

	HP	<u>W1</u>	Vol	SP	M1a
HP					
W1			<p>Multi Collinearity problem.</p>		
<u>Vol</u>					
SP					
M1a					

$$y = \beta_0 + \beta_1 \underline{WT} + \beta_2 \underline{VOL} + \beta_3 SP + \beta_4 HP$$

↙ ↘ → Correlation

$$\beta_1 = 0; \quad \beta_2 = 0;$$

"WT affects MPG"

✓ WT	MPG

$$\Rightarrow R^2 = \underline{0.277}$$

"VOL affects MPG"

✓ VOL	MPG

$$R^2 = \underline{0.28}$$

↔ x ↔

WT	<del>VOL</del>	MPG

Expected  $R^2 = 50\%$

$$R^2 = 0.28$$

Ignore y (MPa).

$R^2 \downarrow$  No complex 21

$R^2 \rightarrow \therefore$  Variance  $\checkmark$   $R^2 \uparrow$ , Hlm  $x, y$

X-Value

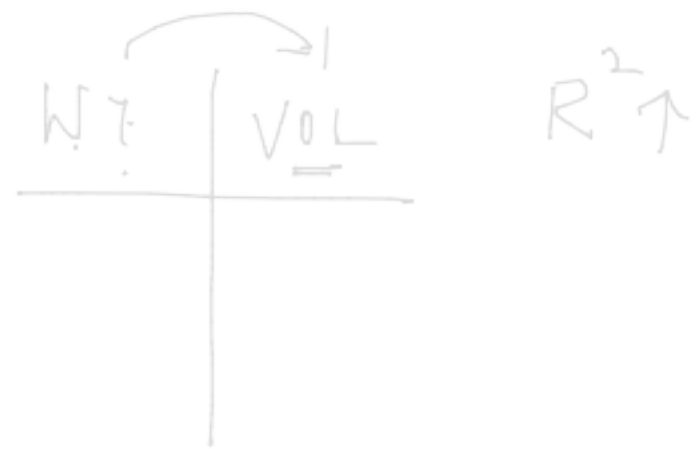
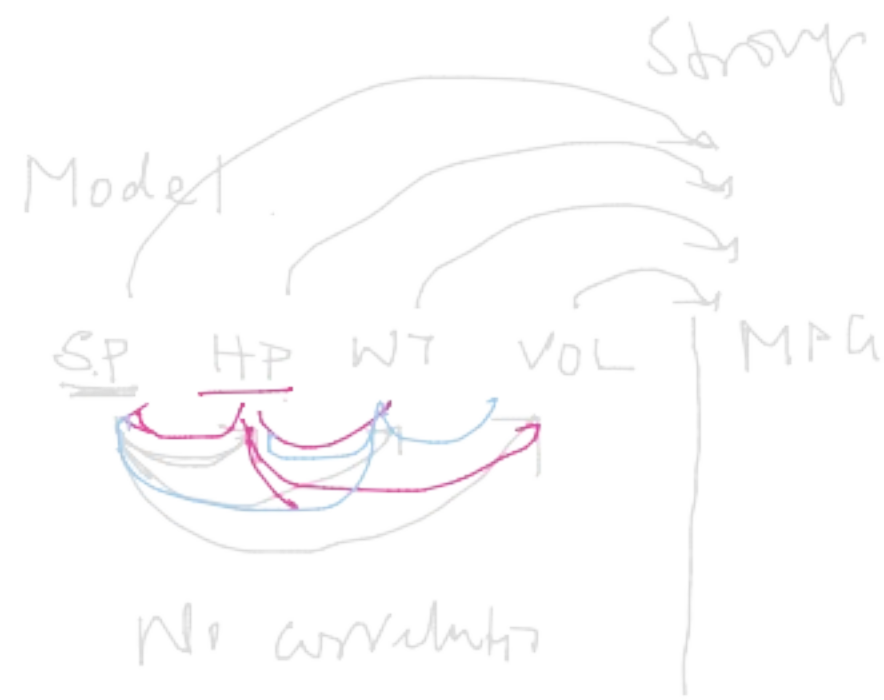
SP | WT VOL SP

↓

Y X

WT | SP HP VOL  $R^2 \downarrow$

VOL | WT SP HP  $R^2$



$y$				
✓ SP	HP	WT	VOL	

$$\frac{R^2 \text{ less}}{VIF \text{ less}}$$

SP - good feature

$y$				
✓ HP	SP	WT	VOL	

$$R^2 = 0$$

HP is good feature

$y$				
✓ WT	SP	VOL	HP	

$$R^2 = 0$$

WT is good

$y$				
✓ VOL	SP	WT	HP	

$$R^2 = 0$$

VOL is good, but

$$\uparrow VIF =$$

$$\frac{1}{(1 - R^2 \uparrow) \downarrow}$$

Large VIF Bad Feature

(1) Decided The best Feature

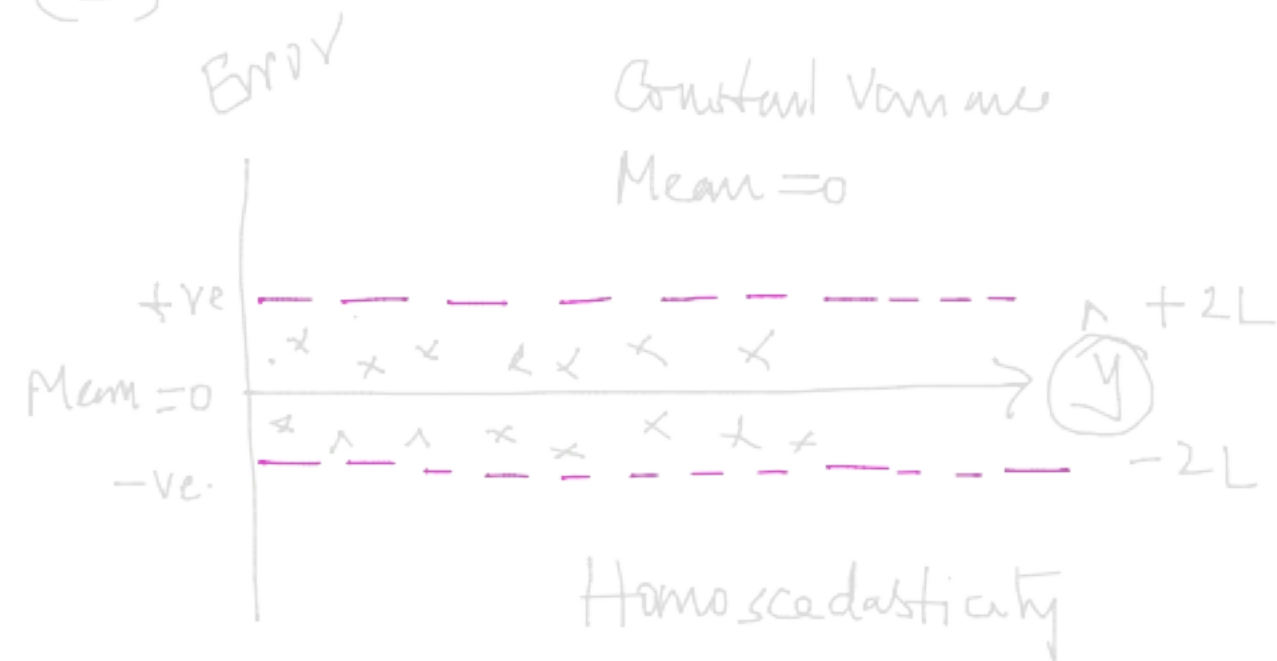
SP	VOL	HP	MPG

✓ 1. Strong correlation betw  $x$  &  $y$

✓ 2. No correlation betw  $x$ -values.

→ removed WT, so now there is no correlation betw  $x$ -values.

(2)



3 Residuals .

✓ Mean = 0 ✓

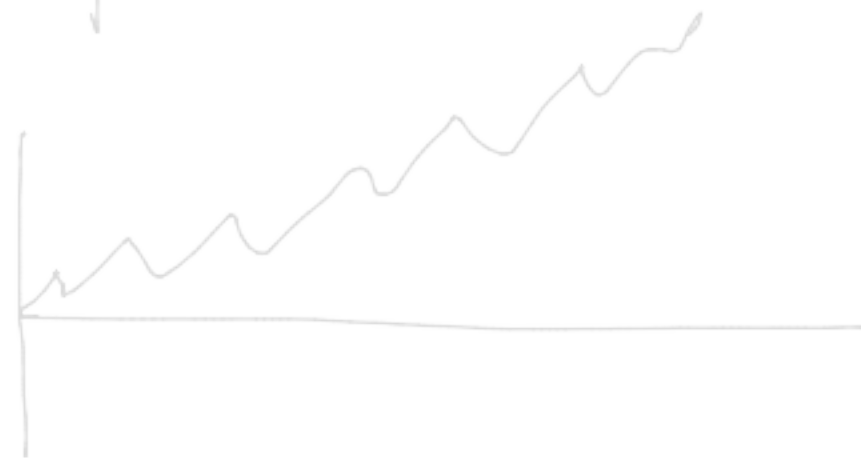
✓ Constant Variance

✓ Normally distributed ✓

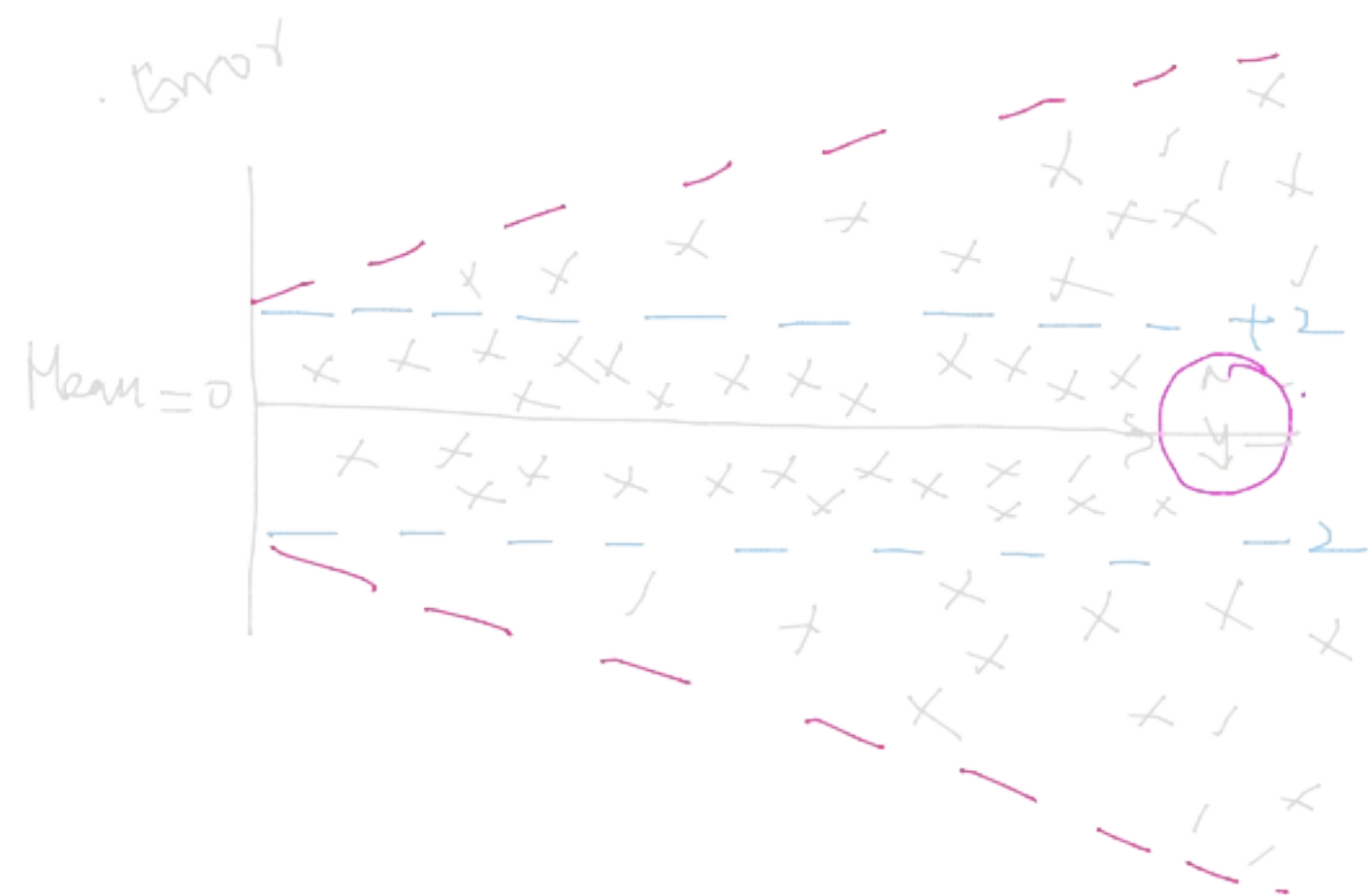
No auto correlation ✓

$y$	$\hat{y}$	$\rightarrow$ Error
50	45.5	+ve
60	43.5	+ve
70	58.1	+ve

Mean = +ve

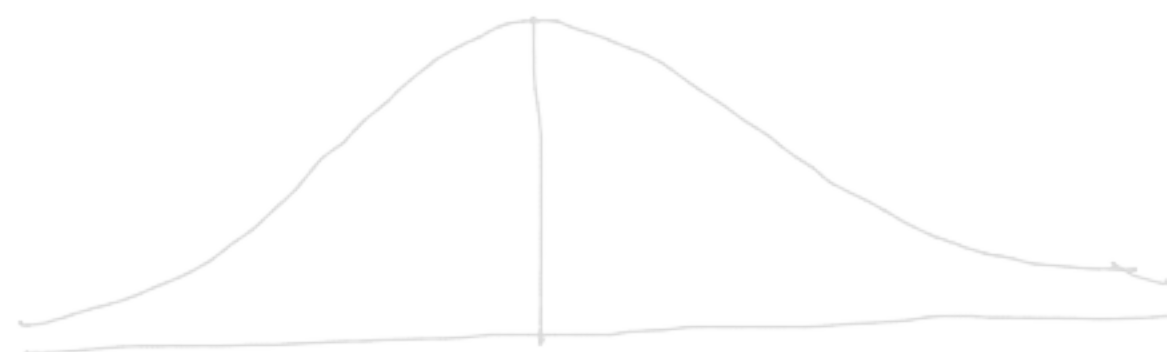


No constant Variance



Funnel  $\rightarrow$  Heteroscedasticity  
 $\downarrow$  Error dist  
 Varying

Random distribution



$3\sigma \rightarrow 2\sigma$

$\pm 2\sigma \rightarrow y$

Homoscedasticity  
 $\downarrow$  Error distribution  
 uniform

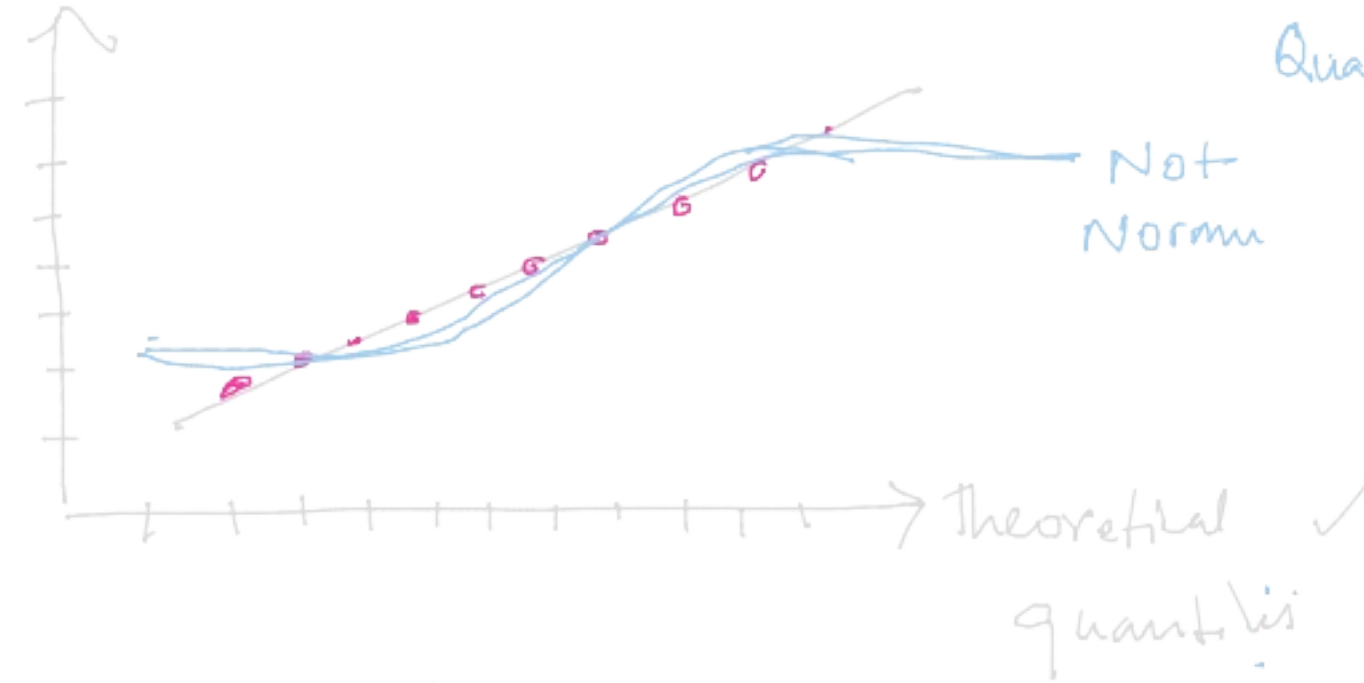
$3\sigma \rightarrow 3\sigma$

# Quantile - Quantile Plot

Q-Q plot

— check if your dataset follows a particular distribution

✓ Error quantiles ✓



Quantile

Not Norm

Theoretical quantiles ✓

↓

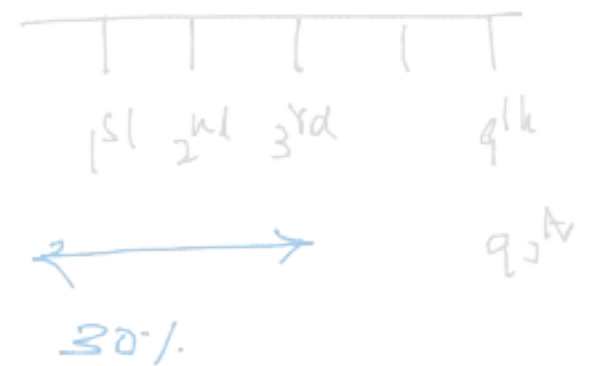
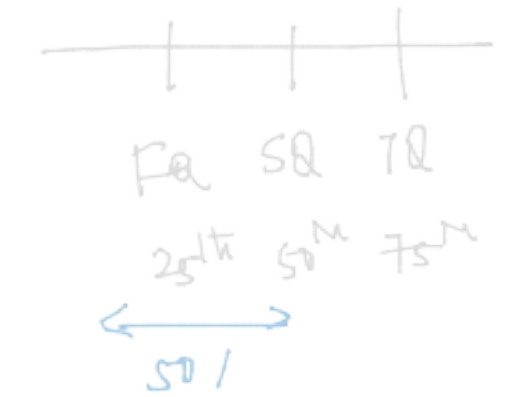
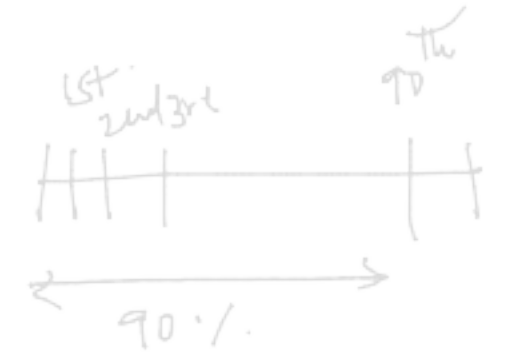
Create a Normal distribution

⇒ Quantile Quantile

Percentiles

Quantile

Deviates



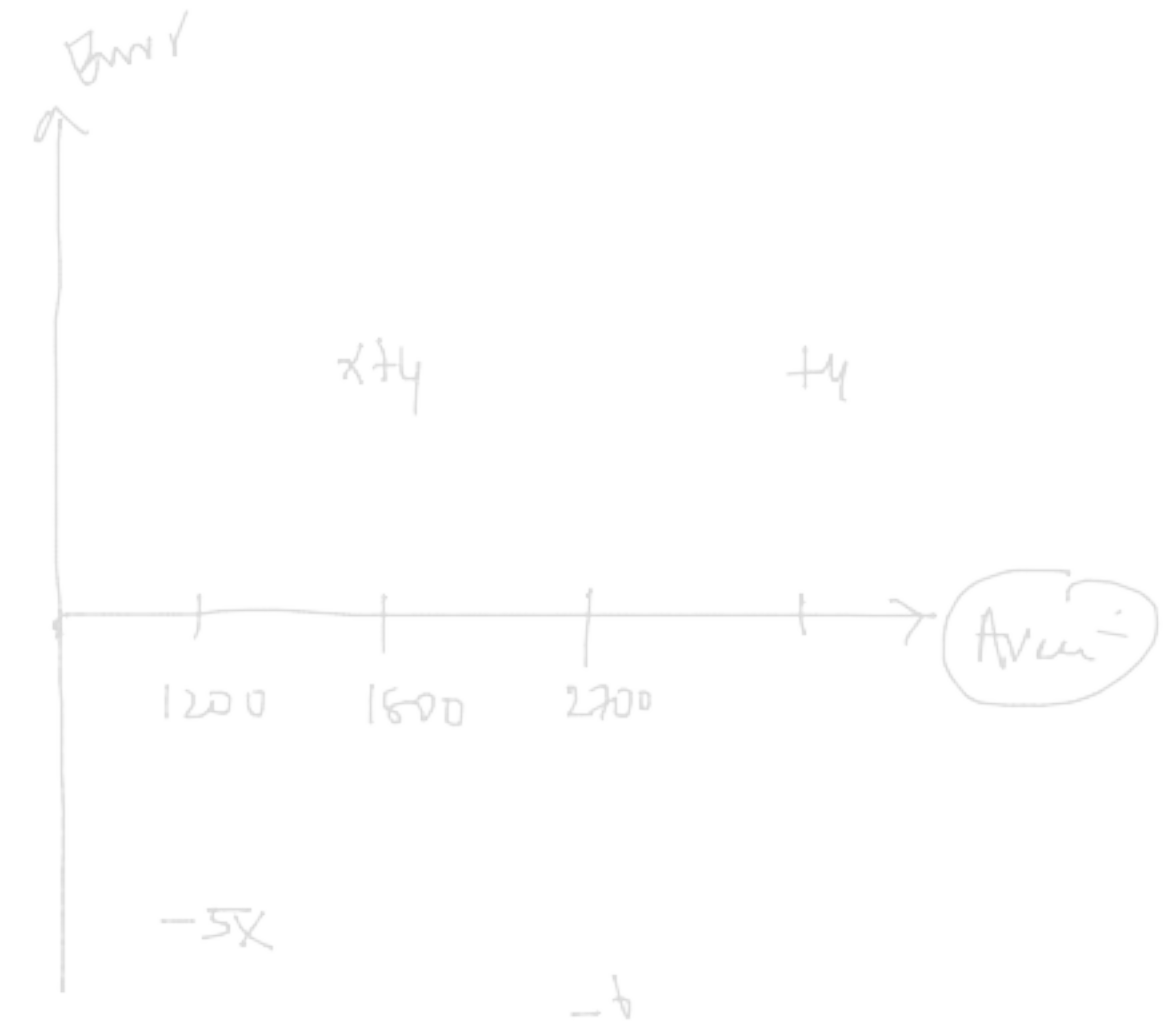
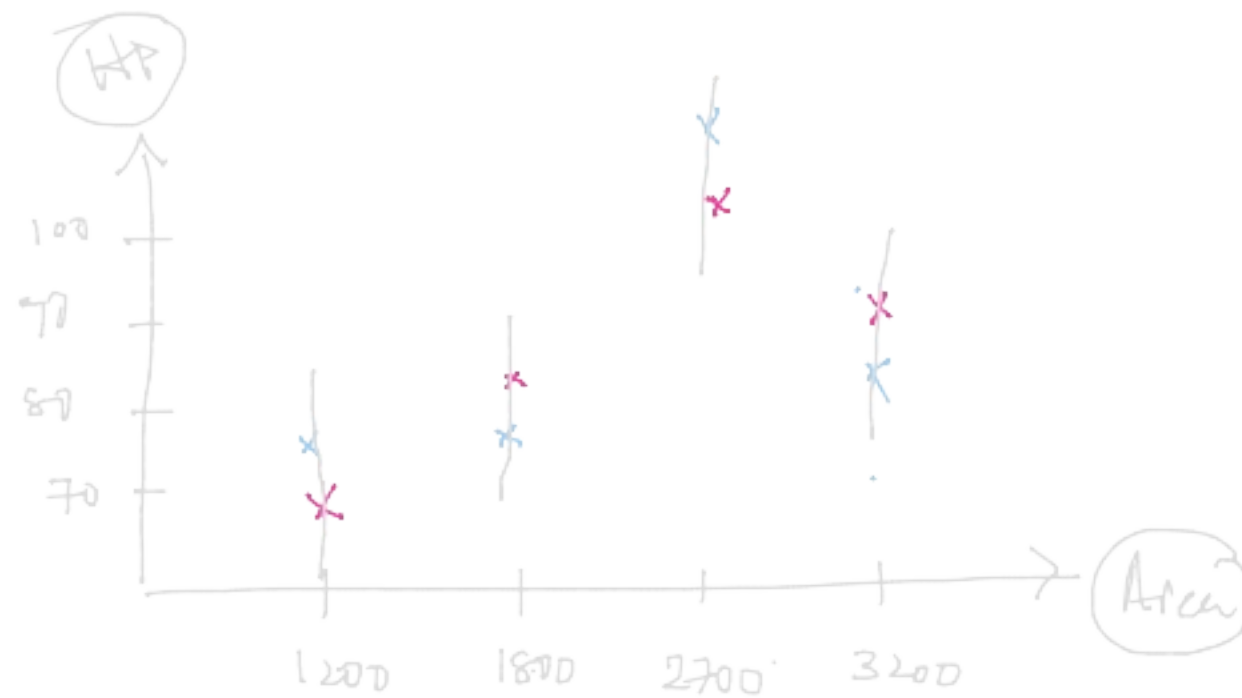
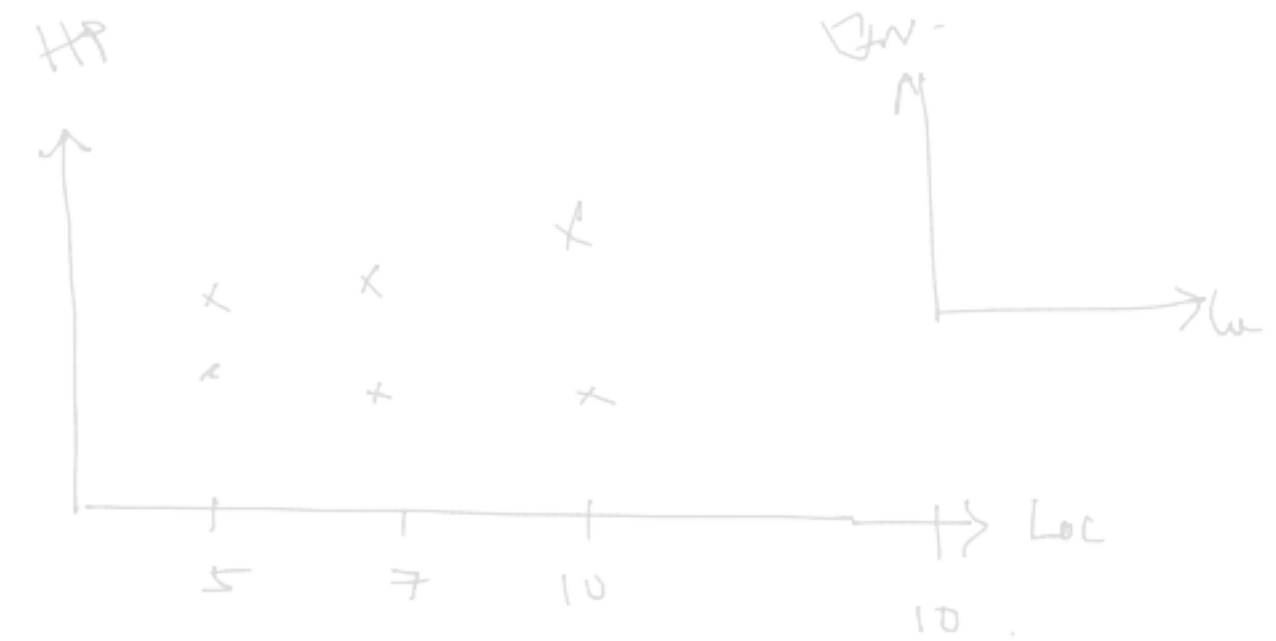
Error	Normal dist
1st	1st
2nd	2nd
3rd	3rd
4th	4th
5th	5th
6th	6th
7th	7th
8th	8th
9th	9th
10th	10th

Percentile Errors	Percentile of ND
1st	1st
2nd	2nd
3rd	3rd
4th	4th
5th	5th
6th	6th
7th	7th
8th	8th
9th	9th
10th	10th





	Area	loc	Bedrooms	HP	HP	Error
→ 1.	1200	5	2	70	75	-5
2	1800	22	2.5	85	81	+4
3	3200	7	3	93	89	+4
4	2700	10	3	102	108	-6



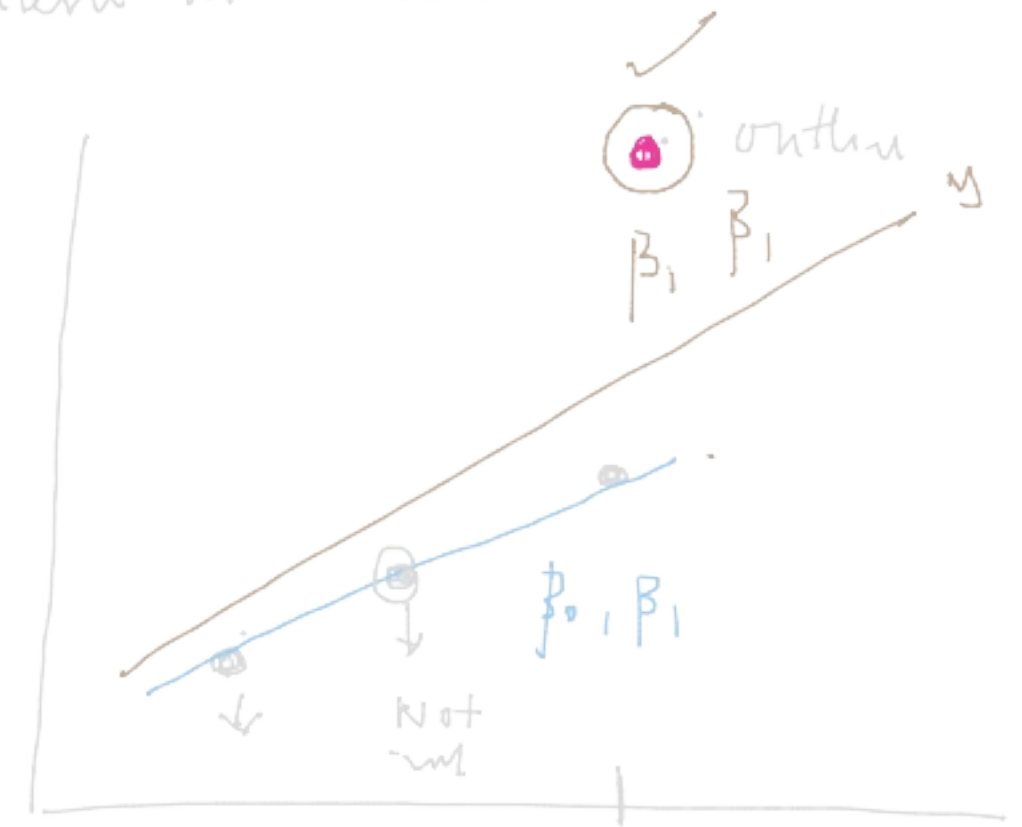
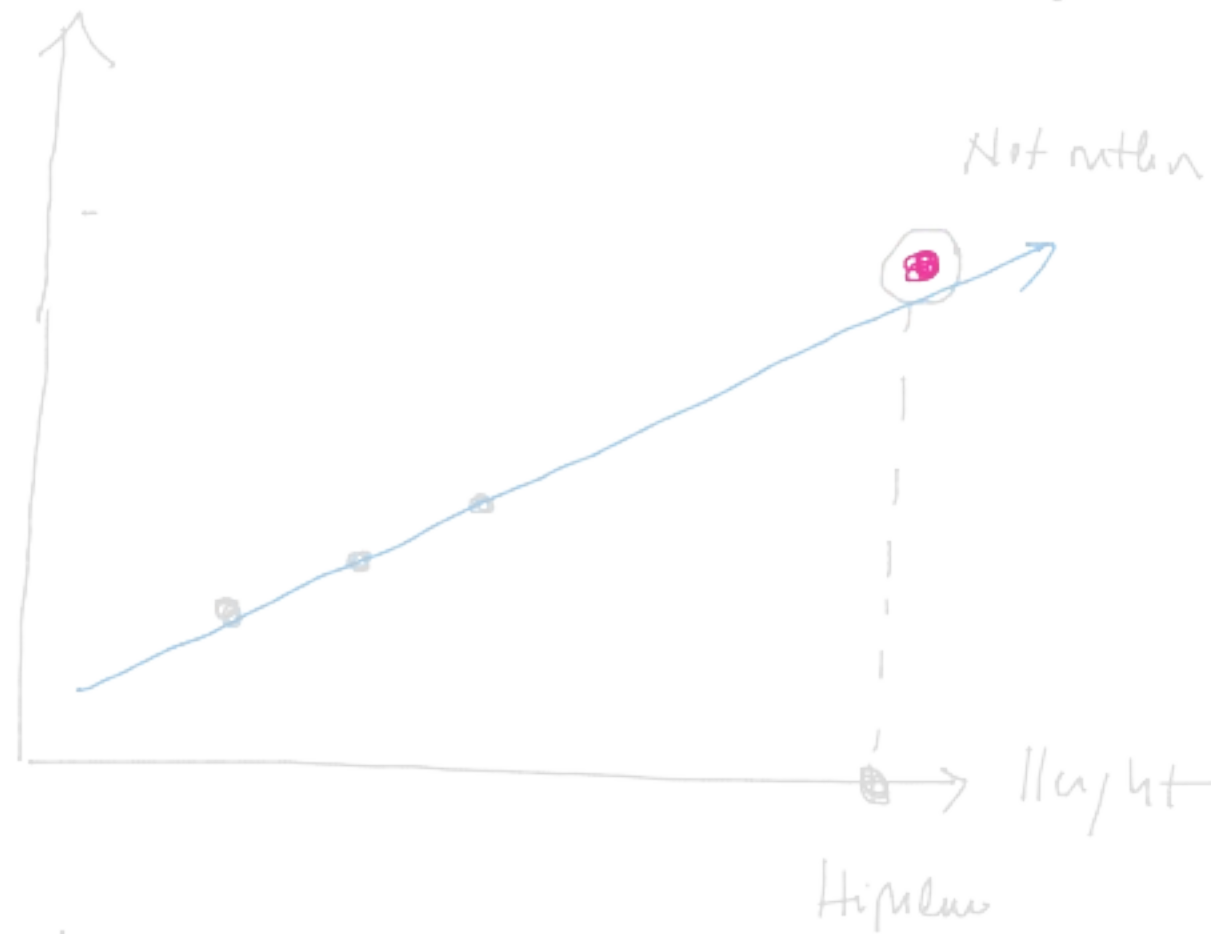
## Outliers

High

Leverage Point

Extreme 'x'-Values

- Extreme Values for 'y'
- It does not follow the pattern in the data



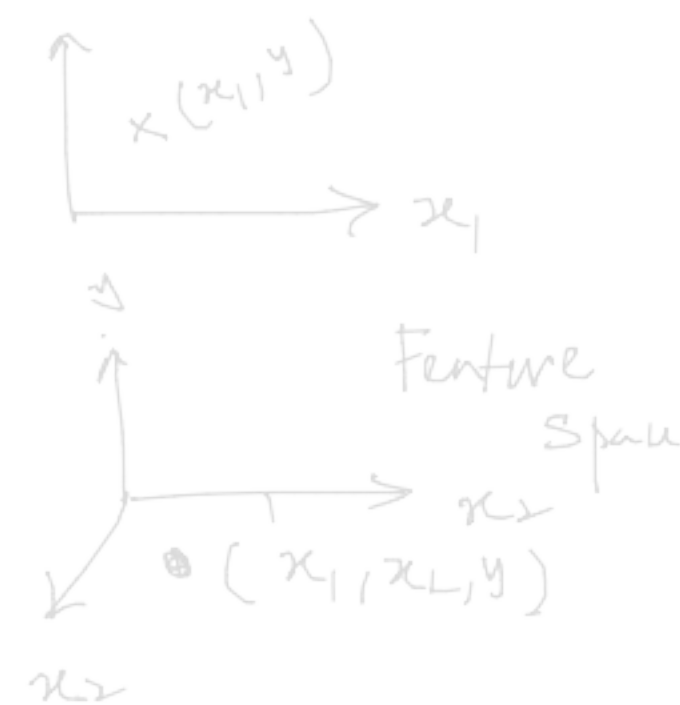
# Highly Influential Points

— Changing the model itself ✓



Include PI find  $\beta_0, \beta_1$   
Exclude PI find  $\beta_0, \beta_1$

Influence of a Point



'Cook's distance'  
 $p_1 \rightarrow > 1.0$  Highly influential

PI  $\rightarrow$   $(\underline{\hspace{2cm}})$   
 $\rightarrow$   $(\underline{\hspace{2cm}})$

Build model with all records  $\rightarrow (\beta_0, \beta_1) \rightarrow$  reference  
Remove  $Pt1$  Build model  $\rightarrow (\beta_0, \beta_1)$   
Remove  $Pt2$  Build the model  $\rightarrow \beta_0, \beta_1$



