# Regression

$$Y = f(x)$$

Dependent variable
(Number)
Continuous

Independent variable
(features)

package (CTC)

# of questions solved

Can't draw 1 line to fit all points
$\Rightarrow$ "Best fit" line

"Linear Regressor"

package (CTC)

#of questions solved

→ New test point

$$y = wx + b$$

co-efficient    intercept (bias : o/p in the absence of input)

| x(# of q solved) | y (CTC) |
|---|---|
| 20 | 6 |
| 35 | 9 |
| 50 | 14 |
| 70 | 19 |
| 100 | ? |

# How do you solve for $w$ & $b$?

① Ordinary Least Squares
  — Closed form solution
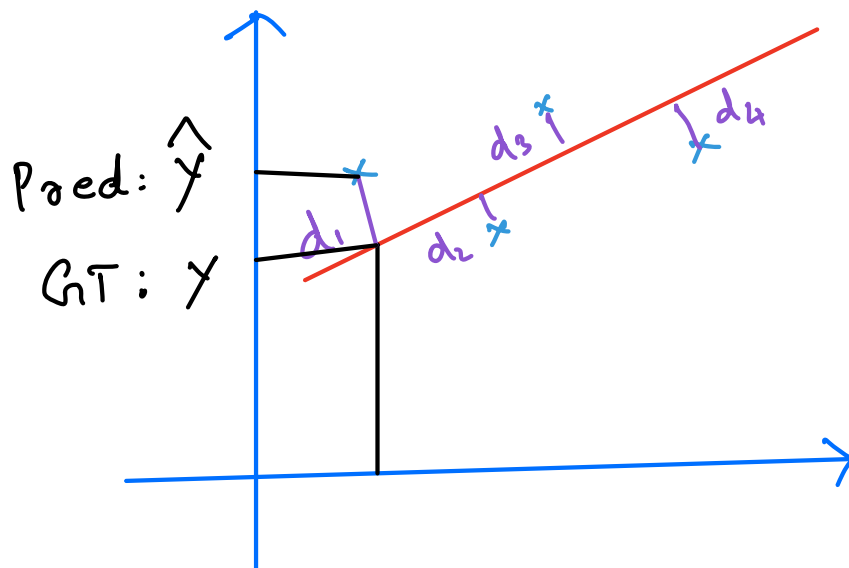


If $d_1 \ldots d_4 \} = 0$
$\Rightarrow$ Perfect fit!

$\downarrow$

Find a line that minimizes this!

Residuals

$$J = (d_1)^2 + d_2^2 + \ldots + d_n^2$$

$$J = \sum_{i=1}^{n} d_i^2 \longrightarrow \text{Find 'm' \& 'b' that min. } J$$

where are $m$ & $b$?

Pred: $\hat{Y}$

GT: $Y$

$$d_j^2 = \left(X_i - \hat{X}_i\right)^2 + \left(Y_i - \hat{Y}_i\right)^2$$

$$\Downarrow$$

$$O$$

$$\Rightarrow d_i^2 = \left(Y_i - \hat{Y}_i\right)^2$$

$$\therefore J = \sum_{i=1}^{n} \left(X_i - \hat{X}_i\right)^2 \longrightarrow \omega X_i + b$$

$$J(\omega, b) = \sum_{i=1}^{\wedge} \left( Y_i - \omega X_i - b \right)^2$$

① If $\omega$ is constant :



② If $b$ is constant



Maxima

Minima
(slope = 0)

$$\frac{\partial J}{\partial \omega} = 0 \qquad \frac{\partial J}{\partial b} = 0$$

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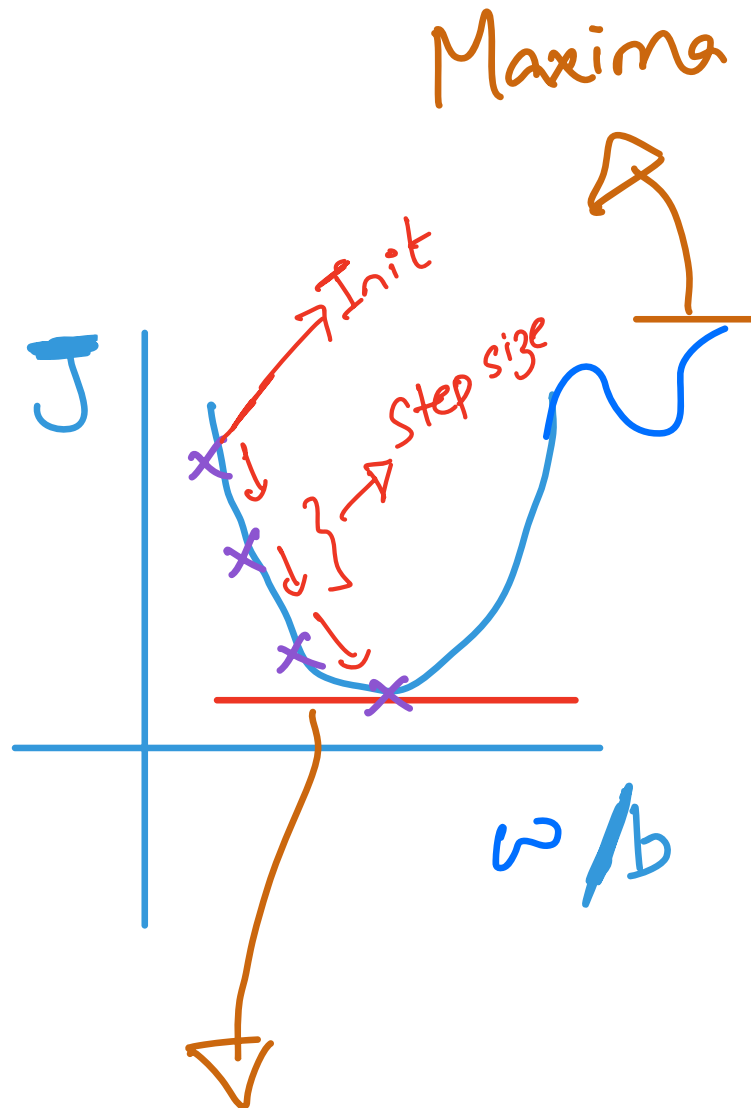
$$b = \bar{y} - \omega \bar{x}$$

$$\bar{x}, \bar{y} : \text{Mean.}$$

② Gradient Descent

— Not closed form solution



- Higher dimensional space
- Approximation algo

# Multiple Linear Regression

$X$ : # of questions solved

College

Cgpa

Company → ordinal data

$y$ : CTC

$$Y = \underbrace{w_1 X_1}_{\#q} + \underbrace{w_2 X_2}_{\substack{small = 0 \\ med = 1 \\ large = 2}} + \underbrace{w_3 X_3}_{\#\,cgpa}$$

$$\left. \begin{array}{l} + w_4 X_4 \\ + w_5 X_5 \\ + w_6 X_6 \end{array} \right\} \longrightarrow$$

$$+ w_7$$

$$\underline{\frac{Company}{one\text{-}hot}}$$

001 → google

010 → MS

100 → apple

→ categorical

# Collinearity

How features are related?

$\#$ questions solved $\propto$ College

"Variance Inflation Factor"

Checks Correlation

1 : no collinearity

1-5 : moderate
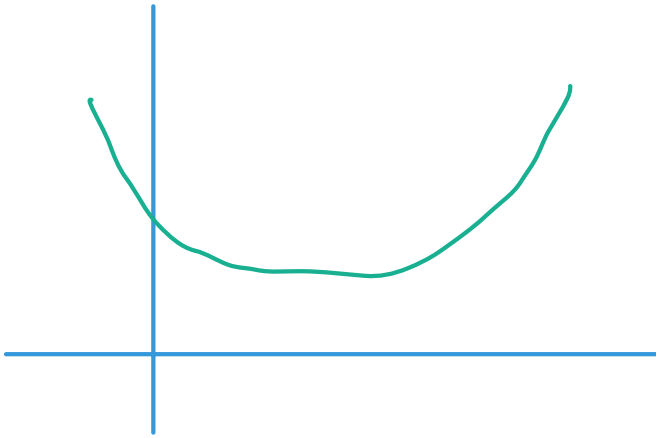
$\geq 5$ : severe

Mitigation strategy

Lowers VIF

① Centering features

$$x = x - \bar{x}$$

$$Y = W_1 x_1 + W_2 x_2 + W_3 x_1 x_2 + W_4 x_1 x_1$$

Correlation

Multiple : Matrices

$X_s$     Calc derivatives to minimize $J$

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

$$MSE = \frac{1}{m} \sum_i (\hat{y}_i - y_i)^2$$

$$= \frac{1}{m} \| \hat{y} - y \|_2^2$$

Mean
Squared
Error

Minimize MSE i.e, gradient/slope=0

$$\nabla MSE = 0$$

$$\Rightarrow \nabla \frac{1}{m} \| \hat{y} - y \|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w \| Xw - y \|_2^2 = 0$$

$$\Rightarrow \frac{1}{m} \nabla_w (Xw - y)^T (Xw - y) = 0$$

$$(x^2 = x \cdot x = x^T x)$$

$$\Rightarrow \boxed{\omega = \left(X^T X\right)^{-1} X^T y}$$

Normal Equations

# Types of Regression

① Polynomial regression

$$Y = \omega_1 x_1 + \omega_2 x_2^2 + \omega_3 x_1 x_2$$

Quadratic in $X_2$

BUT

Linear in $\omega$ !!!

∴ Normal Equations can be used

(2) Non-linear Regression (in $w$)

$$y = w_1 x_1 + w_2^2 x_2 + w_1 w_2 x_3$$

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

# How has the model improved ?

① More features

$\Rightarrow$ More parameters ('w')

② Representational Capacity

(i) Polynomial in X

(ii) Polynomial in W