

STOR 664: FALL 2022
Midterm Exam, October 6, 2022

Open book in-class exam: time limit 75 minutes.

This is a single multi-part question but each part will be graded independently of the other parts. You are allowed to consult course notes (printed or e-read), homework assignments and any personal notes you have made during the course. Other outside materials are not permitted. Computers or ipads may be used *only* for the purpose of accessing pre-stored course notes; they are not to be used for computations during the exam. A hand-held calculator is permitted. Answers should preferably be written in a university examination book (“blue book”). You may consult the teaching assistant (in class) or the instructor (email, text or phone) if the wording is unclear or if you think there might be an error, but the teaching assistant or instructor will not give hints how to solve the exam. The university Honor Code is in effect at all times.

Consider the linear regression model

$$y_i = \beta_0 + x_i^2\beta_1 + x_i\beta_2 + x_i^3\beta_3 + \epsilon_i, i = 1, \dots, n$$

where the ϵ_i are independent normally distributed random variables with mean 0 and a common unknown variance σ^2 . (Please note the order of covariates: x_i^2 , x_i and x_i^3 , in that order.) Defining $S_k = \sum_{i=1}^n x_i^k$, $T_k = \sum_{i=1}^n y_i x_i^k$ for any $k \geq 0$, we assume that the x_i 's are symmetric about 0 to guarantee that $S_k = 0$ for all odd values of k .

- (a) Find explicit expressions for the least squares estimators $\hat{\beta}_0, \dots, \hat{\beta}_3$, and their variances. You should express the answer in terms of the values of S_k and T_k for $k \geq 0$, or any expressions derived from them. [25 points]
- (b) We would like to test the hypotheses $H_0 : \beta_2 = \beta_3 = 0$ versus the alternative H_1 that at least one of β_2 or β_3 is not 0. How would the estimates in (a), and their variances, change under the assumption that H_0 is true? [10 points]
- (c) Now suppose we are setting up the formal F test of H_0 against H_1 . Write SSE_0 and SSE_1 for the residual sum of squares under H_0 and H_1 respectively. Show that

$$SSE_0 - SSE_1 = A\hat{\beta}_2^2 + B\hat{\beta}_2\hat{\beta}_3 + C\hat{\beta}_3^2$$

where A , B and C are constants that you should identify (functions of n , S_2 , S_4 , etc.). Hence write down the formal test of H_0 against H_1 and define the rejection region for a test of significance level 0.01. (You are not expected to make an explicit numerical calculation but describe how to calculate it; for example, you may use R notation to define the needed percentage point of the F distribution, which will be one component of your answer.) [25 points]

Turn the page for the last two parts of the question.

- (d) Now consider the case where x_i goes from -3 to $+3$ in steps of 0.5 (so $n = 13$). You can assume (no need to check this) $S_2 = 45.5$, $S_4 = 284.375$, $S_6 = 2099.094$, the last to three decimal places. Also assume $\sigma^2 = 40$, $\beta_2 = 1$, $\beta_3 = 0.5$. What, in that case, will be the power of the test in part (c)? **[20 points]**

(You should give detailed numerical calculations as far as you are able to go, but the final answer will depend on the non-central F distribution and you should give the formula for calculating that as an R function or any equivalent notation that makes clear how to do the numerical calculation. If you cannot do explicit numerical calculations, at least state the formulas on which they may be based.)

- (e) Now suppose that the real purpose of the experiment is to determine a 95% prediction interval for a new observation taken at a new value $x = x^*$. Show (i) how to calculate such a prediction interval under the assumption H_1 , (ii) how the calculations in (i) would change if the experimenter did indeed assume H_0 to be true. **[20 points]**

(Since it's not possible for you to give a numerical answer here, you should describe precisely the sequence of steps, including any formulas for percentage points of relevant probability distributions. For (ii), you do not need to repeat the full calculation of (i), but indicate at which steps the calculation would change.)

Solutions

(a) $X = \begin{pmatrix} 1 & x_1^2 & x_1 & x_1^3 \\ 1 & x_2^2 & x_2 & x_2^3 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^2 & x_1 & x_1^3 \end{pmatrix}$ and hence $X^T X = \begin{pmatrix} n & S_2 & 0 & 0 \\ S_2 & S_4 & 0 & 0 \\ 0 & 0 & S_2 & S_4 \\ 0 & 0 & S_4 & S_6 \end{pmatrix}$. This is of block diagonal

form, so $(X^T X)^{-1} = \begin{pmatrix} S_4/\Delta_1 & -S_2/\Delta_1 & 0 & 0 \\ -S_2/\Delta_1 & n/\Delta_1 & 0 & 0 \\ 0 & 0 & S_6/\Delta_2 & -S_4/\Delta_2 \\ 0 & 0 & -S_4/\Delta_2 & S_2/\Delta_2 \end{pmatrix}$ where $\Delta_1 = nS_4 - S_2^2$,

$\Delta_2 = S_2S_6 - S_4^2$ (just invert both 2×2 submatrices). We also have $X^T Y = \begin{pmatrix} T_0 \\ T_2 \\ T_1 \\ T_3 \end{pmatrix}$,

so $\hat{\beta}_0 = \frac{S_4T_0 - S_2T_2}{\Delta_1}$, $\hat{\beta}_1 = \frac{-S_2T_0 + nT_2}{\Delta_1}$, $\hat{\beta}_2 = \frac{S_6T_0 - S_4T_2}{\Delta_2}$, $\hat{\beta}_3 = \frac{-S_4T_0 + S_2T_2}{\Delta_2}$, with variances respectively $\frac{S_4\sigma^2}{\Delta_1}$, $\frac{n\sigma^2}{\Delta_1}$, $\frac{S_6\sigma^2}{\Delta_2}$, $\frac{S_2\sigma^2}{\Delta_2}$.

- (b) If the model is refitted under H_0 , the estimates for $\hat{\beta}_0$ and $\hat{\beta}_1$ are the same, and their variances are the same as well. This is *because* of the block diagonal structure of $X^T X$, implying that if you just take the first two rows and columns of $X^T X$, you get the same inverse elements (the result would not be true without this). Of course, under H_0 , we don't consider $\hat{\beta}_2$ and $\hat{\beta}_3$ because these are assumed to be 0.
- (c) Various ways to do this, but I think the following argument is the simplest.

First, we note that $SSE_0 - SSE_1 = SSR_1 - SSR_0$ where SSR denotes the regression sum of squares.

Second, recall the formula (under either H_0 or H_1) that says $\sum(y_i - \bar{y})^2 = \sum(y_i - \hat{y}_i)^2 + \sum(\hat{y}_i - \bar{y})^2$ where the first term is SSE and the second term is SSR . Therefore, using the fact that $\hat{\beta}_0$ and $\hat{\beta}_1$ are the same in both models, we can write $SSR_0 = \sum(\hat{\beta}_0 + \hat{\beta}_1 x_i^2 - \bar{y})^2$ and $SSR_1 = \sum(\hat{\beta}_0 + \hat{\beta}_1 x_i^2 + \hat{\beta}_2 x_i + \hat{\beta}_3 x_i^3 - \bar{y})^2 = \sum(\hat{\beta}_0 + \hat{\beta}_1 x_i^2 - \bar{y})^2 + \sum(\hat{\beta}_2 x_i + \hat{\beta}_3 x_i^3)^2$; the cross-product $\sum(\hat{\beta}_0 + \hat{\beta}_1 x_i^2 - \bar{y})(\hat{\beta}_2 x_i + \hat{\beta}_3 x_i^3)$ is 0 because every term in the cross-product includes S_k for some odd k .

Therefore, $SSE_0 - SSE_1 = SSR_1 - SSR_0 = \sum(\hat{\beta}_2 x_i + \hat{\beta}_3 x_i^3)^2$ which expands to $S_2 \hat{\beta}_2^2 + 2S_4 \hat{\beta}_2 \hat{\beta}_3 + S_6 \hat{\beta}_3^2$. This is of the given form with $A = S_2$, $B = 2S_4$, $C = S_6$.

The F statistic is then

$$F = \frac{(SSE_0 - SSE_1)/2}{SSE_1/(n-4)}$$

and has the distribution $F_{2,n-4}$ under H_0 . The degrees of freedom arise because the original model H_1 has $p = 4$ unknown parameters which H_0 has $p - q = 2$ unknown parameters, therefore, $p = 4$, $q = 2$. The test will reject H_0 when $F > c$, where $c = \text{qf}(0.99, 2, n-4)$ in R notation (any equivalent notation for the F distribution will be accepted but the answer must include an explicit formula).

- (d) When H_1 is true, $F \sim F'(2, n - 4, \lambda)$ where the noncentrality parameter λ is given by the formula $\lambda\sigma^2 = S_2\beta_2^2 + 2S_4\beta_2\beta_3 + S_6\beta_3^2$ which you should be able to reduce to $\lambda = 21.37$ (however, I'll give credit for the correct formula even without the numerical answer). We also have $n = 13$. Using c from part (c), the final answer is given in R notation as either `1-pf(c,2,9,21.37)` or `pf(c,2,9,21.37,lower.tail=F)` or any equivalent notation for the non-central F distribution.

Note 1: Since the students did not have access to F tables during the exam, they were not expected to obtain the numerical answers for the last two quantities, but the actual values are $c = 8.02$ and power 0.73 to two decimal places.

Note 2: The book used δ^2 instead of λ and it would also be acceptable to express the answer this way.

- (e) The estimate is $\hat{y}^* = \hat{\beta}_0 + x^{*2}\hat{\beta}_1 + x^*\hat{\beta}_2 + x^{*3}\hat{\beta}_3$ with variance $(K + 1)\sigma^2$ where

$$\begin{aligned}
 K &= \begin{pmatrix} 1 & x^{*2} & x^* & x^{*3} \end{pmatrix} \begin{pmatrix} S_4/\Delta_1 & -S_2/\Delta_1 & 0 & 0 \\ -S_2/\Delta_1 & n/\Delta_1 & 0 & 0 \\ 0 & 0 & S_6/\Delta_2 & -S_4/\Delta_2 \\ 0 & 0 & -S_4/\Delta_2 & S_2/\Delta_2 \end{pmatrix} \begin{pmatrix} 1 \\ x^{*2} \\ x^* \\ x^{*3} \end{pmatrix} \\
 &= \frac{S_4}{\Delta_1} - 2\frac{S_2}{\Delta_1}x^{*2} + \frac{n}{\Delta_1}x^{*4} + \frac{S_6}{\Delta_2}x^{*2} - 2\frac{S_4}{\Delta_2}x^{*4} + \frac{S_2}{\Delta_2}x^{*6}.
 \end{aligned}$$

(Alternative notations will be accepted here. In the course text, K was written as $\mathbf{c}^T(X^T X)^{-1}\mathbf{c}$ where $\mathbf{c}^T = \begin{pmatrix} 1 & x^{*2} & x^* & x^{*3} \end{pmatrix}$. You can also write it that way so long as \mathbf{c} and $(X^T X)^{-1}$ are correctly defined.)

The 95% prediction interval for y^* is therefore $\hat{y}^* \pm cs\sqrt{K + 1}$ where s is the residual standard deviation and c is the appropriate percentage point of the t_{n-4} distribution, written in R notation as `qt(0.975,n-4)`. (Again, alternative notations will be accepted if completely defined.)

Under model H_0 make the following three changes: (i) omit the $\hat{\beta}_2$ and $\hat{\beta}_3$ terms in \hat{y}^* ; (ii) omit the last three terms in the definition of K ; (iii) the degrees of freedom for the t distribution is $n - 2$ instead of $n - 4$.