

STAT 303 / 1002

Note Title

1/30/2012

Pygmalion Example: Two-way Analysis of Variance

Does mean treatment effect differ with company?
Pygmalion is control

$H_0: \beta_{11} = \dots = \beta_{1q} = 0$ (coef. of interaction terms)
 H_a : at least one not zero

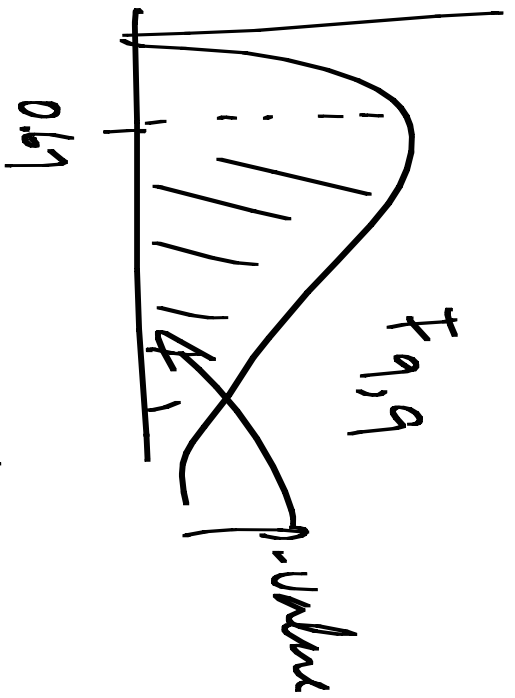
Partial F-test: Given in Type III SS

Type III SS = 311.46 = difference in SSes
 between full and reduced models
 (reduced model doesn't have a interaction term)

MSE for full model = 51.89

Test stat $311.46 / 9 = 51.89 = 0.67$

If H_0 is true this is an observation from a
 F distribution with 9, 9 df for error for full model



p-value is large

Data are consistent with 0 coefficients for interaction terms
 No evidence that treatment effect differs with company.

Next step: Since interaction is not significant, fit additive model so that we can assess "main effects" (Are there differences among companies? or treatments?)

$$\text{Model } Y_i = \beta_0 + \beta_1 I_{PYG,i}$$

ETH
Sampling

- 2
- 3
- 4
- 5
- ...
- 10

	Treatment	Control
Pygmalion	$\beta_0 + \beta_1 + \beta_2$	$\beta_0 + \beta_2$
	$\beta_0 + \beta_1 + \beta_3$	$\beta_0 + \beta_3$
	'	
	'	
	$\beta_0 + \beta_1$	β_0

$$+ \beta_2 I_{COMP1,i} + \beta_3 I_{COMP2,i} + \dots + \beta_{10} I_{COMP9,i} + \epsilon_i$$

Treatment effect
 $\frac{\beta_1}{\beta_1} = \text{control}$

- '
- '
- '
- β_1

Test: Is there a difference in mean score between

Pygmalion and control groups?

$$H_0: \beta_1 = 0 \quad \text{vs} \quad H_a: \beta_1 \neq 0$$

Test stat: $F_{obs} = 7.84$

This is an observation from a F distribution with 1, 18

$$p\text{-value} = .0119$$

We have evidence of a difference in mean score between pygmalion and control (over and above differences among companies)

Differences among companies?

$$H_0: \beta_2 = \dots = \beta_{10} = 0 \quad \text{vs} \quad H_a: \text{at least one}$$

Test stat: $F_{obs} = 1.75$, obs'd from $F(9, 18)$ under H_0

p-value = 0.1484
 No evidence of differences among companies.

Means statement output

Pygmalion
 Control
 N
 10
 19

Mean
 78.7

71.63

Std

7.24

$\sum Y_i$
 (ctrl) —

s.d. of
 19 control
 observations

LS means statement output

Pygmalion
 Control

78.7
 71.48

N ctrl
 Values estimated
 from model.

Here LSMEANS for control is NOT equal to Mean
of the 19 control observations because there
is not an equal number of control observations
per company.
(Company 3 only had 1 control station, the
rest of the companies had 2)

Note: Because of this lack of balance can't use
the Type III SS to assess main effects for a
model that has an interaction
(SAS uses LSMEANS in its calculation of Type III
SS so it is no longer equivalent to the Partial
F-test)

for main effects when model has interaction term

Check model assumptions

Form-Plots - no outliers

- perhaps decreasing variance

- normality OK

Independent observations? OK if we assume platforms
weren't interacting

→ p-values may only be approximate

Conclusion Evidence of a difference in mean score
between Pygmalion and control groups ($p = .01$)
(will consider this only weak evidence since I
have some concern about the variance estimate)
On average, Pygmalion students scored higher
than control students (mean 71.6)

Estimation of Regression Parameters

$$\text{Model: } Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i$$

Least squares find $\hat{\beta}_0, \dots, \hat{\beta}_p$ to minimize

$$RSS = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})^2$$

Because our assumption is $\varepsilon_i \sim N(0, \sigma^2)$
Then $y_i | x_{i1}, \dots, x_{ip} = x_{ip}$

$$\sim N(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}, \sigma^2)$$

So Joint density function of
 $y_1 | x_{11} = x_{11}, \dots, x_{1p} = x_{1p}, \dots, y_n | x_{n1} = x_{n1}, \dots, x_{np} = x_{np}$

(since y_i 's uncorrelated because ε_i 's are uncorrelated)

is $\prod_{i=1}^n f(y_i | x_i)$

$$= \left(\frac{1}{\sqrt{n} \sigma} \right)^n \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2 \right\}$$

Find β 's that

Minimize

β 's that

RSS is equivalent to finding the maximum of the joint density function of the Y 's conditional on the X 's

Least squares is equivalent to maximum likelihood estimation under the assumption of normally distributed errors

New Example

Dinner Party Example 1846 - covered wagon party left Illinois for California

- got stuck in snow in Sierra Nevada mountains in October
- rescued in April 1847

Data: for adults (≥ 15 years old)

- age
- sex
- survived or not

Want to model odds of survival based on age and sex.

Y_i is a binary variable (survived, died)

Model: ^(Binary) Logistic Regression

$$Y_i | X_i = \begin{cases} 1 & \text{if response in category of interest} \\ 0 & \text{otherwise} \end{cases}$$

$$Y_i | X_i \sim \text{Bernoulli}(\pi_i)$$

$$E(Y_i | X_i) = \pi_i$$

$$\text{Var}(Y_i | X_i) = \pi_i(1 - \pi_i)$$

Logistic regression model is an example of a Generalized Linear Model

Generalized Linear Models

Have response: y
+ set of explanatory variables X_1, \dots, X_p

Could have $Y|X \sim$ Bernoulli

or \sim Binomial

or \sim Normal

or \sim Poisson

or \sim Gamma

(regular regression)

Want to model $E(Y)$ as function of X_1, \dots, X_p

often the function is $f(x) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$
"linear"
 $= X\beta$

Key idea of Generalized Linear Models

LINK FUNCTION

function g st.

$$g(E(Y)) = X\beta$$

Some choices for the link function

① Identity link : $g(E(Y)) = E(Y)$

Model: $E(Y) = X\beta$

This is regression (STA 302)

Usual distribution:

$$Y | X \sim \text{Normal}$$