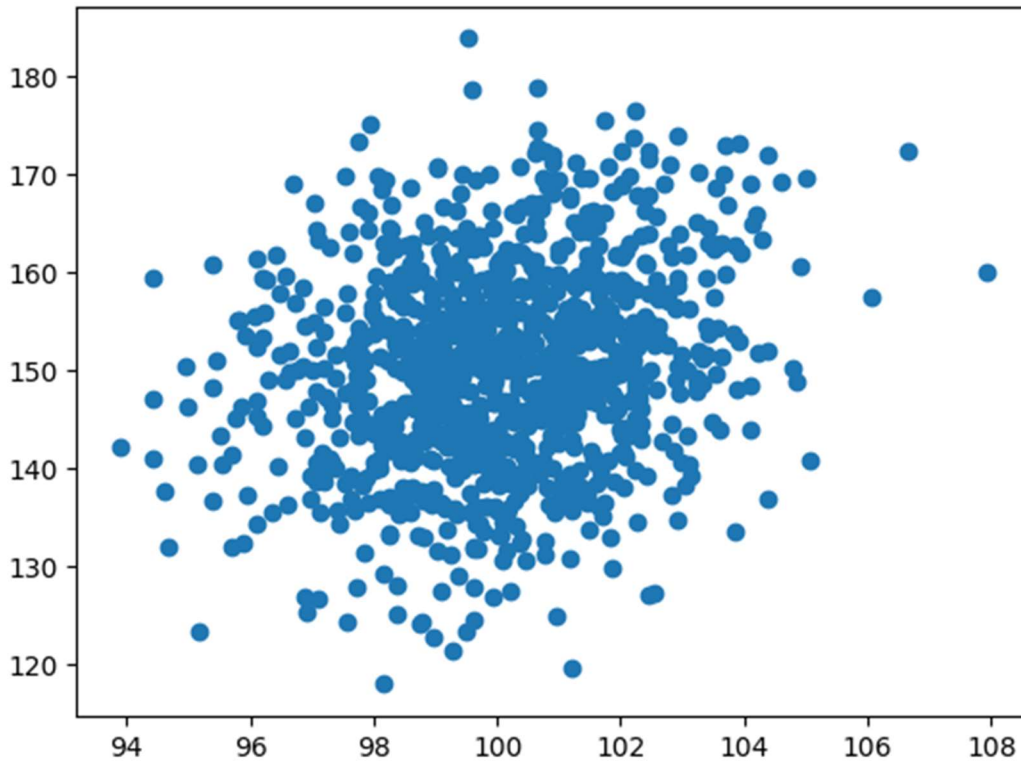
# prepare data
data1 = 2 * randn(1000) + 100
data2 = data1 + (10 * randn(1000) + 50)

# calculate covariance matrix
covariance = cov(data1, data2)

# calculate Pearson's correlation
corr, _ = pearsonr(data1, data2)

# plot
plt.scatter(data1, data2)
plt.show()

# summarize
print('data1: mean=%.3f stdv=%.3f' % (mean(data1), std(data1)))
print('data2: mean=%.3f stdv=%.3f' % (mean(data2), std(data2)))
print('Covariance: %.3f' % covariance[0][1])
print('Pearsons correlation: %.3f' % corr)
```



```
data1: mean=100.078 stdv=1.962  
data2: mean=150.351 stdv=10.528  
Covariance: 4.295  
Pearsons correlation: 0.208
```

Linear Regression

```
import matplotlib.pyplot as plt
from scipy import stats

#Create the arrays that represent the values of the x and y axis
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)

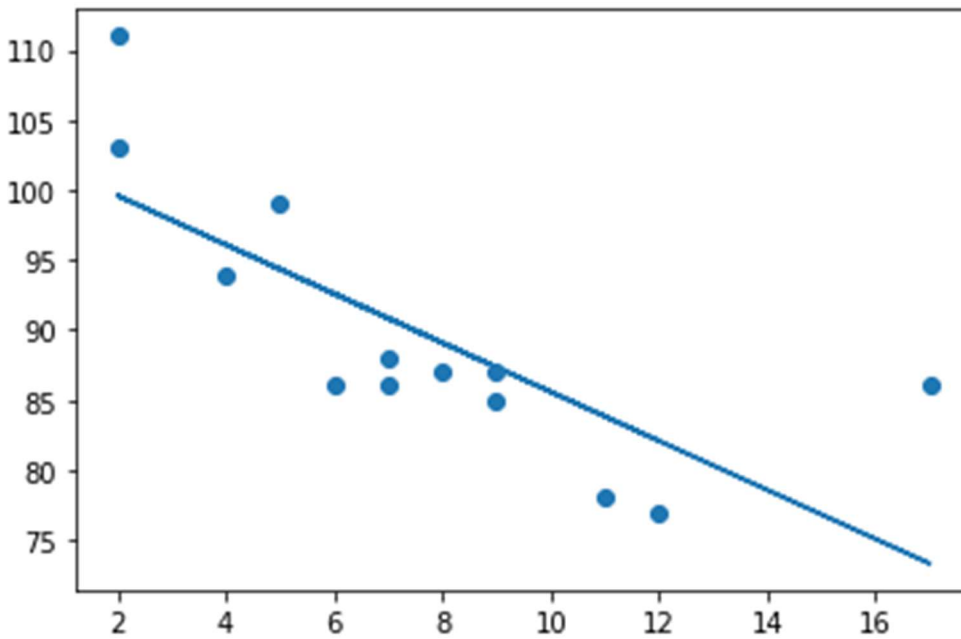
# measure the correlation
corr, _ = stats.pearsonr(x, y)
print('Pearsons correlation: %.3f' % corr)

#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. This will result in a
new array with new values for the y-axis
mymodel = list(map(myfunc, x))

#Draw the original scatter plot & the line of linear regression
plt.scatter(x, y)
plt.plot(x, mymodel)
plt.show()
```

Pearsons correlation: -0.759



To predict future values we use the below formulas

```
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 intercept + slope * x

speed = myfunc(10)

print(speed)
```

85.59308314937454

The above figure of 85.59308314937454 shows that if x were to equal 10 then the predicted figure would be 85.6 to one dp.

```

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 intercept + slope * x

speed = myfunc(15)

print(speed)

```

when the figure is changed to 15 for x then y is predicted to be 76.83664459161147

If x is 20 then y is predicted to be 68.08020603384843

Multiple Linear Regression

```

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]

```

When the weight is 2300kg and the volume is 1300cm³ then the predicted co₂ is 107.2

Coeffecient

The coefficient is a factor that describes the relationship with an unknown variable. In this case, we can ask for the coefficient value of weight against CO₂, and for volume against CO₂. The answer(s) we get tells us what would happen if we increase, or decrease, one of the independent values.

```
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 CO2 emission increases by 0.00755095g.

And if the engine size (Volume) increases by 1 cm3, the CO2 emission increases by 0.00780526 g.

I think that is a fair guess, but let test it!

We have already predicted that if a car with a 1300cm3 engine weighs 2300kg, the CO2 emission will be approximately 107g.

What if we increase the weight with 1000kg (from 2300 to 3300) what will be the CO2 emission?

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

```
predictedCO2 = regr.predict([[3300, 1300]])
```

```
print(predictedCO2)
```

```
[114.75968007]
```


Polynomial Regression

18 cars passing a certain tollboth at different time of the day (x) with different speed (y)

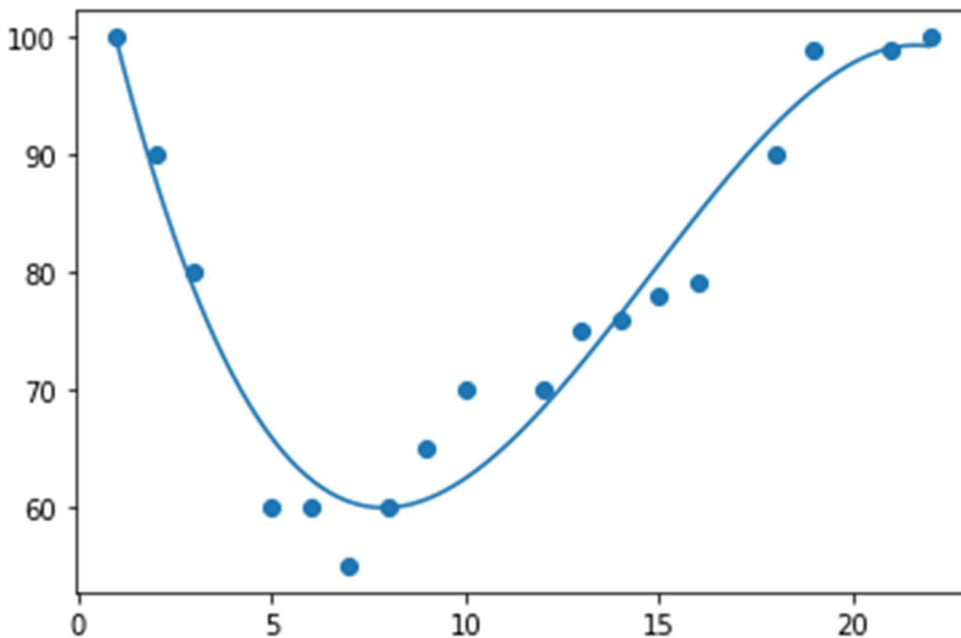
```
import numpy
import matplotlib.pyplot as plt

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]

#NumPy has a method that lets us make a polynomial model
mymodel = numpy.poly1d(numpy.polyfit(x, y, 3))

#specify how the line will display, we start at position 1, and end at
position 22
myline = numpy.linspace(1, 22, 100)

plt.scatter(x, y)
plt.plot(myline, mymodel(myline))
plt.show()
```



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

```
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)))
```

0.9432150416451027

Let us try to predict the speed of a car that passes the tollbooth at around 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)
```

88.87331269697987

At 17:00 the speed would be 88mph, if we change the figure to 20:00 for example then the speed would be 97mph as shown in the code below

```
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(20)
print(speed)
```

97.76768832747528