

## Coefficient of Determination ( $R^2$ )

↳ goodness of fit.  $\in [0, 1]$

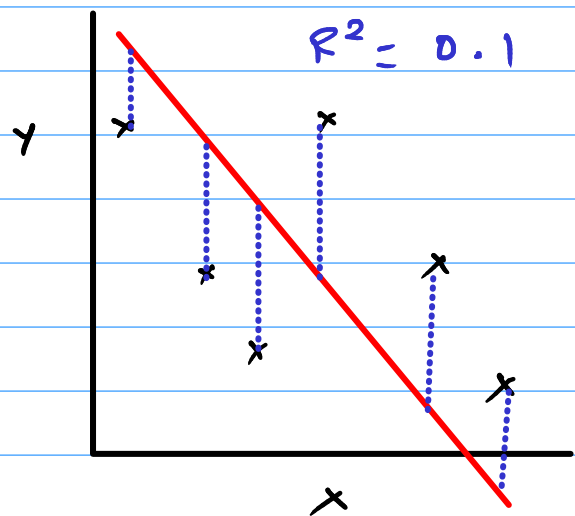
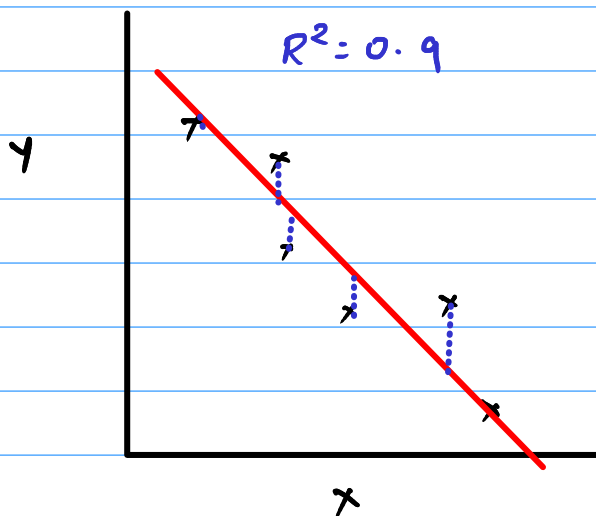
fails to model the data.      perfect fit

$R^2 = 0.2 \Rightarrow 20\%$  of the dependent variable is predicted by the independent variable.

Exam's score(s)  $\rightarrow y$   
time spent in studying  $\rightarrow x$

- $\rightarrow R^2 = 0$  : Can't predict
- $\rightarrow R^2 = 1$  : Perfect prediction
- $\rightarrow R^2 \in (0, 1)$  : partial prediction.

$R^2$  is a goodness of fit.

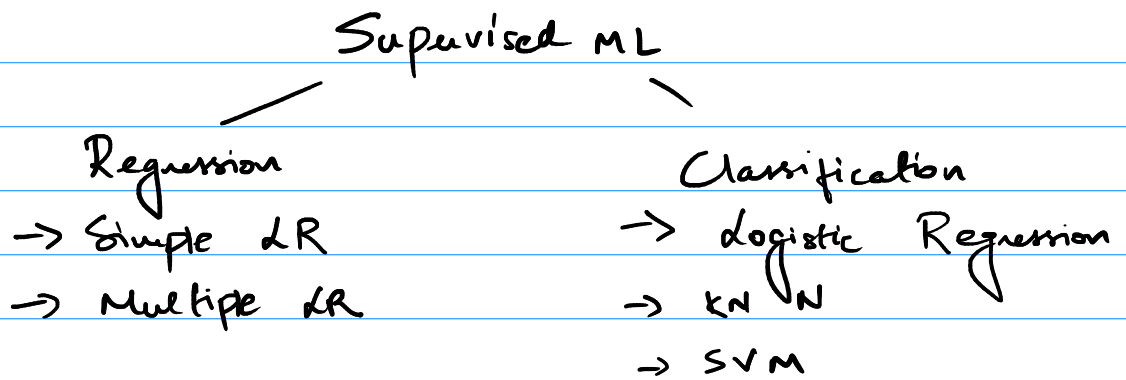


NOTE:  $R^2$  tells you how  $x$  &  $y$  are correlated with each other, but it shouldn't be taken as an evidence

for causality.

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

$$\epsilon_i = \hat{y} - y$$



Logistic Regression

Y (dependent variable) → categorical (Discrete)

Logistic Regression Model →

$$f(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_n x_n)}}$$

where  $e = 2.7$

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \text{predictor variables}$$

$\beta_0$  = Bias

$\beta_1, \dots, \beta_n$  = weights

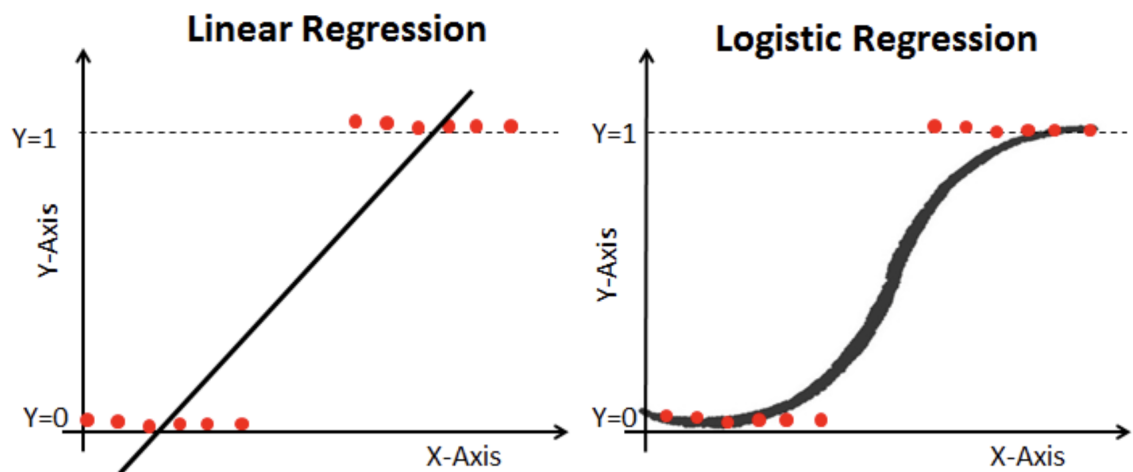
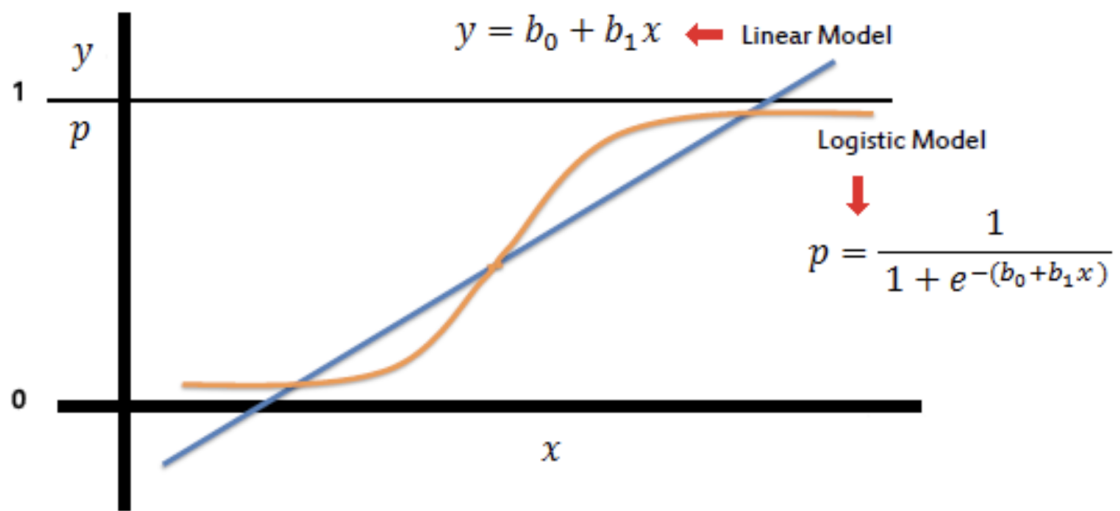
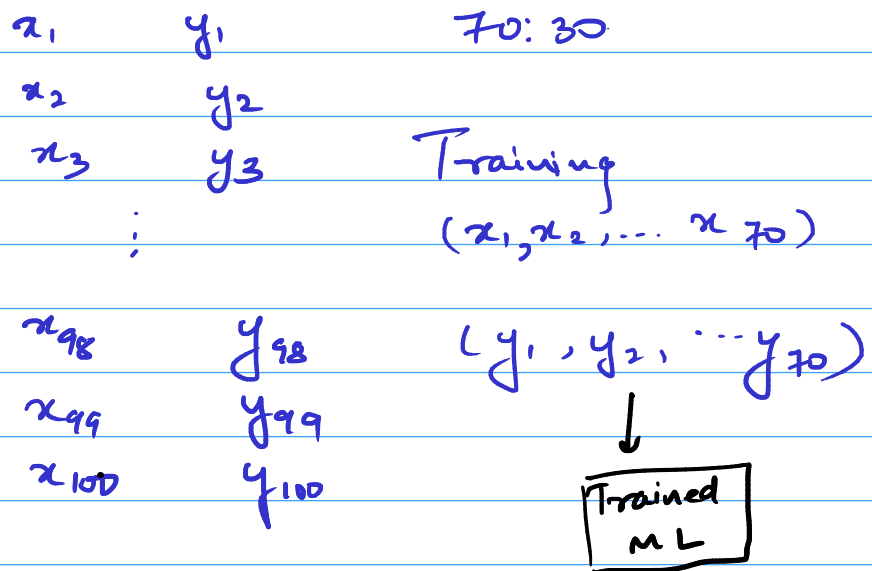


Image Source [www.datacamp.com](http://www.datacamp.com)

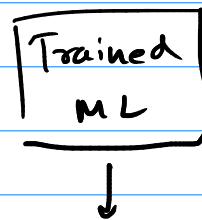
$$p = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

where  $x = \text{Age}$

if  $p > 0.5 \Rightarrow$  The person has CHD or 1  
 else  $p < 0.5 \Rightarrow$  The person don't have CHD or 0.



Testing (30%)  $(x_{71}, x_{72}, \dots, x_{100})$



predicted values:  $(\hat{y}_{71}, \hat{y}_{72}, \dots, \hat{y}_{100})$

Actual values:  $(y_{71}, y_{72}, \dots, y_{100})$

70: 30

80: 20

Training: Testing  
 should be depend on  
 the size of the dataset.

Dataset:

Training :

Testing.

70: 30 \* 10,000,000  
 ✓

7,000,000

3,000,000

9,900,000 (99.9%)

10,000 (0.1%)

✓  
Training set

Validation set

✓  
Test Set.

500

if he fails  $\Rightarrow$  he failed to generalise

(Overfitting)

500

[5-20]  
 $\sim 60\%$   
=

5  
 $\sim 85\%$   
=

Train the  
model  $\rightarrow$

Evaluation within  
the env., with  
the available dataset

$\rightarrow$  Rolling out to  
the world &  
getting the  
feedback.

Random state = k

dataset = [a, b, c, d, e]

splitting: 8:2  $\Rightarrow$  test-size = 0.2

I: Random state = 0

training = [a, b, c, d]

Test = [e]

II Random state = k

training = [a, b, d, e]

test = [c]