Week 9

Linear Regression $\Big\langle$

$\bigcirc{Y}$ : continuous response    Final

$\bigcirc{X}$ : continuous or <u>categorical</u> predictors

# of hrs

Goal $\Big\{$

<span style="color:blue">Inference</span>

(i) identify linear relationship btw $Y \cdot X$

<span style="color:blue">no</span> if $\boxed{\text{significant}}$ → <u>quantify</u>

<span style="color:blue">$H_0: X \cdot Y$ linear rela</span>
<span style="color:blue">$H_a: X \cdot Y$ linea relation</span>

(ii) <u>prediction</u>    $\bigcirc{X} \to \underline{Y} \qquad \bigcirc{\hat{Y}}$

$\bigcirc{Y} - \bigcirc{\hat{Y}}$

simple linear regression

$\boxed{\bigcirc{E(Y)} = \boxed{\beta_0 + \beta_1 x}}$



$\uparrow \quad \bigcirc{Y} \quad \bigcirc{\beta_0 + \beta_1 x} + \bigcirc{\varepsilon}$

<span style="color:blue">$X_2$</span>

$\varepsilon \sim N(0, \sigma^2)$

mean
$[$ Intercept

<span style="color:blue">Error</span>

$\searrow \quad \hat{\beta_0} \quad \hat{\beta_1}$

$\Rightarrow E(Y) =$

$\Big[$

(1) linear relationship btw $X \cdot \bigcirc{Y}$

(2) conditional $\boxed{\text{Normality}}$ of $\underline{Y}$ given $\underline{X = x}$

(3) Equal $\boxed{\text{variance}}$

(4) indep samples

LSE (OLS)

$y - \hat{y}$

$$\hat{y} = \hat{\beta_0} + \hat{\beta_1} x$$

$\beta_1 \; \beta_2 \cdots \; \beta_p$

$p$

simple
$$\boxed{Y} = \beta_0 + \beta_1 x_1 + \varepsilon$$

multiple
$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \varepsilon$$

(1) Model significance

$H_0$: model is not useful
$H_a$: model is useful

$H_0: \beta_1 = 0$
$H_a: \beta_1 \neq 0$

$H_0: \beta_1 = \beta_2 = \cdots = \beta_p = 0$
$H_a$: at least one $\beta \neq 0$

(2) Individual term significance
$H_0: \beta_{alcohol} = 0$
$H_a: \beta_{alcohol} \neq 0$

$H_0: \beta_1 = 0$ vs $H_a: \beta_1 \neq 0$
$\vdots$
$H_0: \beta_p = 0$ vs $H_a: \beta_p \neq 0$

(3) Estimated regression line
$$* \quad \hat{Y} = -5.99 + 1.919 \cdot \text{alcohol}$$
$\hat{\beta_0} \; \hat{\beta_1}$

prediction power

$\hat{y}$  $y$  (4) $R^2 = \dfrac{SS_{model}}{SS_{Total}}$



$R_1^2 > R_2^2$

(5) Model diagnostics

95%

$-2$   $2$
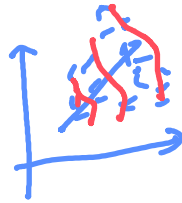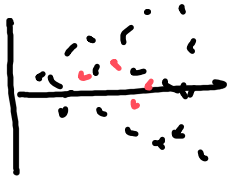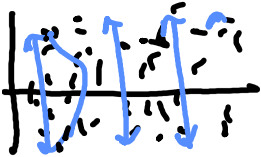
Diagnostics

$Y - \hat{Y}$ = residual

① Model assumption
- Normality —— Normal Q-Q.
  $\sqrt{\text{std re}}$ below $\sqrt{2}$
- Homoscedasticity   residual
- linear relationship   residual

② Influential points check



$Y_1$   $Y_2$   $Y_3$   $Y_4$

Simulation studies
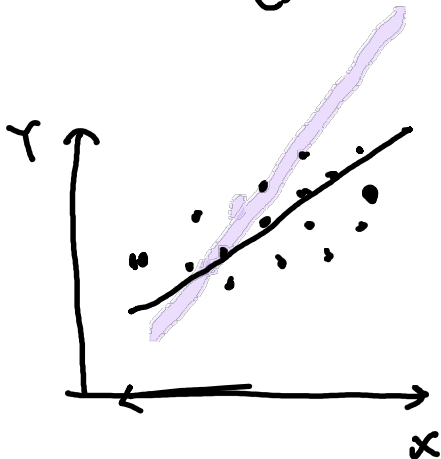
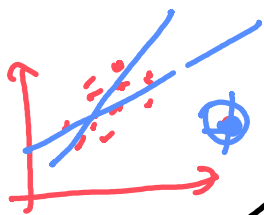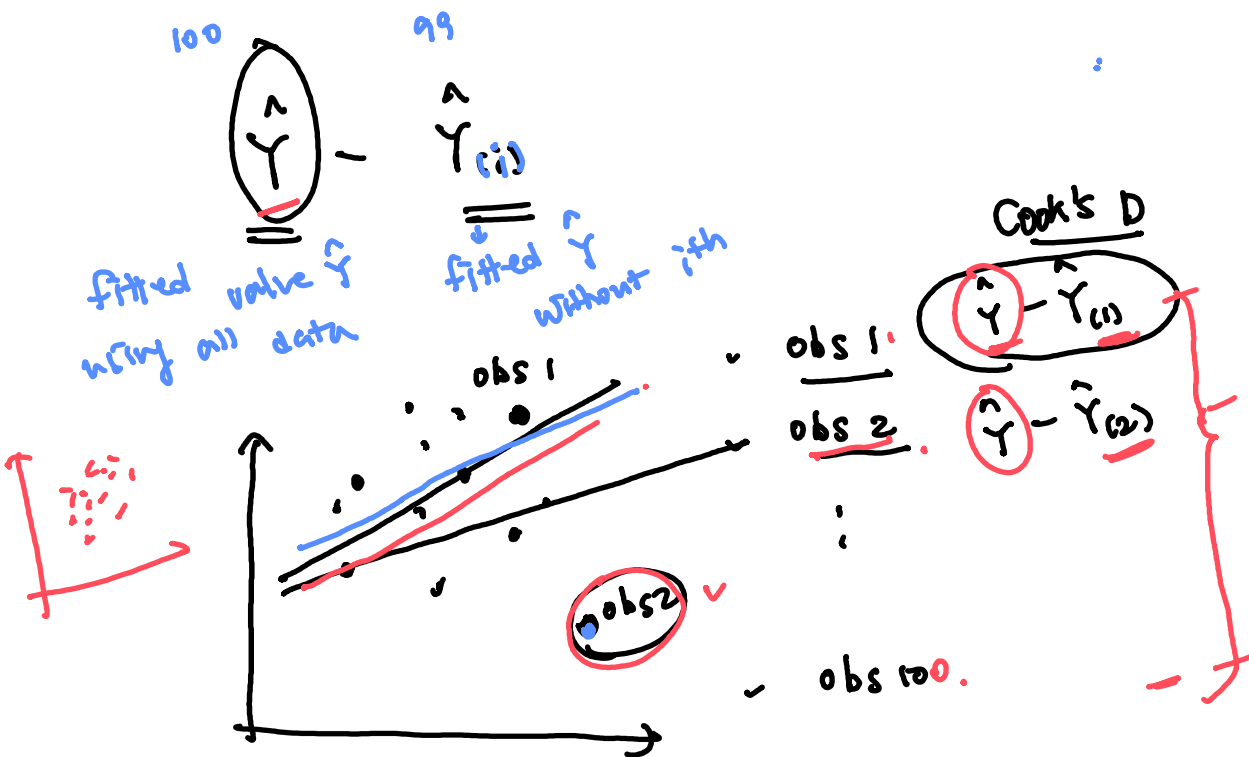| | linearity X-Y | Normality | Homoscedasticity | |
|---|---|---|---|---|
| 1 | O | O | O | lm $(Y_1 \sim x)$ |
| 2 | O | X | O | lm $(Y_2 \sim x)$ |
| 3 | O | O | X | lm $(Y_3 \sim x)$ |
| 4 | X | O | O | lm $(Y_4 \sim x)$ |

## Influential points

⇒ observations that greatly affect the slope of regression line.

$<$ outliers : abnormal beha in $y$
   leverage points : " in $x$

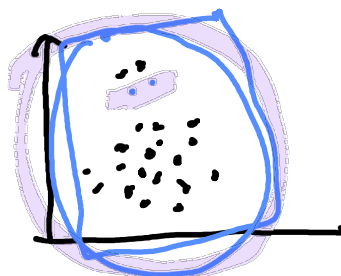# Cook's distance: Influence of individual data points on model fitting

$\widehat{Y}$ ~ $\widehat{Y}_{(i)}$

100    99

fitted value $\widehat{Y}$ using all data

fitted $\widehat{Y}$ without $i$th

## Cook's D

$\widehat{Y} - \widehat{Y}_{(1)}$

$\widehat{Y} - \widehat{Y}_{(2)}$

obs 1

obs 2

obs 100.

obs 1

obs 2

obs 100.

large Cook's D    :    influential points

small Cook's D    :

$\frac{4}{n}$

## NOTE

(Q) several influential obs removed...    all done?

— drinking data example  $\hat{y} = -5.99 + 1.97 * alcohol$  (using all data)

$\downarrow$

without France

$\hat{y} = -3.61 + 1.64 * alcoh$

— crime.csv data

$\hat{y} = -86.20 + \boxed{49.03} * poverty$  ~5%

w/o 51st dc  $\hat{y} = 209 + \boxed{25.45} * po$  $\boxed{13.6\%}$

$\uparrow$ SS_Tota

$\dfrac{SS_{Total}}{S}$

$40 + 50 = 90$

$\boxed{50\%}$  $\underline{\underline{50\%}}$

$35 \rightarrow \boxed{15^{-1}} + 6 \rightarrow 61$

$25 \rightarrow \boxed{25'} \times 40\% \rightarrow +10$

47