

where $S\{\hat{y}_n\} = \sqrt{\text{MSE} \left[\frac{1}{n} + \frac{(x_n - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]}$

HW - 2 : M (1-29)

Quiz - 1 : M (1-29)

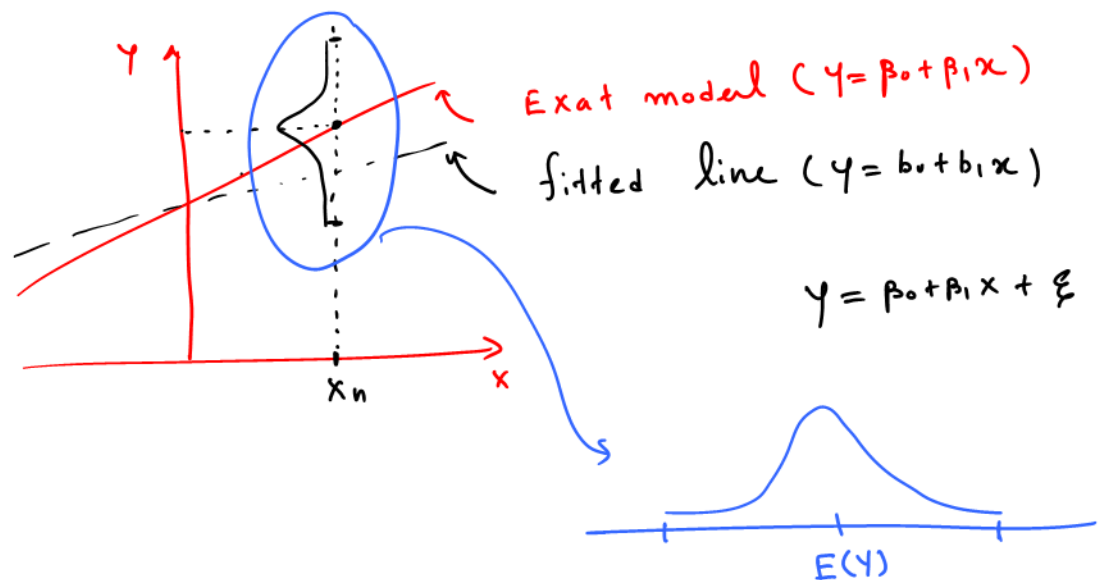
HW - 3 : M (2-05)

Quiz - 2 : M (2-05)

Class-3

Prediction Interval for a new observation ($\hat{y}_{n(\text{new})}$)

Consider a new observation $\hat{y}_{n(\text{new})}$ at $x = x_n$. Suppose $\hat{y}_{n(\text{new})}$ is independent of the observations on which the regression model is based.



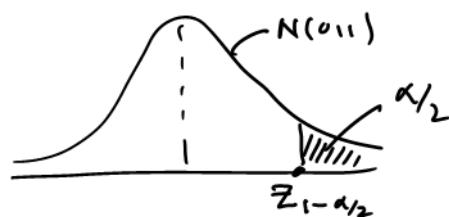
* when model parameters known

$$\hat{y}_{n(\text{new})} \sim N(E(y), \sigma^2)$$

$\therefore (1-\alpha)100\%$ prediction interval for $\hat{y}_{n(\text{new})}$ is

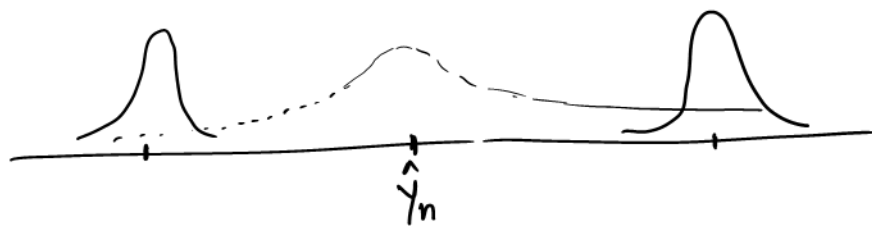
$$E(y_n) \pm Z_{1-\alpha/2} \sigma,$$

where $P(Z > Z_{1-\alpha/2}) = \alpha/2$, and $Z \sim N(0,1)$.



* When model parameters unknown

Now we have to estimate the mean response $E(Y_n)$ too.



So the point estimator for $\hat{y}_{n(new)}$ is \hat{y}_n .

But,

There are two variations to consider

1. Variation of $\hat{y}_{n(new)}$
2. Variation of \hat{y}_n .

Suppose the total variation is $\sigma^2\{\text{Pred}\}$.

Then considering all of the above,

$$\frac{\hat{y}_{n(new)} - \hat{y}_n}{\sigma\{\text{Pred}\}} \sim N(0,1), \text{ when } \sigma^2 \text{ is known.}$$

* when σ^2 is unknown

$$\frac{\hat{y}_{n(new)} - \hat{y}_n}{s\{\text{Pred}\}} \sim t_{n-2}.$$

← not a constant

← constant

$$\frac{b_1 - \beta_1}{s\{b_1\}}$$

$$s^2\{b_1\} = \sigma^2\{b_1 - \beta_1\}$$

Now,

$$\begin{aligned} \sigma^2\{\text{Pred}\} &= \sigma^2\{\hat{y}_{n(new)} - \hat{y}_n\} = \sigma^2\{\hat{y}_{n(new)}\} + \sigma^2\{\hat{y}_n\} \quad (\because \text{independent}) \\ &= \sigma^2 + \sigma^2\left[\frac{1}{n} + \frac{(x_n - \bar{x})^2}{\sum (x_i - \bar{x})^2}\right] \\ &= \sigma^2\left(1 + \frac{1}{n} + \frac{(x_n - \bar{x})^2}{\sum (x_i - \bar{x})^2}\right). \end{aligned}$$

← linear combination of y_1, y_2, \dots, y_n

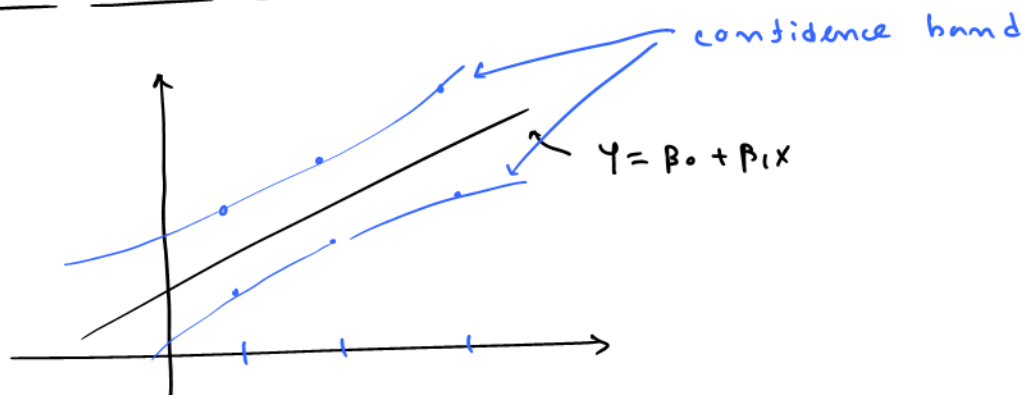
* when σ^2 is unknown,

$$S^2\{\text{pred}\} = \text{MSE} \left[1 + \frac{1}{n} + \frac{(x_n - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

$\therefore (1-\alpha)100\%$ prediction interval for $\hat{y}_{n(\text{new})}$ is

$$\hat{y}_n \pm t_{1-\alpha/2} S\{\text{pred}\}.$$

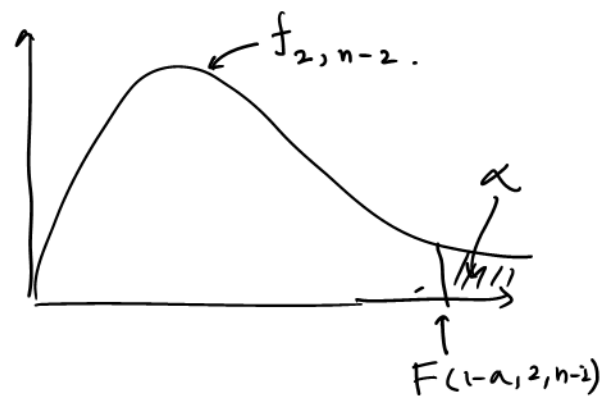
Confidence band for the regression line



Working Hotelling $(1-\alpha)100\%$ confidence level for the regression line for any level of x_n is given by

$$\hat{y}_n \pm w S\{\hat{y}_n\}, \text{ where } w^2 = 2 F(1-\alpha, 2, n-2),$$

F distribution value with degrees of freedom 2 and $n-2$.

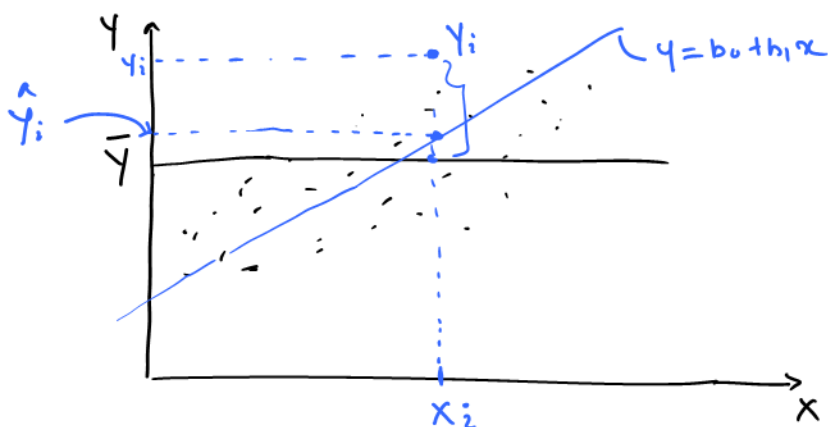


Note:

* W value is larger than $t(1-\alpha/2, n-2)$. So the boundary points of the confidence band for X_n is wider than the corresponding confidence interval.



Analysis of variance approach



$Y = \mu + \epsilon_i$ - center + error model

LSE of $\mu = \hat{\mu} = \bar{y}$

when $\beta_1 = 0$.

Total variation of variable Y (without fitting a regression model) is

$$SSTO = \sum_{i=1}^n (y_i - \bar{y})^2 - \text{total Sum of Squares.}$$

By fitting a regression model we reduce the variation. But still there is a variation due to the error term and it is given by

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 - \text{error Sum of Squares.}$$

Variation reduced by the regression model is

$$SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 - \text{regression Sum of Squares.}$$

From the graph:

$$y_i - \bar{y} = (y_i - \hat{y}_i) + (\hat{y}_i - \bar{y})$$

Note:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$$

Proof: HW

$$\text{ie } \underline{SSTO} = \underline{SSE} + \underline{SSR}$$

* degrees of freedom:

$$\text{df of } SSTO = \sum (y_i - \bar{y})^2 = n-1.$$

$$\text{df of } SSE = \sum (y_i - \hat{y}_i)^2 = n-2.$$

\uparrow
 $b_0 + b_1 x_i$

$$\text{df of } SSR = \sum (\hat{y}_i - \bar{y})^2 = 1.$$

$$\text{ie } df(SSTO) = df(SSE) + df(SSR).$$

* Mean Squares

A Sum Square divided by its degrees of freedom is called the mean square.

$$\text{* Total mean Square} = MSTO = \frac{\sum (y_i - \bar{y})^2}{n-1} \leftarrow \text{Sample variance} = \frac{SSTO}{n-1}$$

$$\text{* Regression mean Square} = MSR = \frac{SSR}{1} = \frac{\sum (\hat{y}_i - \bar{y})^2}{1}$$

$$\text{* Error mean Square} = MSE = \frac{SSE}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2}$$

Note:

Mean Squares are not additive

$$\text{ie } MSTO \neq MSR + MSE.$$

Analysis of Variance (ANOVA) table

All of the above values are summarized in a table called ANOVA table.

Source of Variation	SS	df	MS	$E[MS]$
Regression	$SSR = \sum (y_i - \bar{y})^2$	1	$MSR = \frac{SSR}{1}$	$\sigma^2 + \beta_1^2 \sum (x_i - \bar{x})^2$
Error	$SSE = \sum (y_i - \hat{y}_i)^2$	$n-2$	$MSE = \frac{SSE}{n-2}$	σ^2
Total	$SSTO = \sum (y_i - \bar{y})^2$	$n-1$		

Note:

1. $E[MSE] = \sigma^2$

Proof:

$y_i - \hat{y}_i \sim \text{distributed}$

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

Note: $\frac{SSE}{\sigma^2} \sim \chi_{n-2}^2$ - chi-square distribution with df $n-2$.

$$\Rightarrow E\left[\frac{SSE}{\sigma^2}\right] = n-2 \quad \left[\because E(\chi_n^2) = n \right]$$

$$\Rightarrow E\left[\frac{SSE}{n-2}\right] = \sigma^2$$

$$\Rightarrow E[MSE] = \sigma^2$$

2. $E(MSR) = \sigma^2 + \beta_1^2 \sum (x_i - \bar{x})^2$

Proof:

$$\begin{aligned} SSR &= \sum (\hat{y}_i - \bar{y})^2 = \sum [b_0 + b_1 x_i - (b_0 + b_1 \bar{x})]^2 \\ &= b_1^2 \sum (x_i - \bar{x})^2 \end{aligned}$$

$$\begin{aligned}
 E(SSR) &= E\left[b_1^2 \underbrace{\sum (x_i - \bar{x})^2}_{\square}\right] \\
 &= \sum (x_i - \bar{x})^2 E(b_1^2) \\
 &= \sum (x_i - \bar{x})^2 \left[\text{Var}(b_1) + (E(b_1))^2 \right] \\
 &= \sum (x_i - \bar{x})^2 \left[\frac{\sigma^2}{\sum (x_i - \bar{x})^2} + \beta_1^2 \right] \\
 &= \sigma^2 + \beta_1^2 \sum (x_i - \bar{x})^2_{\square}
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(X) &= E(X^2) - (E(X))^2 \\
 \Rightarrow E(X^2) &= \text{Var}(X) + (E(X))^2
 \end{aligned}$$

F-test

Another hypothesis test for β_1 .

Steps:

1) Hypotheses

$$\begin{array}{ll}
 H_0 : \underbrace{\beta_1 = 0}_{\substack{\text{no relationship} \\ \text{between } X \text{ and } Y}} & \text{vs} \quad H_1 : \underbrace{\beta_1 \neq 0}_{\substack{X \text{ and } Y \text{ are related.}}}
 \end{array}$$

2) Test statistic:

$$F = \frac{MSR}{MSE} \sim f_{1, n-2} \quad \text{— } F\text{-distribution with df 1 and } n-2.$$

3) Calculating the critical value or the p-value.

Critical value: $F(1-\alpha, 1, n-2)$.

p-value (Let f^* be the observed value of F)

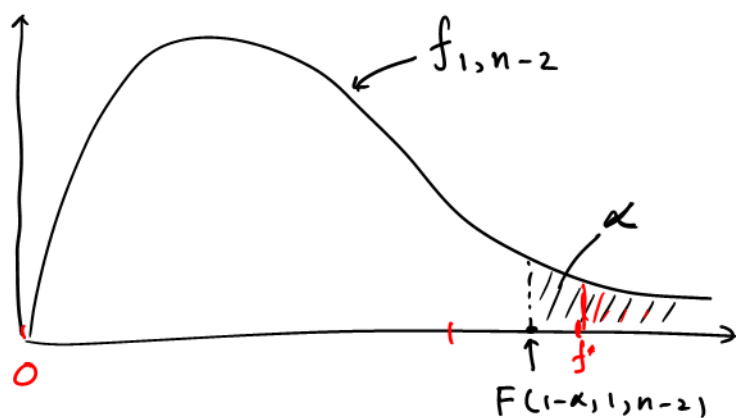
$$\text{p-value} = P(F > f^*)$$

4) Conclusion

Reject H_0 if $p\text{-value} < \alpha$ (usually $\alpha = 0.05$)

OR

Reject H_0 if $f^* > F(1-\alpha, 1, n-2)$.



Note:

$$\frac{MSR}{MSE} = \frac{\frac{MSR}{1}}{\frac{MSE}{n-2}} = \frac{\left(1 \cdot \frac{MSR}{1}\right) / 1}{\left(\frac{(n-2) \frac{MSE}{n-2}}\right) / (n-2)} = \frac{\chi^2_{(1)}/1}{\chi^2_{(n-2)}/(n-2)} \sim f_{1, n-2}.$$

General Linear Test approach

General Linear test is used to compare two linear models.

Steps:

1. Calculate SSE for the full model (unrestricted model)
2. Calculate SSE for the reduced model (restricted model)

Then,

$$F = \frac{SSE(R) - SSE(F)}{df_R - df_F} \div \frac{SSE(F)}{df_F} \sim f_{df_R - df_F, df_F}$$

If we compare the SLR model with $Y_i = \mu + \epsilon_i, i=1,2,\dots,n$,
 model: $\underbrace{Y = \beta_0 + \beta_1 X}_{\text{Full model}} \quad \underbrace{Y_i = \mu + \epsilon_i}_{\text{reduced model}}$

1) H_0 : Restricted model is better (i.e. $\beta_1 = 0$)

H_1 : The full model is better (i.e. $\beta_1 \neq 0$)

$$\begin{aligned} 2) F &= \frac{SSE(R) - SSE(F)}{df_R - df_F} \div \frac{SSE(F)}{df_F} \\ &= \frac{SSTO - SSE}{(n-1) - (n-2)} \div \frac{SSE}{n-2} \\ &= \frac{\frac{SSR}{1}}{\frac{SSE}{n-2}} \sim F_{1, n-2} \end{aligned}$$

* This is same as the test statistic of the F-test.

Note:

If the full model is SLR model, the general linear test and the F-test are both same. But these are different for multiple linear regression models.

Descriptive Measures of linear association between X and Y

1. Coefficient of Determination [R^2]

Coefficient of determination is a measure of reduction of variation as a proportion to the total variation.

$$\begin{aligned} R^2 &= \frac{SSTO - SSE}{SSTO} = \frac{SSR}{SSTO} \\ &= 1 - \frac{SSE}{SSTO} \end{aligned}$$

Note:

* Since $0 \leq SSE \leq SSTO$, $0 \leq R^2 \leq 1$.

* When all the observations fall on the fitted line, $SSE = 0$ and then $R^2 = 1$.



* When the fitted regression line is horizontal, $b_1 = 0$ and then $\hat{y}_i = \bar{y}$ for all $i = 1, 2, \dots, n$.

Then $SSE = SSTO$

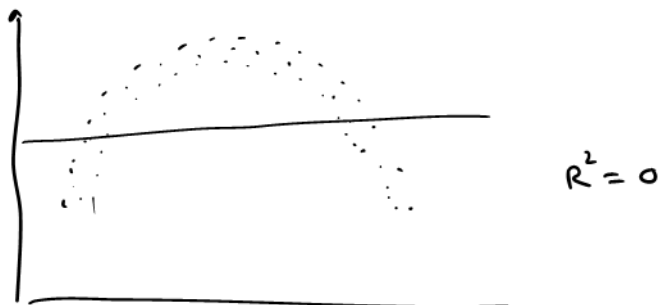
$$\Rightarrow R^2 = 0$$

i.e. regression does nothing.

(i.e. No linear relationship between X and Y).

* There may be another relationship with $R^2 = 0$.

Eg:-



2) Coefficient of correlation (r)

$$r = \pm \sqrt{R^2}$$

The sign of r is decided based on the sign of b_1 .

Chapter-3: Diagnostic and Remedial Measures

Before using the fitted regression model for future predictions we need to check the quality of the fit for data. Here we need to check the assumptions made for the model.

SLR model:

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

↖ constant.

Fitted model:

$$\hat{Y}_i = b_0 + b_1 X_i$$

There are two types of tools to check the quality of a fit. There are graphical tools (i.e. graphs) and numerical tools (tests).

Residuals play the main role of accessing the quality of the fit. The key point is the following.

$$\begin{array}{ccccccc} Y_i & = & \beta_0 & + & \beta_1 X_i & + & \epsilon_i \\ & & \downarrow & & \downarrow & & \downarrow \\ Y_i & = & b_0 & + & b_1 X_i & + & e_i \end{array}$$

Residuals can be used to estimate the error terms. So residuals should show properties (i.e. assumptions) of terms ϵ_i if the fitted model is appropriate.

Properties of residuals:

residuals

$$\underline{e_i} = Y_i - \hat{Y}_i = Y_i - (b_0 + b_1 X_i), \quad i = 1, 2, \dots, n.$$

$$\text{mean: } \bar{e} = \frac{\sum e_i}{n} = 0.$$

Variance:

$$\text{Var}\{e\} = \frac{\sum (e_i - \bar{e})^2}{n-2} = \frac{\sum e_i^2}{n-2} = \frac{\text{SSE}}{n-2} = \text{MSE}.$$

* Nonindependence

Since all e_i s are involved the fitted value \hat{y}_i which is based on the same regression function (both x_i), e_i s are not independent.

But for large sample, dependency among e_i s can be ignored.

* Semistudentized residuals

$$e_i^* = \frac{e_i - \bar{e}}{\sqrt{\text{MSE}}} = \frac{e_i}{\sqrt{\text{MSE}}}$$

* Things to check about the quality of the SLR fitted model

- ① Nonlinearity
- ② Non constancy of the error variance.
- ③ presence of outliers
- ④ Non-independence of error variance
- ⑤ Non-normality of error variance
- ⑥ Omission of important predictor variables.