

Lecture 9

Last class: categorical variables (variables that can take values that fall into several categories). For >2 categories,

- Treat them as numerical (if appropriate)
- Converted into indicator/dummy variables
- In general, k categories will require $k-1$ indicator variables.

ANOVA

- Consider the regression analysis from new perspective called Analysis Of Variance (ANOVA)
- How well does a regression model fit the response variable?
- Idea of ANOVA is to partition the variability in the responses. In particular, the variability in the response is measured by

Total Sums of Squares ($SS(\text{tot})$):

$$SS(\text{tot}) = \sum_{i=1}^n (y_i - \bar{y})^2 : \text{clearly related to sample var. of } y_1, \dots, y_n.$$

(sample variance is $SS(\text{tot})/n-1$).

- Measures the deviation of the responses from the sample mean.
- Greater the $SS(\text{tot})$, greater the variation.

ANOVA decompose $SS(\text{tot})$ as follows:

$$SS(\text{tot}) = SS(\text{res}) + SS(\text{reg})$$

$$(1) \quad SS(\text{res}) = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 = \sum_{i=1}^n e_i^2 \quad (\text{residual sums of squares}).$$

→ Measures the deviation of responses from fitted values.

$$(2) \quad SS(\text{reg}) : \text{regression sums of squares}$$

$$SS(\text{reg}) = \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2 \quad \left(\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \right)$$

→ Measures the deviation of fitted values from sample mean.

* $SS(\text{reg})$: variability that is explained by our model.

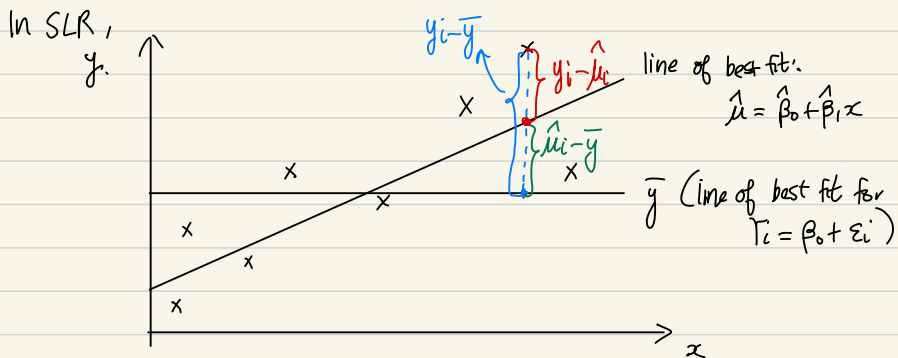
$SS(\text{res})$: variability that is not explained by our model.

$$SS(\text{tot}) = SS(\text{res}) + SS(\text{reg})$$

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

Decomposition of $y_i - \bar{y}$:

$$\underline{y_i - \bar{y}} = \underline{y_i - \hat{\mu}_i} + \underline{\hat{\mu}_i - \bar{y}}$$



- Note
1. When regression model fits the data well, the observations will be close to fitted values $\Rightarrow y_i - \hat{\mu}_i$ may be smaller, $\hat{\mu}_i - \bar{y}$ may be larger.
 2. Line given by \bar{y} is obtained using regression where it's assumed that $\beta_1 = 0$ (no relationship between explanatory var. and response)

Mathematically, we can show the partition of $SS(\text{tot})$ as follows:

$$\begin{aligned} SS(\text{tot}) &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{\mu}_i + \hat{\mu}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 + \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2 + 2 \sum_{i=1}^n \underbrace{e_i}_{\text{blue}} (\hat{\mu}_i - \bar{y}) \\ &= \sum_{i=1}^n (y_i - \hat{\mu}_i)^2 + \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2 + 2 \underbrace{\sum_{i=1}^n e_i \hat{\mu}_i}_{=0} - 2 \bar{y} \underbrace{\sum_{i=1}^n e_i}_{=0} \\ &= SS(\text{res}) + SS(\text{reg}). \end{aligned}$$

$$\therefore \underline{SS(\text{tot})} = \underline{SS(\text{res})} + \underline{SS(\text{reg})}$$

Measure of deviation of responses from sample mean

Measure of deviation of responses from fitted values

Measure of deviation of fitted values from sample mean.

Breakdown of degrees of freedom (df):

- $SS(\text{tot}) = \sum_{i=1}^n (y_i - \bar{y})^2$ has $n-1$ df.
 - 1 df lost due to estimation of sample mean.
 - denominator of sample variance of response.
- $SS(\text{res}) = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2$ has $n-p-1$ df
 - $p+1$ df lost due to estimation of $\vec{\beta}$ for estimating fitted values.
 - denominator of $\hat{\sigma}^2$
- $SS(\text{reg}) = \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2$ has p df.
 - df of $SS(\text{res}) + \text{df of } SS(\text{reg})$ equals df of $SS(\text{tot})$
 - ⇒ df of $SS(\text{tot}) = n-1$, $(n-1) - (n-p-1) = p$
 - p is the number of explanatory variables.

Summaries in ANOVA table

Source	SS	df	Mean Squares ^{MS}	F
Regression	$SS(\text{reg})$	p	$MS(\text{reg}) = SS(\text{reg})/p$	$MS(\text{reg})/MS(\text{res})$
Residual	$SS(\text{res})$	$n-p-1$	$MS(\text{res}) = SS(\text{res})/n-p-1$	//////
Total	$SS(\text{tot})$	$n-1$	//////	//////

Mean Squares:

- $MS(\text{reg}) = SS(\text{reg})/p$
- $MS(\text{res}) = SS(\text{res})/n-p-1 = \hat{\sigma}^2$

F-test: (next class)

- Used to test for the significance of our regression model.

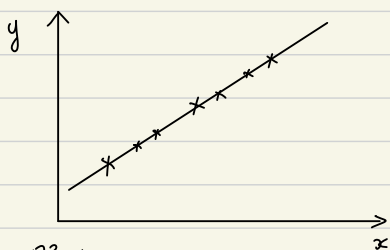
Coefficient of Determination (R^2):

$$R^2 = SS(\text{reg})/SS(\text{tot})$$

- R^2 is the proportion of total variation in the response variable that's explained by regression model.
- Bigger the R^2 value, the better fit. (Fitted values are closer to y_i 's.).
- R^2 will be in $[0, 1]$.

When $R^2=1$, then this implies a perfect fit s.t. $\hat{\mu}_i = y_i$, $\forall i$
and that $SS(\text{res})=0$ and $SS(\text{reg})=SS(\text{tot})$ (generally not the case)

In SLR,



$$R^2=1$$

$$SS(\text{reg})=SS(\text{tot}), SS(\text{res})=0$$

$$\hat{\mu}_i = y_i, \forall i$$



$$R^2=0$$

$$SS(\text{res})=SS(\text{tot}), SS(\text{reg})=0 \quad \frac{\sum (\hat{\mu}_i - \bar{y})^2}{n} = 0$$

$$\hat{\mu}_i = \bar{y}, \forall i \quad (\hat{\beta}_0 = \bar{y}, \hat{\beta}_1 = 0)$$

Note: In SLR, R^2 is equal to squares of sample correlation between x and y (or r)
 $\therefore R^2 = r^2$

$$\text{Recall: } r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \Rightarrow r^2 = \frac{S_{xy}^2}{S_{xx} S_{yy}}$$

$$\begin{aligned} R^2 &= \frac{\sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \rightarrow \sum_{i=1}^n (\hat{\mu}_i - \bar{y})^2 = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y})^2 \\ &= \sum_{i=1}^n (\bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i - \bar{y})^2 \quad (\text{def. of } \hat{\beta}_0) \\ &= \sum_{i=1}^n \left[\frac{S_{xy}}{S_{xx}} (x_i - \bar{x}) \right]^2 \quad (\text{def. of } \hat{\beta}_1) \\ &= \frac{S_{xy}^2}{S_{xx}} \end{aligned}$$

$$\therefore R^2 = \frac{S_{xy}^2}{S_{xx} S_{yy}} \Rightarrow R^2 = r^2 \quad \square$$

Numerical Example Lung Function Data ($n=655$; age, height, gender, smoke).

Source	SS	df	MS	F
Reg	369.88	4	92.47	462
Res	130.67	650	0.20	//////
Total	500.55	654	//////	//////

$$R^2 = SS(\text{reg}) / SS(\text{tot})$$

$$= 0.74.$$

\therefore regression model with age, height, gender and smoking status explains 74% of variation in the response.