

10/04/16 ALSM 1

ANALYSIS OF VARIANCE (ANOVA IN THE DEPENDENT VARIABLE)

Variation in y can be quantified by error:

$$\sum (y_i - \bar{y})^2 = \sum y_i^2 \quad \text{Total Uncorrected sum of Squares} \quad SS(\text{uncorrected})$$

or more usually

$$\sum (y_i - \bar{y})^2 \quad \text{Total Corrected sum of squares} \quad SS(\text{corrected})$$

$$\sum (y_i - \bar{y})^2 = \sum (y_i^2 - 2\bar{y}y_i + \bar{y}^2)$$

$$= \sum y_i^2 - 2\bar{y}\sum y_i + n\bar{y}^2$$

$$= \sum y_i^2 - n\bar{y}^2$$

$$= SS(\text{uncorrected}) - \text{correction}$$

$$\text{Total corrected } \sum (y_i - \bar{y})^2 = \text{Total Uncorrected } \sum y_i^2 - \text{Correction } n\bar{y}^2$$

Write $SS(\text{uncorrected})$ in terms of fitted value and residual: $y_i = \bar{y} + \hat{e}_i$

$$SS(\text{uncorrected}) = \sum y_i^2 = \sum (\bar{y} + \hat{e}_i)^2$$

$$= \sum (\bar{y}^2 + 2\bar{y}\hat{e}_i + \hat{e}_i^2)$$

$$= \sum \bar{y}^2 + \cancel{\sum \bar{y}\hat{e}_i} + \sum \hat{e}_i^2 + 2\sum \bar{y}\hat{e}_i$$

$$\sum \hat{y}_i \hat{e}_i = \sum \hat{y}_i (y_i - \bar{y})$$

$$= \sum \hat{y}_i (y_i - \beta_0 - \beta_1 x_i)$$

$$= \sum \hat{y}_i (y_i - \bar{y} + \beta_1 \bar{x} - \beta_1 x_i)$$

$$= \sum \hat{y}_i [(y_i - \bar{y}) + \beta_1 (x_i - \bar{x})]$$

$$= \sum [\bar{y} + \beta_1 (x_i - \bar{x})] [(y_i - \bar{y}) + \beta_1 (x_i - \bar{x})]$$

$$= \sum [\bar{y}(y_i - \bar{y}) + \bar{y}\beta_1(x_i - \bar{x}) + \beta_1(x_i - \bar{x})(y_i - \bar{y}) + \beta_1^2(x_i - \bar{x})(x_i - \bar{x})]$$

$$= \bar{y} \sum (y_i - \bar{y}) + \bar{y}\beta_1 \sum (x_i - \bar{x}) + \beta_1 \sum (x_i - \bar{x})(y_i - \bar{y}) + \beta_1^2 \sum (x_i - \bar{x})^2$$

$$= \bar{y}(0) + \bar{y}\beta_1(0) + \beta_1 S_{xy} - \beta_1^2 S_{xx}$$

$$\frac{S_{xy}}{S_{xx}} S_{xy} - \frac{S_{xy}^2}{S_{xx}} S_{xx} = \frac{S_{xy}^2 - S_{xy}^2}{S_{xx}} = 0$$

$$\Rightarrow = \sum \hat{y}_i^2 + \sum \hat{e}_i^2$$

$$SS(\text{uncorrected}) = SS(\text{Model}) + SSE$$

$$\begin{aligned} SS(\text{Model}) &= \sum \hat{y}_i^2 = \sum (y_i - \bar{y} + \bar{y})^2 \\ &= \sum [(y_i - \bar{y})^2 + 2\bar{y}(y_i - \bar{y}) + \bar{y}^2] \\ &= \sum (y_i - \bar{y})^2 + 2\bar{y} \sum (y_i - \bar{y}) + n\bar{y}^2 \end{aligned}$$

$$\begin{aligned} \sum (y_i - \bar{y}) &= \sum (\beta_0 + \beta_1 x_i - \bar{y}) \\ &= \sum (y - \beta_0 - \bar{x} - \beta_1 x) \\ &= \sum (x - \bar{x}) \beta_1 \\ &= 0 \end{aligned}$$

$$\begin{aligned} \sum (y_i - \bar{y}) &= \sum (y - \beta_0 - \bar{y}) \\ &= \sum (y - \bar{y}) - \sum (\bar{y} - \bar{y}) \\ &= \sum (y - \bar{y}) \\ &= 0 \end{aligned}$$

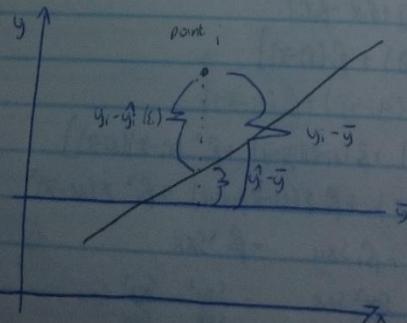
$$\rightarrow SS(\text{Model}) = \sum (y_i - \bar{y})^2 + n\bar{y}^2 = SS(\text{Regression}) + \text{Correction}$$

The regression sum of squares represents the reduction in variation (around the predicted mean) by adding the $\beta_1 x_i$ term to a model containing only \bar{y} .

$$\begin{aligned} SS(\text{uncorrected}) &= SS(\text{Reg}) + n\bar{y}^2 + SSE \\ \sum \hat{y}_i^2 &= \sum (y_i - \bar{y})^2 + n\bar{y}^2 + \sum \varepsilon_i^2 \end{aligned}$$

If we subtract correct $n\bar{y}^2$ from both sides:

$$\begin{aligned} \sum \hat{y}_i^2 - n\bar{y}^2 &= \sum (y_i - \bar{y})^2 + \sum \varepsilon_i^2 \\ SS(\text{corrected}) &= SS(\text{Reg}) + SSE \end{aligned}$$



$$\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{y})^2$$

10/04/16

ALSM I

7

ANVA

$$\begin{aligned}\sum (y_i - \bar{y})^2 &= \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2 \\ \text{SS(uncorrected)} &= \text{SS(Regression)} + \text{SSE} \\ &\quad (\text{model}) \quad (\text{error in general})\end{aligned}$$

Model Comparison

$$M_0: y_i = \beta_0 + \epsilon_i \quad \beta_0 = \bar{y} \quad \text{SSE}_0 = \sum (y_i - \bar{y})^2$$

$$M_1: y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad \beta_1 = \frac{\sum y_i}{\sum x_i} \quad \beta_0 = \bar{y} - \beta_1 \bar{x} \quad \text{SSE}_1 = \sum (y_i - \hat{y}_i)^2 \quad \hat{y}_i = \beta_0 + \beta_1 x_i$$

M₁ less restricted than M₀. SS E₂ < SS E₀SS E₁, SS E₂ measure of how much better M₁ is than M₀

$$\begin{aligned}\text{SSE}_1 - \text{SSE}_0 &= \sum (y_i - \bar{y})^2 - \sum (y_i - \hat{y}_i)^2 \\ &= \text{SS(uncorrected)} - \text{SSE}_2 \\ &= \text{SS(Regression)}\end{aligned}$$

SS(Reg) indicates importance of β_1 term in the model

$$\text{SS(Reg)} = R(\beta_1 | \beta_0)$$

$$M_0: y_i = \epsilon_i \quad \text{SSE}_0 = \sum (y_i - 0)^2 = \sum (y_i)^2$$

$$M_1: y_i = \beta_0 + \epsilon_i \quad \text{SSE}_1 = \sum (y_i - \bar{y})^2$$

$$\Rightarrow \text{SSE}_0 - \text{SSE}_1 \Rightarrow \sum (y_i)^2 - \sum (y_i - \bar{y})^2 \\ = \text{SS(uncorrected)} - \text{SS(corrected)} = n\bar{y}^2 \text{ (correction)}$$

Hence $n\bar{y}^2$ represents importance of β_0 in the model $R(\beta_0) = n\bar{y}^2$

$$M_0: y_i = \epsilon_i \quad \text{SSE} = \sum y_i^2$$

$$M_1: y_i = \beta_0 + \beta_1 x_i + \epsilon_i \quad \text{SSE} = \sum (y_i - \hat{y}_i)^2$$

$$\text{SSE}_0 - \text{SSE}_1 = \sum y_i^2 - \sum (y_i - \hat{y}_i)^2$$

$$= \text{SS(uncorrected)} - \text{SSE}$$

$$= \text{SS(Reg)} + n\bar{y}^2$$

$$= \text{SS(Model)}$$

$$-\text{Importance of } \beta_0 + \beta_1 \text{ in model} : R(\beta_0, \beta_1) = \text{SS}(M_0, M_1) = R(\beta_0) + R(\beta_1 | \beta_0)$$

Degrees of Freedom

$SS(\text{Unorrected})$	$\sum y_i^2$	n
$SS(\text{Corrected})$	$\sum (y_i - \bar{y})^2 \rightarrow \sum y_i^2 - n\bar{y}^2$	$n-1$
Correction	$n\bar{y}^2$	1
$SS(\text{Model}) R(\beta_0, \beta_1)$	$\sum (\hat{y}_i - \bar{y})^2 + n\bar{y}^2$	2 (1 for each param)
$SS(\text{Regression})$	$\sum (\hat{y}_i - \bar{y})^2$	1 (1 for β_1)
SSE	$\sum e_i^2 \Rightarrow \sum (y_i - \hat{y}_i)^2$	$n-2$ d.f.

$$SS(\text{corrected}) = SS(\text{uncorrected}) - \text{correction}$$

$$SS(\text{uncorrected}) = SS(\text{Model}) + SSE$$

$$SS(\text{Model}) = SS(\text{Regression}) + \text{correction}$$

$$SS(\text{Uncorrected}) = SS(\text{Regression}) + \text{correction} + SSE$$

$$SS(\text{Corrected}) = SS(\text{Regression}) + SSE$$

$$\begin{aligned} SS(\text{Reg}) &= \sum (y_i - \bar{y})^2 = \sum s_{xy} \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \\ &= \sum (\hat{\beta}_0 + \hat{\beta}_1 x_i - \bar{y})^2 \\ &= \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i - \bar{y})^2 \\ &= \sum (\hat{\beta}_1 (x_i - \bar{x}))^2 \\ &= \hat{\beta}_1^2 \sum (x_i - \bar{x})^2 \\ &= \hat{\beta}_1^2 S_{xx} \\ &= \frac{\sum x_i^2 - \bar{x}^2}{n-2} = \hat{\beta}_1 S_{xy} \end{aligned}$$

Source	DF	SS	Mean Square	$\frac{SS}{DF}$
Regression on X	1	$\sum (\hat{y}_i - \bar{y})^2$	$MS(\text{Reg})$	
Residual error	$n-2$	$\sum (y_i - \hat{y}_i)^2$	MSE	
Total corrected	$n-1$	$\sum (y_i - \bar{y})^2$		

15/04/16 ASMI SAMPLE PAPER

(Q1) A. $y_i = \mu + \epsilon_i \quad i=1, \dots, N$
 $Q(\mu) = \sum (y_i - \mu)^2 \quad \text{for SLR} \quad \sum (y_i - \mu - \beta_0 - \beta_1 x_i)^2$

B LS estimator of μ

$$\frac{dQ}{d\mu} = \sum_{i=1}^N (-2)(y_i - \mu)$$

$$= -2 \sum (y_i - \mu) = 0$$

$$\sum y_i - n\mu = 0$$

$$n\mu = \sum y_i$$

$$\mu = \bar{y} / N$$

$$\mu = \bar{y}$$

C. SSE: Sum of Squared Errors = $\sum (y_i - \bar{y}_i)^2$ deviations around fitted mean

Usually SSE = $\sum (y_i - \bar{y})^2$ where $y_i = \mu + \epsilon_i \quad \bar{y}_i = \mu \quad \bar{y} = \bar{y}$

MSE = SSE / d.f. Here we will have $n-1$ d.f. because we have only estimated one parameter μ of this model

D. Model with a "day effect"

→ n observations from the first day

→ m observations from the second day

"Day Effect" introduce a variable $x_i = \begin{cases} 0 & 1 \leq i \leq n \\ 1 & n+1 \leq i \leq n+m \end{cases}$

Use the model: $y_i = \mu + \mu_d x_i + \epsilon_i \quad 1 \leq i \leq n+m$

First day: $y_i = \mu + \epsilon_i \quad i=1, \dots, n$

Second day: $y_i = \mu + \mu_d + \epsilon_i \quad n+1 \leq i \leq n+m$

- There will be a day effect if $\mu_d \neq 0$

- We want to test $H_0: \mu_d = 0 \quad v \quad H_1: \mu_d \neq 0$ for a day effect

$$\Sigma = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & n+m & \dots & n+m \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} \quad \text{SVD INVERSE}$$

-H₀ can be carried out using a t-test

$$t = \frac{\bar{D}_{1-0}}{SE(\bar{D}_{1-0})}$$

-All quantities from LS SUR model - compare $|t|$ with $t_{n+m-2, \alpha/2}$ and reject H₀ if
 $|t| > t_{n+m-2, \alpha/2}$

15/04/16 ALSM 1

Exam Questions

Q2 A $E[y_i] = \beta_0 + \beta_1 x_i$

\hat{y}_i at x_i

\hat{y}_i is the predicted mean value of y at x_i

\hat{y}_i is a predictor of $E[y|x=x_i] = \beta_0 + \beta_1 x_i$

It's given by $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$

B. Show $\sum (y_i - \bar{y})^2 = \sum (y_i - \hat{y}_i)^2 + \sum (y_i - \bar{y})^2$

$SS(\text{corrected}) = SS(\text{Reg}) + SSE$

$$\begin{aligned}\sum (y_i - \bar{y})^2 &= \sum (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\ &= \sum ((y_i - \hat{y}_i)^2 + (\hat{y}_i - \bar{y})^2 + 2(y_i - \hat{y}_i)(\hat{y}_i - \bar{y})) \\ &= \sum (y_i - \hat{y}_i)^2 + \sum (\hat{y}_i - \bar{y})^2 + 2 \sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y})\end{aligned}$$

\uparrow need to show that this = 0

$$\begin{aligned}\sum (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum (y_i - b_0 - b_1 x_i)(b_0 + b_1 x_i - \bar{y}) \\ &= \sum (y_i - \bar{y} + b_1 \bar{x} - b_1 x_i)(\bar{y} - b_1 \bar{x} + b_1 x_i - \bar{y}) \\ &= \sum [(y_i - \bar{y}) - b_1(x_i - \bar{x})][b_1(x_i - \bar{x})^2] \\ &= b_1 \sum (x_i - \bar{x})(y_i - \bar{y}) - b_1^2 \sum (x_i - \bar{x})^2 \\ &= b_1 S_{xy} - b_1^2 S_{xx} = \frac{S_{xy} S_{yy}}{S_{xx}} - \frac{S_{xy}^2}{S_{xx}} = 0\end{aligned}$$

$$\Rightarrow \sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$$

C $SS(\text{corrected})$ has $n-1$ d.f. since this is like the SSE of a simple mean model

$SSE = \sum (y_i - b_0 - b_1 x_i)^2$ has $n-2$ d.f., one for each b_0, b_1

The remaining 1 d.f. is for $SS(\text{Reg})$

D Source	on SS	D.F.	MS	F
Regression	$\sum (y_i - \bar{y})^2$	1	$\frac{S_{xy} S_{yy}}{S_{xx}}$	$\frac{SS(\text{Reg})}{SS(\text{Error})} MSe$
Error	$\sum (y_i - \hat{y}_i)^2$	$n-2$	$\frac{S_{yy}}{n-2}$	
Total	$\sum (y_i - \bar{y})^2$	$n-1$		

Ques: Why 1 D.F. for $SS(\text{Res})$?

$$\begin{aligned} \sum (y_i - \bar{y})^2 &= \sum (b_0 + b_1 x_i - \bar{y})^2 \\ &= \sum (y - b_0 \bar{x} + b_1 x - \bar{y})^2 \\ &= b_1^2 \sum (x_i - \bar{x})^2 = S_{xx} b_1^2 \end{aligned}$$

↳ only comes as an end product of other processes and therefore is the single amalgamation of other processes which is why it has 1 d.f.

15/04/16 ALSM 1

EXAM QUESTIONS

Q3 A. One way classification model

Design Matrix

$$J \left\{ \begin{bmatrix} 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ \vdots & 0 & \ddots & \dots \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \right. \left. \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_E \end{bmatrix} \right\} J$$

$$X \quad \mu$$

LS estimates of μ_1, \dots, μ_E

$$Q(\mu_1, \dots, \mu_E) = \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \mu_i)^2$$

$$\frac{\partial Q}{\partial \mu_i} = \sum_{j=1}^J (-2)(y_{ij} - \mu_i) = 0$$

$$\sum_{j=1}^J y_{ij} - J\mu_i = 0$$

$$\bar{\mu}_i = \sum_{j=1}^J y_{ij}/J$$

$$= \bar{y}_{i.}$$

ANOVA for Model?

$$SS(\text{corr}) \text{ for model? } \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.})^2 = \bar{y}_{i.}$$

$$SS(\text{err}) = \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.})^2.$$

Within a population LS estimate of μ_i is $\bar{y}_{i.}$

$$\begin{aligned} \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.})^2 &= \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.} + \bar{y}_{i.} - \bar{y}_{..})^2 \\ &= \sum_{i=1}^I \sum_{j=1}^J [(y_{ij} - \bar{y}_{i.})^2 + (\bar{y}_{i.} - \bar{y}_{..})^2 + \underbrace{2(y_{ij} - \bar{y}_{i.})(\bar{y}_{i.} - \bar{y}_{..})}_{=0}] \\ &= \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i.})^2 + I \sum_{i=1}^I (\bar{y}_{i.} - \bar{y}_{..})^2 \end{aligned}$$

B. $H_0: \mu_i = 0$ I different intercepts, one for each population
 $H_1: \mu_1 - \mu_2 = \mu_k$ implied that $\mu_1 - \mu_2 = 0$ for $i \neq k$

Write L as

$$\begin{bmatrix} 1 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & 0 & -1 \\ \vdots & & & \ddots \end{bmatrix}$$

C. An F-test

$$F = \frac{MS(R_{\text{adj}})}{MSE} = \frac{R(\beta_1, \beta_2 | R) / 2(J-1)}{MSE}$$

If H_0 is true, i.e. $\mu_1 = \mu_2 = 0$ then F follows an F-distribution with I and $2(J-1)$ D.F.

If $F > F_{\alpha/2, 2(J-1)}$ then reject H_0 at $\alpha/2(1-\alpha)$ y. significance level.

15/04/16

ALSM 1

EXAM QUESTIONS

Q4 A

 y $n \times 1$ response X $n \times (p+1)$ Design Matrix β $(p+1) \times 1$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

$$\mathbb{E}[y] = X\beta \quad \text{or} \quad y = X\beta + \varepsilon$$

(where ε is a $n \times 1$ error vector with 0 mean and uncorrelated entries)B The Hat matrix is such that $\hat{y} = Hy$ Here $\hat{y} = x\hat{\beta}$

$$H = X(X^T X)^{-1} X^T y$$

$$H = X(X^T X)^{-1} X^T$$

C Vector of residuals

$$\hat{\varepsilon} = y - \hat{y} = y - H\beta$$

$$= y - x\hat{\beta}$$

$$D SSE = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \hat{\varepsilon}^T \hat{\varepsilon}$$

$$= (y - x\hat{\beta})^T (y - x\hat{\beta})$$

$$= y^T y - y^T x\hat{\beta} - \hat{\beta}^T X^T y + \hat{\beta}^T X^T x\hat{\beta}$$

$$= y^T y - y^T H y - \hat{\beta}^T H^T y + y^T H^T H y \quad H^T H = H$$

$$= y^T y + y^T H y = y^T (I - H) y$$

$$H = X(X^T X)^{-1} X^T$$

$$H^T = (X^T)^{-1} (X^T)^T X^T = X(X^T X)^{-1} = H \Rightarrow \text{Symmetric Matrix}$$

$$H^T H = HH = X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T$$

$$= X(X^T X)^{-1} X^T = H$$

Q5 Algae growth vs concentration of chemical X

P value for $X=0.15$

This is the p-value for testing the hypothesis, size of 10 i.e. $[H_0: \beta=0, H_1: \beta \neq 0]$

Concentration decreases algae growth?

Correlation of sample, correlation between X and Y, of 0.5, large?

Correlation suggests Y and X are positively related, is this actually correct conclusion?

What statistical software are you using?

Does it have default settings for hypothesis testing? i.e. SLR?

- Assuming the default test (relating to X above, $p=0.15$) is whether the slope is equal to zero, a p-value of 0.15 may suggest you make tentative conclusions without further investigation.

- You should carry out regression diagnostics, have you considered standardized residuals or normality plots?

17/05/16

ALSM 1 EXAM NOTES : SLR, ANOVA

SLR

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \quad E[y_i | x_i] = \beta_0 + \beta_1 x_i$$

1. For a fixed value of x , y is a random variable with a finite mean and variance:
 $|E[y_i | x_i]| < \infty \quad \text{Var}[y_i | x_i] < \infty$ (i.e. they exist)

2. The values of y are uncorrelated

3. The model is linear; i.e. $E[y_i | x_i] = \beta_0 + \beta_1 x_i$ $N_{y|x} = \beta_0 + \beta_1 x_i$

4. Homoscedasticity - Variance of y is not dependent on x , i.e. variance is constant

$$\text{Var}[y_i | x_i] = \text{Var}[y_i] = \sigma^2$$

Statement of model:

- Data $(x_1, y_1), \dots, (x_n, y_n)$
- $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$
- $E[\varepsilon_i] = 0 \quad \forall i$
- $\text{Var}[\varepsilon_i] = \sigma^2 \quad \forall i$ (Homoscedastic)
- $\text{Cov}[\varepsilon_i, \varepsilon_j] = 0 \quad \text{if } i \neq j$ y_i 's uncorrelated

Least Squares

$$SS_{\text{Error}} = SSE = \sum \varepsilon_i^2 = \sum (y_i - \hat{y}_i)^2 = Q(\hat{\beta}_0, \hat{\beta}_1)$$

$$\frac{\partial SSE}{\partial \beta_0} = -2 \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\frac{\partial SSE}{\partial \beta_1} = -2 \sum (x_i) (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) \quad \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

$$S_{xy} = \sum (x_i)(y_i - \bar{y}) / \sum (x_i - \bar{x})(y_i - \bar{y}) / \sum (y_i)(x_i - \bar{x}) / \sum x_i y_i - \frac{\sum x_i \sum y_i}{n}$$

$$S_{xx} = \sum (x_i)(x_i - \bar{x}) / \sum (x_i - \bar{x})(x_i - \bar{x}) / \sum (x_i)^2 - \bar{x}^2$$

$$SSE = \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \text{ can be written: } S_{yy} - \frac{S_{xy}^2}{S_{xx}}$$

ANOVA

$$SS(\text{Uncorrected}) = \sum y_i^2$$

$$SS(\text{Corrected}) = \sum (y_i - \bar{y})^2$$

$$SS(\text{Corrected}) = SS(\text{Uncorrected}) - \text{Correction}$$

$$= \sum (y_i^2 - 2y_i\bar{y} + \bar{y}^2)$$

$$= \sum y_i^2 - 2n\bar{y}^2 + n\bar{y}^2$$

$$= \sum y_i^2 - n\bar{y}^2$$

$$SS(\text{Corrected}) = \text{Correction}$$

$$SS(\text{Uncorrected}) = \sum y_i^2 = \sum (g_i + \varepsilon_i)^2$$

$$= \sum g_i^2 + \sum \varepsilon_i^2 + 2 \sum g_i \varepsilon_i$$

$$= \sum y_i^2 + \sum \varepsilon_i^2 \rightarrow 0$$

$$SS(\text{Model}) = SSE$$

$$SS(\text{Model}) = \sum \hat{y}_i^2$$

$$SSE = \sum (y_i - \hat{y}_i)^2$$

$$SS(\text{Model}) = \sum \hat{y}_i^2 = \sum [y_i - (\bar{y} + \beta_0 + \beta_1 x)]^2$$

$$= \sum (y_i - \bar{y})^2 + 2\beta_0 \sum (y_i - \bar{y}) + n\bar{y}^2$$

$$\beta_0$$

$$= \sum (y_i - \bar{y})^2 + n\bar{y}^2$$

$$SS(\text{Regression}) + \text{Error}$$

$$S(\text{Regression}) = \sum (y_i - \hat{y}_i)^2$$

Source	SS	DF
SS(Uncorrected)	$\sum y_i^2$	n
SS(Corrected)	$\sum (y_i - \bar{y})^2 = \sum y_i^2 - n\bar{y}^2$	n-1
Correction	$n\bar{y}$	1
SS(Model) $P(\beta_0, \beta_1)$	$\sum (y_i - \bar{y})^2 + n\bar{y}^2$	2 (1 for each parameter) importance of β_0 and β_1
SS(Regression)	$\sum (\hat{y}_i - \bar{y})^2$	1 for β_1 importance of β_1 term in model
SSE	$\sum \varepsilon_i^2 = \sum (y_i - \hat{y}_i)^2$	n-2

Regression on X	1	$\sum (y_i - \hat{y}_i)^2$	MS(Reg)
Total Error	n-2	$\sum (y_i - \bar{y})^2$	MSF
Total (Uncorrected)	n-1	$\sum (y_i - \bar{y})^2$	

105/16

ALSM 1 EXAM NOTES: ANOVA

$$SS(\text{corrected}) = SS(\text{uncorrected}) - \text{correction}$$

$$R^2 = \frac{SS(\text{reg})}{SS(\text{corrected})}$$

$$SS(\text{Unadjusted}) = SS(\text{reg}) + SSE$$

$$SS(\text{Model}) = SS(\text{reg}) + \text{correction}$$

$$SS(\text{Uncorrected}) = SS(\text{reg}) + (\text{correction} + SSE)$$

$$SS(\text{Corrected}) = SS(\text{reg}) + SSE$$

$$\begin{aligned} \text{Show } \sum(y_i - \bar{y})^2 &= \sum(y_i - \bar{y} + \bar{y} - \bar{y})^2 + \sum(\bar{y} - \bar{y})^2 \\ &= \sum(y_i - \bar{y})^2 + (\bar{y} - \bar{y})^2 + 2(y_i - \bar{y})(\bar{y} - \bar{y}) \\ &= \sum(y_i - \bar{y})^2 + S(\bar{y} - \bar{y})^2 + 2S(y_i - \bar{y})(\bar{y} - \bar{y}) \end{aligned}$$

$$\begin{aligned} &\sum(y_i - \beta_0 - \beta_1 x_i)(\beta_0 + \beta_1 x_i - \bar{y}) \quad b_0 = \bar{y} - b_1 \bar{x} \\ &= \sum(y_i - \bar{y} + b_1 \bar{x} - \beta_1 x_i)(\bar{y} - b_1 \bar{x} + \beta_1 x_i - \bar{y}) \\ &= \sum(y_i - \bar{y})^2 + b_1^2 \sum(x_i - \bar{x})^2 \\ &- b_1 \sum(y_i - \bar{y})(x_i - \bar{x}) + b_1^2 \sum(x_i - \bar{x})^2 \\ &\stackrel{\text{SSE}}{\sum} \frac{Sxy}{Sxx} \rightarrow \frac{Sxy^2}{Sxx^2} = 0 \end{aligned}$$

$$= \sum(y_i - \bar{y})^2 + S(\bar{y} - \bar{y})^2 \checkmark$$

Properties

$$\begin{aligned} E[\beta_1] &= \beta_1 & \text{Var}[\beta_1] &= \sigma^2 / S_{xx} & SE(\beta_1) &= \sqrt{\frac{MSE}{S_{xx}}} \\ E[\beta_0] &= \beta_0 & \text{Var}[\beta_0] &= \sigma^2 \left(\frac{1}{n} + \frac{S_{xx}}{S_{xx}^2} \right) & SE(\beta_0) &= \sqrt{MSE \left(\frac{1}{n} + \frac{S_{xx}}{S_{xx}^2} \right)} \\ E[\beta_0, \beta_1] &= \frac{S_{xx}^2}{S_{xx}} \end{aligned}$$

$$\text{If p-value below 0.05} \Rightarrow \text{reject } H_0 \quad P[-t_{n-2, 1-\alpha/2} \leq t_{n-2} = t_{n-2, \alpha/2}] = 1 - \alpha$$

$$(1) \text{ For } \sigma^2? \quad \frac{(n-2)MSE}{S_{xx}^2} \sim \chi^2_{n-2} \text{ via this } (\chi^2 \text{ has a symmetric distribution})$$

$$\hat{\sigma}^2 = \frac{(n-2)MSE}{2(n-2, \alpha/2)}$$

$$\hat{\sigma}^2 = \frac{(n-2)MSE}{2(n-2, \alpha/2)} \quad \text{Is it introduced?} \quad \checkmark$$

F-test

für direktionale

$$\frac{F = \frac{SS(\text{Reg})/\sigma^2}{1}}{\frac{MSE(n-2)/\sigma^2}{n-2}} = \frac{SS(\text{Reg})}{MSE} \sim F_{1, n-2}$$

Wann H₀ ist LVR ($\beta_1 = 0$) $SS(\text{Reg})/\sigma^2 \sim \chi^2$

If F-value < F-tabelle REJECT H₀

Rohr J. für das Chi-Squared Test

$$F \approx t^2$$

INFERENZ CONCERNING $E[y|x]$

Point estimator: $\hat{y}' = \hat{\beta}_0 + \hat{\beta}_1 x'$ (mein vorher)

$$\text{Var}[y'] = \sigma^2 \left(1 + \frac{(x'-\bar{x})^2}{\sum x_i^2} \right) \quad (\text{Intervall im CI})$$

Prediction

$y - \hat{y}'$ besteht zwischen Prediktiv und actual

$$E[y - \hat{y}'] = 0$$

$$\text{Var}[y - \hat{y}'] = \sigma^2 \left[1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum x_i^2} \right]$$

$$\Rightarrow \frac{(y - \hat{y}') - 0}{\sqrt{MSE \left(1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum x_i^2} \right)}} \sim t_{n-2}$$

18/05/16 ALSM1 EXAM NOTES: PROBLEM SHEET

$$y_i = \mu + \epsilon_i$$

$$(Q(\mu)) = \sum (y_i - \mu)^2$$

$$L.S. = \sum 2(y_i - \mu)$$

$$\sum y_i - n\mu = 0$$

$$\mu = \bar{y}$$

SSE = sum of squared errors = $\sum (y_i - \bar{y})^2 = \sum (y_i - \mu)^2$ which is the contribution of square

MSE = SSE / DF $DF = n-1$ because we have only one fixed parameter

$$\text{Day } t \text{ fresh. } n \text{ first day, } m \text{ second day} \\ \text{include } n \text{ intercepts } \alpha_0, \dots, \alpha_n = \begin{cases} 0 & 1 \leq i \leq n \\ 1 & n+1 \leq i \leq m \end{cases}$$

$$\text{We make } y_i = \alpha_0 + \alpha_1 x_1 + \epsilon_i \quad 1 \leq i \leq m$$

$$\text{First day: } y_i = \mu + \epsilon_i \quad 1 \leq i \leq n \quad \text{Second day: } y_i = \mu + \alpha_0 x_1 + \epsilon_i \quad n+1 \leq i \leq m$$

- trend to left $\alpha_0 = 0$ v $\alpha_0 \neq 0$ for a day effect

- use t-test $t = \bar{y} / \sqrt{SSE/m}$

- compare $|t|$ with $t_{m-2, \alpha/2}$ and decide if $|t| > t_{m-2, \alpha/2}$. accept H₀

$$\begin{aligned} E[\text{MSE}(\text{reg})] &= E[\sum (y_i - \bar{y})^2] = E[\sum (\beta_i^2/m)] \\ &= \sum \alpha_i^2 E[\beta_i^2] \\ &= \sum \alpha_i^2 [Var(\beta_i) + E[\beta_i]^2] \\ &= \sum \alpha_i^2 [\sigma^2/m + \beta_i^2] \\ &= \sigma^2 + \hat{\beta}_1^2 S_N \end{aligned}$$

$$\begin{aligned} Q3. \quad y_i &= \alpha_0 + \alpha_1 (x_i - \bar{x}) + \epsilon_i \\ &= \alpha_0 + \alpha_1 x_i - \alpha_1 \bar{x} \\ \hat{\beta}_1 &= \alpha_0 - \alpha_1 \bar{x} \quad \hat{\alpha}_1 = \alpha_1 \end{aligned}$$

$$\begin{aligned} L) \quad \sum (y_i - \alpha_0 - \alpha_1 x_i) \\ \frac{dL}{d\alpha_0} &= \sum 2(y_i - \alpha_0 - \alpha_1 x_i) = -2 \sum (y_i - \alpha_0 - \alpha_1 x_i) \\ \sum y_i - n\alpha_0 - \alpha_1 \sum x_i &= 0 \\ \alpha_0 &= \bar{y} \end{aligned}$$

$$\frac{dL}{d\alpha_1} = -2 \sum (y_i - \alpha_0 - \alpha_1 x_i) x_i =$$

$$\sum (y_i - \alpha_0 - \alpha_1 x_i) x_i^2 =$$

$$\alpha_1 = S_{xx}^{-1} S_{xy}$$

$$\begin{aligned} \text{Cov}[y_i, \frac{\partial \eta_j}{\partial x_i}] &= \frac{1}{N-1} \text{Cov}[y_i, \varepsilon(x) \cdot \eta_j] \\ &\rightarrow \frac{1}{N-1} \text{Cov}[y_i, (x_0 + b_0 x_i + (x_i - \bar{x}) \eta_j)] \\ &= \frac{1}{N-1} \text{Cov}[(x_i - \bar{x}) \eta_j] \\ &= \frac{\sigma^2 (x_i - \bar{x})}{N-1} \end{aligned}$$

181
182: OS (x_{0i}, y_{0i}) first day (x_{ni}, y_{ni}) , (x_{ni}, y_{ni}) second day
be carry day is day before

$$\text{Day 1: } a_0 + b_0 x_i \quad i = 1, \dots, n$$

$$\text{Day 2: } a_2 + b_2 x_i \quad i = 1, \dots, n$$

$$Q(a_0, a_1, b_1) = \sum y_i^2 = \sum_{i=1}^n \hat{y}_i^2 + \sum_{i=n+1}^m \hat{e}_i^2$$

$$\sum (y_i - a_0 - b_0 x_i)^2 + \sum (y_i - a_1 - b_1 x_i)^2$$

$$\frac{\partial Q}{\partial a_0} = -2 \sum (y_i - a_0 - b_0 x_i)^2 \quad (1) \quad \frac{\partial Q}{\partial a_1} = -2 \sum (y_i - a_1 - b_1 x_i)^2 \quad (2)$$

$$\frac{\partial Q}{\partial b_0} = -2 \sum (x_i)(y_i - a_0 - b_0 x_i) - 2 \sum (x_i)(y_i - a_1 - b_1 x_i) \quad (3)$$

$$0 = \sum y_i - n a_0 - \beta \sum x_i = 0$$

$$a_0 = \bar{y}_1 - \beta \bar{x}_1 \quad \Rightarrow \quad a_1 = \bar{y}_2 - \beta \bar{x}_2$$

$$(1) \quad \sum x_i y_i - a_0 \sum x_i - \beta \sum x_i^2 + \sum x_i (y_i - a_1 - \beta x_i) - \beta \sum x_i^2 = 0$$

$$\sum x_i y_i - (\bar{y}_1 - \beta \bar{x}_1) \sum x_i - \beta \sum x_i^2 + \sum x_i (y_i - \bar{y}_2 - \beta \bar{x}_2) \sum x_i - \beta \sum x_i^2 = 0$$

$$\sum x_i y_i - n_1 \bar{x}_1 \bar{y}_1 - \beta (\sum x_i^2 - n \bar{x}_1^2) + \sum x_i y_i - n_2 \bar{x}_2 \bar{y}_2 - \beta (\sum x_i^2 - n \bar{x}_2^2) = 0$$

$$Sxy^1 - \beta Sx^1 + Sxy^2 - \beta Sx^2 = 0$$

$$\beta (Sxx^1 + Sx^2) = Sxy^1 + Sxy^2$$

$$\beta = \frac{Sxy^1 + Sxy^2}{Sxx^1 + Sx^2}$$

18/05/16

ALG M1 EXAM NOTES: PROBLEM SHEET

$$\begin{aligned} \text{Var}[\beta] &= \text{Var} \left[w_1 \beta_1 + w_2 \beta_2 \right] \\ &= \frac{1}{(w_1 + w_2)^2} \text{Var} [w_1 \beta_1 + w_2 \beta_2] \\ &= \frac{1}{(w_1 + w_2)^2} [w_1^2 \text{Var}[\beta_1] + w_2^2 \text{Var}[\beta_2] + 2w_1 w_2 \text{Cov}[\beta_1, \beta_2]] \end{aligned}$$

$$\begin{aligned} \text{Var}[\beta_1] &= \sigma^2 s_{xx}^{-2} & \text{Var}[\beta_2] &= \sigma^2 s_{xx}^{-2} \\ \text{Cov}[\beta_1, \beta_2] &= \text{Cov} \left[\frac{\sum d_i y_i}{n}, \frac{\sum d_i y_i}{n} \right] \quad \text{all Cov will be 0} \Rightarrow \text{Diagonal} \\ \text{Var}[\beta] &= \frac{1}{w_1^2 + w_2^2} \left[\frac{(w_1)^2 \sigma^2}{s_{xx}^{-2}} + \frac{(w_2)^2 \sigma^2}{s_{xx}^{-2}} \right] \end{aligned}$$

ALSM1 EXAM NOTES: MATRIX FORMULATION

$$\text{Var}[e^T x] = \sigma^2 \leq \sigma^2$$

$$y = x\beta + e$$

Residual vector, Design matrix, Coefficient, Error vector
 $E[y] = x\beta$ $E[e] = 0$ (zero vector)

$$\text{Var}[y] = \text{Var}[e] = \begin{bmatrix} \sigma^2 & & \\ & \ddots & \\ & & \sigma^2 \end{bmatrix} = \sigma^2 I \quad \text{Independent} \neq 0 \text{ in off diagonals}$$

↑ identity matrix

Multiple Regression

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ 1 & x_{21} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$n \times 1$ $n \times (p+1)$ $(p+1) \times 1$ $n \times 1$

Least Squares

$$\text{Minimize } SSE = \sum e^T e$$

$$\frac{\partial SSE}{\partial \beta} = (y - x\beta)^T (y - x\beta) = y^T y - y^T x\beta - y x^T \beta + x^T x\beta$$
$$= 0 - x^T y - x^T y + 2x^T x\beta = 0$$
$$x^T x\beta = x^T y$$

Assuming $x^T x$ is non singular (i.e. invertible) then $(x^T x)^{-1} x^T y = (x^T x)^{-1} x^T y$

$$\hat{\beta} = (x^T x)^{-1} x^T y$$

Predicted mean of y is $\hat{y} = x\hat{\beta} = x(x^T x)^{-1} x^T y = Hy$ $H = x(x^T x)^{-1} x^T$

SSE

$$\begin{aligned} SSE &= \sum e^T e = (y - x\hat{\beta})^T (y - x\hat{\beta}) \\ &= y^T y - y^T x\hat{\beta} - \hat{\beta}^T y^T x + \hat{\beta}^T x^T x\hat{\beta} \\ &= y^T y - \hat{\beta}^T x^T y - \hat{\beta}^T y^T x + \hat{\beta}^T x^T x\hat{\beta} \\ &= y^T y - 2\hat{\beta}^T x^T y + \hat{\beta}^T x^T x[(x^T x)^{-1} x^T y] \\ &= y^T y - 2\hat{\beta}^T x^T y + \hat{\beta}^T x^T y \\ &= y^T y - \hat{\beta}^T x^T y \\ &= y^T y - (x^T x)^{-1} x^T y^T x^T y = y^T (I - H)y \\ &= y^T (I - H)y \quad H = x(x^T x)^{-1} x^T \end{aligned}$$

$y^T y = \sum e^2$

18/05/16

$$\hat{\beta} = (x^T x)^{-1} x^T y$$

$$\text{E}[\hat{\beta}] = \text{E}[C(x^T x)^{-1} x^T y] = (x^T x)^{-1} x^T \text{E}[y]$$

$$= (x^T x)^{-1} x^T [A\beta]$$

$$= \beta \quad (\text{unbiased})$$

$$\text{Var}[\hat{\beta}] = \text{Var}\left[\frac{1}{n} (x^T x)^{-1} x^T y\right]$$

$$= A \text{Var}[y] A^T$$

$$= A \sigma^2 A^T$$

$$= \sigma^2 [x^T x]^T [x^T x]$$

$$= \sigma^2 [x^T x]$$

ANOVA

Source	DF	SS	MS
Regression	p	$\hat{\beta}^T y - \bar{y}^2$	$SSE(\hat{\beta})/p$
Residual	$n-p-1$	$(y - \hat{\beta})^T (y - \hat{\beta})$	$SSE(n-p-1)$
Total (Overall)	$n-1$	$y^T y - \bar{y}^2$	

 $H_0: \beta_0 = \beta_1 = \dots = \beta_p = 0$ vs $H_1: \beta_i \neq 0$ Need estimate for S.E. for $C(\beta - \bar{\beta})$ $\text{Var}[\beta - \bar{\beta}] = \text{Var}[\hat{\beta}] + \text{Var}[\beta] - 2 \text{Cov}[\hat{\beta}, \beta]$ used in one-way classification models

Confidence Intervals

 $\hat{\beta}$ is a linear estimator hence $\hat{\beta} \sim N_{p+1}(\beta, \sigma^2 (x^T x))$

$$\hat{\beta}_j \sim N(\beta_j, c_{jj}, \sigma^2)$$

Where c_{jj} is the $j+1^{\text{th}}$ diagonal entry of $(x^T x)^{-1}$ for $j=0, 1, \dots, p$

$$\Rightarrow \frac{\hat{\beta}_j - \beta_j}{\sqrt{(\text{MSE})(c_{jj})}} \sim t_{n-p-1} \quad \pm \hat{\beta}_j t_{n-p-1, \alpha/2} \sqrt{(\text{MSE})(c_{jj})}$$

Confidence

$$\text{Define } X_0 = \begin{bmatrix} 1 \\ x_0 \end{bmatrix} \times (p+1) \times 1 \text{ column vector}$$

The mean value of y at this point x_0 is $x_0^T \beta$ which is estimated by

$$x_0^T \hat{\beta} = \hat{y}_0$$

18/05/16 ALSM 1 EXAM NOTES: MATRIX FORMULATION

$$\mathbb{E}[g] = \mathbb{E}[x_0^\top \beta] = x_0^\top \beta = \mu_0$$

$$\text{Var}[g] = x_0^\top \text{Var}[\beta] x_0$$

$$= \sigma^2 x_0^\top (X^\top X)^{-1} x_0$$

$$g \sim N(\mu_0, \sigma^2 x_0^\top (X^\top X)^{-1} x_0) \quad \text{Replace } \sigma^2 \text{ by MSE}$$

CE for μ_0 : $\frac{g - \mu_0}{\sqrt{\text{MSE}(x_0^\top (X^\top X)^{-1} x_0)}} \sim N(0, 1)$

$$\frac{g_0 - \mu_0}{\sqrt{\text{MSE}(x_0^\top (X^\top X)^{-1} x_0)}} \sim t_{n-p-1} \text{ Distribution}$$

$$\text{MSE}(1-\alpha) \geq (1-\alpha) \pm t_{n-p-1, \alpha/2} \text{ MIE } x_0^\top (X^\top X)^{-1} x_0$$

Prediction

Pi for a new observation y_0 is the value of a future observation x_0 , estimated by

$$g_0 = x_0^\top \hat{\beta} \text{ and } \alpha \text{ PI is}$$

$$g_0 \pm t_{n-p-1, \alpha/2} \left(\sqrt{\text{MSE}(1 + x_0^\top (X^\top X)^{-1} x_0)} \right) = 1-\alpha$$

Hypothesis Testing

Overall test of Significance given by:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \quad \text{vs} \quad H_1: \beta_j \neq 0 \text{ for certain one.}$$

$$\text{Use } F\text{-test} = \text{MS}(\text{Reg})/\text{MSE}$$

If $F < F_{p, n-p-1, \alpha/2}$ accept H_0

Individual test

$$H_0: \beta_j = 0 \quad \Rightarrow \quad H_1: \beta_j \neq 0 \quad t\text{-test} \quad \sqrt{\frac{\hat{\beta}_j}{\text{MSE}_{jj}}} \quad \text{with } n-p-1 \text{ df}$$

Partial F-Test Extra Sum of Squares

- Create full model (full terms) and reduced model (remove β_j)

- $R(\beta_j) \beta_0, \beta_1, \beta_2, \dots, \beta_{j-1}, \beta_{j+1}, \dots, \beta_p$ = Extra sum of squares

- ESS is the partial SS for X_j and represents the contribution of X_j , adjusted for all other independent variables in the model.

$$F = \frac{R(\beta_1, \dots, \beta_p | \beta_0)}{MSE} \quad DF = 1, n-p-1$$

is the test statistic for $H_0: \beta_1 = \dots = \beta_p = 0$

Called LR partial F-test for β_0

Test for a subset of Regression Coefficients

Reduced model: $E[y_i] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_g x_{ig} \quad g < p$

Full Model: $E[y_i] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$

$$H_0: \beta_{g+1} = \beta_{g+2} = \dots = \beta_p = 0$$

H_a : At least one of those β 's is $\neq 0$

$R(\beta_0, \dots, \beta_p | \beta_0, \dots, \beta_g) = ESS$ (How much more info you explain by including term in model)

Use F-test: $\frac{R(\beta_0, \dots, \beta_p | \beta_0, \dots, \beta_g)}{MSE} / (p-g)$

MSE

If H_0 is true, then F follows $F_{p-g, n-p-1}$ Distribution

Source	DF	Sums	DF	
$R(\beta_0, \dots, \beta_g \beta_0)$	g	or	$R(\beta_0 \beta_0)$	1
$R(\beta_0, \dots, \beta_p \beta_0)$	$n-g$		$R(\beta_1, \dots, \beta_p \beta_0)$	1
$R(\beta_1, \dots, \beta_p \beta_0)$	p		$R(\beta_0, \dots, \beta_p \beta_0)$	1
Residual	$n-p-1$	$R(\beta_0, \dots, \beta_p \beta_0)$	$n-p-1$	
Total	$n-1$	Total	$n-1$	

Sequential SS - Add one predictor at a time and find reduction in SS_{res} not contained in x_1, \dots, x_{g-1}

Partial SS - Find info contained in x_g which is not contained in x_k for $k \neq g$.

$$R(A|P) + R(P|D, A) + R(B_1|P, A, B_2) = R(B_1, B_2, B_3 | P) = SS(\text{Reg})$$

General Linear Hypothesis

$$H_0: L\beta = \underline{c} \quad H_a: L\beta \neq \underline{c}$$

L : $R \times (p+1)$ matrix of coefficients

\underline{c} : $(p+1) \times 1$ vector of constants

MATRIX CALCS

18/05/16 AL5M1 EXAM NOTES: MATRIX FORMULATION

e.g. If $\beta_2 = 0$ & $H \beta_2 \neq 0$
 $\begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} 0 \end{bmatrix}$

b ≠ s

$H_0: \beta_1 = \beta_3 = 0$ $\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

b ≠ s

The extra sum of squares method can be used to test $H_0: L\beta = s$

$F = \frac{\text{ExRSS}/k}{\text{RSS}/(n-p-1)}$

PROBLEM SET 4

18/05/16

ALSM 1 : EXAM NOTES: PROBLEM SHEETS

$$PS 1 \quad X^T X = \begin{bmatrix} 1 & \dots & 1 \\ x_{11} & \dots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nn} \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \\ x_{11} & \dots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{n1} & \dots & x_{nn} \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$$|X^T X| = ad - bc = n \sum x_i^2 - (\sum x_i)^2 = n \lambda_{XX}$$

$$(X^T X)^{-1} = \frac{1}{n \lambda_{XX}} X^T X = \frac{1}{n \lambda_{XX}} \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$$Q3 \quad E[y_t] = x_t^T \beta_1 \quad t=1, \dots, r$$

$$E[y_t] = x_t^T \beta_2 \quad t=r+1, \dots, n$$

$$Q(\beta_1, \beta_2) \quad y_1 = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \quad x_1 = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \end{bmatrix} \quad \text{similarly } y_2 \text{ and } x_2$$

$$\alpha: (y_1, x_1, \beta_1)^T (y_2 - x_2 \beta_2) + (y_2 - x_2 \beta_2)^T (y_1 - x_1 \beta_1)$$

$$\text{LS estimates: } \hat{\beta}_1 = (x_1^T x_1)^{-1} x_1^T y_1,$$

$$\hat{\beta}_2 = (x_2^T x_2)^{-1} x_2^T y_2$$

$$\hat{\sigma}^2 = \text{partial MSE} = \frac{Q(\beta_1, \beta_2)}{n-2} \quad \text{(sses)}$$

B. If variances can be different, would have to use a weighted LSE

$\hat{\beta}_1, \hat{\beta}_2$ are OS bfr

$$\hat{\sigma}^2 = \frac{Q(\beta)}{r-2} \quad \text{use MSE for first segment} \quad \hat{\sigma}^2 = \frac{Q(\beta)}{n-r-2}$$

$$4. H^T H = X(X^T X)^{-1} X^T Y^T (X^T X)^{-1} X = H = X(X^T X)^{-1} X^T = X^T (X^T X)^{-1} X$$

$$5. SSE = \sum e_i^2 = (Y - \hat{Y})^T (Y - \hat{Y}) = (Y - HY)^T (Y - HY)$$

$$= Y^T Y - Y^T Y H - H^T Y^T Y + H^T H Y^T Y \quad H^T H = H$$

$$= Y^T Y - 2 Y^T Y H + Y^T Y H$$

$$= Y^T Y - Y^T Y H$$

$$= Y^T (I - H) Y$$

PS4

$$y_{ij} = \mu_i + \varepsilon_{ij} \quad i = 1, \dots, I \quad \text{I different populations}$$

$$j = 1, \dots, J \quad \text{equal sample size - called balanced design}$$

$$\begin{aligned} LS &= Q(\mu_1, \dots, \mu_I) = \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \mu_i)^2 \\ &= \sum_i (y_{i\cdot}^2 - 2\mu_i y_{i\cdot} + \mu_i^2) \\ &= \sum_i y_{i\cdot}^2 - 2 \sum_i \mu_i y_{i\cdot} + \sum_i \mu_i^2 \end{aligned}$$

$$\frac{\partial Q}{\partial \mu_i} = -2 \sum_j y_{ij} + 2J\mu_i = 0$$

$$\mu_i = \bar{y}_{i\cdot} = \bar{y}_{i\cdot} \quad \text{Sample mean for } i^{\text{th}} \text{ population}$$

dot notation, summing over J.

B Estimate I parameters $SSE = \sum_i (y_{i\cdot} - \bar{y}_{i\cdot})^2$ has DF

$$DF = n-I = IJ-I = I(J-1)$$

$$\sigma^2 = MSE = SSE / I(J-1)$$

C Design Matrix $y = X\mu + \varepsilon$

$$\begin{matrix} & \begin{matrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \\ & \vdots & & \vdots \\ & 1 & 0 & \dots & 0 \end{matrix} & \left[\begin{matrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_I \end{matrix} \right] & IJ \times I \text{ matrix} \end{matrix}$$

D $H_0: L\mu = 0$ vs $H_1: L\mu \neq 0$

$$H_0: \mu_1 - \mu_2 = \dots = \mu_I = 0 \quad \text{vs} \quad H_1: \text{Not } H_0$$

Implies that $\mu_i - \mu_j = 0$ for $i \neq j$

Write L as

$$\begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 1 & \dots & 0 \\ & & & \ddots \\ & & & 1 & -1 \end{bmatrix} \quad L \text{ is } (I-2) \times I \text{ matrix}$$

Use F-test. $\frac{MS(\text{Reg})}{MSE} \sim F_{I-2, I(I-1)}$ compare to F critical value, if smaller reject H_0

DEFINITION SET 4

18/05/16

ALSM1 : EXAM NOTES: EXAM QUESTIONS

$$Q3 \text{ A LS : } Q = \sum_{j=1}^J (y_{ij} - \mu_i)^2$$

$$\frac{\partial Q}{\partial \mu_i} = \sum_{j=1}^J (-2)(y_{ij} - \mu_i) = 0$$

$$\sum y_{ij} - J\bar{\mu}_i = 0$$

$$\bar{\mu}_i = \frac{1}{J} \sum y_{ij}$$

$$\bar{\mu}_i = \bar{y}_{i..} \quad (\text{summing over } j)$$

ANOVA for model?

$$SS(\text{corr}) = \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{i..})^2 = S_{\text{res}}$$

$$SS(\text{tot}) = \sum_{i=1}^I \sum_{j=1}^J (y_{ij} - \bar{y}_{..})^2$$

Within a population LS estimate of μ_i is $\bar{y}_{i..}$

$$\begin{aligned} \sum (y_{ij} - \bar{y}_{i..})^2 &= \sum (y_{ij} - \bar{y}_{i..} + \bar{y}_{i..} - \bar{y}_{..})^2 \\ &= \sum (y_{ij} - \bar{y}_{i..})^2 + (\bar{y}_{i..} - \bar{y}_{..})^2 + 2(y_{ij} - \bar{y}_{i..})(\bar{y}_{i..} - \bar{y}_{..}) \end{aligned}$$

\hookrightarrow

$$F \text{ test } F = MSE / \text{MSE} = R(\beta_1 | \beta_0) / I(J-1)$$

If H_0 is true $\mu_1 = \mu_2 = \dots$ then $F \sim F$ with I and $I(J-1)$ d.f.

If $F \geq F_{2,2}$ then reject H_0 at $\alpha(1-\alpha)/I$ significance level

$$\begin{aligned} \text{Var}[\hat{\beta}] &= \text{Var} \left[\frac{w_1 \beta_1 + w_2 \beta_2}{w_1 + w_2} \right] = \frac{1}{(w_1 + w_2)} \cdot \text{Var}[w_1 \beta_1 + w_2 \beta_2] \\ &= \frac{1}{(w_1 + w_2)^2} [w_1^2 \text{Var}[\beta_1] + w_2^2 \text{Var}[\beta_2] + 2w_1 w_2 \text{Cov}[\beta_1, \beta_2]] \\ \text{Var}[\mu_1] &= \sigma^2/n, \quad \text{Var}[\mu_2] = \sigma^2/n \end{aligned}$$

$$\text{Var}[u] = \frac{1}{(w_1 + w_2)^2} \left[\frac{(w_1)^2 \sigma^2}{n_1} + \frac{(w_2)^2 \sigma^2}{n_2} \right]$$

$$Q4 \quad \mathbf{y} = \mathbf{X}\beta + \varepsilon$$

y $n \times r$ vector of response

X $n \times (p+1)$ Matrix - design matrix

β $(p+1) \times 1$ vector of estimated parameters

$\hat{\beta} = (X^T X)^{-1} X^T y$ estimate of parameters

11/10/15

$$\mathbb{E}[y] = x\beta \quad \text{or} \quad y = x\beta + \varepsilon$$

where ε is a $n \times 1$ error vector with 0 mean and uncorrelated entries

B The hat matrix is such that $\hat{y} = Hy$. Here $\hat{y} = x\hat{\beta}$

$$\hat{y} = (x^T x)^{-1} x^T y$$

$$H = x(x^T x)^{-1} x^T$$

C Vector of residuals

$$\begin{aligned}\hat{\varepsilon} &= y - \hat{y} \\ &= y - x\hat{\beta}\end{aligned}$$

D $SSE = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \hat{\varepsilon}^T \hat{\varepsilon}$

$$= (y - x\beta)^T (y - x\beta)$$

$$= y^T y - y^T x\beta - x\beta^T y + x\beta^T x\beta$$

$$= y^T y - y^T Hy - y^T H\beta + y^T H\beta^T H \quad H^T H = H$$

$$= y^T y + y^T Hy$$

$$= y^T (I - H) y$$

$$H = x(x^T x)^{-1} x^T \quad H^T = x^T (x^T x)^{-1} x \quad (\text{symmetric})$$

$$HTY = x^T (x^T x)^{-1} x^T x^T (x^T x)^{-1} x$$

$$= x^T (x^T x)^{-1} x^T = I \quad \checkmark \quad H^T H = H$$

ST3451-1



Coláiste na Trionóide, Baile Átha Cliath
Trinity College Dublin
Ollscoil Átha Cliath | The University of Dublin

Faculty of Engineering, Mathematics and Science
School of Computer Science & Statistics
Statistics

Trinity Term 2016

Sophister Mathematics

ST3451 : Applied linear statistical methods 1

19/05/2016

Sports Centre

9:30-11:30

Dr. Jason Wyse

Instructions to Candidates:

Full marks will be awarded for complete solutions to any three questions. Each full question is worth 20 marks, and marks for component parts are indicated in brackets. Non-programmable calculators are permitted. Some useful distributions and their properties are provided at the end of the exam.

Materials permitted for this examination:

Non-programmable calculators are permitted for this examination — please indicate the make and model of your calculator on each answer book used.

1. The simple linear regression model is

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, 2, \dots, n$$

where ε_i are the error terms.

(a) What assumptions are made about the error terms? [2 marks]

(b) Derive the least squares estimators $\hat{\beta}_0$ and $\hat{\beta}_1$ of β_0 and β_1 . [6 marks]

(c) Show that

$$\hat{\beta}_1 = \sum_{i=1}^n c_i Y_i \quad \hat{\beta}_0 = \sum_{i=1}^n d_i Y_i$$

giving the explicit expressions for the c_i, d_i . [6 marks]

(d) Using the results of part (c) find $\text{Var}\{\hat{\beta}_1\}$. [6 marks]

2. A one-way classification model assumes that J observations are taken from each of I normal populations, that is

$$Y_{ij} = \mu_i + \varepsilon_{ij} \quad (i = 1, 2, \dots, I; j = 1, 2, \dots, J)$$

where the ε_{ij} are iid $N(0, \sigma^2)$.

(a) Show that the least squares estimates of μ_1, \dots, μ_I are given by $\hat{\mu}_i = \bar{Y}_i = \sum_{j=1}^J Y_{ij}/J$ for $i = 1, \dots, I$. [5 marks]

(b) Explain why

$$\text{SSE} = \sum_{i=1}^I \sum_{j=1}^J (Y_{ij} - \bar{Y}_i)^2.$$

What are its associated degrees of freedom? Explain.

(c) Show that the SSE partitions into two sums of squares as [3 marks]

$$\text{SSE} = \sum_{i=1}^I \sum_{j=1}^J \frac{(Y_{ij} - \bar{y}_{..})^2}{\bar{y}_{..} - \bar{y}_i} + J \sum_{i=1}^I (\bar{Y}_i - \bar{Y}_{..})^2$$

where $\bar{Y}_{..} = \sum_{i=1}^I \sum_{j=1}^J Y_{ij}/(IJ)$ is the overall mean.

(d) Describe how one would carry out the test of hypothesis [7 marks]

$H_0: \mu_1 = \mu_2 = \dots = \mu_I = \mu$ versus $H_A: H_0$ not true.

[5 marks]

3. Consider a model where observations are time ordered

$$(x_1, Y_1), (x_2, Y_2), \dots, (x_t, Y_t), \dots, (x_n, Y_n),$$

where t denotes time. There is a changepoint at time $t = \tau$, so that up to time τ it is reasonable to assume that

$$E\{Y_t\} = \alpha_1 + \beta_1 x_t, \quad t = 1, \dots, \tau$$

while after time τ

$$E\{Y_t\} = \alpha_2 + \beta_2 x_t, \quad t = \tau + 1, \dots, n$$

such that $\alpha_1 \neq \alpha_2$ and $\beta_1 \neq \beta_2$.

- (a) Make a sketch indicating the main features of the model. [2 marks]
- (b) Write this model in the form $E\{\mathbf{Y}\} = \mathbf{X}\theta$ where $\theta = (\alpha_1, \beta_1, \alpha_2, \beta_2)^T$ is a column vector and \mathbf{X} is a design matrix. [6 marks]
- (c) Assuming a constant error variance σ^2 before and after the changepoint τ , suggest an estimator of σ^2 giving reasons for your answer. [6 marks]
- (d) Suggest an approach for determining an estimate of τ if it is not known and must be estimated. [6 marks]

4. Consider the matrix formulation of the simple linear regression model

$E\{\mathbf{Y}\} = \mathbf{X}\beta$, where $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ is an $n \times 1$ column vector of responses, $\beta = (\beta_0, \beta_1)^T$, and \mathbf{X} is the $n \times 2$ design matrix with row i equal to $(1, x_i)$.

- (a) Find $\mathbf{X}^T \mathbf{X}$, $|\mathbf{X}^T \mathbf{X}|$ and $(\mathbf{X}^T \mathbf{X})^{-1}$ in terms of the x_i . [5 marks]
- (b) Compute the least squares estimator $\hat{\beta}$ of β in terms of the x_i, Y_i . You may assume $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$. [5 marks]
- (c) Write down an expression for the vector of model residuals $\hat{\epsilon}$ in terms of \mathbf{Y}, \mathbf{X} and $\hat{\beta}$. [5 marks]
- (d) Show that sum of squared errors (SSE) can be written as

$$\text{SSE} = \mathbf{Y}^T [\mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T] \mathbf{Y}$$

[5 marks]

5. A practitioner (who has little experience with statistics) has asked for a consult to explain some of the terminology associated with regression analysis and regression model diagnostics. They send a list of questions by email before the meeting in order to help you (as the consultant) prepare. The questions follow.

1. What is a studentized residual? I think I understand what a residual is, but I'm not sure if this is the same as a studentized residual.
2. I obtained a 95% confidence interval for b_1 (slope). Does this mean that 95% of my response variable observations will be between these two numbers?
3. What is heteroscedasticity and should I be worried about it?
4. I know the normal curve is the bell curve, but I can't see how this is related to my analysis. How would it be related to my analysis?

Write notes addressing the questions in the email, explaining how you would answer each of them in a non-technical manner (remember, the practitioner is not statistically literate).

[20 marks]