

X = Input variable / predictors / independent variables

y = Output variable / response / dependent variables

N = # of observations

P = # of features

X observation# feature# \Rightarrow

So a single observation can be represented as below,

$$\text{row1} = (x_{11} \quad x_{12} \quad x_{13} \dots \quad x_{1P})$$

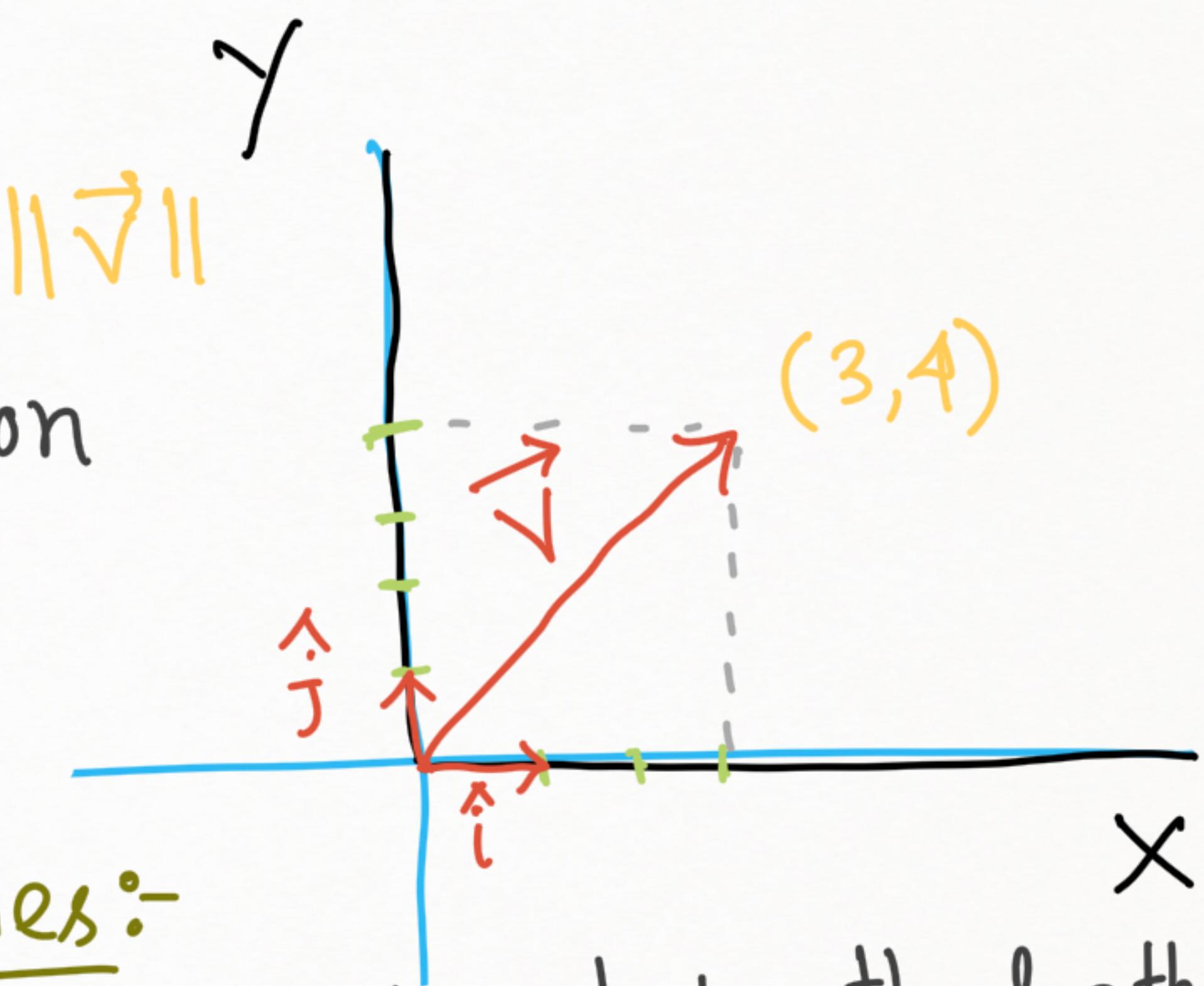
x_{11} = GRE Score of the first observation

x_{12} = TOEFL Score of the first observation

x_{21} = GRE Score of the second observation

$$X = \begin{bmatrix} x_{11} & x_{12} \dots x_{1P} \\ x_{21} & x_{22} \dots x_{2P} \\ \vdots & \vdots \\ x_{N1} & x_{N2} \dots x_{NP} \end{bmatrix}$$

length = $\|\vec{v}\|$
direction



Properties :-

$2\vec{v}$ = A vector twice the length
of \vec{v} but in same direction

$$\vec{a} + \vec{b} \Rightarrow 3\hat{i} + 4\hat{j} = \vec{v}$$

(Column vector) $\begin{bmatrix} 3 \\ 4 \end{bmatrix}$

Transpose :- (row becomes cols &
cols becomes rows)

$$\vec{v} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$\vec{v}^T = [3 \ 4]$$

$$x_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{bmatrix}$$

$$x_i^T = [x_{i1} \ x_{i2} \ \dots \ x_{ip}]$$

(Vector of length P)

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}$$

(Row view of the training inputs)

$$z_j = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{Nj} \end{bmatrix}$$

Vector of length N

Representing each and every feature column as vector z_j

$$X = [z_1 \ z_2 \ \dots \ z_p]$$

(Column view of inputs)

Similarly, the output Y can be represented like below,

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Based on # of features
LR can be sub-divided
in two categories,

LR

Simple/Univariate
LR

Multivariable
LR

For Simple LR \rightarrow

$$X = [z_1] = \begin{bmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{N1} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

So, Basically we can view the whole dataset as,
 $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

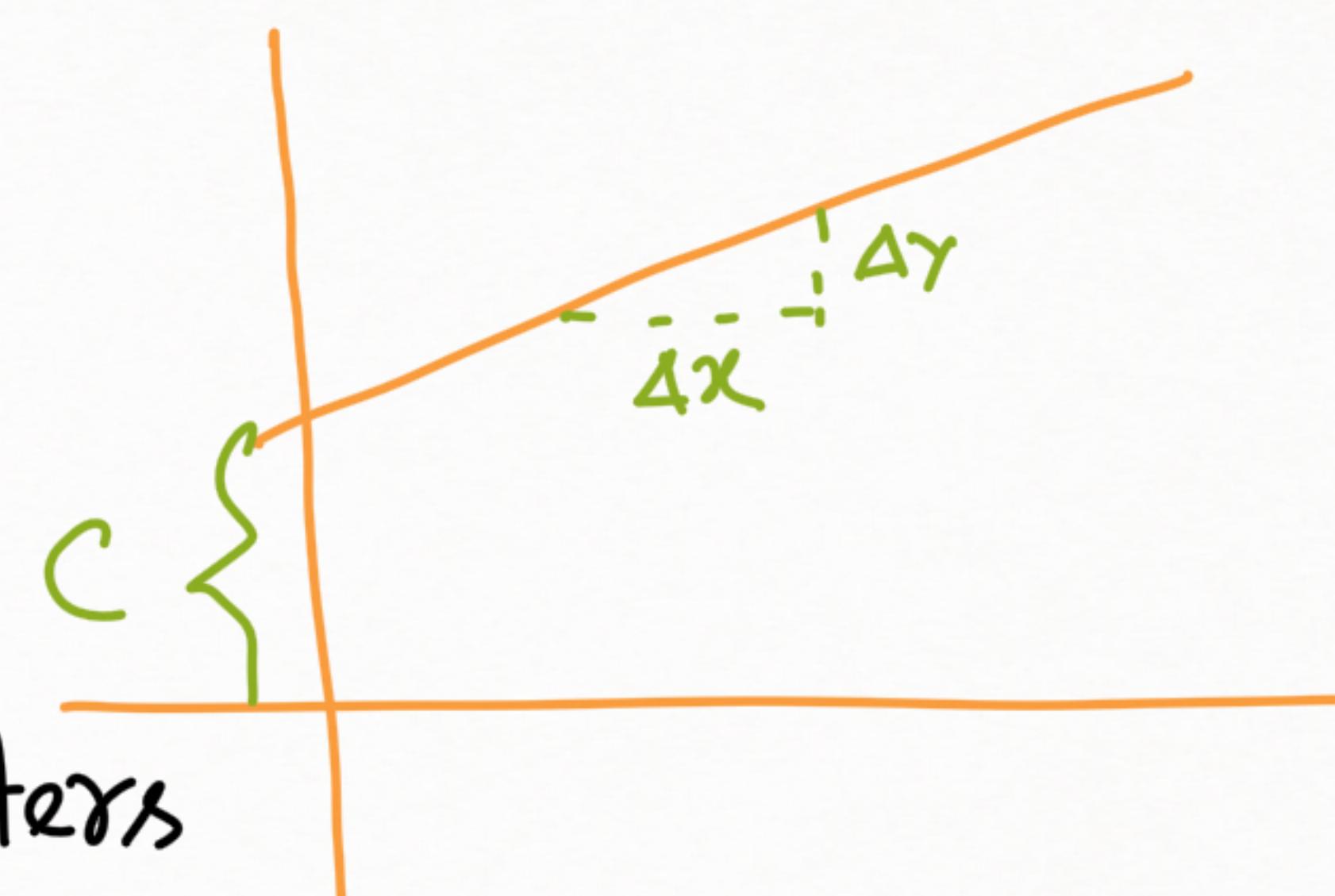
Questions We should ask:-

- 1) Is there a relationship between GRE Score & chance of admit?
- 2) How strong is the relationship? Given a GRE Score can we predict chance of admit?
- 3) Is the relationship linear? ***

Lets do the math 😊

(model)

$$y = mx + c$$



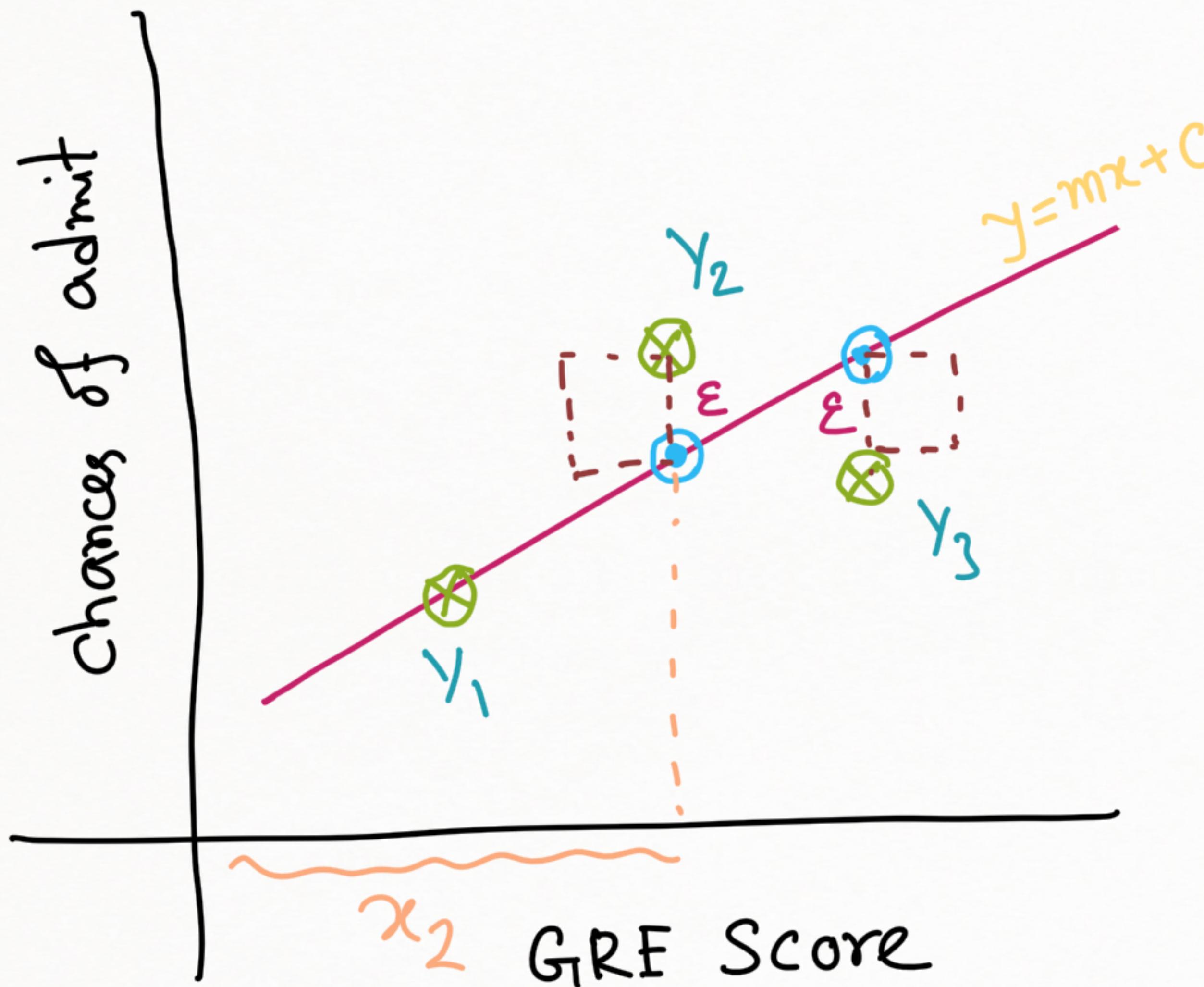
$$m = \frac{\Delta y}{\Delta x}$$

= Slope

C = Y-axis
Intercept

In ML, m & c are called parameters
or co-efficients

(chances of Admit) = m (GRE Score) + C



First think we need to measure how goodness the fit is.

We will define ϵ as a error function for single measuring data point with below rule,

1. If the line pass through that point error (ϵ) is zero.
2. ϵ can't be negative

$$\epsilon = (y - y_i)^2$$

This is error for a single point.

$$\text{Total Error } (E) = \sum_{i=1}^N \varepsilon \quad (\text{Sum of all errors})$$

$$E = \sum (y - y_i)^2 \quad \left[\text{as } \varepsilon = (y - y_i)^2 \text{ and we will assume } \sum = \sum_{i=1}^N \right]$$

$$\Rightarrow E = \sum (mx_i + c - y_i)^2 \quad \left[\text{as the predicted value } y = mx_i + c \right]$$

$$\Rightarrow E = \sum (m^2x_i^2 + c^2 + y_i^2 + 2mx_i c - 2cy_i - 2mx_i y_i) \quad \left[\text{as, } (a+b+c)^2 = a^2 + b^2 + c^2 + 2ab + 2bc + 2ca \right]$$

$$\Rightarrow E = \sum m^2x_i^2 + \sum c^2 + \sum y_i^2 + \sum 2mx_i c - \sum 2cy_i - \sum 2mx_i y_i$$

Let us further simplify it,

Let's say,

$$\alpha = \sum y_i^2$$

$$\beta = \sum x_i^2$$

$$\gamma = \sum x_i y_i$$

$$\mu = \sum y_i$$

$$\theta = \sum x_i$$

$$E = m^2 \sum x_i^2 + c^2 + \sum y_i^2 + 2mc \sum x_i$$

$$- 2c \sum y_i - 2m \sum x_i y_i$$

$$E = m^2 \beta + \sum_{i=1}^N c^2 + \alpha + 2mc\theta - 2c\mu - 2m\gamma$$

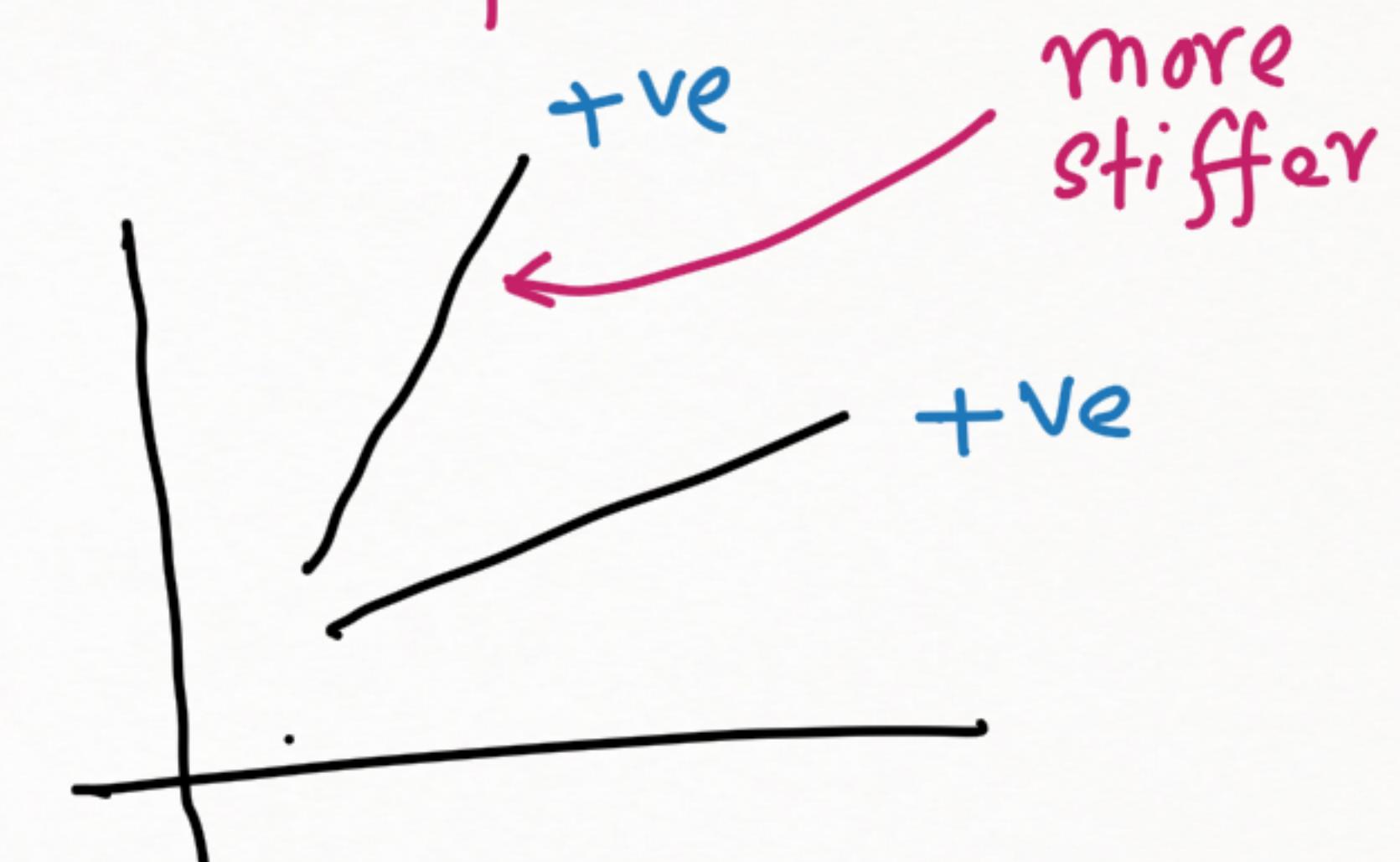
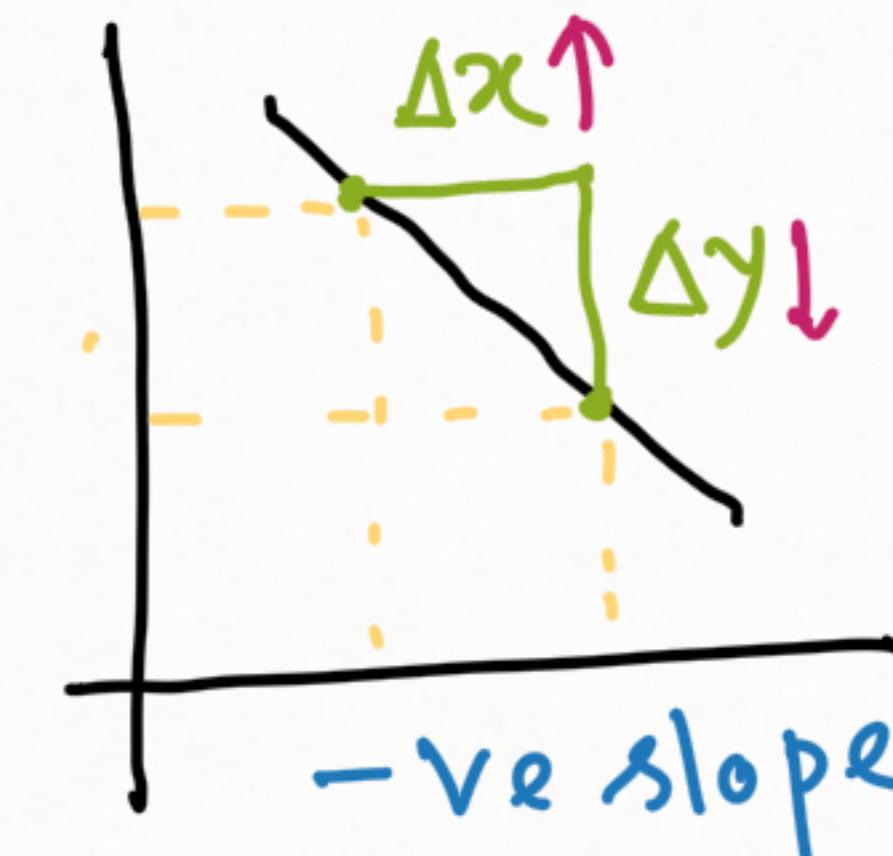
$$E = m^2 \beta + c^2 N + \alpha + 2mc\theta - 2c\mu - 2m\gamma$$

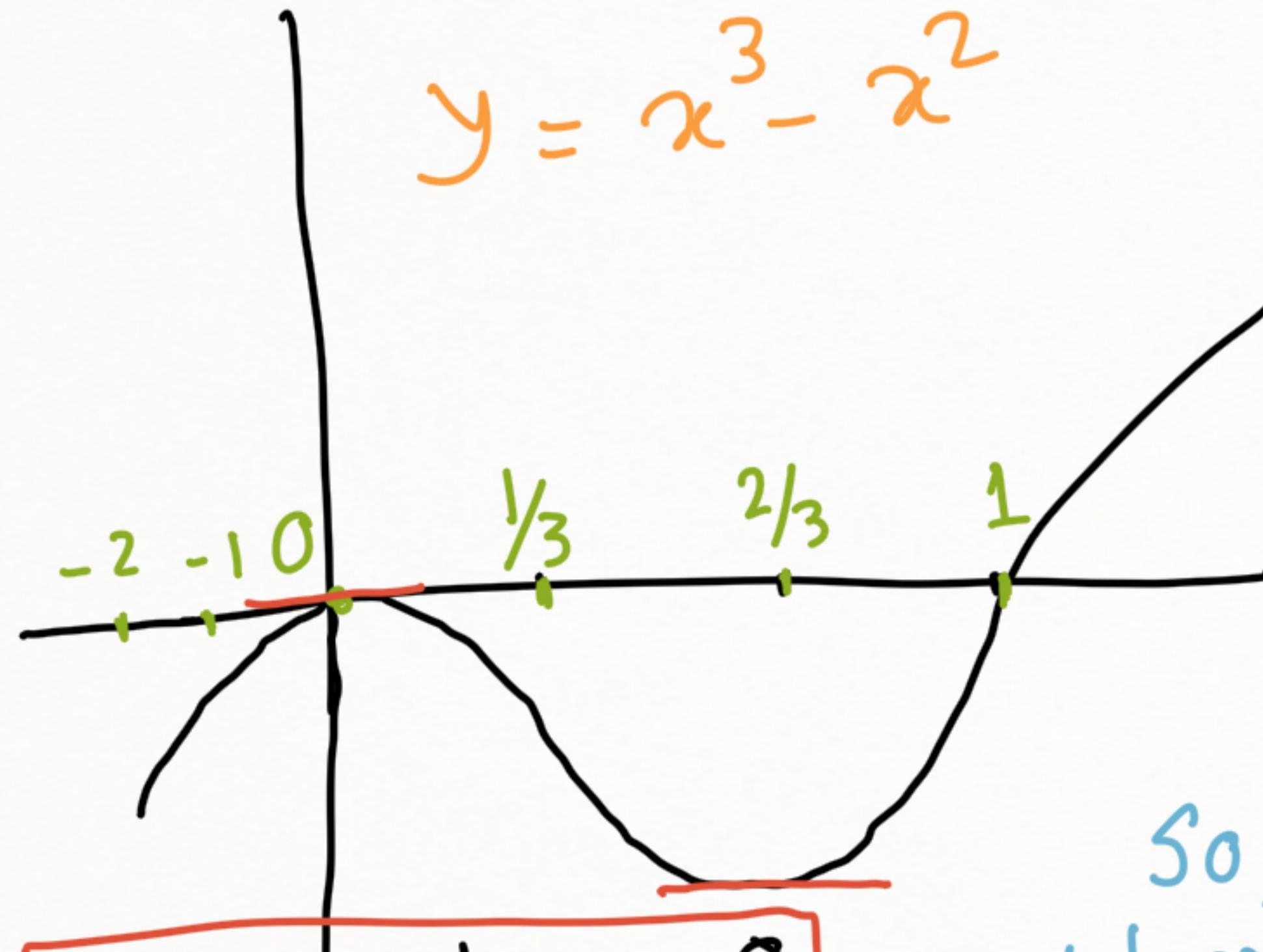
Now our job is to find m & c values for which E is min

Derivative = Slope of tangent line

$$\text{Slope} = \frac{\Delta y}{\Delta x}$$

(when we go from one point to another point)





Slope at $x=0$

$$y' = 0$$

Slope at $x=\frac{1}{3}$

$$y' = -\frac{1}{3}$$

Slope at $x=\frac{2}{3}$

$$y' = 0$$

$$y = x^3 - x^2$$

$$\frac{dy}{dx} = y' = 3x^2 - 2x$$

Slope at $x=-2$

$$y' = 3(-2)^2 - 2(-2) = 16$$

Slope at $x=-1$

$$y' = 3(-1)^2 - 2(-1) = 5$$

So, we are dealing with maxima or minima when derivative = 0. To find whether its maxima or minima we need to take second derivative.

$$\frac{d^2y}{dx^2} = y'' = 6x - 2$$

at $x=0$

$$y'' = -2 < 0$$

(maxima)

at $x=\frac{2}{3}$

$$y'' = 2 > 0$$

(minima)

$$\left[\frac{d}{dx} x^n = n \cdot x^{n-1} \right]$$

Now we reach a point where we need to take derivative of E
 As E is a function of m, c So we need to take partial derivative

$$E = \tilde{m}^2\beta + \tilde{c}^2N + \alpha + 2m\theta - 2c\mu - 2m\gamma$$

$$\frac{\partial E}{\partial m} = 2m\beta + 0 + 0 + 2c\theta - 0 - 2\gamma = 2m\beta + 2c\theta - 2\gamma$$

To get min first we will set this to zero

$$\begin{aligned} \frac{\partial E}{\partial m} = 0 &\Rightarrow 2m\beta + 2c\theta - 2\gamma = 0 \\ &\Rightarrow m = \frac{\gamma - c\theta}{\beta} \quad \dots \quad (\text{i}) \end{aligned}$$

Similarly,

$$\begin{aligned} \frac{\partial E}{\partial c} = 0 + 2cN + 0 + 2m\theta - 2\mu = 0 \\ \Rightarrow c = \frac{\mu - m\theta}{N} \quad \dots \quad (\text{ii}) \end{aligned}$$

$$\frac{\partial^2 E}{\partial m^2} = 2\beta > 0$$

$$\frac{\partial^2 E}{\partial c^2} = 2N > 0$$

$$m = \frac{\gamma - \mu\theta}{\beta}$$

$$\Rightarrow m\beta = \gamma - \left(\frac{\mu - m\theta}{N} \right) \theta$$

$$\Rightarrow m\beta = \gamma - \frac{\mu\theta - m\theta^2}{N}$$

$$\Rightarrow m\beta - \gamma + \frac{\mu\theta - m\theta^2}{N} = 0$$

$$\Rightarrow \frac{m\beta N - \gamma N + \mu\theta - m\theta^2}{N} = 0$$

$$\Rightarrow m(\beta N - \theta^2) = \gamma N - \mu\theta$$

$$\Rightarrow m = \frac{\gamma N - \mu\theta}{\beta N - \theta^2}$$

Putting back $\gamma, \mu, \theta, \beta$ we
get,

$$m = \frac{N \sum x_i y_i - \sum y_i \sum x_i}{N \sum x_i^2 - (\sum x_i)^2}$$

Similarly,

$$C = \frac{\mu - m\theta}{N}$$

$$\Rightarrow CN = \mu - \left(\frac{y_N - \mu\theta}{\beta N - \theta^2} \right) \theta$$

$$\Rightarrow CN = \mu - \frac{y_N\theta - \mu\theta^2}{\beta N - \theta^2}$$

$$\Rightarrow CN = \frac{\mu\beta N - \cancel{\mu\theta} - y_N\theta + \cancel{\mu\theta^2}}{\beta N - \theta^2}$$

$$\Rightarrow C = \frac{\mu\beta N - y_N\theta}{N(\beta N - \theta^2)}$$

Putting back again,

$$C = \frac{\sum y_i \sum x_i^2 N - \sum x_i y_i N \sum x_i}{N (\sum x_i^2 N - (\sum x_i)^2)}$$

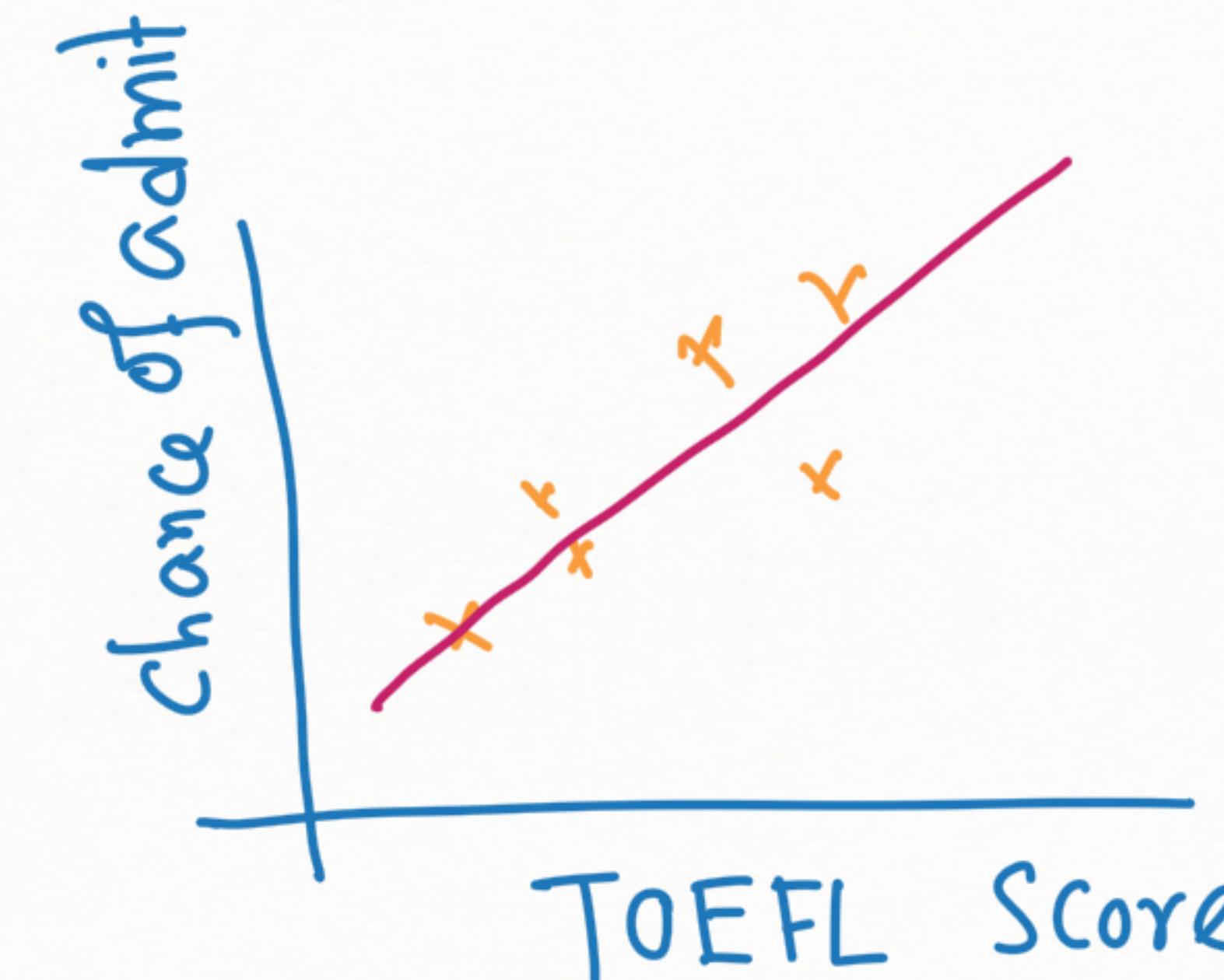
So, Finally we found
the equation of our line

$$y = mx + C$$

Multiple Linear Regression:-

Approach 1

Fit separate Simple linear regression model for each predictor.



Fit separate Simple linear regression model for each predictor.

Cons:-

1. It is unclear how to make a single prediction of Chances of admit given GRE Score & TOEFL score
2. Each of the regression equation will ignore the effect of other predictors while forming the regression Co-efficient, which leads to wrong estimate.

Approach 2

Extend the simple linear regression model so that it can directly accomodate multiple predictors. We can do that by giving each predictor a seperate slope coefficient in a single model.

P = # of features / predictor

$$\hat{y} = \beta_0 + \beta_1 z_1 + \beta_2 z_2 + \dots + \beta_p z_p$$

where,

$$z_i = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{Ni} \end{bmatrix}$$

(N × 1 Vector)

\hat{y} = Predicted output

lets define,

$$Z_0 = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

(N × 1 vector)

$$\hat{y} = \beta_0 z_0 + \beta_1 z_1 + \beta_2 z_2 + \dots + \beta_p z_p$$

$$= \sum_{i=0}^p \beta_i z_i$$

Let's define,

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

($p \times 1$ vector)

We already established,

$$x = [z_0 \ z_1 \ z_2 \ \dots \ z_p]$$

$$x = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ 1 & x_{31} & x_{32} & \dots & x_{3p} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{Np} \end{bmatrix}$$

($N \times p$ matrix)

$$X\beta = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1P} \\ x_{21} & x_{22} & \cdots & x_{2P} \\ x_{31} & x_{32} & \cdots & x_{3P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{NP} \end{bmatrix} \times \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_P \end{bmatrix} \quad (\text{Px1 vector}) = (\text{Nx1 vector})$$

(N x P matrix) x

$$= \begin{bmatrix} \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \cdots + \beta_P x_{1P} \\ \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \cdots + \beta_P x_{2P} \\ \vdots \\ \vdots \\ \beta_0 + \beta_1 x_{N1} + \beta_2 x_{N2} + \cdots + \beta_P x_{NP} \end{bmatrix} \quad (\text{Nx1 vector})$$

Let's define,

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix}$$

($N \times 1$ Vector)

So, we can write,

$$X\beta = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix}$$

$$X\beta = \hat{y}$$

We already established,

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

where, y = Actual output vector
 y_i = Actual output for each observation

\hat{y}_i = predicted output for each observation

In Simple LR we were calculating the error for each observation point of training set and then summing it up their square.
 Here also we will do the same.

Let's say,
 error vector,

$$e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix} = \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_N - \hat{y}_N \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix} = y - \hat{y}$$

Solution for multiple LR because
 we will be working on vectorised
 taking square of a expression of p+1 terms then take derivative
 p times will be exacting task.

Total Error vector =

$$\text{So, } E = e_1^2 + e_2^2 + \dots + e_N^2$$

$$= \sum_{i=1}^N e_i^2$$

$$\begin{bmatrix} e_1^2 \\ e_2^2 \\ \vdots \\ e_N^2 \end{bmatrix}$$

$$E = e^T e$$

$$\Rightarrow E = (y - \hat{y})^T (y - \hat{y})$$

$$\Rightarrow E = (y - X\beta)^T (y - X\beta)$$

Recap

$$\begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \quad (e^T)$$

$$\begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \quad (e)$$

$$= e_1^2 + e_2^2 + e_3^2$$

$$= \sum_{i=1}^3 e_i^2 = e^T e = e e^T$$

Recap

$$\text{Let, } A = \begin{bmatrix} a & b \\ c & d \end{bmatrix} \quad B = \begin{bmatrix} e & f \\ g & h \end{bmatrix}$$

$$A^T = \begin{bmatrix} a & c \\ b & d \end{bmatrix} \quad B^T = \begin{bmatrix} e & g \\ f & h \end{bmatrix}$$

$$A - B = \begin{bmatrix} a-e & b-f \\ c-g & d-h \end{bmatrix}, \quad (A-B)^T = \begin{bmatrix} a-e & c-g \\ b-f & d-h \end{bmatrix}$$

$$A^T - B^T = \begin{bmatrix} a-e & c-g \\ b-f & d-h \end{bmatrix} = (A-B)^T$$

$$E = \left(Y^T - (X\beta)^T \right) (Y - X\beta)$$

~~Recap~~

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$B = \begin{bmatrix} e & f \\ g & h \end{bmatrix}$$

$$A^T = \begin{bmatrix} a & c \\ b & d \end{bmatrix}$$

$$B^T = \begin{bmatrix} e & g \\ f & h \end{bmatrix}$$

$$AB = \begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

$$(AB)^T = \begin{bmatrix} ae + bg & ce + dg \\ af + bh & cf + dh \end{bmatrix}$$

$$B^T A^T = \begin{bmatrix} e & g \\ f & h \end{bmatrix} \begin{bmatrix} a & c \\ b & d \end{bmatrix}$$

$$= \begin{bmatrix} ea + bg & ec + gd \\ fa + bh & fc + hd \end{bmatrix}$$

$$(AB)^T = B^T A^T$$

In general,

$$(ABC)^T = C^T B^T A^T$$

so,

... (a)

$$E = (y^T - \beta^T x^T) (y - x\beta)$$

$$E = y^T y - y^T x \beta - \beta^T x^T y + \beta^T x^T x \beta$$

So, we got our equation for E . Now we need to minimize E . For simple linear Regression we took derivative w.r.t m & c . So, here we will take derivative w.r.t β

So, we are talking about differentiation of a matrix with respect to a vector.

$$\frac{\partial E}{\partial \beta} = 0$$

$$\Rightarrow \frac{\partial (y^T y - y^T X \beta - \beta^T X^T y + \beta^T X^T X \beta)}{\partial \beta} = 0$$

$$\Rightarrow \frac{\partial (y^T y)}{\partial \beta} - \frac{\partial (y^T X \beta)}{\partial \beta} - \frac{\partial (\beta^T X^T y)}{\partial \beta} + \frac{\partial (\beta^T X^T X \beta)}{\partial \beta} = 0$$

- - - (1)

Recap:-

Matrix Derivative

There are 6 common types of matrix derivatives.

Type	Scalar	Vector	Matrix
Scalar	$\frac{\partial y}{\partial x}$	$\frac{\partial \mathbf{y}}{\partial x}$	$\frac{\partial \mathbf{Y}}{\partial x}$
Vector	$\frac{\partial \mathbf{y}}{\partial x}$	$\frac{\partial \mathbf{y}}{\partial x}$	
Matrix	$\frac{\partial \mathbf{y}}{\partial x}$		

Case 1 :- When y & x both are scalar, like $y = x^3 + x^2$

$$y' = \frac{dy}{dx}$$

Case 2 :- When y is vector and x is scalar.

Like, $\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_N}{\partial x} \end{bmatrix}$$

Case 3:-

when y is a matrix and x is scalar.

like

$$y = \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1P} \\ \vdots & \vdots & \ddots & \vdots \\ y_{N1} & y_{N2} & \dots & y_{NP} \end{bmatrix}$$

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_{11}}{\partial x} & \frac{\partial y_{12}}{\partial x} & \dots & \frac{\partial y_{1P}}{\partial x} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_{N1}}{\partial x} & \frac{\partial y_{N2}}{\partial x} & \dots & \frac{\partial y_{NP}}{\partial x} \end{bmatrix}$$

Case 4 :- When y is scalar and x is a vector.

like, $x = [x_1 \ x_2 \ \dots \ x_N] \Rightarrow \frac{\partial y}{\partial x} = \left[\frac{\partial y}{\partial x_1} \ \frac{\partial y}{\partial x_2} \ \dots \ \frac{\partial y}{\partial x_N} \right]$

Case 5 :- When y is vector and x is vector.

like, $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}, \quad x = [x_1 \ x_2 \ \dots \ x_n]$ (Numerator layout notation)

$$\frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & & \vdots \\ \frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix}$$

Case 6 :-

when y is a scalar and X is a Matrix.

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NP} \end{bmatrix}$$

$$\frac{\partial y}{\partial X} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \dots & \frac{\partial y}{\partial x_{N1}} \\ \frac{\partial y}{\partial x_{12}} & \dots & \frac{\partial y}{\partial x_{N2}} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{1P}} & \dots & \frac{\partial y}{\partial x_{NP}} \end{bmatrix}$$

Some derivation

Derivation 1:

Lets say a is a vector x is another vector.
then

$$\frac{\partial a^T x}{\partial x} = \frac{\partial x^T a}{\partial x} = a^T$$

We already saw $a^T x = x^T a = a_1 x_1 + a_2 x_2 + \dots + a_n x_n$

$$\frac{\partial a^T x}{\partial x} = \frac{\partial (a_1 x_1 + a_2 x_2 + \dots + a_n x_n)}{\partial x} \quad (\text{This is Case 4})$$

$$= \left[\frac{\partial (a_1 x_1 + a_2 x_2 + \dots + a_n x_n)}{\partial x_1} \dots \right]$$

$$\left[\frac{\partial (a_1 x_1 + a_2 x_2 + \dots + a_n x_n)}{\partial x_n} \right]$$

$$= [a_1 \ a_2 \ \dots \ a_n] = a^T$$

Derivation 2 :- if A is a Matrix and x is a vector
the

$$\boxed{\frac{\partial Ax}{\partial x} = A}$$

Lets,

$$Ax = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11}x_1 + \dots + a_{1n}x_n \\ \vdots \\ a_{n1}x_1 + \dots + a_{nn}x_n \end{bmatrix}$$

$$\frac{\partial Ax}{\partial x} = \begin{bmatrix} \frac{\partial(a_{11}x_1 + \dots + a_{1n}x_n)}{x_1} & \dots & \frac{\partial(a_{11}x_1 + \dots + a_{1n}x_n)}{x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial(a_{n1}x_1 + \dots + a_{nn}x_n)}{x_1} & \dots & \frac{\partial(a_{n1}x_1 + \dots + a_{nn}x_n)}{x_n} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} = A$$

(Case 5)

Derivation 3:-

Lets say A and x are matrix and a vector
Consicutively, then

$$\boxed{\frac{\partial x^T A x}{\partial x} = x^T (A + A^T)}$$

We need to memorise one product rule here,

$$\boxed{\frac{\partial u^T v}{\partial x} = u^T \frac{\partial v}{\partial x} + v^T \frac{\partial u}{\partial x}}$$

Applying the above rule, $u = x, v = Ax$

$$\begin{aligned}\frac{\partial x^T A x}{\partial x} &= x^T \frac{\partial Ax}{\partial x} + (Ax)^T \frac{\partial x}{\partial x} \\ &= x^T A + (Ax)^T. \quad [\text{from derivation 2}]\end{aligned}$$

$$\begin{aligned}\frac{\partial (x^T A x)}{\partial x} &= x^T A + x^T A^T \\ &= x^T (A + A^T) \quad [\text{distributive property}]\end{aligned}$$

now, lets apply whatever we learn so far to compute the derivative of E. From equation (1) we get,

$$\frac{\partial E}{\partial \beta} = \frac{\partial (y^T y)}{\partial \beta} - \frac{\partial (y^T x \beta)}{\partial \beta} - \frac{\partial (\beta^T x^T y)}{\partial \beta} + \frac{\partial (\beta^T x^T x \beta)}{\partial \beta}$$

now,

$\frac{\partial (y^T y)}{\partial \beta} = 0$

(There is no β term)

$$\frac{\partial (y^T x \beta)}{\partial \beta} = \frac{\partial A \beta}{\partial \beta}$$

[$A = y^T x$ = a matrix]

\Rightarrow $\frac{\partial (y^T x \beta)}{\partial \beta} = y^T x$ [Derivation 2]

$$\frac{\partial (\beta^T x^T y)}{\partial \beta} = \frac{\partial (\beta^T a)}{\partial \beta}$$

[$a = x^T y$]

\Rightarrow $\frac{\partial (\beta^T x^T y)}{\partial \beta} = (x^T y)^T$ [Derivation 1]

$$\begin{aligned}
 \frac{\partial(\beta^T X^T X \beta)}{\partial \beta} &= \frac{\partial(\beta^T A \beta)}{\partial \beta} \quad [A = X^T X] \\
 &= \beta^T (A + A^T) \quad [\text{Derivation 3}] \\
 &= \beta^T (X^T X + (X^T X)^T) \\
 &= \beta^T (X^T X + X^T X)
 \end{aligned}$$

\Rightarrow

$$\boxed{
 \begin{aligned}
 \frac{\partial(\beta^T X^T X \beta)}{\partial \beta} &= 2 \cdot \beta^T X^T X
 \end{aligned}}$$

$$\text{So, } \frac{\partial E}{\partial \beta} = 0 - y^T x - (x^T y)^T + 2\beta^T x^T x$$

$$= -y^T x - y^T x + 2\beta^T x^T x$$

$$= 2\beta^T x^T x - 2y^T x$$

Setting $\frac{\partial E}{\partial \beta} = 0$

$$\Rightarrow 2\beta^T x^T x - 2y^T x = 0$$

$$\Rightarrow \beta^T x^T x = y^T x$$

$$\Rightarrow x^T x \beta = x^T y \quad [\text{taking transpose both side}]$$

$$\Rightarrow \boxed{\beta = (x^T x)^{-1} x^T y}$$

Our equation to find
co-efficients.