

Supervised **Regression** ML-Algo  $\rightarrow$

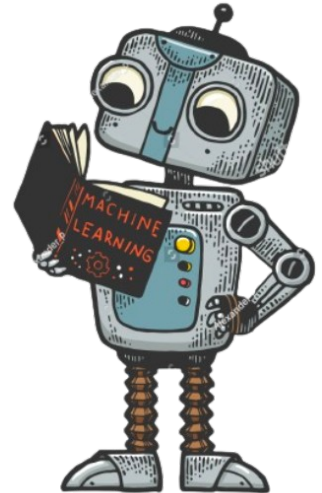
- ① L.R
- ② D.T.R
- ③ R.F.R
- ④ K.N.N.R

# { ML: Linear Regression-1 }

Supervised ML-Algo

we may use this  
algo when the data  
comes with the target  
feature

Target feature  $\rightarrow$  continuous  
data-type



# Summary

## Branches of ML

Supervised  
(If data comes with the target)

unsupervised  
(If data comes only with features)

Regression  
(continuous)

Classification  
(categorical)

⑤

$x_1$	$x_2$	$x_3$	# of cars	Buy/not Buy
age	income			Y
				7
				N
				6
				Y
				9
				N
				2

Task → predict if a customer will purchase the product or not?

target → categorical

ML- Algo  
supervised classification algo

ML- Model

predictions

$$y = f(x)$$

Yes/No



age = 27  
income = 70k  
# of cars = 1

Classification

## Introduction to machine Learning

→ Experience (Historical Data)

→ Task (objective)

→ performance (Evaluation)

Independent var (x)

Dependent variable (y) target

age	Exp	Sex	Salary

Task → predict "Salary"

continuous

"Historical Data"

Regression

ML- Algo

(Learn the relationship between X's and y)

Model

$$y = f(x)$$

predictions

predicted Salary

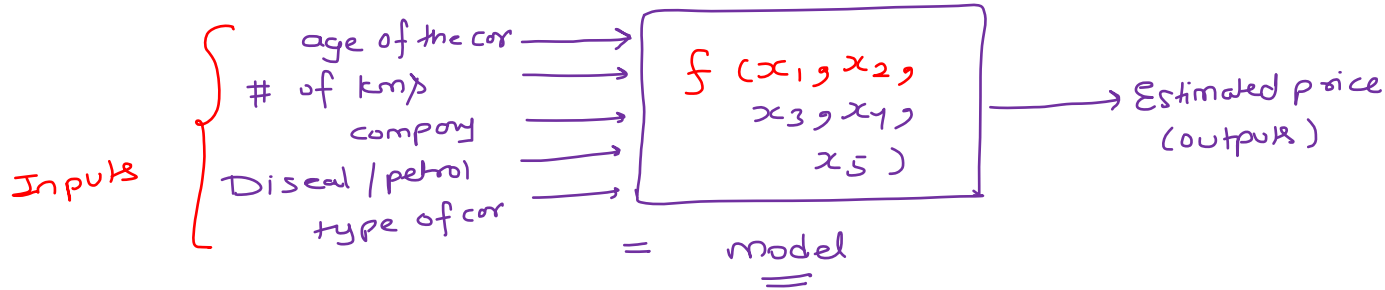
(mathematical Eq. which connect X's with y)

age = 27  
Exp = 7  
Sex = male  
(new employee)

# Motivation

{ Can you think of the factors on which the price of a used car depends on?

- ① Age of the car
- ② # of kms
- ③ company
- ④ Diesel / petrol
- ⑤ Type of car (Hatchback / sedan / SUV)



# Which ML Algorithm To Choose?

Experts

	$x_1$ age	$x_2$ # of kms	$x_3$ Model	$x_4$ P/D	$y$ Price
$c_1$					
$c_2$					
$c_3$					
$\vdots$					
$\vdots$					

(Historical)

→ continuous

2019  
Cross 24 [Last 8 years] →

ML-Algo  
 Linear-Reg

model

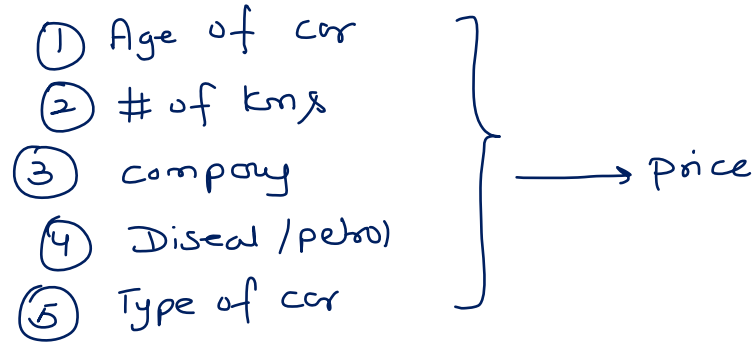
$$\text{price} = f(\text{age}, \# \text{ of kms}, \text{model}, \text{P/D})$$

Data of new car → Estimated price

$\{ \{x_1, x_2, x_3, x_4\} \rightarrow y \}$  → multiple linear regression  
 Independent feature

# Mental Model

## Multiple Linear Regression



$$\text{Price} = \underbrace{(W_1)}_{0.8}(\text{age}) + \underbrace{(W_2)}_{0.2}(\text{\# of kms}) + \underbrace{(W_3)}_{0.01}(\text{company}) + \underbrace{(W_4)}_{0.001}(\text{Diesel/petrol}) + \underbrace{(W_5)}_{0.005}(\text{Type of car})$$

info → age = 2  
→ # of kms = 20k  
→ company = Tata (EV)  
→ car = EV  
→ SUV

3L  
2.9L  
3.25L  
3.25L  
4L

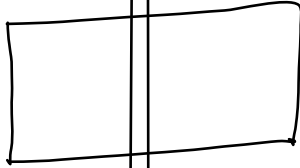
} Equation of the line }

# Simple Linear Regression

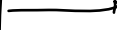
one independent variable  $\rightarrow$  (y)  
(x) — (y)

What we want to predict  $\rightarrow$  Sales

marketing  
budget



Expected  
Sales



$$(y = \underline{w_0} + \underline{w_1}x)$$

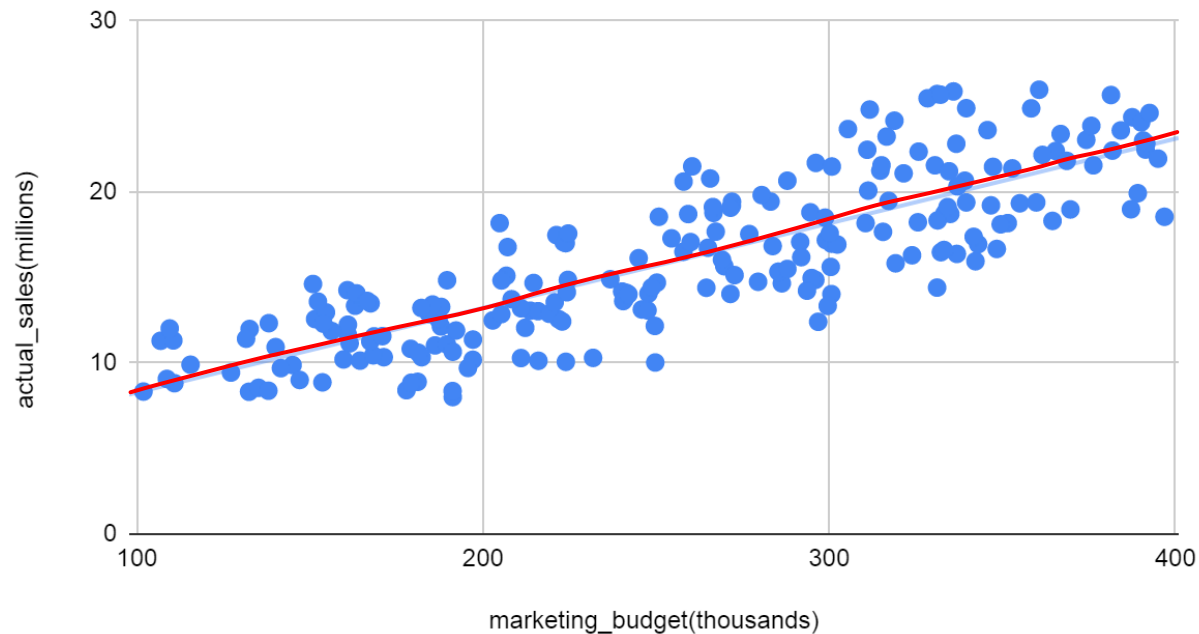
$$[ \text{Sales} = \underline{w_0} + \underline{w_1} (\text{marketing budget}) ]$$

Best value of  $w_0, w_1$

marketing_bud get(thousands)	actual_sale s(millions)
187.86	12.14
138.13	12.33
177.89	8.41
181.13	8.9
151.5	12.57
106.87	11.3
160.94	11.66
140.09	10.93
132.58	11.99

# Simple Linear Regression

actual\_sales(millions) vs. marketing\_budget(thousands)



Best fit line  
 $w_0, \bar{w}$

# Simple Linear Regression

$(w_0, w_1) \rightarrow \text{coefficient}$

$y(\text{dependent})$

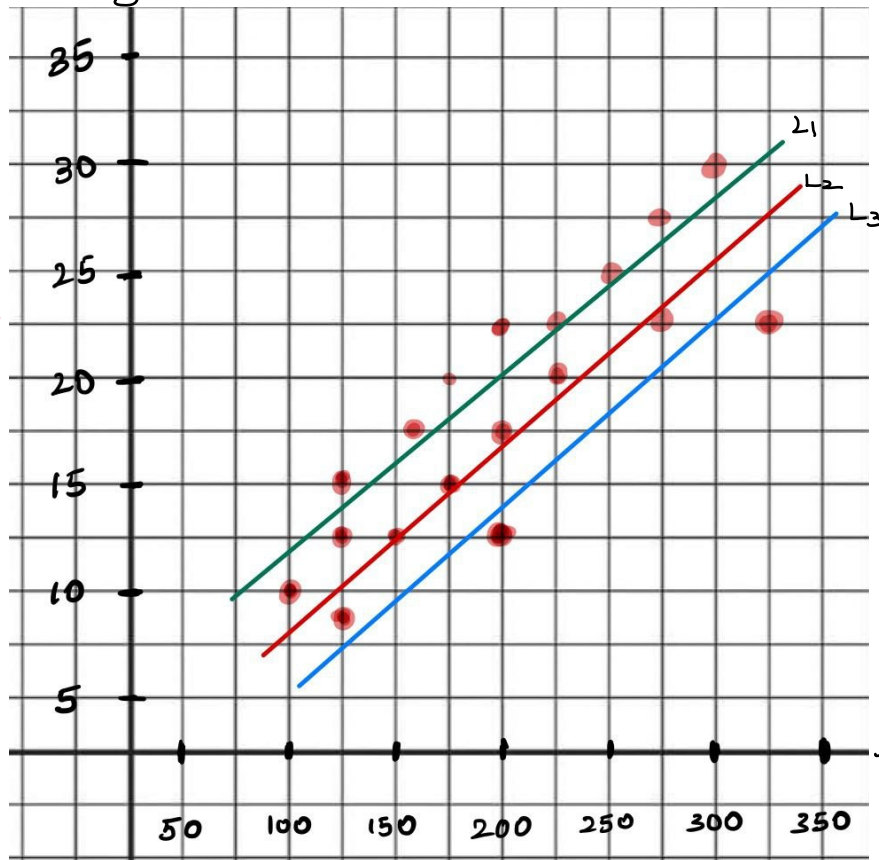
Task  $\rightarrow$  choose the best value of  $w_0, \bar{w}$

Question  $\rightarrow$  on what basis?

$$y_a \approx y_p$$

or

$$(y_a - y_p) \approx 0$$



$$L_1: y = w_0 + w_1(x) \text{ RSS} = 500$$

$$L_2: y = w'_0 + w'_1(x) \text{ RSS} = 200$$

$$L_3: y = w''_0 + w''_1(x) \text{ RSS} = 700$$

OLS (Ordinary Least Square method)

out of all the possible lines (coefficient) that can pass through the data choose that line (coefficient) for which the error (RSS) is minimum



# Simple Linear Regression

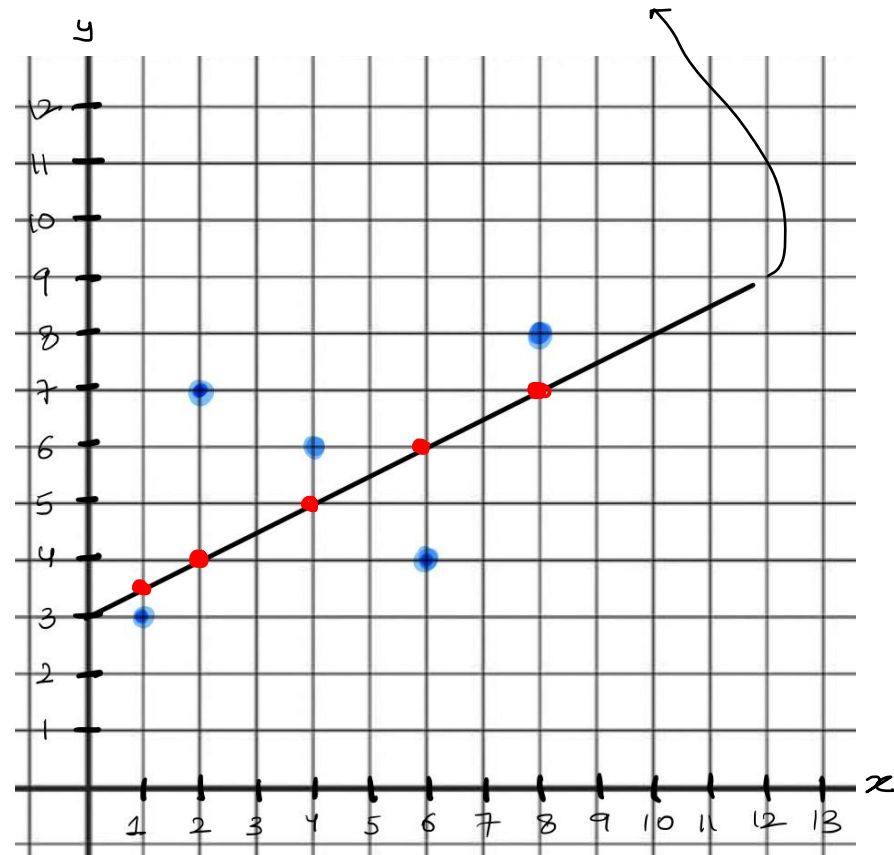
$$y_p = y_{\text{(predicted)}} = \hat{y} \text{ (y estimate)}$$

$$y = w_0 + w_1(x)$$

$x$	$y_a$	$y_p$	$\epsilon = (y_a - y_p)$	$\epsilon^2$
1	3	3.5	-0.5	0.25
2	7	4	3	9
4	6	5	1	1
6	4	6	-2	4
8	8	7	1	1

Total

15.25 RSS  
(Residual sum of Squares)



# Simple Linear Regression

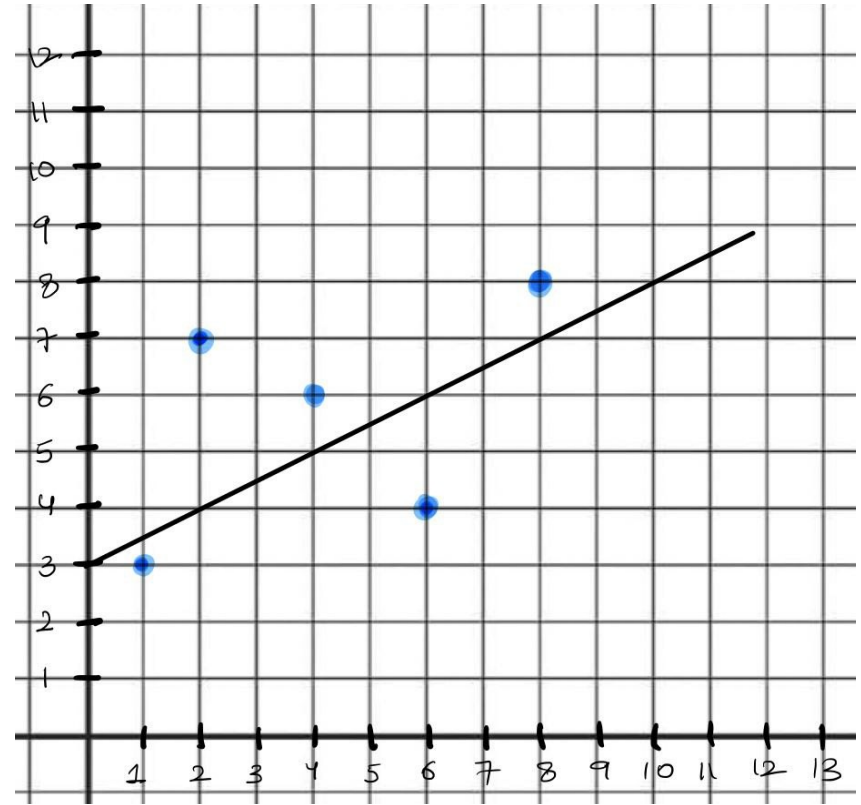
Initialisation  $\rightarrow$  Random value of  $w_0, \bar{w}$

① Random line

$$w_0 + w_1(x) \quad [w_0 = 0, w_1 = 1]$$

[RSS]

② Update  $w_0, \bar{w}$  (gradient Descent)



# Simple Linear Regression

objective  $\rightarrow$  minimize (RSS)

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

$$\hat{y} = w_0 + w_1 x$$

$$\mathcal{L}(w_0, w_1) = \underset{w_0, w_1}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - (w_0 + w_1 x))^2 \right\} \quad \text{Error minimisation}$$

objective : optimise the above function and find out that value of  $w_0, w_1$  for which the function yields the minimum value

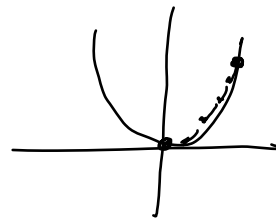
$\rightarrow$  Gradient Descent

Slope of loss function w.r.t  $w_0$

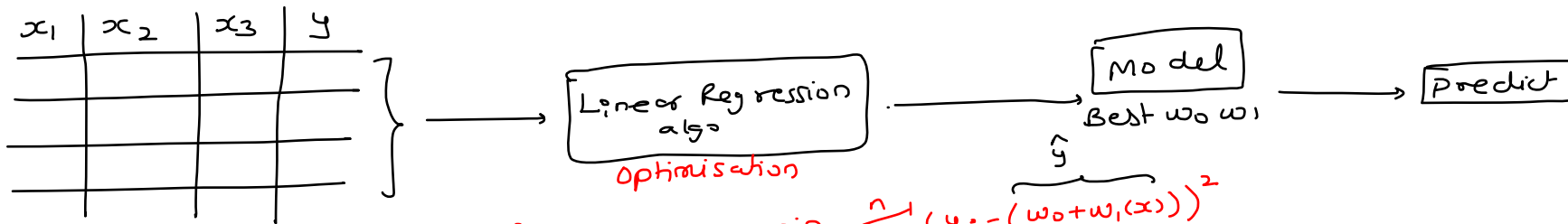
① initialise  $(w_0, w_1)$

$$\left. \begin{aligned} \text{② update } w_0 &= w_0 - \eta (\nabla_{w_0} \mathcal{L}) \\ w_1 &= w_1 - \eta (\nabla_{w_1} \mathcal{L}) \end{aligned} \right\}$$

$\downarrow$   
Step size



# Simple Linear Regression



Loss function  
cost function

$$\mathcal{L}(w_0, w_1) = \arg \min_{w_0, w_1} \sum_{i=1}^n (y_i - (w_0 + w_1(x)))^2$$

↓  
gradient Descent

objective: find that value of  $w_0 w_1$   
corresponding to which the  
loss function is minimum

$$\hat{y} = w_0 + w_1 x$$

( $\hat{y}$ )

$x=2$

## Error Functions (Various Versions)

KR

(y) K		
Income	$\hat{y}$	$E^2$
10	5	25
20	15	25
25	25	0
15	15	0
20	20	0

50 = RSS

RS

Income	$\hat{y}$	$E^2$
10,000	5000	
20,000	15000	
25,000	25000	
15,000	10000	
20,000	20,000	

RSS = 5 M

Problem with using Error function as Evaluation parameter

ERROR IS Dependent on the Scale of target variable

$$MSE (\text{mean Square error}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$


$$MAE (\text{mean-absolute error}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$


$$RMSE (\text{Root mean square error}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

Error

we can't use them for evaluating the goodness of the model

# Motivation For R-2 Score

 (Sharma ji ka ladka)  
("sharma's son")

  
(Sumit)

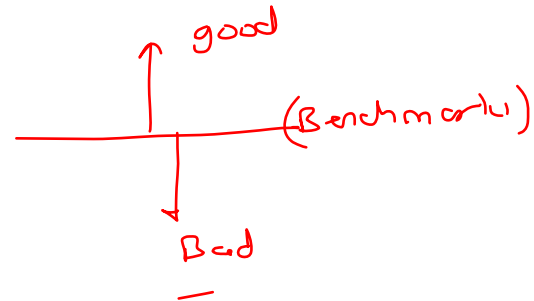
  
(Sumit's father)

compare your marks (good student)

↓  
Reference / Benchmark

Compare (performance)

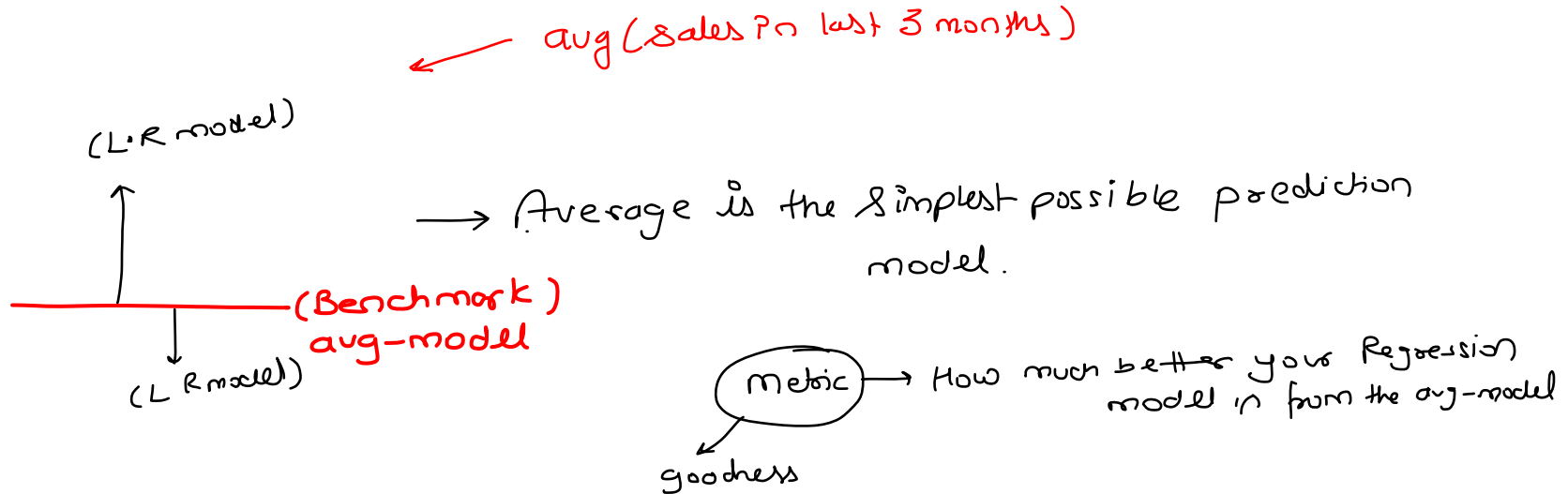
↓  
Benchmark



# Simple Linear Regression

Assume that you work for a Marketing Firm, and your CEO requires you to provide her with a rough estimate of the sales that will occur next month. You already have the data (present on the right). Using the given data, try to determine the approximate sales for the upcoming month.

PS: Please remember that as of now, you do not have any knowledge of Linear Regression or complex Machine Learning algorithms. Try to solve this task using your logical understanding and basic mathematics.



# Simple Linear Regression

$RSS = 15.25$  (Squared Error by the regression model)  
 $TSS = 17.20$  (Squared Error by the avg-model)

$x$	$y_a$	$\hat{y}$	$e^2$	$\bar{y}$	$E^2 = (y_a - \bar{y})^2$
1	3	3.5	0.25	5.6	$(3-5.6)^2 = 6.76$
2	7	4	9	5.6	$(7-5.6)^2 = 1.96$
4	6	5	1	5.6	$(6-5.6)^2 = 0.16$
6	4	6	4	5.6	$(4-5.6)^2 = 2.56$
8	8	7	1	5.6	$(8-5.6)^2 = 5.76$

5.6

15.25

17.2

$$R^2 = 1 - \left( \frac{RSS}{TSS} \right) = 1 - \left( \frac{15.25}{17.2} \right)$$

0.11

Industry  
( $R^2 = 85\%$ )

Regression model is 11% better than the avg-model.

