

Date: _____

u MM allu. M
u u u u u / u

→ When we run too many test type I error increases

Type I error

rejecting the null hypothesis when it's actually true

→ The scenario of testing many pairs of groups is called multiple comparisons

→ The Bonferroni correction suggests that a more stringent significance α is more appropriate for these tests

$$\alpha^2 = \alpha / k$$

where k is the no. of comparison being considered

→ If there are k groups, then usually all possible pairs are compared and

$$k = \frac{k(k-1)}{2}$$

Date: _____

Difference in two means: after ANOVA

$$SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}} \approx \sqrt{\frac{MSE}{n_1} + \frac{MSE}{n_2}}$$

$$T_{d(E)} = \frac{(\bar{X}_1 - \bar{X}_2)}{\sqrt{\frac{MSE}{n_1} + \frac{MSE}{n_2}}}$$

Residuals

Diff b/w
observed (y_i) and
predicted (\hat{y}_i)

$$e_i = y_i - \hat{y}_i$$

Quantifying the relationship

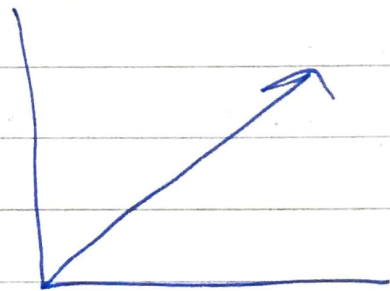
→ Correlation describes the strength of the linear association b/w two variable

→ It takes values b/w -1 (perfectly neg) and +1 (perfectly pos)

→ A value of 0 indicates no linear association

Fitting a line by least
Sq regression

→ Minimize the sum of mag
(abs values) of residuals



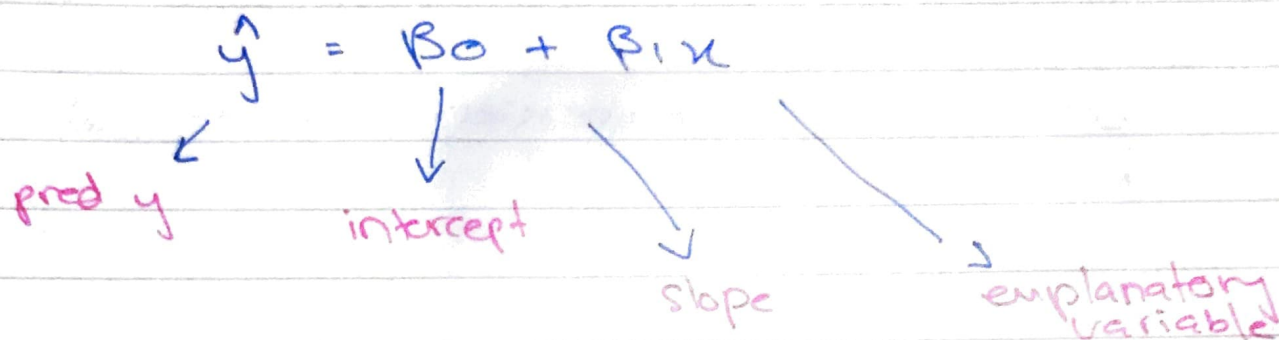
$$|e_1| + |e_2| + \dots + |e_n|$$

Date: _____

→ Minimize the sum of sq of residuals least squares

$$e_1^2 + e_2^2 + \dots + e_n^2$$

The Least Sq line



Intercept

Slope

- Parameter : β_0
- Point estimate : b_0

- Parameter : β_1
- Point estimate : b_1

Conditions

1) Linearity

The relation b/w explanatory and the response variable should be linear

2) Nearly Normal residuals

The residuals should be nearly normal.

3) Constant Variability

The variability of points around

Date: _____

the least sq line should be roughly const. This implies that the variability of residuals around the O line should be roughly const. also called homoscedasticity

Slope

The slope of a regression can be calculated as

$$b_1 = \frac{S_y}{S_x} R$$

Intercept

The intercept is where the regression line intersects the y-axis.

$$b_0 = \bar{y} - b_1 \bar{x}$$

→ when $x=0$, y is expected to equal the intercept

→ For each unit in x , y is expected to increase / decrease on average by the slope

Date: _____

Prediction

→ using linear model to pred the value of response variable for a given value of the explanatory variable is called prediction simply plugging in the x value in the linear model eq

Extrapolation

→ Applying a model estimate to values outside of the realm of the original data is called extrapolation

R^2

The strength of the fit of a linear model is commonly evaluated using R^2

→ R^2 is calculated as the sq of correlation coefficient

least sq regression Line Formula

$$\hat{\beta}_1 = \frac{SS_{xy}}{SS_{xx}}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$SS_{xx} = \sum x^2 - \frac{1}{n} \left(\sum x \right)^2$$

Date: _____

$$SS_{xy} = \sum xy - \frac{1}{n} (\sum x)(\sum y)$$

$$\bar{y} = \hat{\beta}_1 x + \hat{\beta}_0$$

Example

verify that it fits the data better than the line $\hat{y} = \frac{1}{2}x - 1$

x	y
2	0
2	1
6	2
8	3
10	3

x	y	x ²	x	xy	y ²
2	0	4	2	0	0
2	1	4	2	2	1
6	2	36	6	12	4
8	3	64	8	24	9
10	3	100	10	30	9
Total	28	9	208	68	23

$$\boxed{SS_{xx}} = \sum x^2 - \frac{1}{n} (\sum x)^2 = 208 - \frac{1}{5} (28)^2 = 51.2$$

Date: _____

$$\begin{aligned}\boxed{SS_{xy}} &= \sum xy - \frac{1}{n} (\sum x)(\sum y) \\ &= 68 - \frac{1}{5} (28)(9) \\ &= 17.6\end{aligned}$$

$$\begin{aligned}\boxed{\bar{x}} &= \sum x / n \\ &= 28 / 5 \\ &= 5.6\end{aligned}$$

$$\begin{aligned}\boxed{\bar{y}} &= \sum y / n \\ &= 9 / 5 \\ &= 1.8\end{aligned}$$

$$\begin{aligned}\boxed{SS_{yy}} &= \sum y^2 - \frac{1}{n} (\sum y)^2 \\ &= 23 - \frac{1}{5} (9)^2 \\ &= 23 - \frac{1}{5} (81) \\ &= \frac{115 - 81}{5} = \frac{34}{5} = 6.8\end{aligned}$$

$$\boxed{\hat{\beta}_1} = \frac{SS_{xy}}{SS_{xx}} = \frac{17.6}{51.2} = 0.34375$$

$$\begin{aligned}\boxed{\hat{\beta}_0} &= \bar{y} - \hat{\beta}_1 \bar{x} \\ &= 1.8 - 0.34375(5.6) \\ &= -0.125\end{aligned}$$

$$\boxed{\hat{y}} = 0.34375x - 0.125$$

$$\begin{aligned}\boxed{r} &= \frac{SS_{xy}}{\sqrt{SS_{xx} \cdot SS_{yy}}} = \frac{17.6}{\sqrt{51.2 \times 6.8}} = \frac{17.6}{\sqrt{348.16}} = \frac{17.6}{18.66} = 0.9432\end{aligned}$$

Date: _____

$$\begin{aligned}\boxed{SSE} &= SS_{yy} - \hat{B}_1 SS_{xy} \\ &= \overset{6.8}{\cancel{20.41}} - 0.34375 \times 17.6 \\ &= \overset{6.8}{\cancel{20.41}} - 6.05 \\ &= \cancel{14.36} \quad 0.75\end{aligned}$$

$$\begin{aligned}\boxed{\text{Slope}} &= \frac{\sqrt{SS_{yy}/n} * R}{\sqrt{SS_{xx}/n-1}} \\ &= \frac{\sqrt{\cancel{6.8/4} * \cancel{0.94321}}}{\sqrt{\cancel{51.2/2}}} \\ &= \frac{\sqrt{\frac{6.8}{5-1}} * 0.9432}{\sqrt{\frac{51.2}{5-1}}} \\ &= \frac{\sqrt{\frac{6.8}{4}} * 0.9432}{\sqrt{\frac{51.2}{4}}} \\ &= \frac{1.30 * 0.9432}{3.58} \\ &= 0.36 * 0.9432 \\ &= 0.34\end{aligned}$$

Date: _____

Types of Outliers

- Outliers are points that lie away from the cloud of points
- Outliers that lie horizontally away from the center of the cloud are called ~~high~~ high leverage points
- High leverage points that actually influence the slope of the regression line are called influential points

ANOVA Example

u_1	u_2	u_3
Group 1	Group 2	Group 3
2	3	12
8	4	15
6	11	2
10	13	5
20	17	8
$\sum \bar{u}_1 = 9.2$	$\sum \bar{u}_2 = 9.6$	$\sum \bar{u}_3 = 8.4$

$$\bar{X} = \frac{\sum \bar{u}_1 + \sum \bar{u}_2 + \sum \bar{u}_3}{3} = \frac{9.2 + 9.6 + 8.4}{3} = 9.07$$

Date: _____

$$\begin{aligned} df_g &= k-1 \\ &= 3-1 \\ &= 2 \end{aligned}$$

$$\begin{aligned} df_t &= n-1 \\ &= 15-1 \\ &= 14 \end{aligned}$$

$$\begin{aligned} df_E &= 14-2 \\ &= 12 \end{aligned}$$

$$\begin{aligned} SS_G &= \left[(5 \times (9.2 - 9.07)^2) + (5 \times (9.64 - 9.07)^2) \right. \\ &\quad \left. + (5 \times (8.4 - 9.07)^2) \right] \\ &= \left[5 \times (0.13)^2 + 5 \times (0.53)^2 + 5 \times (-0.67)^2 \right] \\ &= (5 \times 0.0169) + (5 \times 0.2809) + (5 \times 0.4489) \\ &= 0.0845 + 1.4045 + 2.2445 \\ &= 3.7335 \end{aligned}$$

$$\begin{aligned} SST &= |2 - 9.07|^2 + |3 - 9.07|^2 + |12 - 9.07|^2 + |8 - 9.07|^2 \\ &\quad + |4 - 9.07|^2 + |15 - 9.07|^2 + |6 - 9.07|^2 + \\ &\quad + |11 - 9.07|^2 + |2 - 9.07|^2 + |10 - 9.07|^2 + |13 - 9.07|^2 \\ &\quad + |5 - 9.07|^2 + |20 - 9.07|^2 + |17 - 9.07|^2 \\ &\quad + |8 - 9.07|^2 \end{aligned}$$

Date: _____

$$\begin{aligned} &= 49.9849 + 36.8449 + 8.5849 + 1.1449 \\ &+ 25.7049 + 35.1649 + 9.4249 + 3.7249 + \\ &44.9849 + 0.8649 + 15.4449 + 16.5649 \\ &+ 119.4649 + 62.8849 + 1.1449 \end{aligned}$$

$$= 372.9078 + 62.8849 + 1.1449$$

$$= 436.9076$$

$$\begin{aligned} SSE &= 436.9076 - 3.7335 \\ &= 433.1741 \end{aligned}$$

$$\begin{aligned} MSG &= SS_G / df_G \\ &= 3.7335 / 2 \\ &= 1.86675 \end{aligned}$$

$$\begin{aligned} MSE &= SSE / df_E \\ &= 433.1741 / 12 \\ &= 36.10 \end{aligned}$$

$$F = \frac{MSG}{MSE} = \frac{1.86675}{36.10} = 0.05$$