

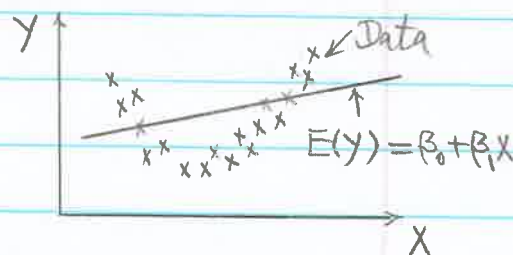
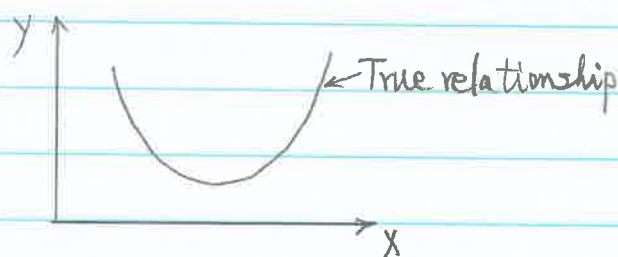
①

Chapter 2 Assessing the Straight Line Fit

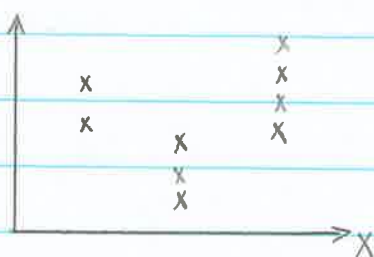
1. Testing for lack of fit

In Chapter 1, we introduced SLR model: $Y = \beta_0 + \beta_1 X + \varepsilon \Leftrightarrow E(Y) = \beta_0 + \beta_1 X$.

Question: Is the model adequate?



Replication:



At X_i , have $Y_{i1}, Y_{i2}, \dots, Y_{in_i}, i=1, \dots, m$, i.e.,

X_1	X_2	\dots	X_m
Y_{11}	Y_{21}	\dots	Y_{m1}
Y_{12}	Y_{22}	\dots	Y_{m2}
\vdots	\vdots		\vdots
Y_{1n_1}	Y_{2n_2}	\dots	Y_{mn_m}

total observations: $n = \sum_{i=1}^m n_i$, $m = \# \text{distinct } X\text{'s}$, $\text{Var}(Y_{ij}) = \sigma^2$.

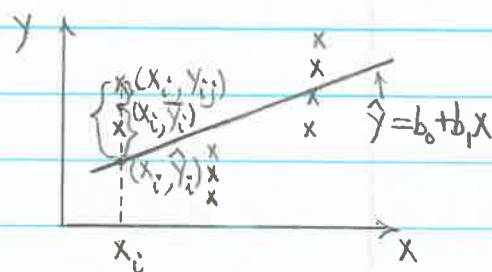
(2)

Deviation about the fitted line can be decomposed into:

(1) deviation of x_{ij} from \bar{y}_i , plus

(2) deviation of \bar{y}_i from $\hat{y}_{ij} = \hat{y}_i$, i.e.,

$$\begin{aligned} e_{ij} &= x_{ij} - \hat{y}_{ij} = x_{ij} - \hat{y}_i \\ &= (x_{ij} - \bar{y}_i) + (\bar{y}_i - \hat{y}_i). \end{aligned}$$



— Square both sides

$$(x_{ij} - \hat{y}_i)^2 = (x_{ij} - \bar{y}_i)^2 + (\bar{y}_i - \hat{y}_i)^2 + 2(x_{ij} - \bar{y}_i)(\bar{y}_i - \hat{y}_i).$$

— Sum over i and j

$$\begin{aligned} \sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij} - \hat{y}_i)^2 &= \sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij} - \bar{y}_i)^2 + \sum_{i=1}^m \sum_{j=1}^{n_i} (\bar{y}_i - \hat{y}_i)^2 + 2 \sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij} - \bar{y}_i)(\bar{y}_i - \hat{y}_i) \\ &= \sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij} - \bar{y}_i)^2 + \sum_{i=1}^m n_i (\bar{y}_i - \hat{y}_i)^2. \end{aligned}$$

Note: $\sum_{i=1}^m \sum_{j=1}^{n_i} (x_{ij} - \bar{y}_i)(\bar{y}_i - \hat{y}_i) = 0.$

Residual sum of squares = Pure error sum of squares
+ Lack of fit sum of squares.

$$RSS = SS(pe) + SS(lof).$$

df of $\sum_{j=1}^{n_i} (x_{ij} - \bar{y}_i)^2$ is $n_i - 1$, $i = 1, 2, \dots, m$.

(3)

$$\Rightarrow (df \text{ of } SS(pe)) = \sum_{i=1}^m (n_i - 1) = \sum_{i=1}^m n_i - m = n - m.$$

$$\text{Then, } (df \text{ of } SS(lof)) = (df \text{ of } RSS) - (df \text{ of } SS(pe))$$

$$= (n - 2) - (n - m) = m - 2.$$

Notes: (1) $SS(pe) = \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$ does not depend on the model.

(2) $SS(lof) = \sum_{i=1}^m n_i (\bar{y}_i - \hat{y}_i)^2$ depends on the model fitted and measures model lack of fit.

Modify ANOVA table to "break out" pure error and lack of fit SS.

<u>Source of variation</u>	<u>df</u>	<u>SS</u>	<u>MS</u>
Regression	1	SS_{reg}	$SS_{reg}/1 = SS_{reg}$
<u>Residual</u>	<u>$n - 2$</u>	<u>RSS</u>	<u>$RSS/(n - 2)$</u>
{ Lack of fit	$m - 2$	$SS(lof)$	$SS(lof)/(m - 2)$
{ Pure error	$n - m$	$SS(pe)$	$SS(pe)/(n - m)$
Total (corrected)	$n - 1$	TSS	

Testing $H_0: Y = \beta_0 + \beta_1 X + E$ (A SLR model is adequate) vs.

$H_a: Y \neq \beta_0 + \beta_1 X + E$ (A SLR model is inadequate).

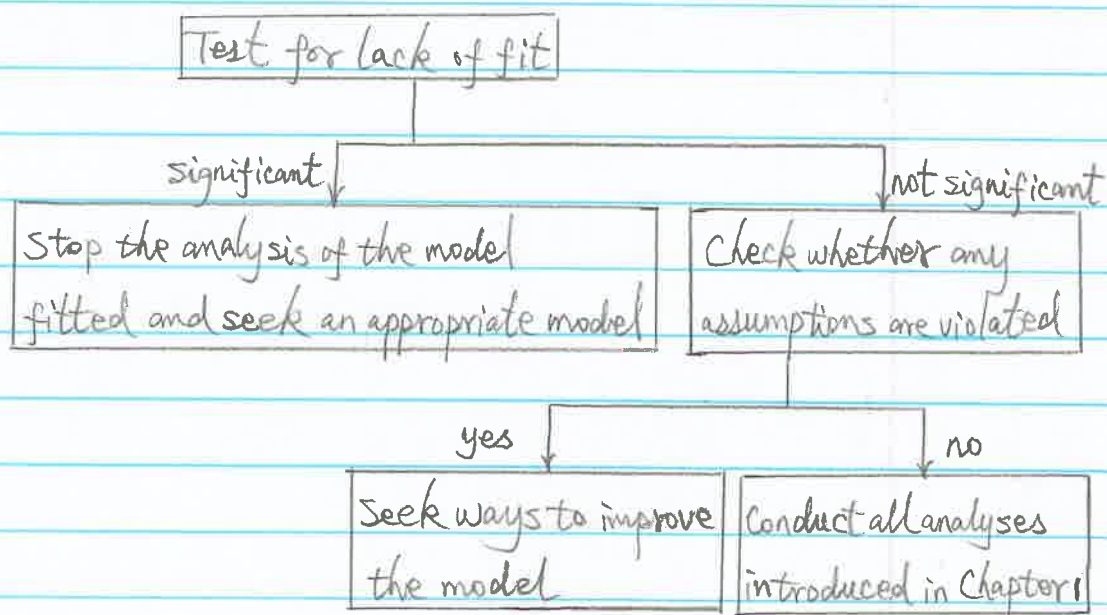
Test statistic: $F = \frac{MS(lof)}{MS(pe)}$ Under H_0 $F_{m-2, n-m}$.

(4)

Reject H_0 if $p\text{-value} = P(F_{m-2, n-m} \geq F_{obs}) \leq \alpha$ or $F_{obs} \geq F_{m-2, n-m}(1-\alpha)$, where $F_{m-2, n-m}(1-\alpha)$ is the $(1-\alpha)^{th}$ quantile of $F_{m-2, n-m}$.

Notes, (1) For the lack of fit test, we need at least 3 distinct X values and replication at at least one X value.

(2) Analysis steps for data with at least 3 distinct X values and replication at at least one X value:



2. Check Model Assumptions

Model: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, i = 1, 2, \dots, n.$

Model assumptions:

- (1) $E(\epsilon_i) = 0$,
- (2) $\text{Var}(\epsilon_i) = \sigma^2$ — constant variance,
- (3) $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are independent,
- (4) ϵ_i is normally distributed.

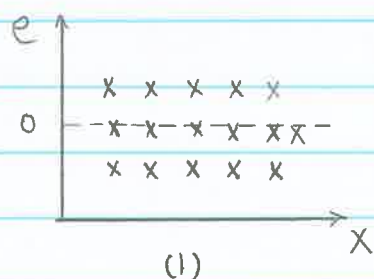
(5)

$$\Leftrightarrow \varepsilon_1, \varepsilon_2, \dots, \varepsilon_n \stackrel{iid}{\sim} N(0, \sigma^2).$$

Since ε_i is estimated by $e_i = y_i - \hat{y}_i$, we check the assumptions e_1, e_2, \dots, e_n — residuals.

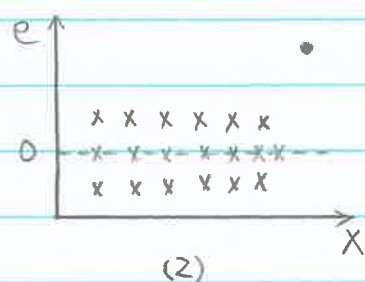
Residual plots: Plot e_i vs. X_i (on X axis)
or e_i vs. $\hat{y}_i (= b_0 + b_1 X_i)$.

Note: Since $\hat{y}_i = b_0 + b_1 X_i$ is simply a linear function of X_i , the only real difference between the two types of residual plots is the scale on the horizontal axis. The pattern of points in the residual plots will be the same, and it is the pattern of points that is important, not the scale.



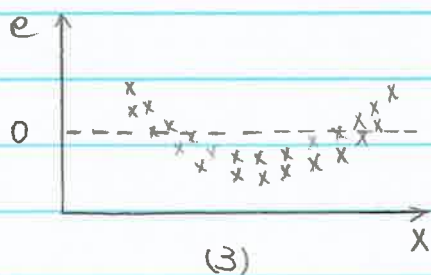
If model is correct, no pattern.

— residuals are roughly uniformly distributed about the horizontal line $e=0$.

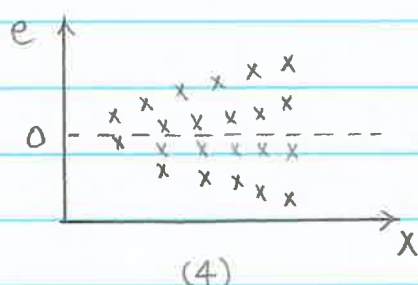


— a possible outlier.

⑥



— Curvature indicates a curved rather than linear relationship between X and Y . (Note: Sometimes it is too subtle to pick up in the scatter plot of y_i vs. x_i .)



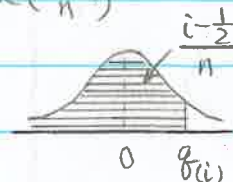
— "Funnel" shape indicates non-constant variance (here variance increases with X).

Note: Read question and answer 1, 2, 3 on pages 63-67.

Assessing normality

Normal plot (Quantile-Quantile plot, i.e., Q-Q plot)

- (1) Order residuals e_1, e_2, \dots, e_n from the smallest to the largest to get the ordered residuals: $e_{(1)}, e_{(2)}, \dots, e_{(n)}$.
- (2) Calculate the normal scores: $g_{(i)} = \Phi^{-1}\left(\frac{i-1/2}{n}\right)$, $i = 1, 2, \dots, n$, where $\Phi^{-1}(\cdot)$ is the inverse function of the standard normal distribution function $\Phi(x) = P(Z \leq x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz$ (i.e., the $\left(\frac{i-1/2}{n}\right)$ th quantile of $N(0,1)$). Thus, $\Phi(g_{(i)}) = P(Z \leq g_{(i)}) = \frac{i-1/2}{n}$.
- (3) Plot the pairs $(g_{(1)}, e_{(1)}), (g_{(2)}, e_{(2)}), \dots, (g_{(n)}, e_{(n)})$.



A substantial linear pattern in a normal plot indicates that the

⑦

normality assumption remains tenable. Otherwise, normality is suspect.

Note: The judgement as to whether a normal plot does or does not show a substantial linear pattern is somewhat subjective.

Test for normality

The Shapiro and Wilk test

$H_0: e_1, e_2, \dots, e_n$ come from a normal population vs.
(the normality assumption is tenable.)

$H_a: e_1, e_2, \dots, e_n$ come from a non-normal population.
(the normality assumption is violated.)

Test statistic: $W = \frac{[\sum_{i=1}^n a_i e_{(i)}]^2}{\sum_{i=1}^n e_i^2}$.

Reject H_0 if $p\text{-value} \leq \alpha$.

Notes: (1) Residuals e_1, e_2, \dots, e_n are not independent. However, the effect of correlations between residuals is usually negligible.

(2) When intercept is in model, $\sum_{i=1}^n e_i = 0 \iff \bar{e} = \frac{\sum_{i=1}^n e_i}{n} = 0$,
Thus, there is need to check $\bar{e} = 0$.

Q: What if the SLR is inadequate or some assumptions are violated?

A: Transform data — Transformations, which are the focuses of chapter 12 & 13.

Example 2.1: For the data given in problem F on page 99 (i.e., the data for HW #1),

X	4.7	5	5.2	5.2	5.9	4.7	5.9	5.2	5.3	5.9	5.6	5
Y	3	3	4	5	10	2	9	3	7	6	6	4

- (a) Test to determine whether it is adequate to fit the data by a straight line at $\alpha = 0.05$.
 (b) Conduct Shapiro and Wilk test for normality based on residuals using $\alpha = 0.05$.

Example 2.2: In a study, 15 American females aged 30-39 are randomly selected. The following data on X = height (in) and Y = weight (lb) are observed.

X	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72
Y	113	115	118	121	124	128	131	134	137	141	145	150	153	159	164

Is the relationship between X and Y linear?