

1. Mb/1 Let $F = \hat{\beta}_1^2 \sum (x_i - \bar{x})^2 / s^2$, the F-statistic for testing $H: \beta_1 = 0$ for a straight line.

using the notation of Section 6.1.5, prove that $\tilde{x}_0 - \bar{x} = \frac{F}{F + (n-2)} (\hat{x}_0 - \bar{x})$

$$\text{Pf: } \hat{x}_0 = \frac{y_0 - \hat{\beta}_0}{\hat{\beta}_1} = \bar{x} + \frac{y_0 - \bar{y}}{\hat{\beta}_1} \quad , \quad \hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

\tilde{x}_0 : inverse estimate $X \sim d_0 + \gamma_{d_1} + \varepsilon$

$$\text{then: } \tilde{d}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (y_i - \bar{y})^2}, \quad \tilde{d}_0 = \bar{x} - \bar{y} \cdot \tilde{d}_1$$

$$\Rightarrow \tilde{x}_0 = \tilde{d}_0 + \tilde{d}_1 y_0$$

$$= \bar{x} - \bar{y} \tilde{d}_1 + y_0 \tilde{d}_1$$

$$= \bar{x} + \tilde{d}_1 (y_0 - \bar{y})$$

$$\Rightarrow \tilde{x}_0 - \bar{x} = (y_0 - \bar{y}) \cdot \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (y_i - \bar{y})^2}$$

$$\text{and } \hat{x}_0 - \bar{x} = (y_0 - \bar{y}) / \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\Rightarrow \frac{(\hat{x}_0 - \bar{x})}{(\tilde{x}_0 - \bar{x})} = \frac{(\sum (x_i - \bar{x})(y_i - \bar{y}))^2}{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}$$

$$\text{that is the prove: } \frac{F}{F + (n-2)} = \frac{(\sum (x_i - \bar{x})(y_i - \bar{y}))^2}{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}$$

$$\begin{aligned} \text{LHS} &= \frac{\hat{\beta}_1^2 \sum (x_i - \bar{x})^2}{\hat{\beta}_1^2 \sum (x_i - \bar{x})^2 + (n-2) \cdot s^2} \\ &= \frac{(\sum (y_i - \bar{y})(x_i - \bar{x}))^2}{(\sum (y_i - \bar{y})(x_i - \bar{x}))^2 + (n-2) \cdot s^2 \cdot \sum (x_i - \bar{x})^2} \\ &= \frac{(\sum (y_i - \bar{y})(x_i - \bar{x}))^2}{(\sum (y_i - \bar{y})(x_i - \bar{x}))^2 + \sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2} \\ &= \frac{(\sum (y_i - \bar{y})(x_i - \bar{x}))^2}{(\sum (y_i - \bar{y})(x_i - \bar{x}))^2 + \sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2 - (\sum (y_i - \bar{y})(x_i - \bar{x}))^2} \\ &= \frac{(\sum (y_i - \bar{y})(x_i - \bar{x}))^2}{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2} \\ &= \text{RHS} \\ \Leftrightarrow \tilde{x}_0 - \bar{x} &= (\hat{x}_0 - \bar{x}) \cdot \frac{F}{F + (n-2)} \quad \# \end{aligned}$$

- 100(1-2\alpha)%
2. Derive the Feller's interval estimate for the ratio of means of two correlated random variables with finite variances; using the notation from the wikipedia page

Notation: $(m_L, m_U) = \left(\frac{1}{(1-g)} \left(\frac{a}{b} - \frac{gV_{12}}{V_{22}} \pm \frac{\text{tr}_r s}{b} \left(V_{11} - 2\frac{a}{b}V_{12} + \frac{a^2}{b^2}V_{22} - g(V_{11} - \frac{V_{12}}{V_{22}}) \right)^{1/2} \right), \text{ for } \frac{Ma}{M_b} \right)$, for $\frac{Ma}{M_b}$

where $A \sim f(\mu_A, V_{11}\sigma^2)$, $B \sim f(\mu_B, V_{22}\sigma^2)$, $\text{cov}(A, B) = \sigma^2 V_{12}$

$a = \bar{A}$, $b = \bar{B}$, $g = \frac{\text{tr}_r s^2 V_{22}}{b^2}$, s^2 is the unbiased estimator for σ^2 with $\text{df} = r$

Pf: let $\theta = \frac{Ma}{Mb}$, $w = A - \theta B$

$$\text{then: } E(w) = E(A) - E(\theta B) = Ma - Ma = 0$$

$$\text{Var}(w) = \text{Var}(A) + \theta^2 \text{Var}(B) - 2\theta \text{cov}(A, B) = \sigma^2 (V_{11} + \theta^2 V_{22} - 2\theta V_{12})$$

$$\text{thus: } w = A - \theta B \sim N(0, \sigma^2 (V_{11} + \theta^2 V_{22} - 2\theta V_{12}))$$

$$\text{and: } \frac{rs^2}{\sigma^2} \sim \chi^2_r$$

$$\Rightarrow \frac{(A - \theta B) / \sigma (V_{11} + \theta^2 V_{22} - 2\theta V_{12})^{1/2}}{\sqrt{r} s / \sigma} \sim \text{tr}$$

$$\frac{A - \theta B}{(V_{11} + \theta^2 V_{22} - 2\theta V_{12})^{1/2} \cdot s} \sim \text{tr}$$

therefore: the $(1-2\alpha)$ -CI for $(A - \theta B)^2$ is $(\text{tr}_r a \cdot s^2 (V_{11} + V_{22}\theta^2 - 2\theta V_{12}))$

plugging in estimation $a = \bar{A}$, $b = \bar{B}$

$$(a - \theta b)^2 = (\text{tr}_r a \cdot s^2 (V_{11} + V_{22}\theta^2 - 2\theta V_{12}))$$

$$\theta^2 - \frac{2a}{b} \theta + \frac{a^2}{b^2} = \frac{(\text{tr}_r a \cdot s^2)^2}{b^2} (V_{11} + V_{22}\theta^2 - 2\theta V_{12})$$

$$\theta^2 - \frac{2a}{b} \theta + \frac{a^2}{b^2} = \frac{g V_{11}}{V_{22}} + g \theta^2 - 2 \frac{g V_{12}}{V_{22}} \theta$$

$$(1-g)\theta + 2\left(\frac{g V_{12}}{V_{22}} - \frac{a}{b}\right)\theta + \left(\frac{a^2}{b^2} - \frac{g V_{11}}{V_{22}}\right) = 0$$

$$\theta = \frac{1}{2(1-g)} \left[-2 \frac{g V_{12}}{V_{22}} + \frac{a}{b} \pm \left(\left(\frac{g V_{12}}{V_{22}} - \frac{a}{b} \right)^2 - 4(1-g)\left(\frac{a^2}{b^2} - \frac{g V_{11}}{V_{22}}\right) \right)^{1/2} \right]$$

$$= \frac{1}{(1-g)} \left[\frac{a}{b} - \frac{g V_{12}}{V_{22}} \pm \frac{1}{2} \left(\frac{4g^2 V_{12}^2}{V_{22}^2} + \frac{4a^2}{b^2} - 8 \frac{g V_{12} \cdot a}{V_{22} \cdot b} - 4(1-g) \frac{a^2}{b^2} + 4(g-g^2) \frac{V_{11}}{V_{22}} \right)^{1/2} \right]$$

$$= \frac{1}{(1-g)} \left[\frac{a}{b} - \frac{g V_{12}}{V_{22}} \pm \frac{1}{2} (1g - g^2) \frac{V_{11}}{V_{22}} - \frac{a}{b} \cdot \frac{g V_{12}}{V_{22}} + \frac{a^2}{b^2} g + \frac{g^2 V_{12}^2}{V_{22}^2} \right]$$

$$\text{since } g = \frac{\text{tr}_r s^2 V_{22}}{b^2} \Rightarrow \frac{g}{V_{22}} = \left(\frac{\text{tr}_r a s}{b} \right)^2$$

$$\text{therefore: } CI = \frac{1}{(1-q)} \left[\frac{q}{b} - \frac{qV_{12}}{V_{22}} \pm \left(\frac{\text{triaS}}{b} \right) \left((1-q)V_{11} - 2\frac{q}{b}V_{12} + \frac{q^2}{b^2}V_{22} + q\frac{V_{12}^2}{V_{22}} \right)^{1/2} \right] \\ = \frac{1}{(1-q)} \left[\frac{q}{b} - \frac{qV_{12}}{V_{22}} \pm \left(\frac{\text{triaS}}{b} \right) \left(V_{11} - 2\frac{q}{b}V_{12} + \frac{q^2}{b^2}V_{22} - q(V_{11} - \frac{V_{12}^2}{V_{22}}) \right)^{1/2} \right]$$

The V_{ij} is exactly the same as that from wikipedia

3. Find a 95% confidence interval for the unknown change point in a two straight line regression.

$$\hat{Y} = -\frac{\hat{\alpha}_1 - \hat{\alpha}_2}{\hat{\beta}_1 - \hat{\beta}_2}$$

$$\text{let } U = (\hat{\alpha}_1 - \hat{\alpha}_2) + Y(\hat{\beta}_1 - \hat{\beta}_2)$$

$$E(U) = 0$$

$$\text{Var}(U) = \sigma^2 \left\{ \frac{1}{n_1} + \frac{(\bar{x}_1 - Y)^2}{\sum(x_{1i} - \bar{x}_1)^2} + \frac{1}{n_2} + \frac{(\bar{x}_2 - Y)^2}{\sum(x_{2i} - \bar{x}_2)^2} \right\} \\ \cong \sigma^2 \cdot W$$

then: 95% CI for $U^2: F_{1, n-q}^{0.95} S^2 W$

$$(\hat{\alpha}_1 - \hat{\alpha}_2 + Y(\hat{\beta}_1 - \hat{\beta}_2))^2 = F_{1, n-q}^{0.05} S^2 W \\ \Rightarrow [(\hat{\beta}_1 - \hat{\beta}_2)^2 - F_{1, n-q}^{0.05} S^2 (\sum(x_{1i} - \bar{x}_1)^2)^{-1} + (\sum(x_{2i} - \bar{x}_2)^2)^{-1}] Y^2 \\ + 2[(\hat{\alpha}_1 - \hat{\alpha}_2)(\hat{\beta}_1 - \hat{\beta}_2) + F_{1, n-q}^{0.05} S^2 \left(\frac{\bar{x}_1}{\sum(x_{1i} - \bar{x}_1)^2} + \frac{\bar{x}_2}{\sum(x_{2i} - \bar{x}_2)^2} \right)] \cdot Y \\ + (-F_{1, n-q}^{0.05} S^2) \left(\frac{\bar{x}_1^2}{\sum(x_{1i} - \bar{x}_1)^2} + \frac{\bar{x}_2^2}{\sum(x_{2i} - \bar{x}_2)^2} \right) + (\hat{\alpha}_1 - \hat{\alpha}_2)^2 = 0$$

for simple notation:

Write as $CY^2 + 2DY + E = 0$

then 95% CI for Y is: $(-\frac{D}{C} \pm \frac{\sqrt{D^2 - CE}}{C})$

4. Verify two results from the climate change paper.

Notation: Model: $Y_{it} = \alpha_0 + \beta_0 t + \beta_1 (t - c) \cdot I(t > c) + \epsilon_{it}$ (2)

$$b = b_1 - b_0$$

① Likelihood ratio test statistic: $U = \frac{(S_0 - S)/3}{S/(n-4)}$, S_0 : residual sum of square from $H_0: Y_{it} = \alpha_0 + \beta_0 t + \epsilon_{it}$

② CI for c is the co sit. $(S' - S)/(S/(n-4)) \leq F_{1, n-4}(1-\alpha)$, S' : residual sum of square from fitting (2) at $c=c_0$

reference: Inference in Two-Phase Regression, 1971, Hinkley (a)

Inference about the Intersection in Two-Phase Regression (b)

Pf: ① write model as:

$$Y_{it|t} = \begin{cases} \alpha_0 + \beta_0 t + \epsilon_{it}, & t \leq c, 1 \leq i \leq I \\ \alpha_1 + \beta_1 t + \epsilon_{it}, & t > c, 1 \leq i \leq N \end{cases} \quad \alpha_1 = \alpha_0 + \beta_0 c - \beta_1 c$$

likelihood $L(Y_{it}; \alpha_0, \beta_0, \alpha_1, \beta_1, c, \sigma^2)$

$$= (2\pi\sigma^2)^{-\frac{N}{2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{j=1}^I (y_{ij} - \alpha_0 - \beta_0 t_j)^2 - \frac{1}{2\sigma^2} \sum_{j=I+1}^N (y_{ij} - \alpha_1 - \beta_1 t_j)^2 \right\}$$

$$L(c) = L_I(c) \quad (t_I \leq c \leq t_{I+1}; I=2, \dots, N-2)$$

$L_I(x) = (2\pi\sigma^2)^{-\frac{N}{2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} S_I^2(x) \right\}$, where $S_I^2(y)$ is the residual sum of squares for two regression lines constrained to meet at $t=x$

For the residual sum of squares of a single regression:

$$S_0^2 = \sum_{j=1}^N (y_{ij} - \bar{y}_N)^2 = \left(\sum_{j=1}^N (t_j - \bar{t})(y_{ij} - \bar{y}) \right)^2 / \sum_{j=1}^N (t_j - \bar{t})$$

let $Z_I^2(x) = S_0^2 - S_I^2(x)$, then the overall maximum, $\tilde{Z}_I^2(c)$ is the likelihood ratio statistic for $H_0: \beta_0 = \beta_1$ (b) / 1.1.2.

test statistic: from MR theories: $S = S_{\perp}^2(c)$ follows $S(\sigma^2) \sim \chi^2(n-4)$

from strong empirical evidence, (though not theoretically proved),

$$S_0 - S = \tilde{Z}_I^2(c) \sim \chi^2(3) \quad (a) / 2.3$$

$$\Rightarrow \text{Test statistic } U = \frac{(S_0 - S)/3}{S/(n-4)} \sim f_{2, n-4} \text{ under } H_0: \beta_0 = \beta_1$$

② $H_0: C = c_0$.

$$\begin{aligned}\text{likelihood ratio test statistic: } \Lambda &= -2 \log \left\{ L(c_0) / L(c) \right\} \rightarrow \chi^2(1) \text{ (Wilk's thm)} \\ &= -\frac{1}{8^2} \left\{ \bar{Z}_{\hat{\Sigma}_0}^2(w) - \bar{Z}_{\hat{\Sigma}}^2(w) \right\} \\ &= -\frac{1}{8^2} (S_0 - S' - S_0 + S) \\ &= \frac{1}{8^2} (S' - S) \xrightarrow{\text{asymptotically}} \chi^2(1)\end{aligned}$$

meanwhile: $S/\delta^2 \sim \chi^2(n-4)$

therefore: under $H_0: C = c_0$: $\frac{(S' - S)}{S/(n-4)} \sim F_{1, n-4}$

the asymptotic $100(1-\alpha)\%$ CI for c_0 is $\left\{ c : \frac{1(S' - S)}{S/(n-4)} \leq F_{1, n-4}(1-\alpha) \right\}$

5.

- a. Verify that SSPE = 3.036669

```
lof <- data.frame(x = c(1,1,2,3.3,3.3,4,4,4,4.7,5,5.6,5.6,5.6,6,6,6.5,6.9),
                    y = c(2.3,1.8,2.8,1.8,3.7,2.6,2.6,2.2,3.2,2,3.5,2.8,2.1,3.4,3.2,3.4,5))
```

$$SSPE = \sum_j \sum_i (Y_{ij} - \bar{Y}_j)^2$$

```
# calculate group mean
mean_y <- lof %>%
  group_by(x) %>%
  summarise(meany = mean(y), n = n())
# calculate SSPE
(lof %>%
  left_join(mean_y, by = c("x" = "x")) %>%
  mutate(err2 = (y - meany)^2) %>%
  summarise(sum(err2)))
```

```
##   sum(err2)
## 1  3.036667
```

- b. Would a cubic polynomial fit well for the same data set? What about a quintic polynomial?

```
lof <- lof %>%
  mutate(x2 = x^2) %>%
  mutate(x3 = x^3) %>%
  mutate(x4 = x^4) %>%
  mutate(x5 = x^5)

mod1 <- lm(y ~ x, data = lof)
mod3 <- lm(y ~ x + x2 + x3, data = lof)
mod5 <- lm(y ~ x + x2 + x3 + x4 + x5, data = lof)
```

```
anova(mod1, mod3)
```

```
## Analysis of Variance Table
##
## Model 1: y ~ x
## Model 2: y ~ x + x2 + x3
##   Res.Df   RSS Df Sum of Sq    F Pr(>F)
## 1      15 7.4150
## 2      13 4.2043  2     3.2107 4.9639 0.02502 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
anova(mod1, mod5)
```

```
## Analysis of Variance Table
##
## Model 1: y ~ x
## Model 2: y ~ x + x2 + x3 + x4 + x5
##   Res.Df   RSS Df Sum of Sq    F Pr(>F)
## 1      15 7.4150
## 2      11 4.1306  4     3.2843 2.1866 0.1376
```

```
anova(mod3, mod5)

## Analysis of Variance Table
##
## Model 1: y ~ x + x2 + x3
## Model 2: y ~ x + x2 + x3 + x4 + x5
##   Res.Df   RSS Df Sum of Sq    F Pr(>F)
## 1     13 4.2043
## 2     11 4.1306  2  0.073615 0.098 0.9074
```

The cubic polynomial has a significant better fit than the model with only linear terms. While for the quintic polynomial model, it has neither significantly better fit than the linear model nor than the cubic polynomial model.

b. $Z_1, \dots, Z_k \sim N(0,1)$, $U \sim \chi^2_m(0)$, $U \perp (Z_1, \dots, Z_k)$. $M \triangleq \max_{1 \leq i \leq k} \frac{|Z_i|}{\sqrt{U/m}}$

M ~ studentized maximum modulus distribution ; $M \sim M_{k,m}$

- (i) use studentized maximum modulus distribution to find a simultaneous CI for the set of all $\mu_i = \theta + \alpha_i$ for one way ANOVA, k groups, n obs/group.

under H_0 : $\alpha \sim N(\mu_\alpha, \Sigma)$, $\Sigma = \frac{1}{n} \sigma^2 I_k$.

let $\bar{x}_1, \dots, \bar{x}_k$ be the group sample mean

then $[(\bar{x}_1 - \mu_1) - (\bar{x}_2 - \mu_2)], \dots, [(\bar{x}_1 - \mu_1) - (\bar{x}_k - \mu_k)] \dots [(\bar{x}_{k-1} - \mu_{k-1}) - (\bar{x}_k - \mu_k)] \stackrel{iid}{\sim} N(0, \sigma^2 (\frac{1}{n} + \frac{1}{n}))$

have $k^* = \frac{k(k-1)}{2}$ such pair of group mean difference.

S^2 ; unbiased estimator for σ^2 , $(nk-k)S^2/\sigma^2 \sim \chi^2_{nk-k}$ (2)

$$\text{therefore: } M = \max_{1 \leq i < j \leq k} \frac{|(\bar{x}_i - \mu_i) - (\bar{x}_j - \mu_j)| / \sqrt{\frac{S^2}{n}}}{S} \sim M_{(k^*, nk-k)}$$

Let $m_{k,m}^\alpha$

be the upper α quantile

- (ii) a_1, \dots, a_k : a set of numbers; Is it true that $\max|a_i| \leq c$ iff $|\sum_{i=1}^k d_i a_i| \leq c \cdot \sum_{i=1}^k |d_i|$ for all numbers d_1, d_2, \dots, d_k ?

$$\max|a_i| \leq c \Leftrightarrow |\sum_{i=1}^k d_i a_i| \leq c \cdot \sum_{i=1}^k |d_i|$$

$$\Rightarrow |\sum_{i=1}^k d_i a_i| \leq \sum_{i=1}^k |d_i a_i| \leq \sum_{i=1}^k |a_i| \cdot |d_i| \leq \sum_{i=1}^k c \cdot |a_i| = c \cdot \sum_{i=1}^k |a_i|$$

\Leftarrow suppose $|\sum_{i=1}^k d_i a_i| \leq c \cdot \sum_{i=1}^k |d_i|$ is true for any $\vec{d} = (d_1, \dots, d_k)$

the it holds for any $\vec{d} = \vec{e}_i = (0, \dots, \overset{i}{\underset{i \neq 0}{\dots}}, \dots, 0)$, $i = 1, 2, \dots, k$

$$\Rightarrow |\sum_{i=1}^k d_i a_i| = |d_i a_i| = |a_i|$$

$$c \cdot \sum_{i=1}^k |d_i| = c \cdot |d_i| = c$$

$$\Rightarrow |a_i| \leq c, \quad \forall i = 1, 2, \dots, k \Rightarrow \max|a_i| \leq c$$

Yes, this conclusion is true.

- (iii) Utilize (ii) and studentized maximum modulus distribution, find a simultaneous set of confidence interval for $\sum_{i=1}^k d_i \mu_i$.

Solns:

$$\bar{x}_i - \mu_i \stackrel{iid}{\sim} N(0, \sigma^2/n) , \quad (nk-k)\frac{s^2}{\sigma^2} \sim \chi^2_{(nk-k)}$$

$$\Rightarrow \max_{1 \leq i \leq k} \frac{|\bar{x}_i - \mu_i| \sqrt{n}}{s} \sim M_{k, nk-k}$$

then: $(1-\alpha) \cdot 100\%$ CI for $\max |\bar{x}_i - \mu_i|$ is $\max |\bar{x}_i - \mu_i| \leq M_{k, nk-k}^{\alpha} \cdot \frac{s}{\sqrt{n}}$
where $M_{k, nk-k}^{\alpha}$ is the upper α quantile

$$\max |\bar{x}_i - \mu_i| \leq c \Leftrightarrow \left| \sum_{i=1}^k d_i (\bar{x}_i - \mu_i) \right| \leq c \cdot \sum_{i=1}^k |d_i|$$

$$\Leftrightarrow \left| \sum_{i=1}^k d_i \bar{x}_i - \sum_{i=1}^k d_i \mu_i \right| \leq c \cdot \sum_{i=1}^k |d_i|$$

$$\Rightarrow (100(1-\alpha)\%) \text{ CI for } \sum_{i=1}^k d_i \bar{x}_i \text{ is } \sum_{i=1}^k d_i \bar{x}_i \pm \frac{s}{\sqrt{n}} M_{k, nk-k}^{\alpha} \cdot \sum_{i=1}^k |d_i|$$