

# STAC67: Regression Analysis

## Lecture 11

Sohee Kang

Feb. 24, 2021

# Outline of the Lecture

- Write the CI and PI in terms of matrix notation
- Re-write the Sum of Squares in ANOVA with quadratic forms
- Prove that MSE is the unbiased estimator of  $\sigma^2$
- Introduce the adjusted  $R^2$

# Properties of a prediction $\hat{Y}_{pred_0}$

- The prediction of a new observation

$$Y_0 = \underline{x}_0' \underline{\beta} + \epsilon_0$$

for a vector of values of independent variables  $\underline{x}_0$  is

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}$$

$$\hat{Y}_{pred_0} = \underline{x}_0' \hat{\underline{\beta}}$$

- The mean and variance of  $\hat{Y}_{pred_0}$  are

$$E(\hat{Y}_{pred}) = \underline{x}_0' \underline{\beta} \quad \text{Var}(\hat{Y}_{pred}) = \left( 1 + \underline{x}_0' (\underline{X}'\underline{X})^{-1} \underline{x}_0 \right) \sigma^2$$

accommodating for  $\epsilon_0$

Moreover, when  $\underline{\epsilon} \sim N(\underline{0}, \sigma^2 \underline{I})$  then

$$\hat{Y}_{pred} \sim N(\underline{x}_0' \underline{\beta}, (1 + \underline{x}_0' (\underline{X}'\underline{X})^{-1} \underline{x}_0) \sigma^2)$$

# Properties of a prediction $\hat{Y}_{pred_0}$

- The variance of a prediction for a new observation is larger than the variance of the estimator of the mean response even though the point estimate is the same. That is, for a vector of values of predictor variables  $\underline{x}_0$

- Prediction:

$$\hat{Y}_{pred_0} = \underline{x}_0' \hat{\beta} \quad \text{with variance} \quad [1 + \underline{x}_0' (\mathbf{X}' \mathbf{X})^{-1} \underline{x}_0] \sigma^2$$

- Mean response:

$$\hat{Y}_0 = \underline{x}_0' \hat{\beta} \quad \text{with variance} \quad \underline{x}_0' (\mathbf{X}' \mathbf{X})^{-1} \underline{x}_0 \sigma^2$$

$\sim N(\underline{x}_0' \hat{\beta}, \underline{x}_0' (\mathbf{X}' \mathbf{X})^{-1} \underline{x}_0 \sigma^2)$

# Estimates and precision: summary

Consider the model  $\underline{Y} = \underline{X}\underline{\beta} + \epsilon$ , with  $E(\underline{\epsilon}) = \underline{0}$ ,  $Var(\underline{\epsilon}) = I\sigma^2$

Quantity	Estimator	Variance of the estimator
$\underline{\beta}$	$\hat{\underline{\beta}}$	$(\underline{X}'\underline{X})^{-1} \sigma^2$
$E(\underline{Y})$	$\hat{\underline{Y}} = H\underline{Y}$	$\sigma^2 H$
$\underline{\epsilon}$	$\underline{e} = (I - H)\underline{Y}$	$(I - H) \sigma^2$
$Y_0$	$\underline{X}_0' \hat{\underline{\beta}}$	$(1 + \underline{X}_0' (\underline{X}'\underline{X})^{-1} \underline{X}_0) \sigma^2$

# Analysis of Variance and Quadratic Forms

- Quadratic forms of  $\underline{Y}$ :  $\underline{Y}'\mathbf{A}\underline{Y}$ , where  $\mathbf{A}$  is a symmetric matrix of coefficients called **defining matrix**
- **Next section**: Study the properties of residual, regression and total sum of squares and sum of squares used in inference
- They are all quadratic forms of  $\underline{Y}$

# Partitioning of total sum of squares

- We know that

$$\underline{Y} = \underline{\hat{Y}} + \underline{e}$$

- We will generalize the partitioning of the total sum of squares that we had for simple linear regression, i.e.

$$SST = SSR + SSE$$

$$\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

to multiple linear regression.

# Total sum of squares

- $SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$  in matrix notation:

## Exercise

Show that  $SST = \underline{Y}'\underline{Y} - \frac{1}{n}\underline{Y}'\underline{J}\underline{Y}$ , where  $\underline{J}$  is the  $n \times n$  square matrix with all elements equal to 1.

$$\begin{aligned} SST &= \sum Y_i^2 - n\bar{Y}^2 \\ &= \sum Y_i^2 - n\left(\frac{1}{n}\sum Y_i\right)^2 \\ &= \underline{Y}'\underline{Y} - \frac{1}{n}\underline{Y}'\underline{J}\underline{Y} \end{aligned}$$

- $SST$  is a quadratic form of  $\underline{Y}$  because

$$\underline{Y}' \underbrace{\left( \underline{I} - \frac{1}{n}\underline{J} \right)}_{\text{is symmetric}} \underline{Y}$$

- The defining matrix associated is

$$\underline{I} - \frac{1}{n}\underline{J}$$



# Residual sum of squares

- $SSE = \sum_{i=1}^n e_i^2$  in matrix notation:

## Exercise

Show that  $SSE = \underline{Y}'\underline{Y} - \hat{\underline{\beta}}'\underline{X}'\underline{Y}$

$$\begin{aligned} &= \sum e_i^2 \\ &= \underline{e}'\underline{e} \\ &= [\underline{Y} - \underline{X}\hat{\underline{\beta}}]'\underline{[Y - X\hat{\beta}]} \\ &= \underline{Y}'\underline{Y} - \underline{Y}'\underline{X}\hat{\underline{\beta}} - \hat{\underline{\beta}}'\underline{X}'\underline{Y} + \hat{\underline{\beta}}'\underline{X}'\underline{X}\hat{\underline{\beta}} = \underline{Y}'\underline{Y} - \hat{\underline{\beta}}'\underline{X}'\underline{Y} \end{aligned}$$

$\underline{Y}'\underline{X}\hat{\underline{\beta}}$  (show by  $\hat{\underline{\beta}}' = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{Y}$ )

- $SSE$  is a quadratic forms of  $\underline{Y}$  because

$$SSE = \underline{Y}'(\underline{I} - \underline{H})\underline{Y}$$

- The defining matrix is  $\underline{I} - \underline{H}$ .

# Regression sum of squares

- $SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$  in matrix notation:

## Exercise

Show that  $SSR = \hat{\beta}' \mathbf{X}' \mathbf{Y} - \frac{1}{n} \mathbf{Y}' \mathbf{J} \mathbf{Y}$

$$\begin{aligned} &= \sum \hat{Y}_i^2 - n \bar{Y}^2 \\ &= \hat{\mathbf{Y}}' \hat{\mathbf{Y}} - \frac{1}{n} (\sum Y_i)^2 \\ &= \hat{\mathbf{B}}' \mathbf{X}' \mathbf{X} \hat{\mathbf{B}} - \frac{1}{n} \mathbf{Y}' \mathbf{J} \mathbf{Y} \\ &= \hat{\mathbf{B}}' \mathbf{X}' \cancel{(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'} \mathbf{Y} - \frac{1}{n} \mathbf{Y}' \mathbf{J} \mathbf{Y} \\ &= \hat{\mathbf{B}}' \mathbf{X}' \mathbf{Y} - \frac{1}{n} \mathbf{Y}' \mathbf{J} \mathbf{Y} \end{aligned}$$

# Regression sum of squares

- $SSR$  is a quadratic forms of  $Y$  because

$$Y'X(X'X)^{-1}X'Y - \frac{1}{n} \tilde{Y}'JY = Y'(H - \frac{1}{n}J)Y$$

- The defining matrix associated is:  $H - \frac{1}{n}J$

## Exercise

Check that

$$\begin{aligned} SST &= SSR + SSE \\ Y'(I - \frac{1}{n}J)Y &= Y'(H - \frac{1}{n}J)Y + Y'(I - H)Y \\ &= Y'(I - \frac{1}{n}J)Y \end{aligned}$$

# Degrees of freedom

- For now: the number of values in the calculation of a statistic that can freely vary.
- $SST$  has  $n - 1$  degrees of freedom
- $SSE$  has  $n - (p + 1)$  degrees of freedom
- $SSR$  has  $(p + 1) - 1$  degrees of freedom.

## Mean squares

- **Mean squares:** sum of squares divided by its associated degrees of freedom
- **Regression mean squares:**

$$MSR = \frac{SSR}{p}$$

- **Residual mean squares:**

$$MSE = \frac{SSE}{n - (p + 1)}$$

# Analysis of variance table

- Analysis of variance (ANOVA) table to display the sum of squares and degrees of freedom

Source of variation	Sum of squares	df	Mean squares
Regression	SSR	$(p+1)-1$	$MSR = \frac{SSR}{(p+1)-1}$
Residual	SSE	$n-(p+1)$	$MSE = \frac{SSE}{n-(p+1)}$
Total	SST	$n-1$	

$$F = \frac{MSR}{MSE}$$

$$\sim F((p+1)-1, n-(p+1))$$

- The results in the ANOVA table will be used to construct a global test for the regression coefficients.

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$$

$$H_1: \text{at least one } \beta_i \neq 0$$

# Properties of a quadratic form of a random vector

$$\text{Var}(e) = (I - H) \sigma^2$$

$$E(e) = 0$$

$$HY = XB$$

- Consider a quadratic form

$$U = \underline{\underline{Z}}' \mathbf{A} \underline{\underline{Z}}$$

of a random vector  $\underline{\underline{Z}}$  where  $\mathbf{A}$  is a symmetric matrix (the defining matrix). We have

$$E(\underline{\underline{Z}}' \mathbf{A} \underline{\underline{Z}}) = \text{tr} [A \text{Var}(z)] + E(z)' A E(z)$$

$$\begin{aligned} \text{Var}(\underline{\underline{Z}}' \mathbf{A} \underline{\underline{Z}}) = & 2 \text{tr} (A \text{Var}(z) A \text{Var}(z)) \\ & + 4 E(z)' A \text{Var}(z) A E(z) \end{aligned}$$

# Unbiased estimator of $\sigma^2$

**Exercise:** Show that  $E(\underline{e}'\underline{e}) = (n - p')\sigma^2$ .

Hint: use that  $\text{tr}(\underline{H}) = p'$  (without proof) and the quadratic formulation of  $\underline{e}'\underline{e}$ .

$$\begin{aligned} E(\underline{e}'\underline{e}) &= \text{tr}(\text{Var}(\underline{e})) + E(\underline{e})'E(\underline{e}) \\ &= \text{tr}((\underline{I} - \underline{H})\sigma^2) \quad \text{to} \\ &= \sigma^2(\text{tr}(\underline{I}) - \text{tr}(\underline{H})) = \sigma^2(n - p') \quad p' = p+1 \end{aligned}$$

**Exercise:** Show that the estimator

$$MSE = s^2 = \frac{\underline{e}'\underline{e}}{n - p'}$$

is an unbiased estimator of  $\sigma^2$ .

$$E\left(\underline{e}'\underline{e} / (n - p')\right) = \frac{1}{n - p'} \cdot (n - p')\sigma^2 = \sigma^2$$

Thus is an unbiased estimator.

# Coefficient of multiple determination

$$R^2 = \frac{SSR}{SST}$$

- Fraction of the variation in  $Y$  explained by the model (i.e. its linear relationship with  $X_1, \dots, X_p$ )
- We have  $0 \leq R^2 \leq 1$

**Exercise** Show that

- 1  $R^2 = 0$  when  $\hat{\beta}_k = 0$  for each  $k = 1, \dots, p$ .
- 2  $R^2 = 1$  when  $\hat{Y}_i = Y_i$  for each  $i = 1, \dots, n$ , i.e. when all the observations fall on the fitted regression surface.



# Adjusted coefficient of multiple determination}

- Adding more  $X$  to the model  $R^2$ .
- Adjusted  $R^2$ : modified measure that accounts for the number of variables in the model.
- Adjusted coefficient of multiple determination:

$$R_{adj}^2 = 1 - \frac{\frac{SSE}{n-p'}}{\frac{SST}{n-1}} = 1 - \left( \frac{n-1}{n-p'} \right) \frac{SSE}{SST}$$

- $R_{adj}^2$  does not have the same interpretation as  $R^2$ .
- $R_{adj}^2$  may decrease when we add a new variable because
- $R_{adj}^2$  useful for selecting explanatory variables.