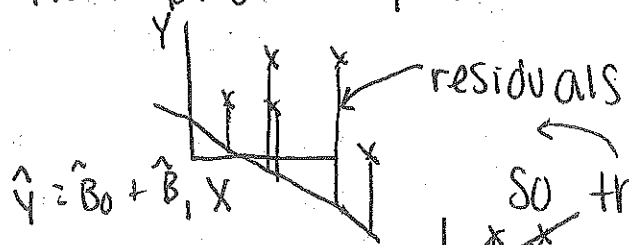


AMST  
Lecture 13  
7/18/16



- Galton 1890
- 1: reg line for predicting x from y
  - 2: SD line to capture trend
  - 3: reg line for predicting y from x

How predict y from x?



So this is a bad line

A better line →



← least squares line

$$\sum_{i=1}^n [y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i)]^2$$

$\hat{y}_i$

math fact: regression line = least squares line

ex. tail + wing length example cont.

L-248

Inf. Summary

unknown estimate	$\beta_1 =$ pop slope for predicting TL from WL $\hat{\beta}_1 = .77 \text{ cm TL/cm WL}$
------------------	--

math facts: 1.  $E_{IID}(\hat{\beta}_1) = \beta_1$  given

$$2. \hat{SE}_{IID}(\hat{\beta}_1) = \frac{S_y \sqrt{1-r^2}}{S_x \sqrt{n-2}}$$

where  $S_{y|x} = S_y \sqrt{1-r^2} \cdot \sqrt{\frac{n-1}{n-2}}$   
 residual SD = "root mean squared error" RMSE

To judge slope  $\hat{\beta}_1$  is large in practical terms, use same reasoning as w/ sample corr. r

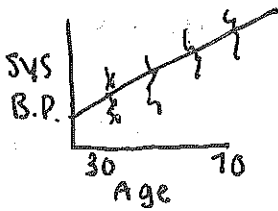
L-253

Biological Meaning of Y-int.

For WL/TL, it makes sense that  
as WL  $\rightarrow 0$  TL  $\rightarrow 0$ .

$\beta = -0.67, .77$

L-254



unobservable

$$y_i = (\beta_0 + \beta_1 x_i) + e_i$$

obs = truth + error

observable

$$\hat{y}_i = (\hat{\beta}_0 + \beta_1 x_i) + \hat{e}_i$$

obs = predicted + residual

We define  $\hat{e}_i$  as  $y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) = \text{obs} - \text{predicted}$

$\sigma$  = residual SD

$$\sigma_{y|x}^2 = \sqrt{\frac{1}{n-2} \sum_{i=1}^n \hat{e}_i^2}$$

RMSE

$\hat{\sigma}_{y|x}$  Represents the typical amt

$$\hat{\sigma}_{y|x} = S_{y|x} = S_y \sqrt{1-r^2} \cdot \sqrt{\frac{n-1}{n-2}}$$

math fact

Is the reg. practically useful?

①  $r^2 = \frac{S_y^2}{S_y^2} = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$

fact:  $v(\hat{e}) = (1-r^2) v(y)$

$r^2 = \frac{v(\hat{y})}{v(y)}$  = % of variance in y "explained by" res. of y on x

$r^2$  is called the coefficient of determination

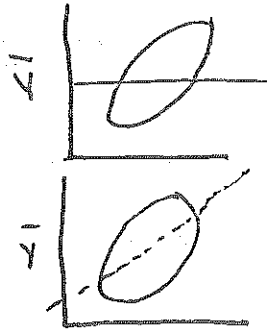
[0-1]

"associated with" is better  $\downarrow$

For complete formula list + deriving equations: See L notes + the reader!  $\ddot{\smile}$

How useful is the regression?

① Predict  $y$  ignoring  $x$   
 $\hat{y}_{no\ x} = \bar{y}$

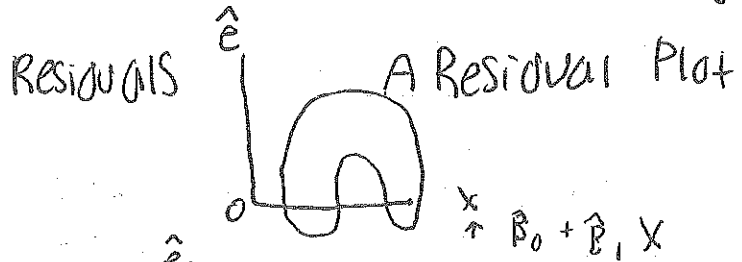


② Predict  $y$  using  $x$   
 $\hat{y}_{no\ x} = \hat{\beta}_0 + \hat{\beta}_1 x$

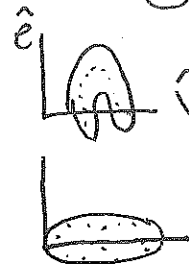
$$w/ \hat{SE}(\hat{y}_{use\ x}) = s_y \sqrt{1-r^2}$$

↑ this is smaller than  $SE(\hat{y}_{no\ x})$

Residual Plots



Official Res Plot



unhealthy: non linear curvature

healthy: no trend / pattern

one column: univariate sample  
 2 or more: multivariate sample

$$y_i = \beta_0 + \beta_1 x_i + e_i \quad \text{simple linear reg, } 1 \times 1$$

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} \dots \text{multiple linear reg}$$

$k > 1$   $x$  variables

Can generalize least squares to get estimate

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_k x_{ik}$$

Q: IS the mult. reg useful?

A:  $R^2$  = multiple  $R^2$  = coefficient of det.  
↑ want big

### Section 7: One-way Analysis of Variance

L-270

case study: Trees w/ 4 treatment groups  
models for each treatment

$M_1, M_2, \dots, M_4$   
 $\sigma$   
 $\sigma$   
 $\sigma$  } 4

remember basic model assumes all  $\sigma_{t,j} = \sigma$   
(I) how many  $I = 4$