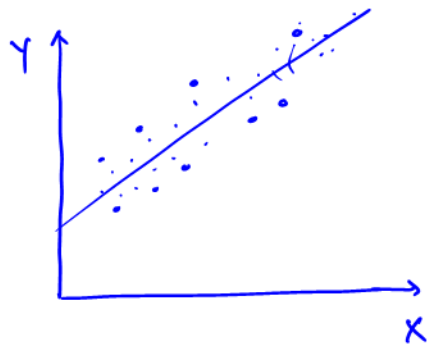


Residual plots (plots of residuals or Semi Studentized residuals) can be used to identify above departures from the model.

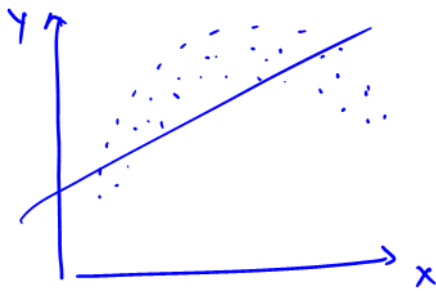
1. Non-linearity of regression function

There are two types of plots

a) Scatter plot (Y vs X)



→ linear \Rightarrow good fit

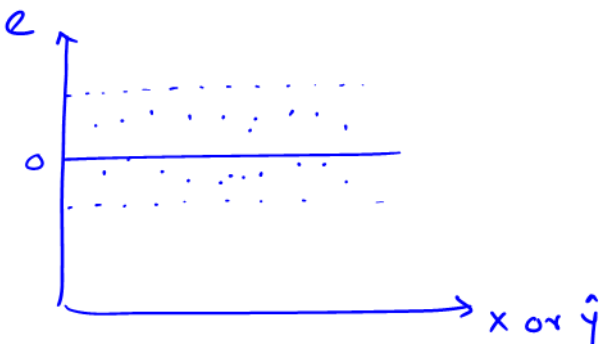


→ not linear \Rightarrow not a good fit.

b) Residual plots (residuals vs X or residuals vs \hat{Y})

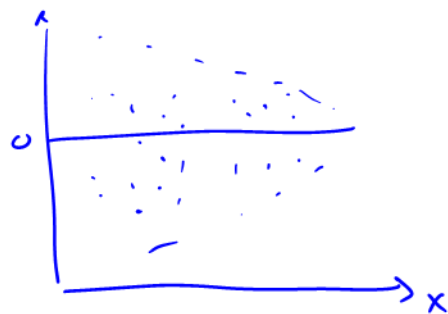
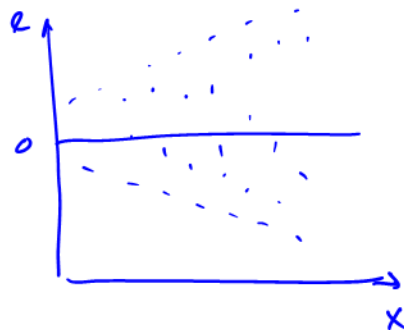
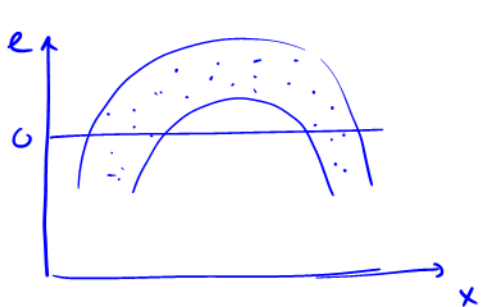
For SLR model plots for residual vs X and residual vs \hat{Y}

Show the same pattern.



All the residuals fall within a horizontal band around 0.
 \Rightarrow good fit.

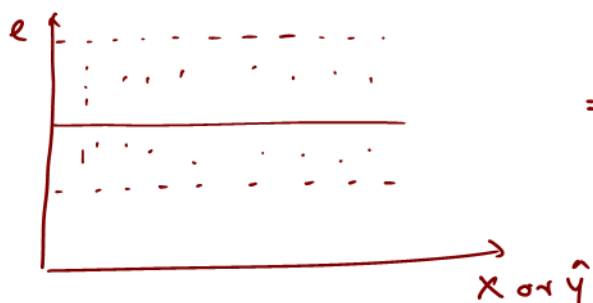
Departures:



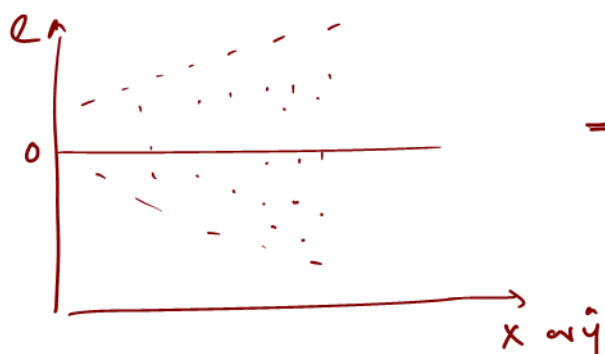
\Rightarrow not a good fit.

② Non-constancy of error variance

a) Residual plots (vs x or \hat{y})



\Rightarrow constant variance.

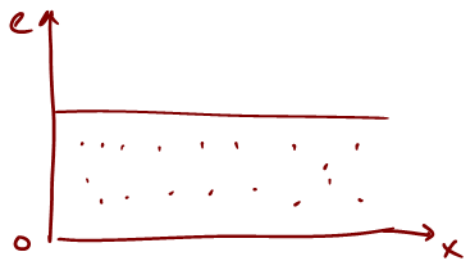


\Rightarrow variation increases with x
(not a good fit)

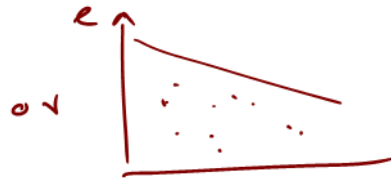
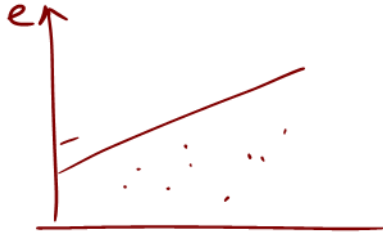


\Rightarrow variation decreases with x
(not a good fit).

b) Absolute residuals vs X (or \hat{y})



\Rightarrow good fit



\Rightarrow not a good fit.

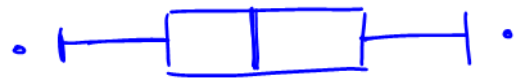
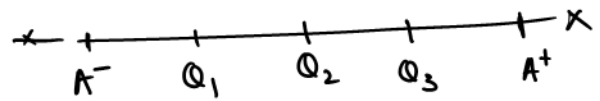
This is useful when the Sample Size is small.

③ Presence of outliers

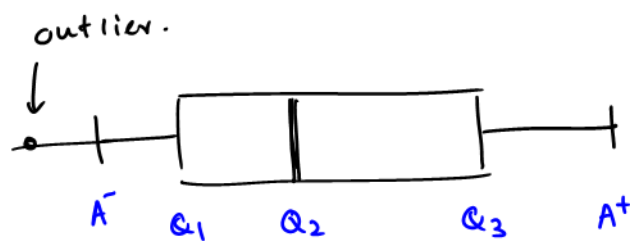
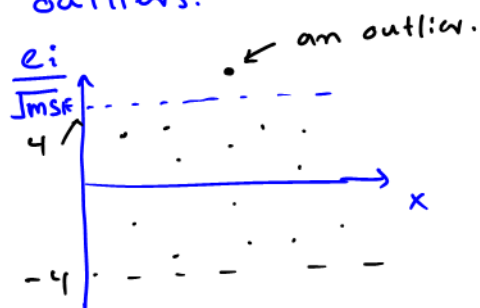
$$A^-, Q_1, Q_2, Q_3, A^+$$

$$A^- = Q_1 - 1.5(Q_3 - Q_1)$$

$$A^+ = Q_3 + 1.5(Q_3 - Q_1)$$

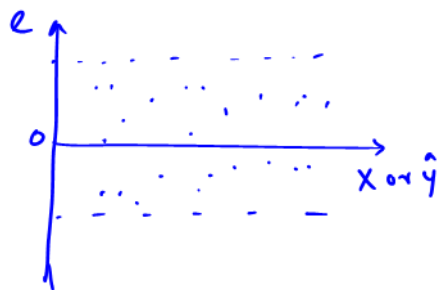


If $\left| \frac{e_i}{\overline{\text{Jmse}}} \right| \geq 4$, then the i^{th} data point is an outlier. Residual plots or box plot of residuals can be used to identify outliers.



④ Non-independence of error terms

Residual plots (e vs x or e vs \hat{y}).



\Rightarrow No pattern \Rightarrow independent

\Rightarrow Some pattern \Rightarrow not independent.

⑤ Non-normality of error terms

Assumption $\epsilon_i \sim N(0, \sigma^2)$

Box-plot or histogram can be used to check the shape of the distribution.

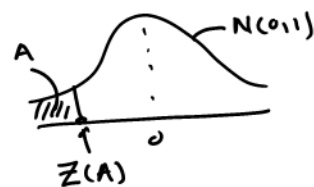
Normal probability plot can be used to check the normality.

Normal probability plot

plot of residuals vs their expected values under normality.

Calculating expected values:

$$\text{Expected value for } e_i = \sqrt{\text{MSE}} \left(Z \left(\frac{k - 0.375}{n + 0.25} \right) \right),$$



where $Z(A)$ is the A -th percentile for the standard normal distribution and k is the rank of the residual e_i . If there is a tie, take the average of the ranks.

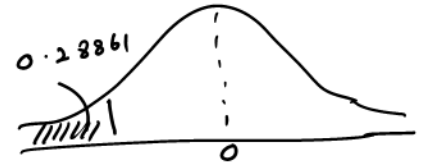
Eg: $-2.43, -2.31, 0.01, 0.01, 0.5, 0.67$

	1	2	3	4	5	6
rank	1	2	3.5	3.5	5	6

Eg: plastic hardness example:

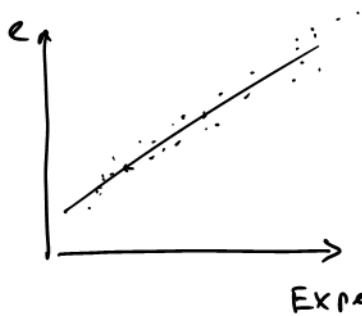
$$\text{Residual } (e_1) = -2.150 \quad \text{rank } (e_1) = K = 5, n = 16$$

$$\text{So } \frac{K - 0.375}{n + 0.25} = \frac{5 - 0.375}{16 + 0.25} = 0.28861$$



$$\text{MSE} = (3.23403)^2, \quad Z(0.28861) = -0.55748$$

$$\begin{aligned} \text{Expected value} &= (3.23403) \times (-0.55748) \\ &= -1.792903. \end{aligned}$$



linear \Rightarrow normality.

non linear \Rightarrow non-normality.

Test for Lack of fit (Numerical tools)

* Correlation test for normality

Here calculate the coefficient of correlation residuals and their expected values under normality. Then compare that with the corresponding value in table B6 (in the text). Values in table B6 are the percentiles of the distribution of the correlation with normally distributed error terms.

If the observed coefficient $>$ table value

\Rightarrow error terms are normal.

Note:

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\sigma(X) \cdot \sigma(Y)}} = \frac{E(XY) - E(X) \cdot E(Y)}{\sqrt{\sigma(X) \cdot \sigma(Y)}} \text{ — parameter.}$$

Sample term:

$$r = \frac{S_{xy}^2}{\sqrt{S_x^2 \cdot S_y^2}} \text{ — Sample correlation,}$$

$$\text{where, } S_{xy}^2 = \sum(XY) - \frac{\sum(X) \cdot \sum(Y)}{n},$$

$$S_x^2 = \sum X^2 - \frac{(\sum X)^2}{n}$$

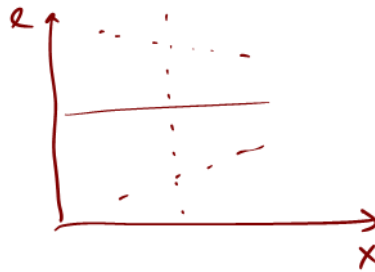
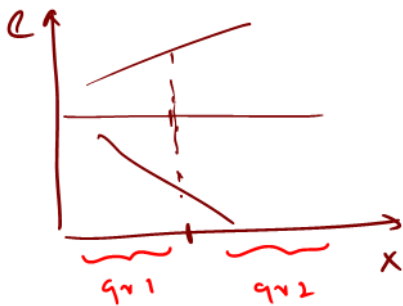
$$S_y^2 = \sum Y^2 - \frac{(\sum Y)^2}{n}$$

* Test for Constant Error Variance

Brown - Forsythe test

Idea:

When error variance is not constant.



Steps:

1) divide residuals into two groups based on the level of X.

group 1: n_1 values (lower X values)

group 2: n_2 values (higher X values)

2) calculate

$$d_{i1} = |e_{i1} - \tilde{e}_1|$$

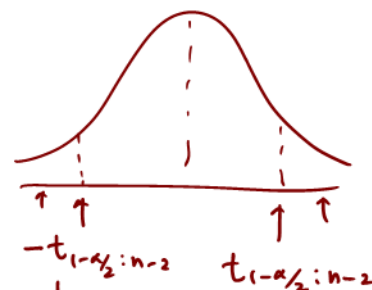
$$d_{i2} = |e_{i2} - \tilde{e}_2|,$$

where e_{ik} — i^{th} residual for the k^{th} group ($k=1,2$),
 \tilde{e}_k — median of the residuals for the k^{th} group ($k=1,2$).

3) calculate

$$T = \frac{\bar{d}_1 - \bar{d}_2}{S \sqrt{1/n_1 + 1/n_2}} \sim t_{n-2},$$

where
$$S^2 = \frac{\sum (d_{i1} - \bar{d}_1)^2 + \sum (d_{i2} - \bar{d}_2)^2}{n-2}$$



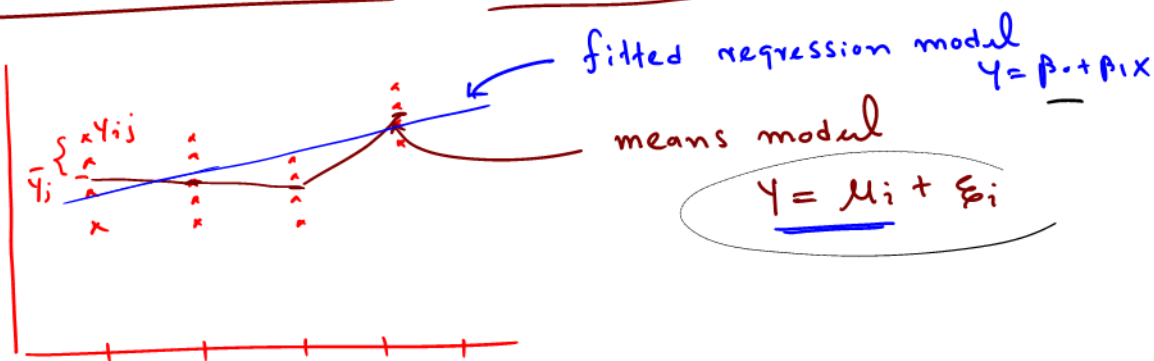
4) Suppose t^* be the observed value of T .

If $|t^*| \leq t_{1-\alpha/2; n-2} \Rightarrow$ error variance is constant

$|t^*| > t_{1-\alpha/2; n-2} \Rightarrow$ error variance is not constant.

* Test for non-linearity (general linear test)

Idea:



Assumptions:

* $y_i \stackrel{\text{ind}}{\sim} N(E(y), \sigma^2)$

↑ constant

* There are repeat observations at one or more X -levels.

Note:

Repeat observations for the same X -level are called replicates.

Notation:

Data should be arranged by level of X and replicate number.

level of X	(X_1) $j=1$	(X_2) $j=2$	$j=3$...	$j=c$
replicate					
$i=1$	y_{11}	y_{12}	y_{13}		y_{1c}
$i=2$	y_{21}	y_{22}	y_{23}		y_{2c}
\vdots	\vdots	\vdots	\vdots		\vdots
\vdots	y_{n1}	y_{n2}	y_{n3}		y_{nc}
	\bar{y}_1	\bar{y}_2	\bar{y}_3	- ... -	\bar{y}_c

Hypotheses:

$H_0: \underline{E(Y)} = \beta_0 + \beta_1 X$ (i.e. Regression function is linear)

$H_1: E(Y) \neq \beta_0 + \beta_1 X$ (i.e. Regression function is not linear)

Test statistic:

* Full model:

If a regression model is not fitted (i.e. predictor variable is not used) the model is

$$y_{ij} = \mu_j + \varepsilon_{ij}, \text{ means model (ANOVA model)}$$

$$i = 1, 2, \dots, n_j$$

$$j = 1, 2, \dots, c.$$

$$\varepsilon_{ij} \stackrel{i.i.d.}{\sim} N(0, \sigma^2),$$

$$E(y_{ij}) = \mu_j$$

* least square estimate of $\mu_j = \hat{\mu}_j = \bar{y}_j \leftarrow$ sample mean of the j^{th} group.

$$* SSE(F) = \sum_j \sum_i \underbrace{\left(y_{ij} - \frac{1}{n_j} \sum_i y_{ij} \right)^2}_{\text{pure error}}$$

This is also called as pure error sum of squares (SSPE).

$$* \text{df of } SSE(F) = \sum_{i=1}^c (n_i - 1) = n - c.$$

* Reduced model: (SLR model)

model: $Y_{ij} = \beta_0 + \beta_1 X_j + \epsilon_{ij}$, $\epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$

$E[Y_{ij}] = \beta_0 + \beta_1 X_j$ — does depend on X .

* least square estimator $E[Y_{ij}] = b_0 + b_1 X_j$

* error Sum of Squares = $SSE(R) = \sum_i \sum_{j=1}^c (Y_{ij} - b_0 - b_1 X_j)^2$
 $= \sum_i \sum_{j=1}^c (Y_{ij} - \hat{Y}_{ij})^2 = SSE$

df or $SSE(R) = n - 2$.

Test statistic

$$F = \frac{SSE(R) - SSE(F)}{df_R - df_F} \div \frac{SSE(F)}{df_F}$$
$$= \frac{SSE - SSPE}{n - 2 - (n - c)} \div \frac{SSPE}{n - c}$$

But $SSE - SSPE = SSLF$ — lack of fit Sum of Squares.

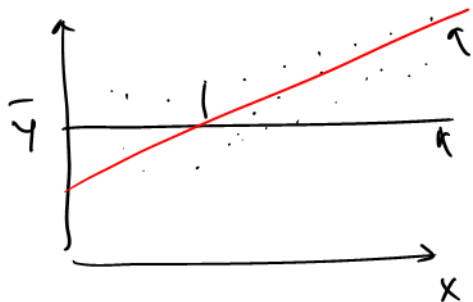
Now $F = \frac{SSLF}{c - 2} \div \frac{SSPE}{n - c} = \frac{MSLF}{MSPE}$

← lack of fit mean square
← pure error mean square

Decision: Let f^* be observed value of F .

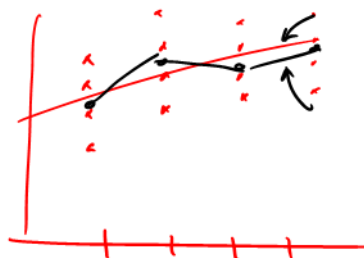
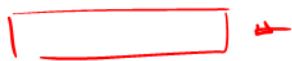
* $f^* \leq F(1 - \alpha, c - 2, n - c) \Rightarrow$ conclude H_0 (ie regression function is linear)

$f^* > F(1 - \alpha, c - 2, n - c) \Rightarrow$ conclude H_1 (ie regression function is not linear).



$$y_i = \mu + \varepsilon_i$$

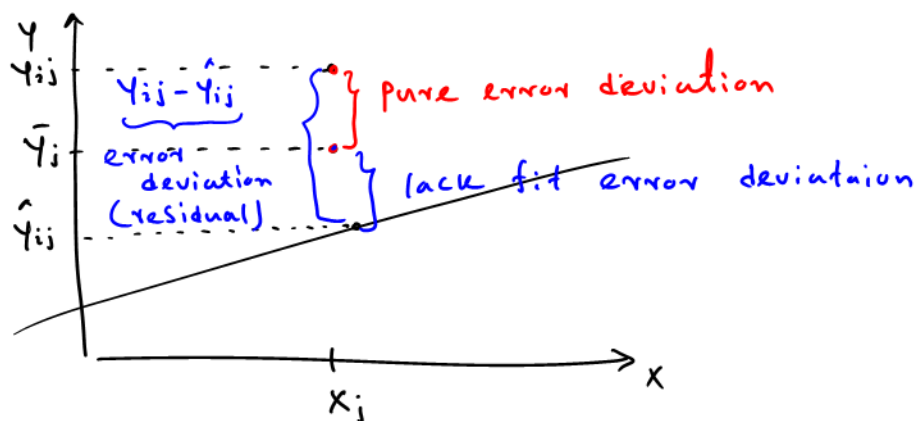
\uparrow
 \bar{y}



$$y_{ij} = \mu_j + \varepsilon_{ij}$$

\uparrow
 \bar{y}_j

Note:



From the graph

$$y_{ij} - \hat{y}_{ij} = y_{ij} - \bar{y}_j + \bar{y}_j - \hat{y}_{ij}$$

We also can show.

$$\underbrace{\sum (y_{ij} - \hat{y}_{ij})^2}_{SSE} = \underbrace{\sum (y_{ij} - \bar{y}_j)^2}_{SSPE} + \underbrace{\sum (\bar{y}_j - \hat{y}_{ij})^2}_{SSLF}$$

Proof: HW.

All of above values are Summarized in a table called general ANOVA table.

General ANOVA table:

Source of variation	SS	df	MS
Regression	$SSR = \sum \sum (Y_{ij} - \bar{Y})^2$	1	$MSR = \frac{SSR}{1}$
<u>error</u>	$SSE = \sum \sum (Y_{ij} - \hat{Y}_{ij})^2$	$n-2$	$MSE = \frac{SSE}{n-2}$
{ Lack of fit Pure error	$SSLF = \sum \sum (\bar{Y}_j - \hat{Y}_{ij})^2$	$c-2$	$MSLF = \frac{SSLF}{c-2}$
	$SSPE = \sum \sum (Y_{ij} - \bar{Y}_j)^2$	$n-c$	$MSPE = \frac{SSPE}{n-c}$
Total	$SSTO = \sum \sum (Y_{ij} - \bar{Y})^2$	$n-1$	