

2009

$$Y = \beta X + \epsilon$$

STATISTICS

2019

$$Y = \beta X + \epsilon$$

MACHINE LEARNING

✖ 10 YEARS CHALLENGE

LINEAR
REGRESSION

Linear Regression

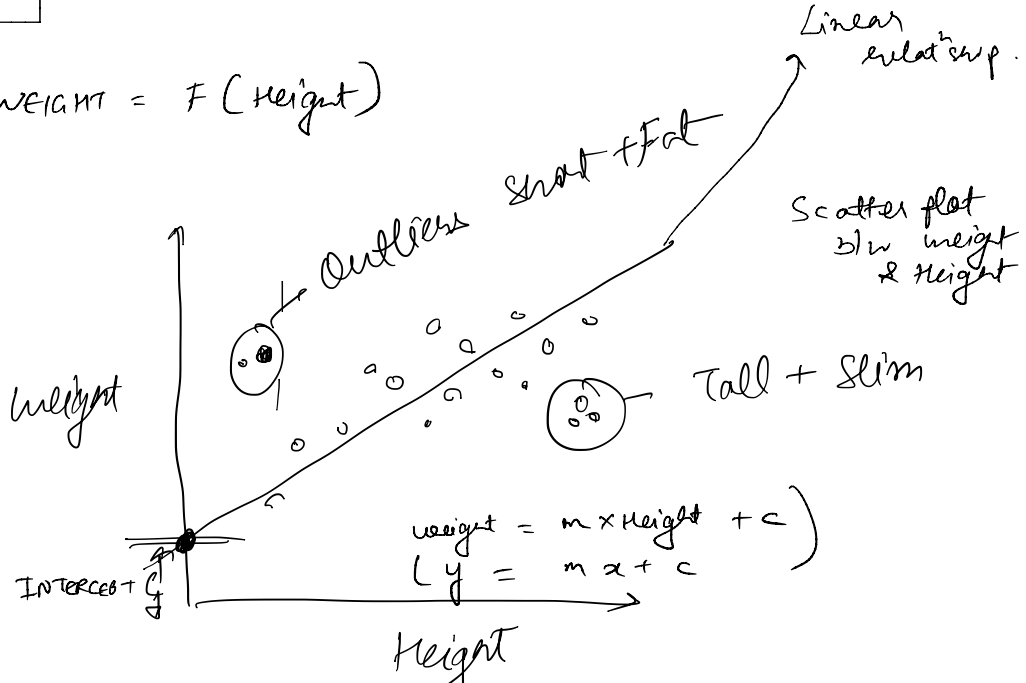
→ iff it continuous in nature

$$F = m \times a$$
$$V = u + a(t)$$

Some linear systems

TASK: PREDICT WEIGHT = $F(\text{height})$

height	weight
3	25
5	8
⋮	⋮
1	⋮



WRITING THE EXPRESSION IN MATRIX FORM

$$\text{weight}_i = 1 \times \theta_0 + \text{height}_i * \theta_1$$

$$\text{weight}_1 = 1 \times \theta_0 + \text{height}_1 * \theta_1$$

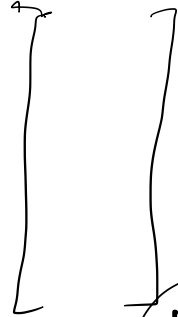
$$\text{weight}_2$$

⋮

$$\text{weight}_N$$

$$1 \times \theta_0 + \text{height}_N * \theta_1$$

weight

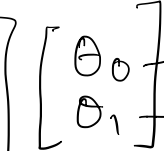


$N \times 1$



height₁

height_N



2×1

Bias

Intercept

N examples

$N \times 2$

Previous example was $y = F(x)$ where x is one-dimensional.

Examples in multiple dimensions

IIT GV water demand = $F(\# \text{ occupants}, \text{Temperature})$

$$\text{Demand} = \text{Base demand} + k_1 \times \# \text{ occupants} + k_2 \times \text{Temp.}$$

we expect

Demand \uparrow if $\# \text{ occupants} \uparrow \Rightarrow k_1$ likely positive

Demand \uparrow if Temp $\uparrow \Rightarrow k_2$ likely positive

Base demand is demand independent of temp. and $\# \text{ occupants}$.

Bias

MORE GENERALLY

M features

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1} = \begin{bmatrix} 1 \\ \vdots \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \dots \\ x_{21} & \dots & \dots \\ \vdots & \vdots & \vdots \\ x_{N1} & \dots & \dots \end{bmatrix} \begin{bmatrix} \theta_0 \\ \vdots \\ \vdots \\ \theta_M \end{bmatrix}_{(M+1) \times 1}$$

$(N \times (M+1))$

$$\underbrace{y}_{N \times 1} \approx \underbrace{X}_{N \times (M+1)} \underbrace{\theta}_{(M+1) \times 1}$$

UNKNOWN

KNOWN

N knowns
M unknowns

(or N equations)
(or M parameters/variables)

TRIVIAL CASE (Back to the weight example)

$$y_{N \times 1} \approx X \theta$$

$N \times (m+1)$ $(m+1) \times 1$

$m=1$

$$\therefore y = \theta_1 x + \theta_0$$

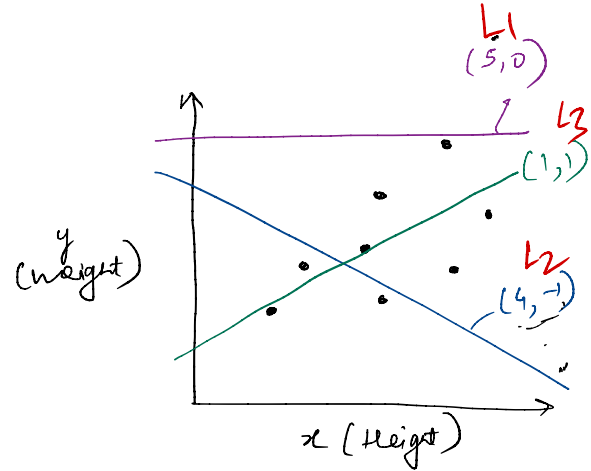
(Same form as $y = mx + c$)

For different θ_0 & θ_1 , different relationship can be learnt.

3 examples

θ_0	θ_1
5	0
4	-1
1	1

} which is best?

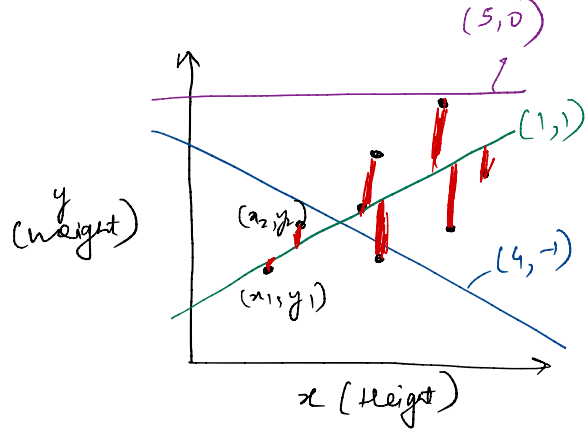


For different θ_0 & θ_1 , different relationship can be learnt.

3 (of ∞ possible parameters)

$$\begin{matrix} \theta_0 & \theta_1 \\ 5 & 0 \\ 4 & -1 \\ 1 & 1 \end{matrix}$$

} which is best?



CHOOSE θ_0, θ_1 s.t.

$|e_{y_i}|$ are reduced, where $e_{y_i}^2$

$$e_{y_i} = y_i - \hat{y}_i$$

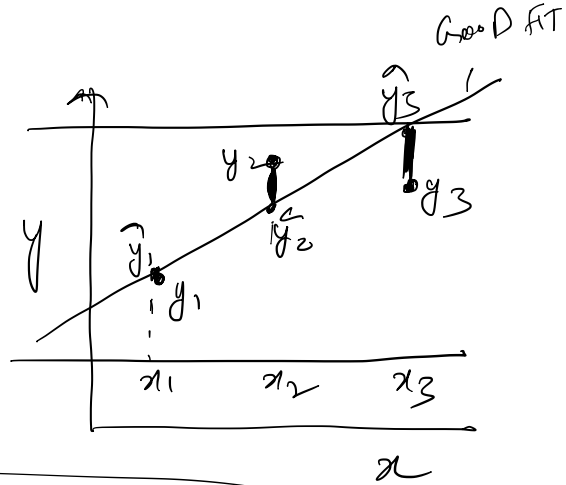
Residual G.T. Pred.

$$e_{y_i} = y_i - (\theta_0 + \theta_1 * x_i)$$

θ_0, θ_1 : Same value for all examples
 e_{y_i} : Variable residual for example

For good fit

$|e_{y1}|, |e_{y2}|, |e_{y3}|, \dots$
all should be
small

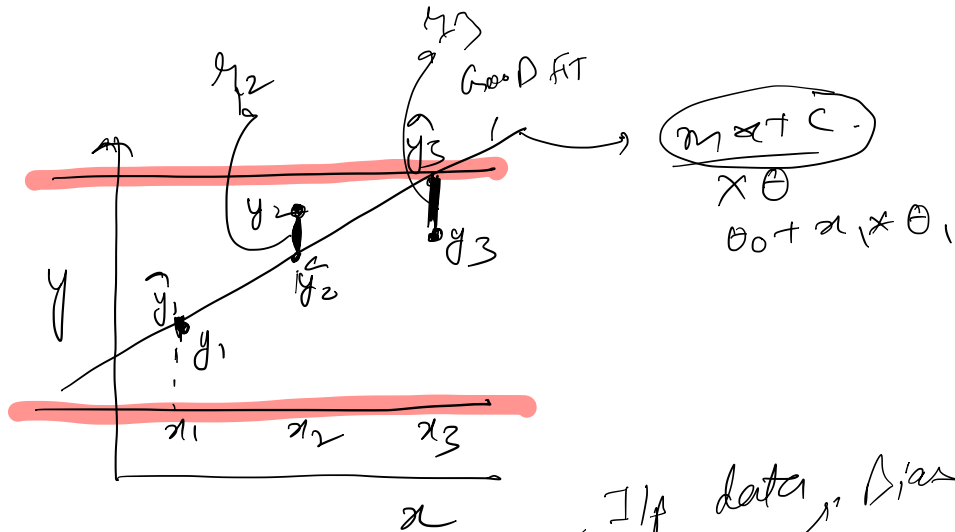


Minimize $(e_{y1})^2 + (e_{y2})^2 + \dots + (e_{yn})^2$ L_2

or $|e_{y1}| + |e_{y2}| + \dots + |e_{yn}|$ L_1

TRAIN | DATA

x_1 | y_1
 x_2 | y_2
 x_3 | y_3

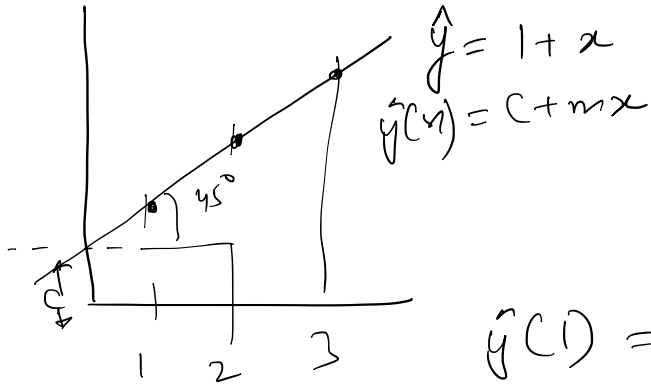


$$y = X\theta + \epsilon_y$$

If data, bias term

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ \vdots \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix} + \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_{y1} \\ \epsilon_{y2} \\ \epsilon_{y3} \end{bmatrix}$$

$\hat{y} = (M \times (C+1))$



$$\hat{y}(1) = 2 = y(1)$$

$$\hat{y}(2) = 3 = y(2)$$

$$e_{y1} = \hat{y}(1) - y(1) = 0$$

NORMAL ϵ_0^N

$$y \approx X \Theta (+ \epsilon_y)$$

 vector
 matrix

To learn: Θ

Objective: Minimize $\epsilon_{y_1}^2 + \dots + \epsilon_{y_N}^2$

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

$$\epsilon_y^T = [\epsilon_{y_1} \dots \epsilon_{y_N}]$$

Equivalent to .

Minimize $\epsilon_y^T \epsilon_y$

$$X: N \times (M+1)$$

$$\theta: (M+1) \times 1$$

$$e_y = y - X\theta$$

$$e_y^T = (y - X\theta)^T = y^T - \theta^T X^T$$

$$e_y^T e_y = (y^T - \theta^T X^T)(y - X\theta)$$

$$e_y^T e_y = \underset{1 \times N}{y^T} \underset{N \times 1}{y} - \underset{1 \times (M+1)}{\theta^T} \underset{(M+1) \times N}{X^T} \underset{N \times 1}{y} - \underset{1 \times N}{y^T} \underset{N \times 1}{X\theta} + \theta^T X^T X \theta$$

Both these are the same

$$= y^T y - 2y^T X \theta + \theta^T X^T X \theta$$

Objective minimize $e_y^T e_y = y^T y - 2y^T X \theta + \theta^T X^T X \theta$

$$\text{Objective Minimize } e_y^T e_y = \overset{\textcircled{1}}{y^T y} - 2 \overset{\textcircled{2}}{y^T x} \theta + \theta^T x^T x \theta$$

w.r.t. θ

$$\frac{\partial e_y^T e_y}{\partial \theta} = 0$$

$$\textcircled{1} \frac{\partial y^T y}{\partial \theta} = \vec{0}$$

$$\textcircled{2} \frac{\partial (-2 \underbrace{y^T x}_A \theta)}{\partial \theta} = A^T = (-2 y^T x)^T = -2 x^T y$$

$$(3) \quad \frac{\partial}{\partial \theta} (\theta^T \underbrace{X^T X} \theta) = 2 X^T X \theta$$

$$0 = -2 X^T y + 2 X^T X \theta$$

$$X^T y = X^T X \theta \dots$$

$$(X^T X)^{-1} X^T y = \theta$$

NORMAL \mathbb{R}^N

Relationship b/w # variables (M) & # examples (N)

(b) $N < M$

example

$$\text{weight}_i = \theta_0 + \text{height}_i \times \theta_1 + \text{Age}_i \times \theta_2$$

$$N=2$$

$$M=3$$

$$\begin{bmatrix} 30 \\ 40 \end{bmatrix} = \begin{bmatrix} 1 & 6 & 30 \\ 1 & 5 & 20 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}$$

$$30 = \theta_0 + 6\theta_1 + 70\theta_2$$

$$40 = \theta_0 + 5\theta_1 + 20\theta_2$$

$$-10 = -\theta_1 - 10\theta_2$$

How many solⁿ? \rightarrow Infinite!!

under-determined system: $e_i = 0 \forall i$

(ii)

$$N > M$$

over-determined

Sum of squared residuals > 0

Question 1

we had

$$X^T y = X^T X \theta$$

can we do?

$$\underbrace{(X^T)^{-1}} X^T y = \underbrace{(X^T)^{-1}} X^T X \theta$$

$$y = X \theta$$

$$X^{-1} y = \theta$$

X may not be a square matrix

X^{-1} may not exist.

Left inverse

$$HX = I_{m+1}$$

$$X_{N \times (m+1)}$$

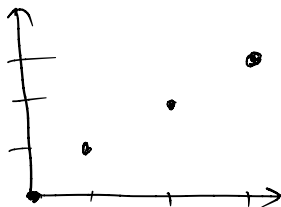
Right inverse

$$XH = I_N$$

WORKED OUT

x	y
0	0
1	1
2	2
3	3

EXAMPLE



Find θ_0 & θ_1

$$X = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 4 & 6 \\ 6 & 14 \end{bmatrix}$$

$$X^T Y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 6 \\ 14 \end{bmatrix}$$

$$|X^T X| = 14 \times 4 - 6 \times 6 = 20$$

$$(X^T X)^{-1} = \frac{\text{adj}(X^T X)}{|X^T X|} = \frac{1}{20} \begin{bmatrix} 14 & -6 \\ -6 & 4 \end{bmatrix}$$

$$\text{adj}(X^T X) = \begin{bmatrix} 4 & -6 \\ -6 & 4 \end{bmatrix}$$

$$\therefore (X^T X)^{-1} = \frac{1}{20} \begin{bmatrix} 14 & -6 \\ -6 & 4 \end{bmatrix}$$

$$X^T y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 6 \\ 14 \end{bmatrix}$$

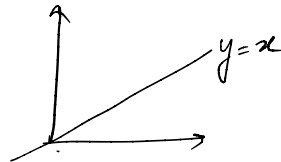
2×4 4×1

$$\theta = (X^T X)^{-1} (X^T y) = \frac{1}{20} \begin{bmatrix} 14 & -6 \\ -6 & 4 \end{bmatrix} \begin{bmatrix} 6 \\ 14 \end{bmatrix}$$

$$= \frac{1}{20} \begin{bmatrix} 14 \times 6 - 6 \times 14 \\ -6 \times 6 + 4 \times 14 \end{bmatrix} = \frac{1}{20} \begin{bmatrix} 0 \\ 20 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

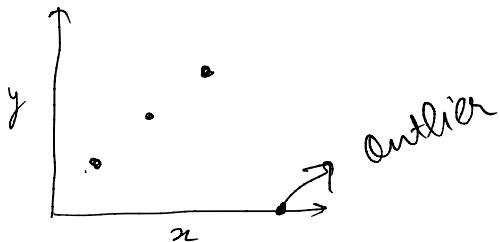
$$= \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$\Rightarrow y = 0 + x$$



CLASS EXERCISE : learn θ_0 & θ_1 (EFFECT OF OUTLIER)

$$\left\{ \begin{array}{l} x \\ y \end{array} \right. \begin{array}{l} 1 \\ 1 \\ 2 \\ 3 \\ 3 \\ 4 \\ 0 \end{array}$$



$$X = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 2 & 2 \\ 3 & 3 \\ 4 & 4 \end{bmatrix} = \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix}$$

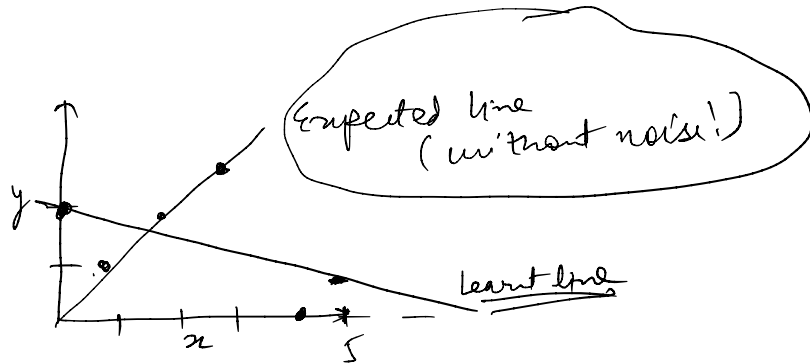
$$|(X^T X)| = 20 \quad (X^T X)^{-1} = \frac{1}{20} \begin{bmatrix} 30 & -10 \\ -10 & 4 \end{bmatrix}$$

$$X^T Y = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 0 \end{bmatrix} = \begin{bmatrix} 6 \\ 14 \end{bmatrix}$$

$$(X^T X)^{-1} (X^T y) = \frac{1}{20} \begin{bmatrix} 30 & -10 \\ -10 & 4 \end{bmatrix} \begin{bmatrix} 6 \\ 14 \end{bmatrix} = \frac{1}{20} \begin{bmatrix} 180 - 140 \\ 56 - 60 \end{bmatrix} = \frac{1}{20} \begin{bmatrix} 40 \\ -4 \end{bmatrix}$$

$$= \begin{bmatrix} 2 \\ -1/5 \end{bmatrix}$$

$$y = 2 - \frac{1}{5}x$$



NEED A WAY TO HANDLE NOISE / OUTLIERS?

CLASS EXERCISE: SOLVE LINEAR REGRESSION

$$\begin{matrix} x_1 & x_2 & y \\ \left[\begin{array}{cc|c} 1 & 2 & 4 \\ 2 & 4 & 6 \\ 3 & 6 & 8 \end{array} \right] \end{matrix}$$

$$\left. \begin{array}{l} y = 2x_1 + 2 \\ y = x_2 + 2 \\ y = \frac{x_2}{2} + x_1 + 2 \end{array} \right\}$$

$$X = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 4 \\ 1 & 3 & 6 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 2 & 4 & 6 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 2 & 3 \\ 2 & 4 & 6 \end{bmatrix} \begin{bmatrix} 1 & 1 & 2 \\ 1 & 2 & 4 \\ 1 & 3 & 6 \end{bmatrix} = \begin{bmatrix} 3 & 6 & 12 \\ 6 & 14 & 28 \\ 12 & 28 & 56 \end{bmatrix}$$

$$\begin{aligned} |X^T X| &= 3(14 \times 56 - 28 \times 28) - 6(6 \times 56 - 28 \times 12) + 12(4 \times 28 - 14 \times 12) \\ &= 0 \end{aligned}$$

$(X^T X)^{-1}$ does not exist

⇒ Matrix X is not full rank

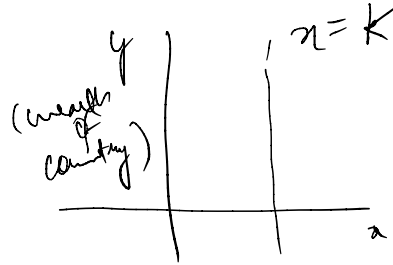
Multi-collinearity: one ^(or more) predictor variable / feature
in X can be expressed as linear combinations of
others

SOLVE (Many ways)

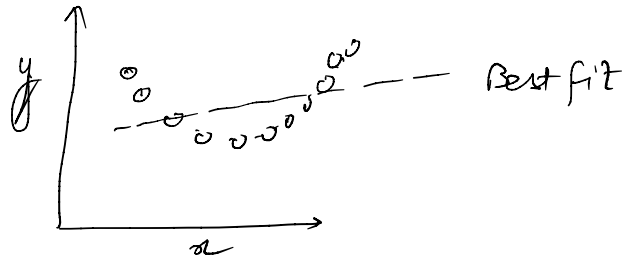
- ① Regularize
- ② Drop variables
- ③ Use different subsets of data
- ④ Avoid dummy variable trap

SOME THINGS TO KNOW ABOUT LINEAR REGRESSION

1) when $x = k$



2) when relationship is non-linear



MODELING

NON-LINEARITIES

Dataset

t	S
0	0
1	6
2	14
3	24
4	36

TRANSFORM

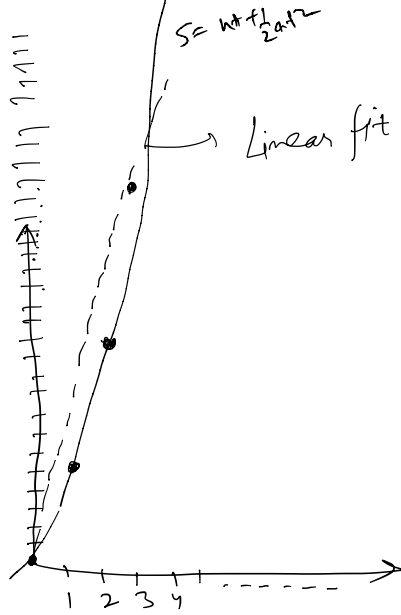


t	t ²	S
0	0	0
1	1	6
2	4	14
3	9	24
4	16	36

X =

$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \\ 1 & 4 & 16 \end{bmatrix}$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 0 & 1 & 2 & 3 & 4 \\ 0 & 1 & 4 & 9 & 16 \end{bmatrix}$$



MODELING INTERACⁿ

$$y = \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m$$

If x_1 increases by 1 unit, y increases by θ_1 units
irrespective of x_2, \dots, x_m

$$y = \theta_0 + \theta_1 x_1 + \underbrace{\theta_2 x_1 x_2}_{\text{Interact term}} + \theta_3 x_2$$

$$= \theta_0 + \theta_1 x_1 + \theta_2 x_2' + \theta_3 x_2 \quad (\text{still linear})$$

DUMMY VARIABLES



$$\text{POLLUTION in DELHI} = \theta_0 + \theta_1 \times \# \text{VEHICLES} + \theta_2 \times \text{WIND SPEED} + \theta_3 \times \text{WIND DIRECTION}$$

(N, E, W, S)

Can we encode wind direction as: $\{N: 0, E: 1, W: 2, S: 3\}$?

No! Implies $S > W > E > N$
(gives some ordering b/w them)

	Is it N?	Is it E?	Is it W?
N	1	0	0
E	0	1	0
W	0	0	1
S	0	0	0

ENCODING
 (USES
 $N-1$ variables for
 N classes)

Why not?

	Is N	Is E	Is W	Is S
N	1	0	0	0
E	0	1	0	0
W	0	0	1	0
S	0	0	0	1

$$x_5 = \begin{bmatrix} x_1 \\ \vdots \\ x_4 \end{bmatrix}$$

Dataset

wind	dir ⁿ	cell ⁿ
N (000)	1	1
E (0100)	2	2
W (0010)	3	3
S (0001)	4	4

$$X = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

Multi-collinearity!!

WHY DUMMY IS ONE-HOT ENCODING & NOT BINARY

BINARY

N	-00
E	-01
W	-10
S	-11

VIS

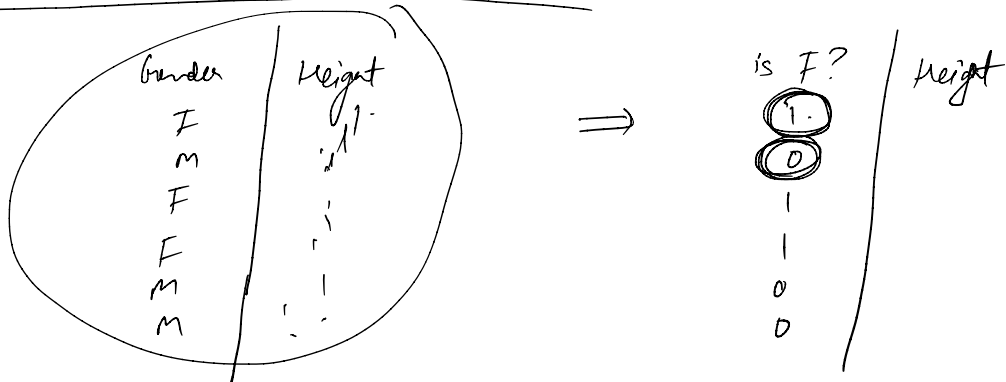
N	-	1	0	0	0
E	-	0	1	0	0
W	-	0	0	1	0
S	-	0	0	0	1

W & S are related by
high bit

THIS INTRODUCES DEPENDENCY
BETWEEN THEM

CAN CONFUSE
CLASSIFIERS ETC...

INTERPRETING DUMMY VARIABLES



$$h_i = \theta_0 + \theta_1 \cdot (\text{is F?}) + \epsilon_i$$

$$h_i = \begin{cases} \theta_0 + \theta_1 + \epsilon_i & \text{if female} \\ \theta_0 + \epsilon_i & \text{if male} \end{cases}$$

- $\theta_0 = \text{Avg. height of male}$
- $\theta_0 + \theta_1 = \text{Avg. height of female}$
- $\theta_1 = \text{Difference in avg. of female and male heights}$

	Is it N?	Is it E?	Is it W?
N	1	0	0
E	0	1	0
W	0	0	1
S	0	0	0

pollⁿ
 {
 |
 |
 |
 |

$$\text{Pollution} = \theta_0 + (\theta_1 (\text{Is it N})) + \theta_2 (\text{Is it E}) + \theta_3 (\text{Is it W})$$

$$P_i = \begin{cases} \theta_0 + \epsilon_i & \text{if wind = South} \\ \theta_0 + \theta_1 & \text{if wind = North} \\ \vdots & \end{cases}$$

θ_0 : Avg. ^{Arg.} pollⁿ South

θ_1 : Difference in arg. of north - south...

$\theta_2 = \dots$
 $\theta_3 = \dots$

Alternative Parameter Estimation (for Linear Eq. in 2 variables)

$$y_i \approx \theta_0 + \theta_1 x_i \quad \text{--- (1)}$$

$$e_i = y_i - \hat{y}_i = \frac{y_i - \theta_0 - \theta_1 x_i}{1}$$

$$\underline{\underline{\sum e_i^2}} = \sum_{i=1}^N (y_i - \theta_0 - \theta_1 x_i)^2$$

$$\text{Derivative } \frac{\partial \sum e_i^2}{\partial \theta_0} = 2 \sum_{i=1}^N (y_i - \theta_0 - \theta_1 x_i) (-1) = 0$$

$$\Rightarrow \sum_{i=1}^N y_i - \theta_1 \sum_{i=1}^N x_i - N\theta_0 = 0$$

$$\Rightarrow \theta_0 = \frac{\sum_{i=1}^N y_i}{N} - \theta_1 \frac{\sum_{i=1}^N x_i}{N} = \bar{y} - \theta_1 \bar{x}$$

$$\sum h_{y_i}^2 = \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2$$

Differentiell: $\frac{\partial \sum h_{y_i}^2}{\partial \theta_1} = 0 \Rightarrow 2 \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i) (-x_i) = 0$

$$\Rightarrow \sum (x_i y_i - \theta_0 x_i - \theta_1 x_i^2) = 0$$

$$\Rightarrow \sum \theta_1 x_i^2 = \sum x_i y_i - \sum \theta_0 x_i$$

$$\Rightarrow \sum \theta_1 x_i^2 = \sum x_i y_i - \sum x_i (\bar{y} - \theta_1 \bar{x})$$

$$\Rightarrow \sum \theta_1 x_i^2 = \sum x_i y_i - \bar{y} \sum x_i + \theta_1 \bar{x} \sum x_i$$

$$\Rightarrow \sum x_i y_i - \sum x_i \bar{y} = \theta_1 (-\bar{x} \sum x_i + \sum x_i^2)$$

$$\Rightarrow \sum x_i y_i - \sum x_i \bar{y} = \theta_1 (-\bar{x} \sum x_i + \sum x_i^2)$$

$$\theta_1 = \frac{\sum x_i y_i - \sum x_i \bar{y}}{\sum x_i^2 - \bar{x} \sum x_i}$$

$$\theta_1 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\text{Cov}(x, y)}{\text{Variance}(x)}$$

