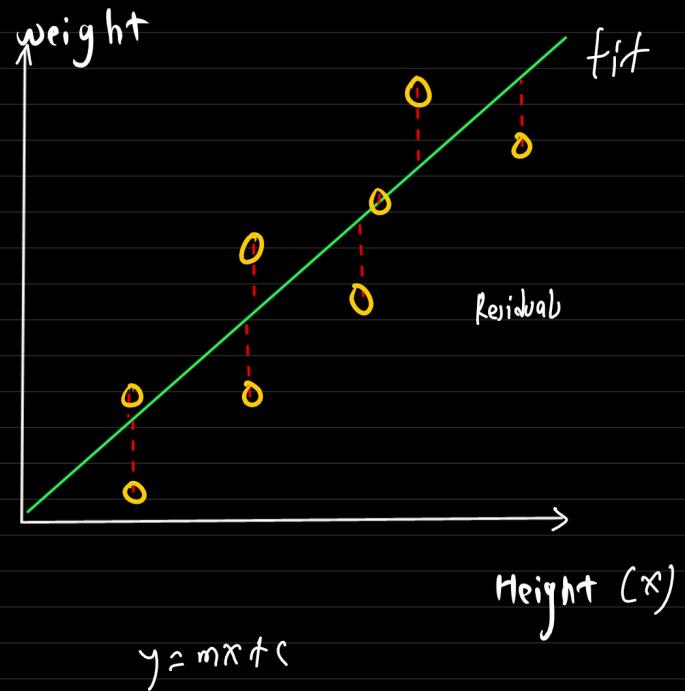
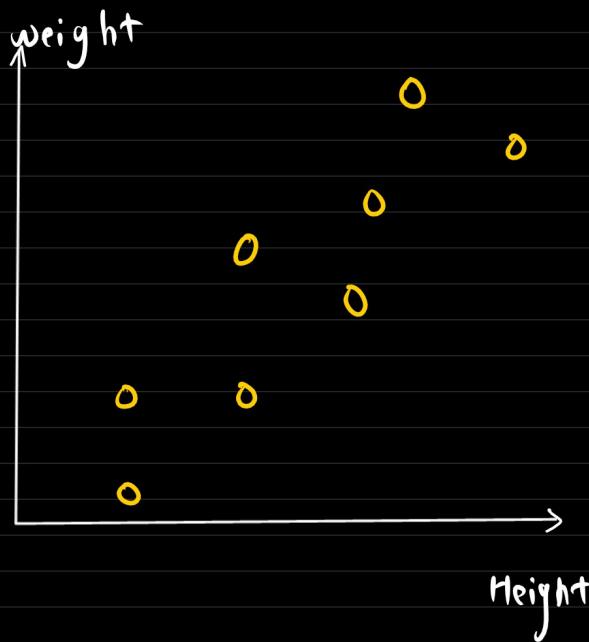


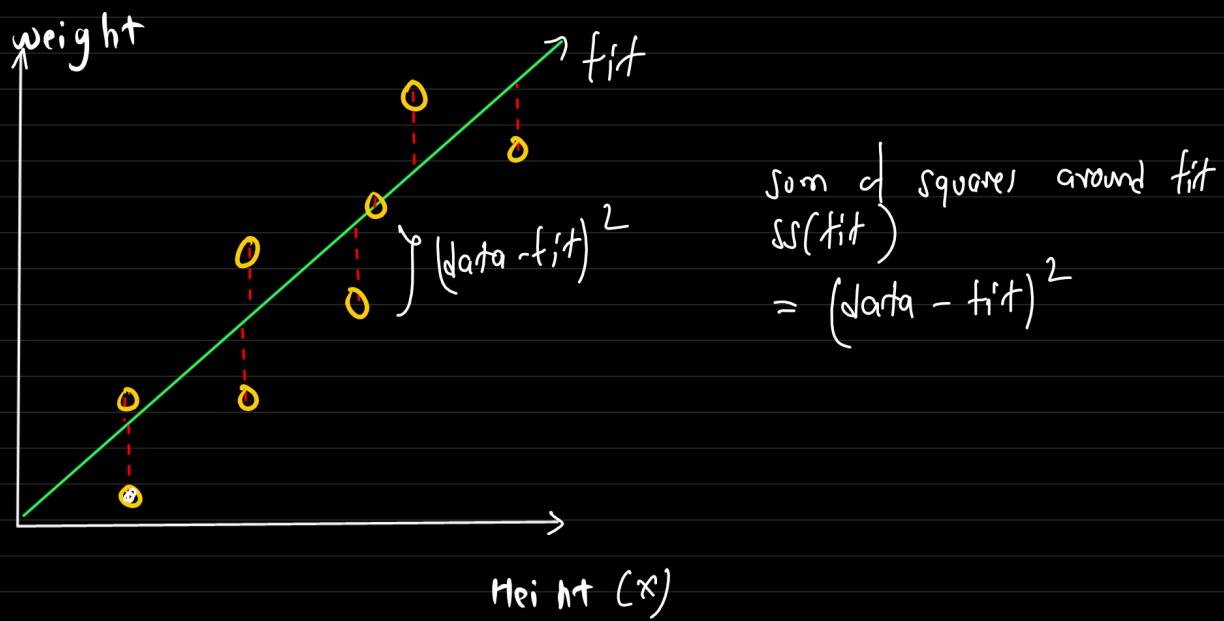
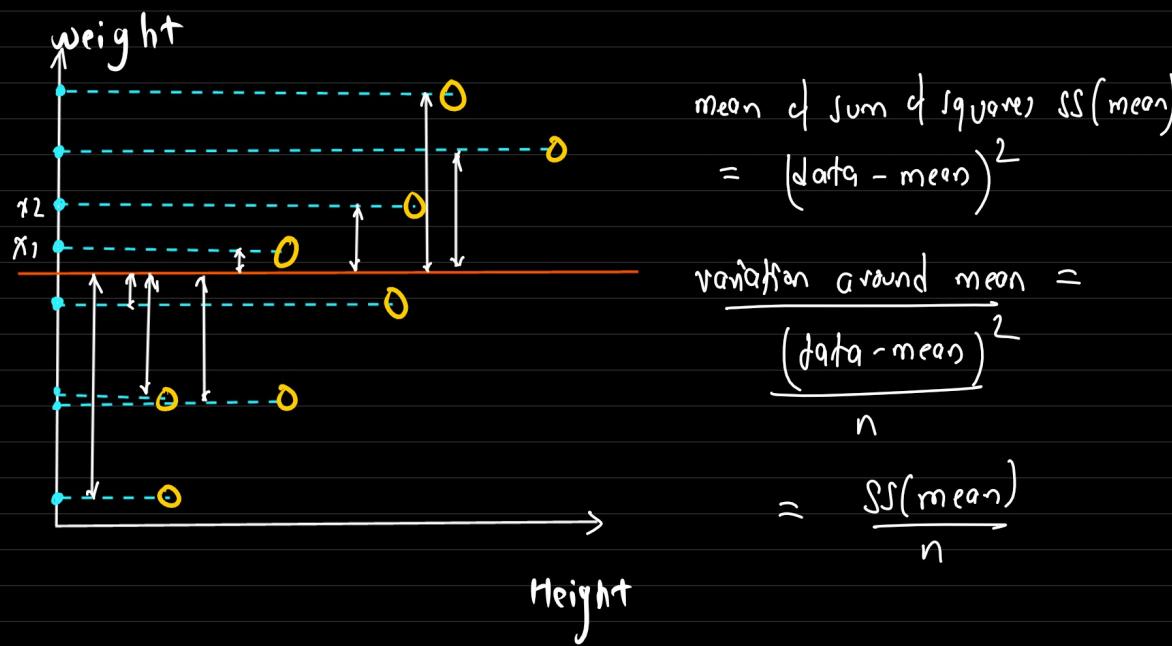
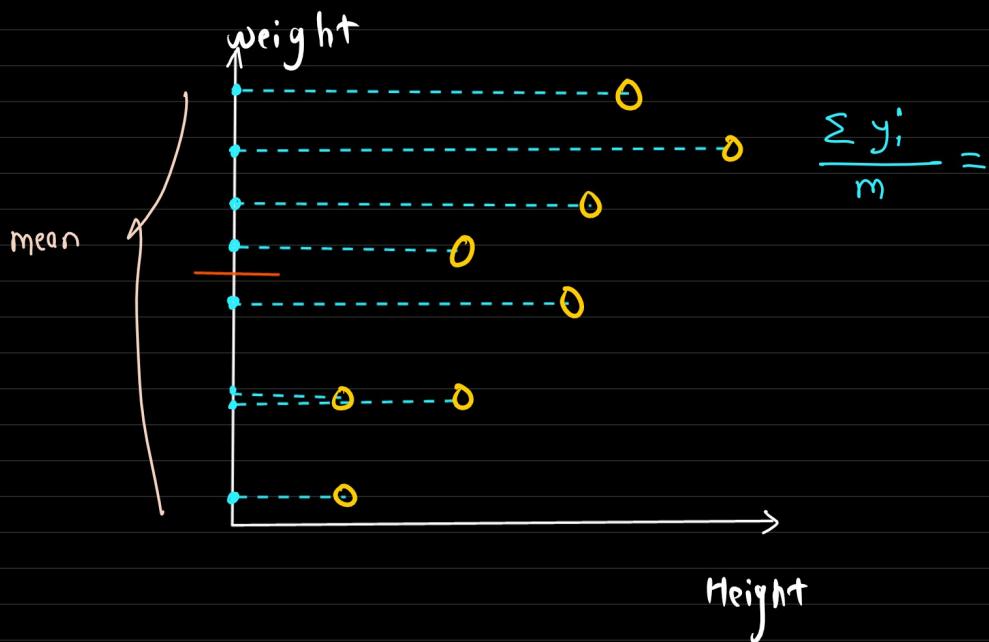
05-10-2025

Agenda:

- $R^2$  score
- Bias & variance
- L1 & L2 regularization
- practicals
- Logistic Regression (Depth)

F-squared score





$$R^2 = \frac{\text{Var}(\text{mean}) - \text{Var}(\text{fit})}{\text{Var}(\text{mean})} = \frac{\text{Var}(\text{weight}) - \text{Var}(\text{fit})}{\text{Var}(\text{weight})}$$

how much of the variation in person weight we can explain by taking person height into account.

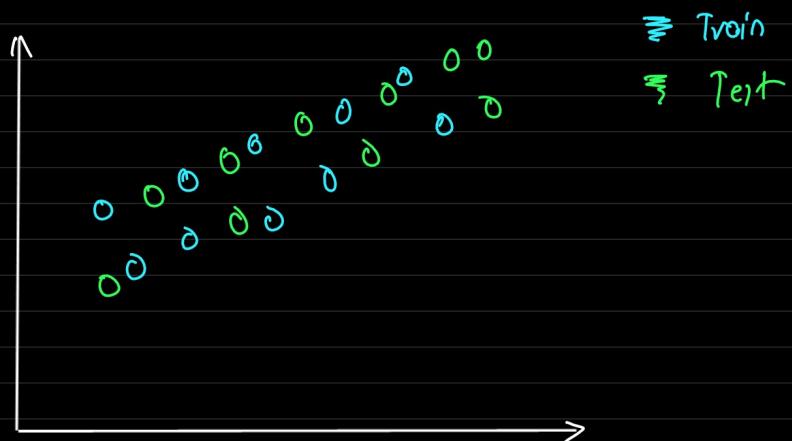
$$\begin{aligned}\text{Var}(\text{mean}) &= 72 \\ \text{Var}(\text{fit}) &= 24 \\ R^2 &= \frac{72 - 24}{72} = 0.66 \approx 60\%.\end{aligned}$$

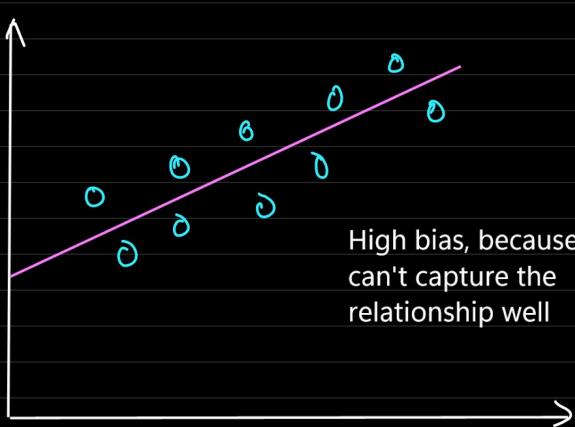
there is 60% reduction in variance when we take person height in account.

person height explains 60% of the variation in person weight.

$R^2$  tells us how much of the variation in the target variable is explained by the model. Higher is better, but context matters.

Big 8 Variable





High bias, because it can't capture the relationship well

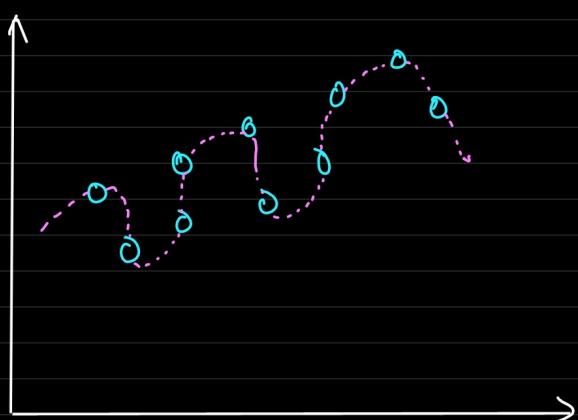


low variance, because of similar fits

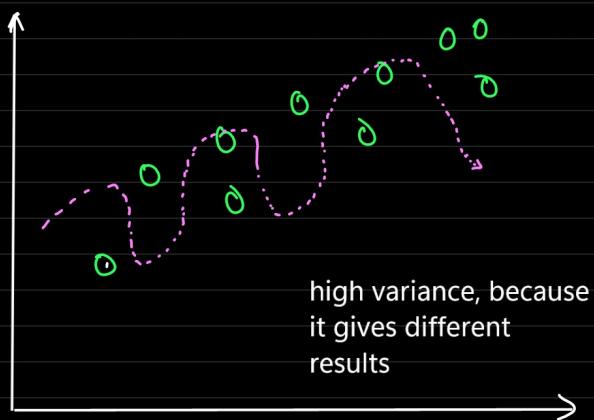
$$\text{loss}(\text{Train}) : 100 \longleftrightarrow \text{loss}(\text{Test}) : 98$$

good fit

The inability for a ML model like linear regression to capture true relation is called as bias.



Low bias, because it is more flexible



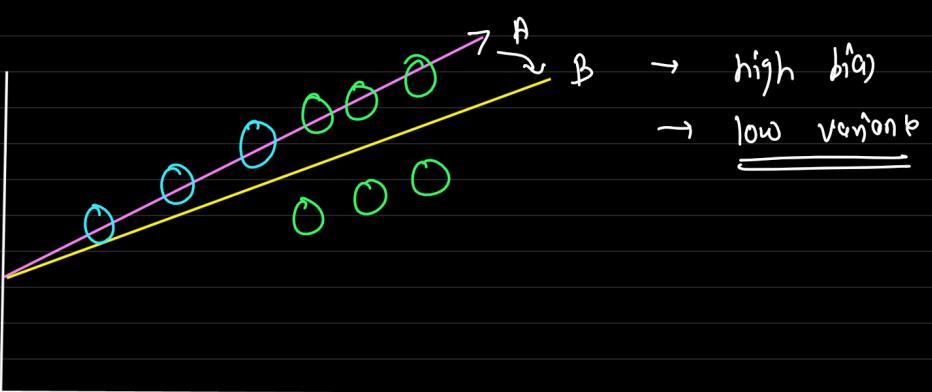
high variance, because it gives different results

$$\text{loss}(\text{Train}) : 0 \longleftrightarrow \text{loss}(\text{Test}) : 70$$

The difference in fits (train, test) is called as variance. The bigger the difference, the higher the variance.

→ Bias & variance tradeoff

→ loss bias & low variance (ideal)



(0))  $(T_{train}) = 0$  : low bias

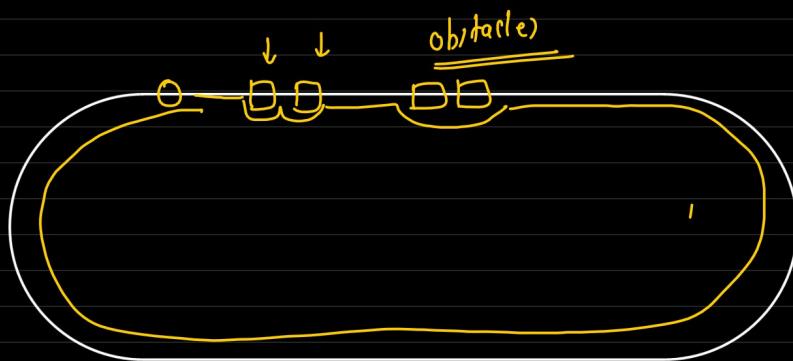
(0))  $(T_{test}) = 100$  : high variance

### Regularization :

#### Driving Test

→ control, Drive  $\infty$ ,  $\bigcirc$ , Traffic sign

→ Traffic,



Don't just minimize error, keep the coefficients controlled.

## The Problems with Vanilla Linear Regression

OLS/SGD works perfectly if:

- You have plenty of data compared to the number of features (columns).
- Features are not strongly correlated (no multicollinearity).
- The model is not too complex (not too many parameters).

Don't just minimize error, keep the coefficients controlled and to not overfit the training data.

L1 ( $(\alpha)_{10}$ ):

$$\text{Loss}_{\text{new}} = \underline{\text{Modified MSE}} + \lambda \sum |\beta_i|$$

X → 11 feature

m → 11 → coefficient / parameter → 100, 200, 300

c → 1

$$y = m_1 x_1 + m_2 x_2 + m_3 x_3 + c$$

$$y = \underline{M}^T \underline{x} + c$$

$$\beta \rightarrow \beta_1 + \beta_2 + \beta_3$$

$$= \underbrace{\frac{1}{2m} \sum_{i=1}^n (y - \hat{y})^2}_{\text{Loss}} + \lambda \sum |\beta_i|$$

X → 11 features, 11 →  $\beta_1 \rightarrow \beta_{11}$

$$\boxed{\beta_1} = 1$$

$$\boxed{\beta_2} = 200$$

$$\beta_3 = 150$$

$$\beta_4 = 2$$

$$\begin{aligned} \textcircled{100} &= 400 + \textcircled{353} \\ \textcircled{100} &= \underline{753} + \boxed{ } \end{aligned}$$

$y = 100 \rightarrow$  reduce coefficient value

$$\begin{matrix} m_1 & x_1 \\ 0 & \end{matrix}$$

: multicollinearity (feature correlated with another feature)

: shrink coefficient towards 0, but never 0

L2 (ridge) : penalty : square of coefficients

$$\text{loss}_{\text{new}} = \text{Modified MSE} + \lambda \sum \beta_i^2$$

: feature selection (pick only the most important feature)

: coefficient shrink all the way to exact 0

Linear Regression (vanilla) :

$$\text{Loss(OJS)} : \text{SSE} \rightarrow \mathbb{P}^2 / F$$

$$\text{Loss(OJS)} : \text{SSE} + L1 \\ \text{SSE} + \lambda \sum |\beta_i| \rightarrow \mathbb{P}^2 / F$$

$$\text{Loss(OJS)} : \text{SSE} + L2 \\ \text{SSE} + \lambda \sum (\beta_i)^2 \rightarrow \mathbb{P}^2 / F$$

$$\text{Loss(OD)} : \text{M.MSE} \rightarrow \mathbb{P}^2 / F$$

$$\text{Loss(OD)} : \text{M.MSE} + L1 \\ \text{M.MSE} + \lambda \sum |\beta_i| \rightarrow \mathbb{P}^2 / F$$

$$\text{Loss(OD)} : \text{M.MSE} + L2 \\ \text{M.MSE} + \lambda \sum (\beta_i)^2 \rightarrow \mathbb{P}^2 / F$$

feature  $\rightarrow 100$   $\rightarrow \text{OJS}$  low - medium  
 $\rightarrow 1 \text{ million data points}$   $\rightarrow \text{OD}$   
iter  $\leftarrow$  low - high

→  $\ell^2$  score (evaluation function/method)

→ bias vs variable tradeoff

→ regularization

→ L1 → add penalty (sum of absolute values of coeff) to loss

→ L2 → add penalty (sum of square values of coeff) to loss

## Logistic Regression

Used: classification

BMI = height, weight

$y = x_1, x_2$

$\hat{y} = x_1 m_1 + x_2 m_2$

regression  $\leftarrow \infty$   $\hat{y} = \begin{cases} 70 & |7| \\ 70.1 & |71.4| \\ 70.001 & |76.8| \\ \infty & \infty \end{cases}$

Temp.mt sick ( $y_e$ ) / no  $\rightarrow$  classification.

$\infty$   $y_e$  / no

$\downarrow$   $\hat{y} (y_e)$  / no

$x_1, x_2 \rightarrow \text{ML algo } (m_1 x_1 + m_2 x_2) \rightarrow \hat{y}$

$m_1, m_2$

$x_1, x_2 \rightarrow \text{ML algo } (m_1 x_1 + m_2 x_2) \rightarrow z \rightarrow \text{function} \rightarrow \hat{y}$

$m_1, m_2$

function  $\rightarrow$  sigmoid  $\rightarrow$  logistic function

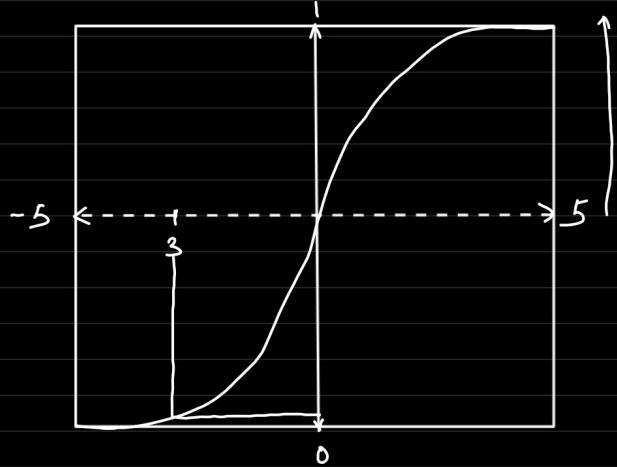
$$\text{Sigmoid} = \frac{1}{1 + e^{-z}}$$

logistic regression

transformation of  
(z) into (p) b/w 0-1

$$p = \frac{1}{1 + e^{-z}}$$

$\rightarrow$  This is sigmoid or logistic function.



$$\begin{array}{l} \text{+ve} \rightarrow 1 \\ \text{-ve} \rightarrow 0 \end{array}$$

Sigmoid (0 - 1)  
Sigmoid (-3)  $\rightarrow$  0.1

$\rightarrow 0.000 | 0.04 | 0.3 | 0.8 | 0.9 | 1$

Temp  $\stackrel{14/10}{\uparrow}$   
 $y$   $\uparrow$

39.9 | 0.8  $\rightarrow$  80% chance person is sick

$\text{Sigmoid}(x) > 0.5 \rightarrow 1 \text{ else } 0$