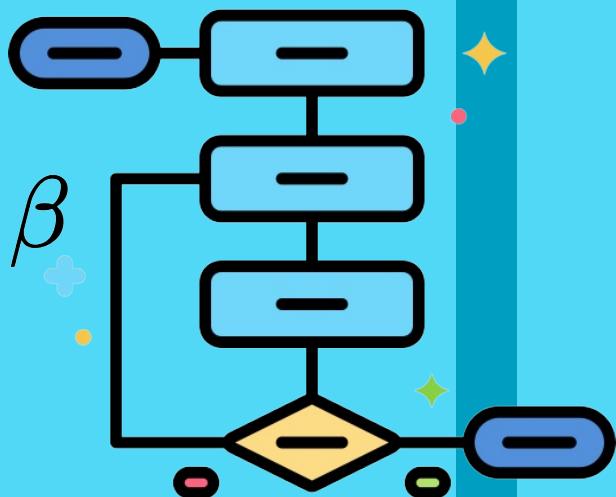


Simple Linear Regression Algorithm from scratch.

$$y = \alpha * x + \beta$$

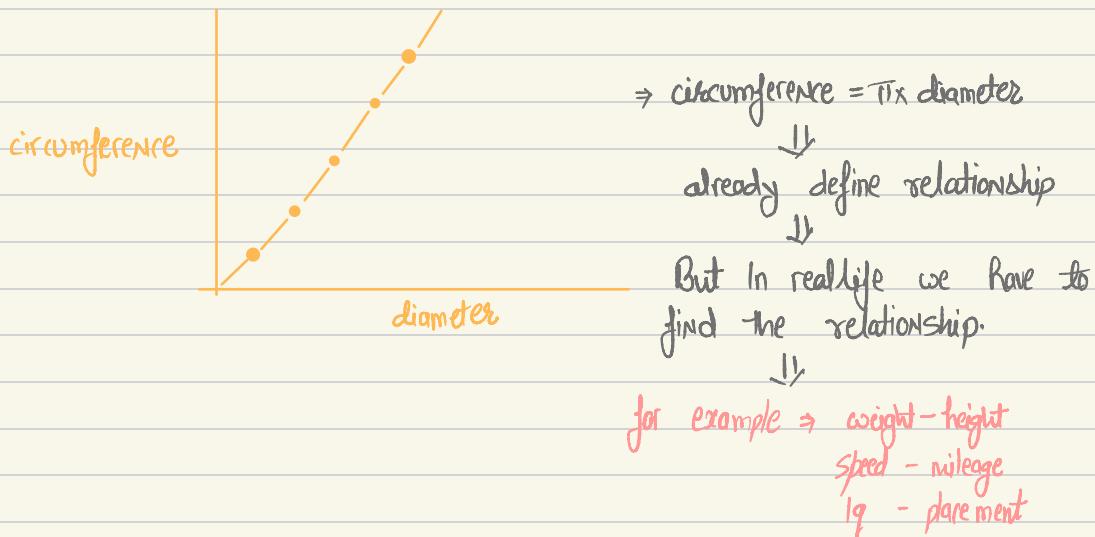


What is Simple Linear Regression?

Simple linear regression is a statistical method that helps us to find out the mathematical relationship between two continuous variables.

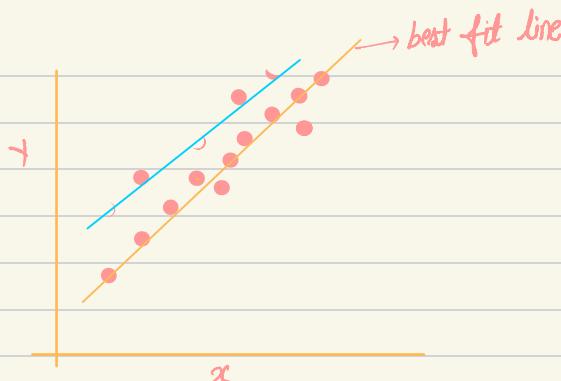
variable $x \rightarrow$ Independent variable

variable $y \rightarrow$ Dependent variable



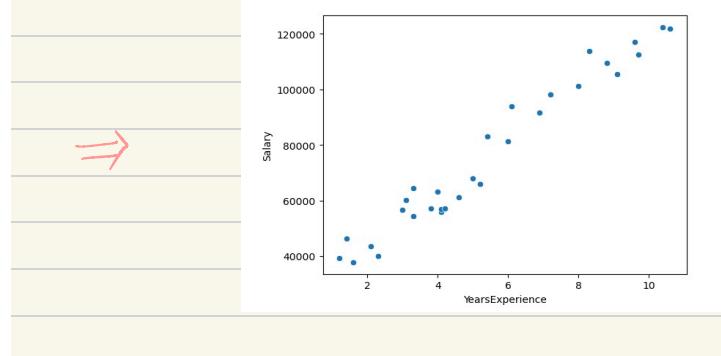
Here, we have to find a best fit line using slope-intercept concept.
 $y = mx + b \Rightarrow m = \text{slope}, b = \text{intercept}$

Best fit line \Rightarrow A line which have lowest error on it. which means that this line cover most of the points on the plan.



For better understanding, I am taking a simple example.

	YearsExperience	Salary
24	24	8.8 109432.0
28	28	10.4 122392.0
25	25	9.1 105583.0
16	16	5.2 66030.0
3	3	2.1 43526.0
1	1	1.4 46206.0
8	8	3.3 64446.0
21	21	7.2 98274.0
7	7	3.3 54446.0
18	18	6.0 81364.0



Now for given data we will find Best fit Line

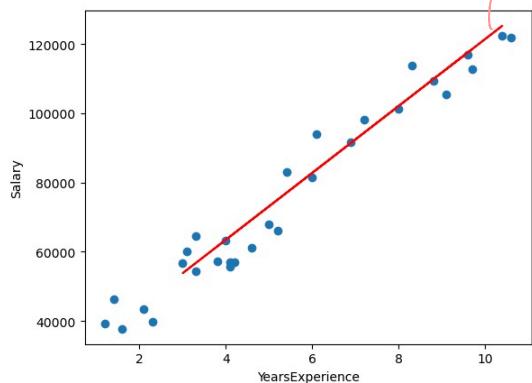
Train the model \Rightarrow

```

: 1 from sklearn.model_selection import train_test_split
: 1 test, y_train, y_test = train_test_split(df[['YearsExperience']], df[['Salary']], test_size=0.25, random_state=10)
: 1 from sklearn.linear_model import LinearRegression
: 1 lr = LinearRegression()  $\Rightarrow$  our model
: 1 lr.fit(x_train,y_train)  $\Rightarrow$  Train the model
: * LinearRegression
LinearRegression()

```

Text(0, 0.5, 'Salary')



Best fit Line

$$\text{Salary} = m(\text{YearsExperience}) + b$$

by give code we find out the value
of Intercept(b) and slope(m)



```
1 lr.coef_
array([[9657.22661273]])
```

→ Slope

```
1 lr.intercept_
array([24857.13238724])
```

→ Intercept

So equation of Best fit line \Rightarrow

$$\text{Salary} = (9657.22)x + 24857.13$$

If $\text{YearsExperience} = 0$ then $\text{Salary} = 24857.13$,

$\therefore b$ is also known as offset.

m, b are called weights

Now, one question should be in your mind that how this function has find these values. and the true maths comes in the picture.

How to find m and b?

There are two main methods to find out the values or weights.

- ① Closed-form method \Rightarrow create a formula using any logical formula with NO free variables. Can not use limit, differentiation or integration. Linear Regression class
(ordinary least square) ←
- ② Non Closed-form method \Rightarrow using Gradient descent, can use limit, differentiation or integration SGD Regression class.

For small features OLS approach is effective but if we have higher dimension of data then Gradient descent have better performance.

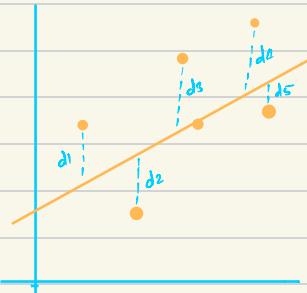
OLS method to find m and b \Rightarrow

$$b = \bar{y} - m \bar{x}$$
$$m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

all the points we have they are sort of linear. and we want to draw a best fit line.

↓

try to cover maximum points.
and reduce the error.



$\text{Error}(E) = d_1 + d_2 + d_3 + \dots + d_n \Rightarrow$ in this equation, value of d may be +ve, and -ve and that may not give correct output, therefore we are taking square of each and every distance to make positive.

$$E = d_1^2 + d_2^2 + d_3^2 + \dots + d_n^2$$

$$E = \sum_{i=1}^n d_i^2 \Rightarrow \text{Error function } (J)$$

Want the value of m and b which have minimum error (E).

Why not made?

$$E = |d_1| + |d_2| + \dots + |d_n|$$

Reason ① \Rightarrow we want to penalize the outliers.

Reason ② \Rightarrow differentiation of mode is NOT possible at origin.

$$E = \sum_{i=1}^n d_i^2$$

(this is always in terms of y)

$$\therefore d_i = y_i - \hat{y}_i$$

here y_i = actual value

\hat{y}_i = Predicted value

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

$$\hat{y}_i = mx_i + b$$

$$E(m, b) = \sum_{i=1}^n (y_i - mx_i - b)^2$$

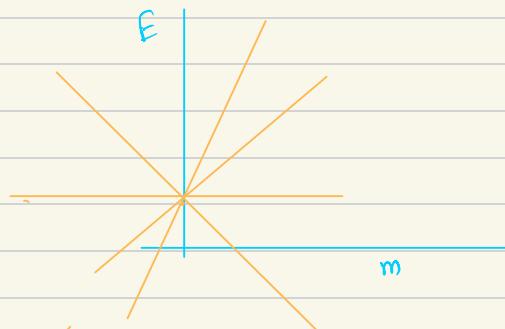
$y_i, x_i \Rightarrow$ constant \Rightarrow cannot be change, we can change m, b to modify the error.

Our error function is depends on m and b . ($E(m, b)$)

① Assume $b=0$

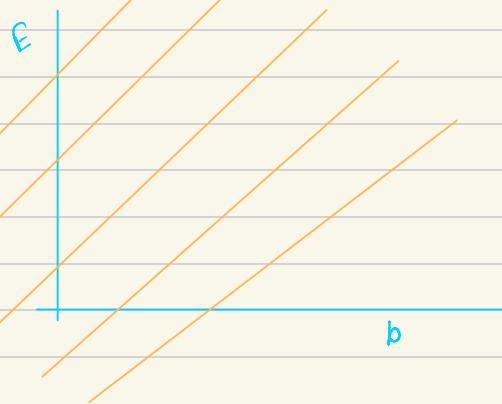
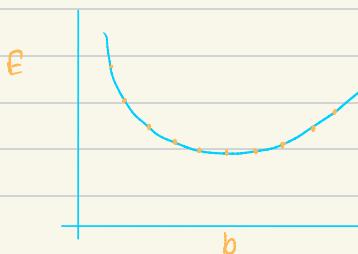
$$E(m) = \sum_{i=1}^n (y_i - mx_i)^2$$

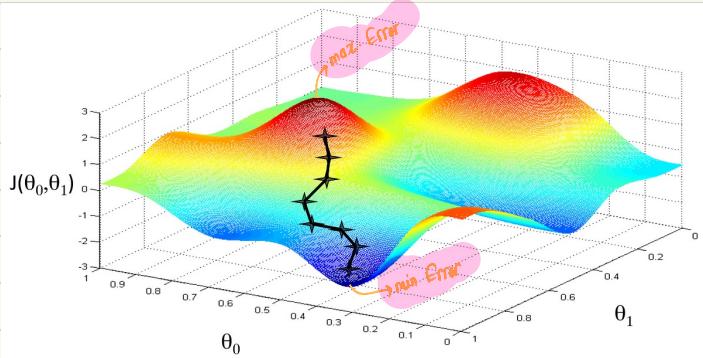
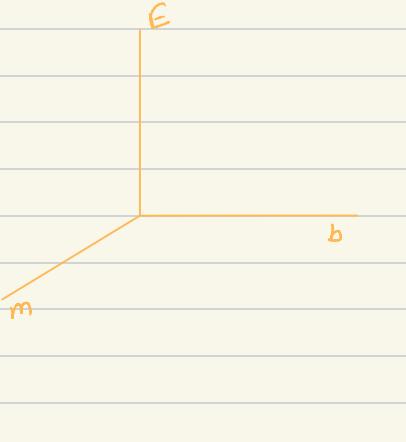
if we create a plot between m and E , then we will get a curve
plot it may be look like \Rightarrow



② Assume $m=1$

$$E(b) = \sum_{i=1}^n (y_i - x_i - b)^2$$





we can find minimum point using maxima-minima.



find out the derivative of error function with respect to m and b.

$$\text{Now } \Rightarrow \frac{\partial E}{\partial b} \Rightarrow \frac{\partial}{\partial b} \left(\sum_{i=1}^n (y_i - mx_i - b)^2 \right) = 0$$

$$\sum \frac{\partial}{\partial b} (y_i - mx_i - b)^2 = 0$$

$$\sum 2(y_i - mx_i - b)(-1) = 0$$

$$\sum (y_i - mx_i - b) = 0$$

$$\frac{\sum y_i}{n} - \frac{\sum mx_i}{n} - \frac{\sum b}{n} = 0$$

$$\bar{y} - m\bar{x} - b = 0$$

$$b = \bar{y} - m\bar{x}$$

$$E = \sum (y_i - mx_i - \bar{y} + m\bar{x})^2$$

$$\frac{\partial E}{\partial m} = \sum \frac{\partial}{\partial m} (y_i - mx_i - \bar{y} + m\bar{x})^2 = 0$$

$$\Rightarrow - \sum 2(y_i - mx_i - \bar{y} + m\bar{x})^2(x_i - \bar{x}) = 0$$

$$\Rightarrow \sum [y_i - mx_i - \bar{y} + m\bar{x}] (x_i - \bar{x}) = 0$$

$$\Rightarrow \sum [(y_i - \bar{y}) - m(x_i - \bar{x})] (x_i - \bar{x}) = 0$$

$$\Rightarrow \sum (y_i - \bar{y})(x_i - \bar{x}) - m \sum (x_i - \bar{x})^2 = 0$$

$$\Rightarrow \sum (y_i - \bar{y})(x_i - \bar{x}) = m \sum (x_i - \bar{x})^2$$

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

```

1 class LinearClass:
2     """
3         This class implements the linear regression algorithm
4     """
5     def __init__(self):
6         """
7             Initializes the model parameters
8         """
9         self.slope = None # Slope of the regression line(m)
10        self.intercept = None # Intercept of the regression line(b)
11
12    def fit(self, X_train,y_train):
13        """
14            Fits the linear regression model to the training data
15        """
16        Args:
17            X_train (Numpy.ndarray) : the training features
18            y_train (Numpy.ndarray) : the training labels
19
20        # Have to calculate the numerator and denominator for the slope
21        numerator = 0
22        denominator = 0
23
24        for i in range(X_train.shape[0]):
25            numerator += (X_train[i] - X_train.mean()) * (y_train[i] - y_train.mean())
26            denominator += (X_train[i] - X_train.mean()) ** 2
27
28        # Calculate the slope and intercept
29        self.slope = numerator/denominator
30        self.intercept = y_train.mean() - (self.slope * X_train.mean())
31
32        # Print the slope and intercept for information purposes
33        print("Slope m", self.slope)
34        print("Intercept b", self.intercept)
35
36    def predict(self, X_test):
37        """
38            Predicts the labels for new data points
39        """
40        Args:
41            X_test (Numpy.ndarray) : The test features
42            ...
43            Numpy.ndarray : the predicted labels
44
45        return self.slope*X_test + self.intercept

```

```

1 lr.fit(x_train['YearsExperience'].values, y_train['Salary'].values)
Slope m 9657.226612725688
Intercept b 24857.132387239224

```