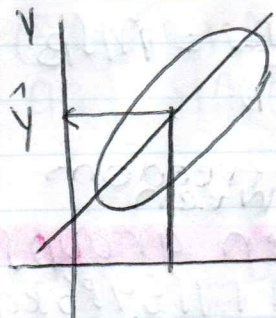


## Lecture #18: ANOVA



regression line for predicting

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

$SE(\hat{y})$  = uncertainty in using  $\hat{y}$  to predict y

$SE(\hat{y}) = s_{y|x} = RMSE = \text{residual SD}$

L-261

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_k x_{ik} + \epsilon_i$$

multiple linear regression (k=1 → simple linear regression)

### 1-way ANOVA

- Fisher (1925)
- L-275 case study
- Data values have 2 subscripts ex)  $y_{ij}$   
 ↑ group trees

• **Basic model:**  $\sigma_1 = \sigma_2 = \sigma_3 = \sigma_4 = \sigma$

• null:  $\mu_1 = \mu_2 = \dots = \mu_k = \mu$

• Alt: not so

- fail or alt / reject null if distance between actual data and expected data (if null true)

is large

- if null is true  $\bar{y}_i$  would be close to

$$\frac{n_1(\bar{y}_1 - \bar{y})^2 + n_2(\bar{y}_2 - \bar{y})^2 + \dots + n_I(\bar{y}_I - \bar{y})^2}{n_1 + n_2 + \dots + n_I - 1} = \sum_{i=1}^I n_i (\bar{y}_i - \bar{y})^2$$

= SS<sub>B</sub> sum of squares

$$\frac{1}{I-1} \sum_{i=1}^I n_i (\bar{y}_i - \bar{y})^2 = \frac{SS_B}{DF_B} = MS_B$$

degrees freedom  
between

**MS<sub>B</sub> = mean square between groups**  
 $= \frac{3.346455 \text{ kg}^2}{3} = 1.115485 \text{ kg}^2$

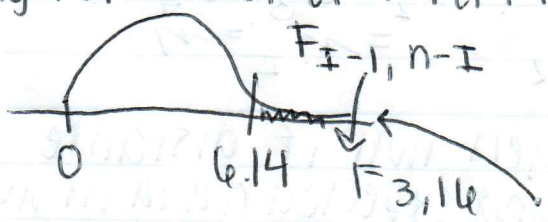
- can make **SS<sub>B</sub>** bigger or smaller by changing units
- Need estimate of noise in kg<sup>2</sup>

$$\frac{s_1^2 \quad s_2^2 \quad \dots \quad s_I^2}{(n_1-1)s_1^2 + (n_2-1)s_2^2 + \dots + (n_I-1)s_I^2}$$

(n<sub>1</sub>-1) + ... + (n<sub>I</sub>-1)  
 ↳ n - I

**signal / noise =  $\frac{MS_B}{MS_W} = F \text{ ratio} = 6.14127$**  in this case

long run hist of F ratio if null true



**statsig!**  
 $P = 0.0056 = 0.6\%$

- reject null if  $p \leq 5\%$ . **statsig ✓**
- Which groups differ?
  - Multiple comparisons: Many methods developed, look at simplest, **Bonferroni for pairwise comparisons**

- ① Decide # comparisons to make = k
- ② Decide level confidence =  $100(1-\alpha)\%$

③ pairwise comparisons in form  $(\bar{y}_i - \bar{y}_j)$  and standard error:

$$\hat{SE}(\bar{y}_i - \bar{y}_j) = \hat{\sigma} \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$$

where  $\hat{\sigma}$  = root mean squared error  
 $= \sqrt{MSW}$  ( $\hat{\sigma}^2 = MSW =$  **pooled variance estimate**)

~~④ there are multiple comparisons prob.~~

④ use a bigger  $t$  number to compensate for multiple comparisons being made (this is from Bonferroni) so make intervals wider

$$t_{1 - \frac{\alpha}{k}}_{n-I}$$

- always be given  $t$  number in problems bc given  $t$  values in table aren't enough