

## Remedial Measures

If the SLR model is not appropriate for data, there are two basic choices.

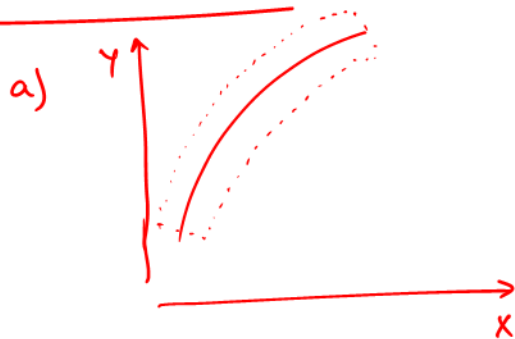
- 1) Using a more general regression model (multiple linear, Polynomial) (Later)
- 2) Using some transformations on  $X$  or  $Y$  or both such that the SLR model is appropriate for transform data.

## Transformations

### \* For non-linear relations

When the relation between  $X$  and  $Y$  is not linear (but the distribution of the error term is normal with constant error variance),  $X$  should be transformed.

### Basic Rule: pattern

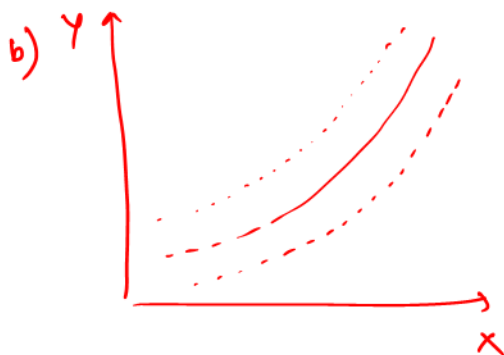


### Transformation

$$x' = \log_{10} X \quad \text{OR}$$

$$x' = \sqrt{X} \quad \text{OR}$$

$$x' = \ln(X).$$



$$x' = x^2$$

$$x' = \exp(X)$$

$$Y = \beta_0 + \beta_1 x'$$

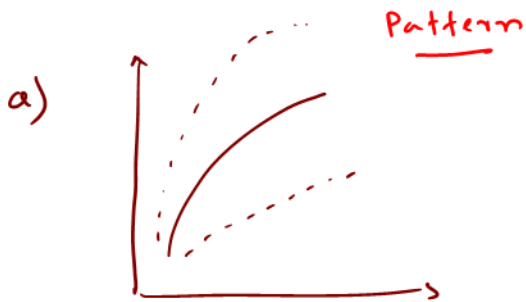


$$x' = 1/x$$

$$x' = \exp(-x)$$

### \* Non-normality and unequal error variance

If error term is not normal and variance is not constant (usually appear together), the response  $y$  should be transformed.



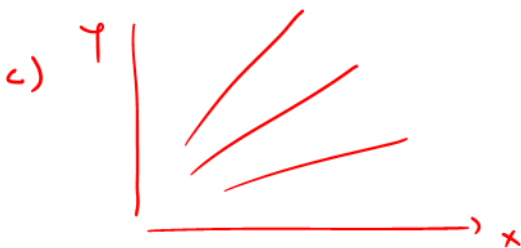
transformation

$$y = \beta_0 + \beta_1 x + \xi$$

$$y' = \sqrt{y}$$



$$y' = \log_e y$$



$$y' = 1/y$$

Note:

- \* Some times we may have to transform both  $X$  and  $Y$ .
- \* Several transformation should be tried, then choose the better one comparing scatter plot and the residual plots.
- \* Some times it is difficult to determine transformation from the plots. Box-cox transformation can be used for those cases.

### Box-cox transformation

Transformation  $Y' = Y^\lambda$ , where  $\lambda$  is a parameter to be determined from data.

The following are transformations based on the value of  $\lambda$ .

<u><math>\lambda</math></u>	<u>transformation</u>
$\lambda = 2$	$Y' = Y^2$
$\lambda = 0.5$	$Y' = \sqrt{Y}$
$\lambda = 0$	$Y' = \log_e Y = \ln Y$ ← (by the definition)
$\lambda = -0.5$	$Y' = \frac{1}{\sqrt{Y}}$
$\lambda = -1$	$Y' = \frac{1}{Y}$

New regression model

$$Y_i^\lambda = \beta_0 + \beta_1 X_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

How estimate  $\lambda$ ?

The Box-cox transformation procedure uses MLE of  $\lambda$  as well as other parameters  $\beta_0, \beta_1$  and  $\sigma^2$ .

\* Simple way to find  $\hat{\lambda}_{MLE}$ .

Steps:

1) Choose the <sup>several</sup> values for  $\lambda$  in the interval  $(-2, 2)$

Eg:  $-2, -1.75, -1.5, \dots, 2$

2) For each  $\lambda$ , standardize  $y_i^\lambda$  observations as follows,

$$W_i = \begin{cases} K_1 (y_i^\lambda - 1) & \lambda \neq 0, \\ K_2 (\log y_i) & \lambda = 0, \end{cases}$$

where  $K_1 = \left( \prod_{i=1}^n y_i \right)^{1/n}$  and  $K_2 = \frac{1}{\lambda K_1^{\lambda-1}}$ .

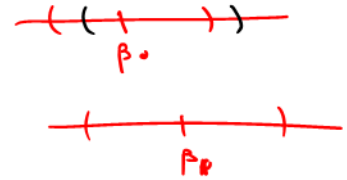
3) For each  $\lambda$ , fit a SLR model for  $W$  and  $X$  and calculate SSE.

(Here the value of SSE does not depend on  $\lambda$ , because we use standardized responses).

4) Choose the value of  $\lambda$  for which SSE is minimum. That is the MLE of  $\lambda$ .

Chapter - 4: Simultaneous Inferences and other topics in Regression Analysis.

Joint estimation of  $\beta_0$  and  $\beta_1$



Joint estimation is needed when we want  $1-\alpha$  (eg-95%) confidence that the conclusions for both  $\beta_0$  and  $\beta_1$ , are correct.

idea:

consider 95% c.i:s for  $\beta_0$  and  $\beta_1$  (separately)

Sample #	95% c.i: for $\beta_0$	95% c.i: for $\beta_1$	
1			✓
2			✗
3			✗
⋮			
100			✓

So if we consider both  $\beta_0$  and  $\beta_1$  together, the confidence coefficient may be less than 95%.

This is called the family confidence coefficient.

Problem: How to find c.i.s for both  $\beta_0$  and  $\beta_1$  with family confidence interval  $1-\alpha$ .

Answer: use Bonferroni method.

\* Bonferroni Method

Idea:

consider  $1-\alpha$  c.i. for  $\beta_0$  and  $\beta_1$  (separately)

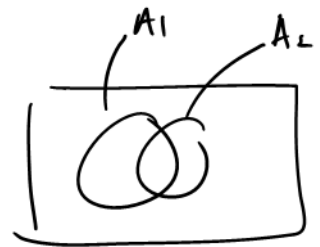
$$b_0 \pm t_{1-\alpha/2; n-2} S\{b_0\} \rightarrow \textcircled{1}$$

$$b_1 \pm t_{1-\alpha/2; n-2} S\{b_1\} \rightarrow \textcircled{2}$$

then,

$$P(\underbrace{\textcircled{1} \text{ is not correct}}_{A_1}) = P(A_1) = \alpha$$

$$P(\underbrace{\textcircled{2} \text{ is not correct}}_{A_2}) = P(A_2) = \alpha$$



$$P(\text{at least one of } \textcircled{1} \text{ OR } \textcircled{2} \text{ is not correct}) = P(A_1 \cup A_2) \\ = P(A_1) + P(A_2) - P(A_1 \cap A_2).$$

$$\Rightarrow P(\text{both correct}) = P(\bar{A}_1 \cap \bar{A}_2) \\ = P(\overline{A_1 \cup A_2}) \quad (\text{De Morgan's Law}) \\ = 1 - \underbrace{P(A_1 \cup A_2)}$$

$$\Rightarrow P(\bar{A}_1 \cap \bar{A}_2) = 1 - P(A_1) - P(A_2) + \underbrace{P(A_1 \cap A_2)}$$

But  $P(A_1 \cap A_2) \geq 0$ ,

$$P(\bar{A}_1 \cap \bar{A}_2) \geq 1 - P(A_1) - P(A_2) \quad (\text{Bonferroni Inequality})$$

$$\Rightarrow P(\text{Both correct}) \geq 1 - P(A_1) - P(A_2) \\ = 1 - \alpha - \alpha = 1 - 2\alpha.$$

(lower bound for the joint (family) confidence coefficient).

So if the family confidence limit is  $1-\alpha$ ,

$$1 - \underbrace{1-\alpha}_{\text{c: limit}} = \alpha$$

$\alpha/2$   
 $\alpha/2 + \beta$   
 $-$

$\alpha/2$   
 $\alpha/2 - \beta$   
 $-$

So Bonferroni c: I: s for  $\beta_0$  and  $\beta_1$  are the c: I: s with confidence coefficient " $1 - \alpha/2$ ".

$$b_0 \pm t_{1-\alpha/4: n-2} S\{b_0\} \quad \text{and}$$

$$b_1 \pm t_{1-\alpha/4: n-2} S\{b_1\}.$$

Note:

more generally  $1-\alpha$  family confidence intervals for  $\beta_0$  and  $\beta_1$  can be obtained by dividing  $\alpha$  into " $\alpha/2 - \beta$ " and " $\alpha/2 + \beta$ " ( $0 \leq \beta \leq \alpha/2$ ).

$$b_0 \pm t_{1-\frac{(\alpha/2 - \beta)}{2}: n-2} S\{b_0\} \quad \text{and}$$

$$b_1 \pm t_{1-\frac{(\alpha/2 + \beta)}{2}: n-2} S\{b_1\}$$

Eg: Suppose family c: c: is 90%. then,

$$1 - 0.9 = 0.1$$

$\beta_0$   
 $\alpha/2 + \beta = 0.08$

$\beta_1$   
 $\alpha/2 - \beta = 0.02$

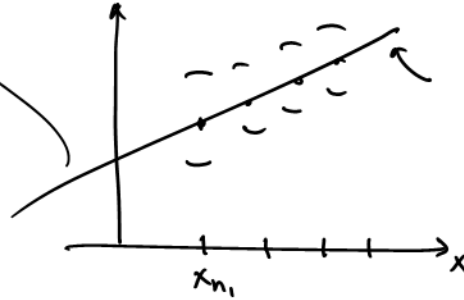
$\beta = 0.03$  (3%).

c: c for  $\beta_0$  is 92%.      c: c for  $\beta_1$  is 98%.

\* This method can be applied for more general regression models.  
 Suppose there are "g" parameters to estimate. If the family confidence coefficient is  $1-\alpha$ , then confidence coefficient for each parameter is " $1-\alpha/g$ ".

### Simultaneous Estimation for Mean Response

$$E[Y] = \beta_0 + \beta_1 X$$



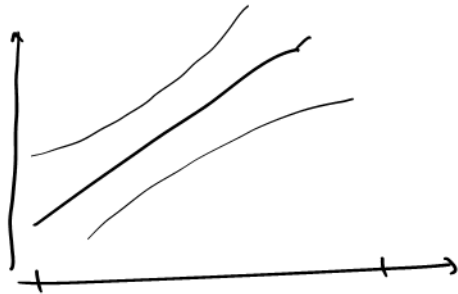
Some times we have to estimate the mean response  $E(Y)$  at a number of  $X$ -levels.

Let family confidence coefficient is  $1-\alpha$ .

There are two methods:

#### ① Working - Hotelling Procedure

This is a more general method and this gives a confidence band for the regression line. (This gives confidence intervals for all the  $X$  values in the range with family c.i.:  $1-\alpha$ ).



c.i.: is given by

$$\hat{y}_n \pm WS\{\hat{y}_n\} \text{ for all } X,$$

where  $w = \sqrt{2F(1-\alpha, 2, n-2)}.$

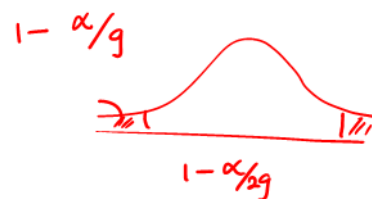


## ② Bonferroni Procedure

$1-\alpha$ , Bonferroni confidence limits for "g" levels are

$$\hat{Y}_n \pm B S\{\hat{Y}_n\},$$

where  $B = t(1-\alpha/2g; n-2)$ .



Note:

\* Working-Hotelling confidence limits do not change with number of intervals (g). But Bonferroni gets wider with g.

\* Both methods provide lower bound for the family confidence coefficients.

\* Given a data set, calculate c.i.s using both methods (i.e. calculate W and B), and then choose the most efficient one.

## Simultaneous Prediction Intervals

Consider the simultaneous prediction limits for "g" new observations at g different x-levels with family c.i.c:  $1-\alpha$ .

There are two methods:

### ① Scheffe Procedure

$$\hat{Y}_n \pm \mathcal{S} S\{\text{pred}\},$$

where  $\mathcal{S} = \sqrt{g F(1-\alpha; g; n-2)}$ .

## ② Bonferroni Procedure

$$\hat{y}_n \pm B S\{\text{Pred}\}$$

$$\text{where } B = t_{1-\alpha/2g; n-2,}$$

\* calculate both prediction intervals and then choose the most efficient one.

## Regression through the Origin

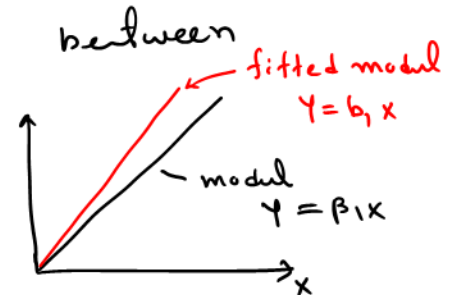
Why important?

There are some applications.

$$Y = \beta_0 + \beta_1 X$$

$$Y = \beta_1 X$$

Eg:- Suppose we need to model the association between  
 $X$  - # of units (output) and  
 $Y$  - Variable cost.



Model:

$$Y_i = \beta_1 X_i + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

$$\Rightarrow E[Y_i] = \beta_1 X_i$$

Least Square Estimator of  $\beta_1$

$$Q = \sum (Y_i - \beta_1 X_i)^2$$

$$\Rightarrow \frac{\partial Q}{\partial \beta_1} = \sum X_i (Y_i - \beta_1 X_i) \stackrel{\text{Set}}{=} 0$$

$$\Rightarrow b_1 = \frac{\sum X_i Y_i}{\sum X_i^2}$$

\*  $i^{\text{th}}$  fitted value =  $\hat{y}_i = b_1 x_i$

\*  $i^{\text{th}}$  residual =  $e_i = y_i - \hat{y}_i = y_i - b_1 x_i$

\* unbiased estimator of error variance

$$S^2 = \text{MSE} = \frac{\sum (y_i - \hat{y}_i)^2}{n-1} = \frac{\sum (y_i - \hat{y}_i)^2}{n-1}$$

Note:

\*  $E(b_1) = \beta_1$

\*  $\text{Var}(b_1) = \frac{\sigma^2}{\sum x_i^2}$  and  $S^2\{b_1\} = \frac{\text{MSE}}{\sum x_i^2}$

\*  $E(\hat{y}_n) = E(y_n)$

\*  $\sigma^2\{\hat{y}_n\} = \frac{x_n^2 \sigma^2}{\sum x_i^2}$  and  $S^2\{\hat{y}_n\} = \frac{x_n^2 \text{MSE}}{\sum x_i^2}$

\*  $\sigma^2\{\text{pred}\} = \sigma^2\left\{1 + \frac{x_n^2}{\sum x_i^2}\right\}$ , and  $S^2\{\hat{y}_n\} = \text{MSE}\left(1 + \frac{x_n^2}{\sum x_i^2}\right)$ .

Proof - HW.

Interval Estimators:

Parameters

$\beta_1$

$E[y_n]$

$\hat{y}_{n(\text{new})}$

Confidence Limit

$b_1 \pm t S\{b_1\}$

$\hat{y}_n \pm t S\{\hat{y}_n\}$

$\hat{y}_n \pm t S\{\text{pred}\},$

where  $t = t_{1-\alpha/2; n-1}$ .

Note:

1)  $\sum_{i=1}^n e_i \neq 0$  (in general).

2)  $\sum_{i=1}^n x_i e_i = 0$

3)  $SSE = \sum e_i^2 \overset{\text{may be}}{\geq} SSTO = \sum (y_i - \bar{y})^2$

4)  $R^2 = 1 - \frac{SSE}{SSTO}$  may be negative.

So residuals can not be used to check the quality of the fit.  
We always try to avoid using regression through the origin.