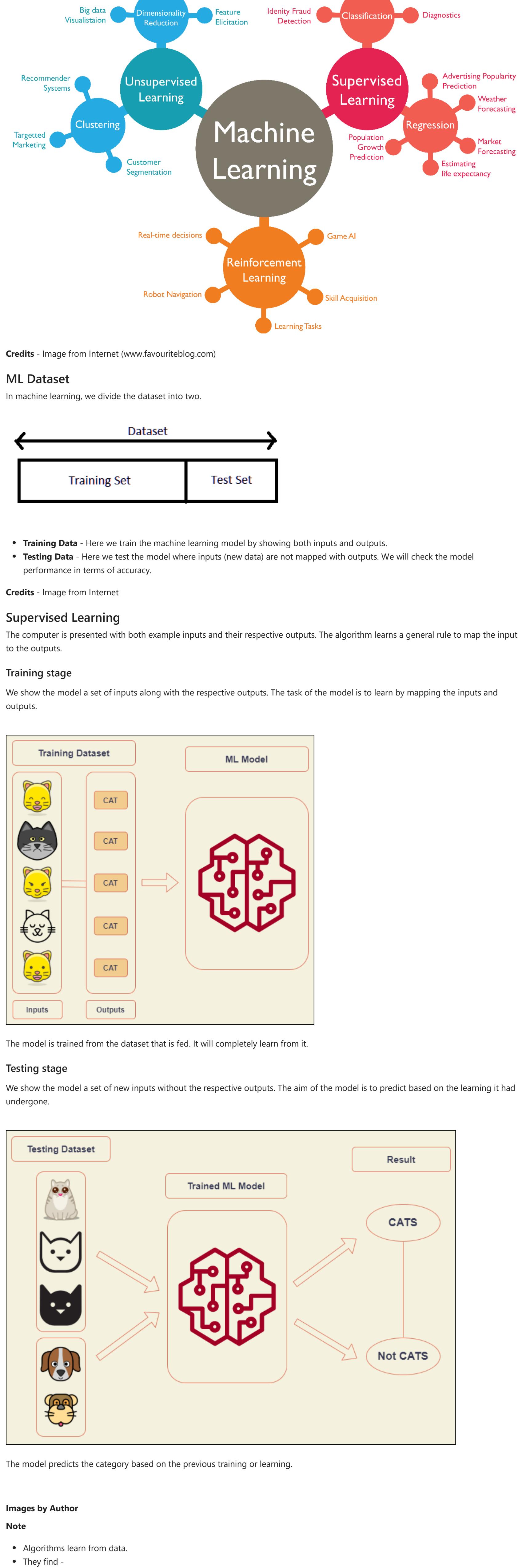


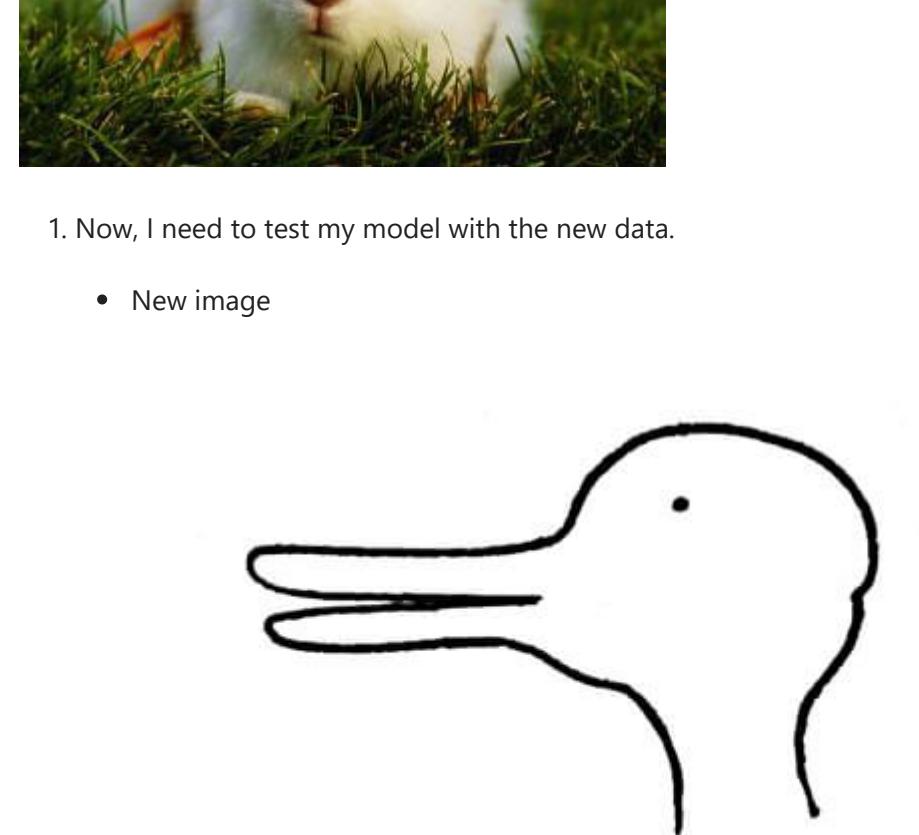
## Machine Learning



Credits - Image from Internet ([www.favouriteblog.com](http://www.favouriteblog.com))

### ML Dataset

In machine learning, we divide the dataset into two.



- Training Data** - Here we train the machine learning model by showing both inputs and outputs.
- Testing Data** - Here we test the model where inputs (new data) are not mapped with outputs. We will check the model performance in terms of accuracy.

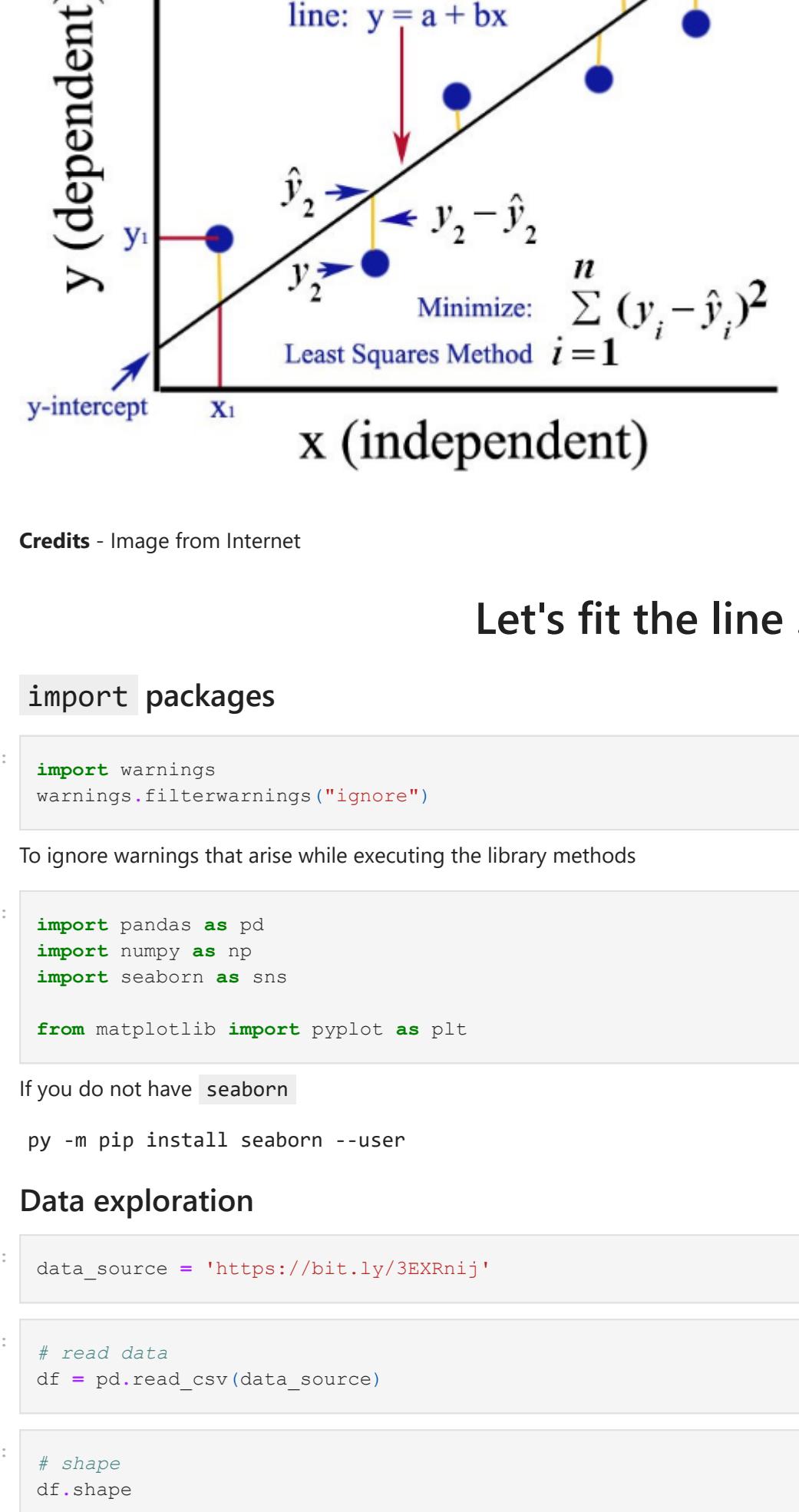
Credits - Image from Internet

### Supervised Learning

The computer is presented with both example inputs and their respective outputs. The algorithm learns a general rule to map the inputs to the outputs.

#### Training stage

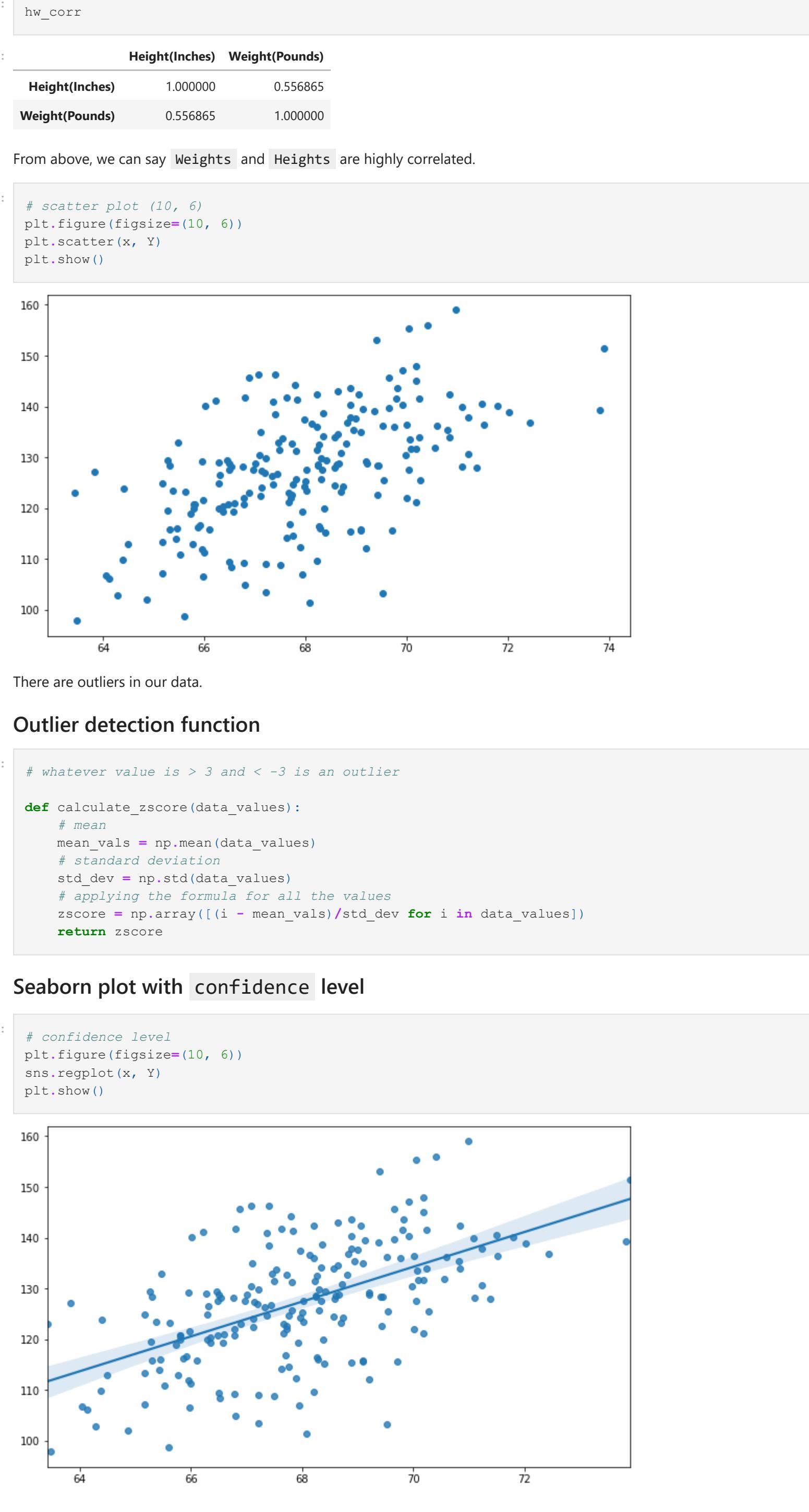
We show the model a set of inputs along with the respective outputs. The task of the model is to learn by mapping the inputs and outputs.



The model is trained from the dataset that is fed. It will completely learn from it.

#### Testing stage

We show the model a set of new inputs without the respective outputs. The aim of the model is to predict based on the learning it had undergone.



The model predicts the category based on the previous training or learning.

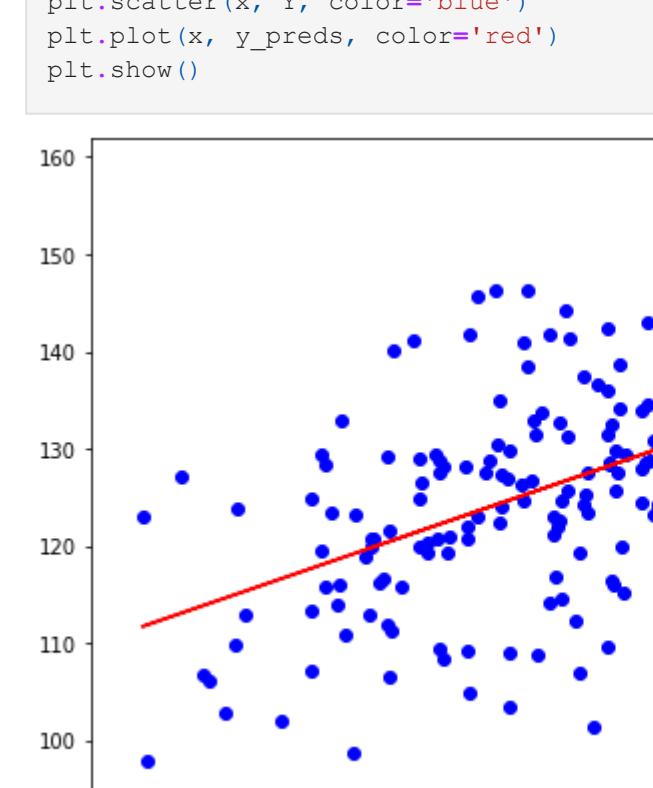
### Images by Author

#### Note

- Algorithms learn from data.
- They find -
  - relationships
  - develop understanding
  - make decisions
  - evaluate their confidence from the training data they are given.
- The better the training data is, the better the model performs.

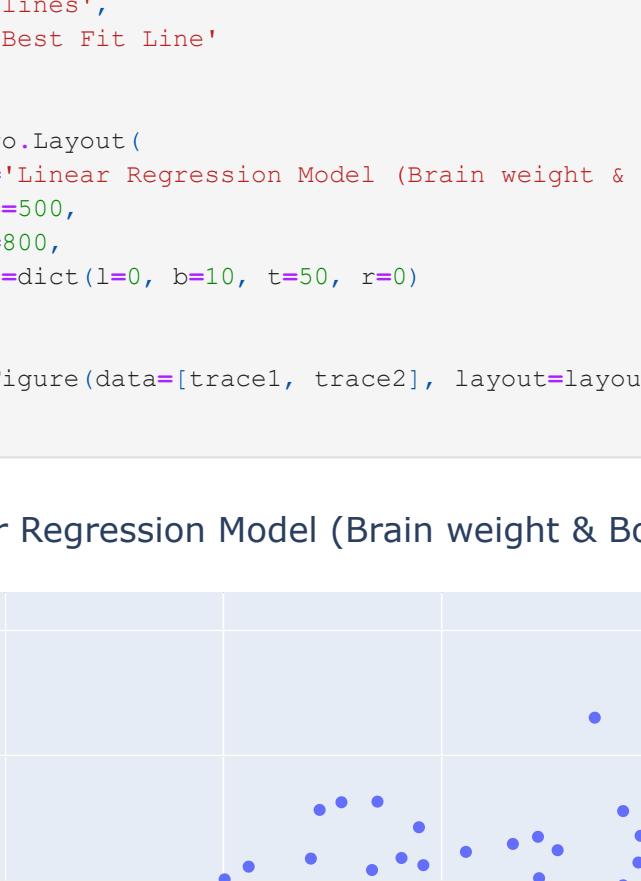
### Exception case for the above example

1. Suppose I have trained my model to identify/separate `duck` and `rabbit` images from the large dataset.

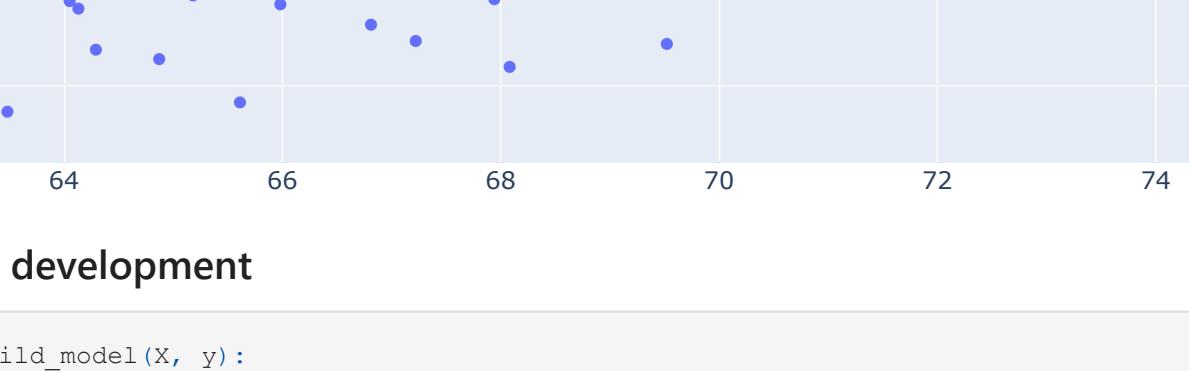


1. Now, I need to test my model with the new data.

- New image



- Is it duck or rabbit?



2. What can you say about this?

Credits - Images from Internet

### Regression

- Regression is a process of predicting the dependant variable based on the independant variable.
- Dependant variable is always considered as  $y$ .
- Independant variable is always considered as  $x$ .
- For doing regression analysis, there needs to be a strong relationship between  $x$  and  $y$ .

Example - Predicting the expenses of an employee based on his/her income.

- Here,
  - $y$  → dependant variable
  - $x$  → independant variable

How can we predict  $y$  from  $x$ ?

To predict  $y$  from  $x$ , we follow a mathematical model -

$$Y = bx + a$$

- Here,

- $b$  → Co-efficient parameter
- $a$  → Bias parameter

- Other terms

- $\hat{y}$  (Y hat) → Predicted Y value
- $y$  → Actual value

#### What is Error?

- The distance between  $\hat{y}$  and  $y$  is called Error.
- Our aim while building the Regression model is to minimize the Error.



Credits - Image from Internet

**Let's fit the line ... !**

**import packages**

In [1]:  
import warnings  
warnings.filterwarnings("ignore")

To ignore warnings that arise while executing the library methods

In [2]:  
import pandas as pd  
import numpy as np  
import seaborn as sns

from matplotlib import pyplot as plt

If you do not have seaborn

py -m pip install seaborn --user

**Data exploration**

In [3]:  
data\_source = 'https://bit.ly/3EXRnij'

In [4]:  
# read data  
df = pd.read\_csv(data\_source)

In [5]:  
# shape  
df.shape

Out[5]: (200, 2)

In [6]:  
# head  
df.head()

Out[6]: Height(inches) Weight(Pounds)

0 65.78 112.99

1 71.52 136.49

2 69.40 153.03

3 68.22 142.34

4 67.79 144.30

- Dependant variable ( $y$ ) → Weight
- Independant variable ( $x$ ) → Height

In [7]:  
# Create Y and x (array)  
Y = np.array(df['Weight(Pounds)'].to\_list())

x = np.array(df['Height(inches)'].to\_list())

How  $x$  is related to  $y$ ?

In [8]:  
# correlation  
hw\_corr = df.corr()

Out[9]:

Height(inches) Weight(Pounds)

Height(inches) 1.000000 0.556865

Weight(Pounds) 0.556865 1.000000

From above, we can say Weights and Heights are highly correlated.

In [10]:  
# scatter plot (10, 6)  
plt.figure(figsize=(10, 6))  
plt.scatter(x, Y, color='blue')  
plt.show()

There are outliers in our data.

**Outlier detection function**

In [11]:  
# whatever value is > 3 and < -3 is an outlier

def calculate\_zscore(data\_values):  
 # mean  
 mean\_vals = np.mean(data\_values)  
 # standard deviation  
 std\_dev = np.std(data\_values)

# applying the formula for all the values  
 zscore = np.array([(i - mean\_vals)/std\_dev for i in data\_values])  
 return zscore

Seaborn plot with confidence level

In [12]:  
# confidence level  
plt.figure(figsize=(10, 6))  
sns.replot(x=x, y=Y)  
plt.show()

160

150

140

130

120

110

100

64 66 68 70 72 74

What about Multiple Linear Regression model?

Dependant variable is depends on mutiple Independant variables.

**Model**

In [13]:  
def build\_model(X, y):  
 x\_dim = 1 if (len(X.shape) == 1) else X.shape[1]

b, a = np.polyfit(x=X, y=y, deg=x\_dim)

return b, a

Out[14]:

3.43267612971628

In [15]:  
# b  
b

Out[15]:

# a

Out[15]: 106.0270644078133

**Broadcasting**

In [16]:  
f = [1, 2, 3, 4]  
s = f \* 3  
print(s)

[1, 2, 3, 4, 1, 2, 3, 4, 1, 2, 3, 4]

In [17]:  
g = np.array([1, 2, 3, 4])  
h = g \* 3  
print(h)

[ 3 6 9 12 ]

In [18]:  
h \* 4

Out[18]: array([ 7, 10, 13, 16])

In [19]:  
type(x)

Out[19]: numpy.ndarray

In [20]:  
# y\_preds = bx + a  
y\_preds = (b \* x) + a

Out[20]:

# y\_preds  
# print(y\_preds)

Out[20]:

Best fit line

In [21]:  
# scatter and line  
plt.figure(figsize=(10, 6))  
plt.scatter(x, Y, color='blue')  
plt.plot(x, y, color='red')

160

150

140

130

120

110

100

64 66 68 70 72 74

Linear Regression Model (Brain weight & Body weight)

160

150

140

130

120

110

100

64 66 68 70 72 74

Data Values

Best Fit Line

160

150

140

130

120

110

100

64 66 68 70 72 74

Data Values

Best Fit Line

160

150

140

130

120

- Implement the above model with both `seaborn` and `numpy` for the dataset of **Brain and Body Relationship**
  - Data source -- <https://bit.ly/349nDxL>