

First Year Examination
Department of Statistics, University of Florida
August 18, 2006, 8:00 am - 12:00 noon

Instructions:

1. Write your **number** on every page that you plan to submit.
2. **Do not** write your name anywhere on any of the pages that you plan to submit.
3. You have four hours to answer questions in this examination.
4. You must show your work to receive credit.
5. **Write only on one side of the paper, and start each question on a new page.**
6. There are 10 problems of which you must answer 8.
7. Only your first 8 problems will be graded.
8. While the 10 questions are equally weighted, some problems are more difficult than others.
9. The parts within a given question are not necessarily equally weighted.
10. You are allowed to use a calculator.

The following abbreviations and terminology are used throughout:

- ANOVA = analysis of variance
- cdf = cumulative distribution function
- SS = sums of squares
- iid = independent and identically distributed
- LRT = likelihood ratio test
- mgf = moment generating function
- MSE = mean squared error
- ML = maximum likelihood
- pdf = probability density function
- α = specified probability of Type I error
- $N(\mu, \sigma^2)$ = normal distribution with mean μ and variance σ^2

You may use the following facts/formulas without proof:

Fact about mgfs: If X has mgf $M_X(t)$ and a and b are constants, then the mgf of $aX + b$ is $e^{bt}M_X(at)$.

Linear Combinations of Independent Normals: Let X_1, X_2, \dots, X_n be independent random variables with $X_i \sim N(\mu_i, \sigma_i^2)$ for $i = 1, \dots, n$. If a_1, \dots, a_n are constants, then the random variable $\sum_{i=1}^n a_i X_i$ has a normal distribution.

Gamma Density: $X \sim \text{Gamma}(\alpha, \beta)$ means X has pdf

$$f(x; \alpha, \beta) = \frac{1}{\Gamma(\alpha) \beta^\alpha} x^{\alpha-1} e^{-x/\beta} I_{(0, \infty)}(x)$$

where $\alpha > 0$ and $\beta > 0$.

Students t Density: $X \sim t_\nu$ means X has pdf

$$f(x; \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right) \sqrt{\nu\pi}} \frac{1}{\left(1 + \frac{x^2}{\nu}\right)^{(\nu+1)/2}}$$

where $\nu > 0$.

Poisson moments: If $X \sim \text{Poisson}(\lambda)$, then $EX = \text{Var}X = \lambda$.

1. Consider the linear model $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ where \mathbf{Y} is the $n \times 1$ vector of dependent variables, \mathbf{X} is an $n \times (p + 1)$ matrix with full column rank and $\boldsymbol{\beta}$ is the vector of regression parameters. Suppose, contrary to the usual least squares assumptions, that the error vector $\boldsymbol{\epsilon}$ has mean vector zero and variance-covariance matrix $\sigma^2\mathbf{V}$, where \mathbf{V} is a known matrix not equal to the identity matrix.

- Find the mean vector and variance-covariance matrix of the ordinary least squares estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$.
- Suppose $\mathbf{V} = \mathbf{T}\mathbf{T}'$, where \mathbf{T} is invertible. Find a square matrix \mathbf{A} such that if $\mathbf{Y}^* = \mathbf{A}\mathbf{Y}$, $\mathbf{X}^* = \mathbf{A}\mathbf{X}$ and $\boldsymbol{\epsilon}^* = \mathbf{A}\boldsymbol{\epsilon}$ then the transformed model $\mathbf{Y}^* = \mathbf{X}^*\boldsymbol{\beta} + \boldsymbol{\epsilon}^*$ satisfies the ordinary least squares assumptions.
- Write the ordinary least squares estimator of $\boldsymbol{\beta}$ for the *transformed* model of the previous part in terms of \mathbf{X} , \mathbf{V} , and \mathbf{Y} . Also find its mean vector and variance-covariance matrix.
- Let \mathbf{e}^* be the vector of ordinary least squares residuals for the *transformed* model (based on the estimator of the previous part). Find the matrix \mathbf{P}^* such that the variance-covariance matrix of \mathbf{e}^* is $\sigma^2\mathbf{P}^*$. Show that \mathbf{P}^* is a projection matrix.
- Find an unbiased estimator of σ^2 .

2. A movie production company pre-releases its films in two test markets, Los Angeles and New York City, before the general release. The total attendances at these test screenings in Los Angeles (X_1 , in thousands) and New York City (X_2 , in thousands) are used as predictors of eventual total box office revenue (Y , in millions of \$) in the linear model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon.$$

Least-squares fitting of this model to data for 15 recent films yields (approximately)

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} 2.80 \\ 6.28 \\ 13.59 \end{bmatrix} \quad (\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} 0.68 & -0.20 & -0.14 \\ -0.20 & 0.21 & -0.14 \\ -0.14 & -0.14 & 0.27 \end{bmatrix} \quad \text{SS(Res)} = 840$$

where \mathbf{X} is the usual matrix of independent variables conforming to $\boldsymbol{\beta} = (\beta_0, \beta_1, \beta_2)'$, and SS(Res) is the residual (error) sum of squares. Assume that the model is adequate and that the errors ϵ are independent and identically distributed with mean zero and the same variance.

- Estimate the variance-covariance matrix of $\hat{\boldsymbol{\beta}}$.
- Give an unbiased prediction \hat{Y} of the total box office revenue (millions of \$) of a new film with test market attendances of 1.5 thousand and 2.2 thousand in Los Angeles and New York City, respectively. Give an unbiased estimate of the variance of \hat{Y} .
- Test the hypothesis that β_1 and β_2 are equal: Determine \mathbf{K} such that the null hypothesis may be written in the form $\mathbf{K}'\boldsymbol{\beta} = \mathbf{0}$, then perform an appropriate test. Use $\alpha = 0.05$.
- Compute 95% simultaneous two-sided confidence intervals for β_0 , β_1 , and β_2 , using the Bonferroni method.

3. A balanced one-factor experiment with t factor levels and r replications at each level yields responses y_{ij} for replication j at treatment level i . Consider the following two alternative models for the data:

$$\text{Model I: } y_{ij} = \mu + \tau_i + \epsilon_{ij}, \quad \sum_{i=1}^t \tau_i = 0$$

$$\text{Model II: } y_{ij} = \mu + a_i + \epsilon_{ij}, \quad a_1, \dots, a_t \sim \text{iid } N(0, \sigma_a^2)$$

where the terms ϵ_{ij} are independent and identically distributed as $N(0, \sigma_e^2)$ (with $\sigma_e^2 > 0$) and are independent of all a_i in Model II.

- For each model, write out the null hypothesis and the alternative hypothesis for the test of whether or not there are any factor effects.
- In terms of the data values y_{ij} , write out the sum of squares for factor effect, $SS(\text{Factor})$, and the sum of squares for error, $SS(\text{Error})$. Also give expressions for their corresponding degrees of freedom.
- Write an expression for the F -statistic (in terms of $SS(\text{Factor})$ and $SS(\text{Error})$) for testing the hypotheses in part (a). What is its distribution under each null hypothesis of part (a)?
- For each model, find the *correlation* between two different responses that have the same treatment level.

4. A balanced two-factor experiment yields responses y_{ijk} for replication k at level i of Factor A and level j of Factor B, for $i = 1, 2$, $j = 1, 2$, $k = 1, 2, 3, 4$. The data are analyzed using the model equation

$$y_{ijk} = \mu + \alpha_i + \beta_j + \alpha\beta_{ij} + \epsilon_{ijk},$$

under the conditions

$$\alpha_1 + \alpha_2 = \beta_1 + \beta_2 = 0, \quad \alpha\beta_{11} + \alpha\beta_{12} = \alpha\beta_{21} + \alpha\beta_{22} = \alpha\beta_{11} + \alpha\beta_{21} = \alpha\beta_{12} + \alpha\beta_{22} = 0$$

and the terms ϵ_{ijk} are $\text{iid } N(0, \sigma^2)$.

Suppose the least squares estimates of the model parameters are

$$\hat{\mu} = 10 \quad \hat{\alpha}_1 = 6 \quad \hat{\alpha}_2 = ? \quad \hat{\beta}_1 = 2 \quad \hat{\beta}_2 = ? \quad \hat{\alpha}\hat{\beta}_{11} = 1 \quad \hat{\alpha}\hat{\beta}_{12} = ? \quad \hat{\alpha}\hat{\beta}_{21} = ? \quad \hat{\alpha}\hat{\beta}_{22} = ?$$

and the sum of squares for error is 100.

- Find $\hat{\alpha}_2$, $\hat{\beta}_2$, $\hat{\alpha}\hat{\beta}_{12}$, $\hat{\alpha}\hat{\beta}_{21}$, and $\hat{\alpha}\hat{\beta}_{22}$.
- Write out an ANOVA table that includes the sum of squares and degrees of freedom for all sources (Factor A, Factor B, Interaction, and Error).
- Test whether there is any interaction between factors A and B. Test whether there is a main effect of Factor A. Test whether there is a main effect of Factor B. Perform all tests individually at level $\alpha = 0.05$.
- Is $\mu + \alpha_1$ a contrast? (Answer yes or no.) Give an unbiased estimate of $\mu + \alpha_1$ and an unbiased estimate of the variance of your estimate.
- Is $\alpha_1 - \alpha_2$ a contrast? (Answer yes or no.) Give an unbiased estimate of $\alpha_1 - \alpha_2$ and an unbiased estimate of the variance of your estimate.

5. An experiment is conducted in a completely randomized design with three treatment groups. In addition to the measured response Y , there is also a measured covariate X . The data are as follows:

Response Y	10	14	2	7	5	4
Treatment Group	1	1	2	2	3	3
Covariate X	-1	2	-2	1	0	0

- Estimate the coefficients in the simple linear regression of Y on X , *completely ignoring* the treatment group. Also find the residual (error) sum of squares.
- Perform an F -test for treatment effects using a one-way ANOVA model that *completely ignores* the covariate X . (Use $\alpha = 0.05$.)
- If both the treatment effects and a linear term in X are included in the model for the response Y , the residual (error) sum of squares is 0.75. Perform an analysis-of-covariance F -test for treatment effects (i.e. a test for treatment effects after adjusting for a linear effect in the covariate X .) (Use $\alpha = 0.05$.)

6. Toss n fair dice. Pick up only those *not* showing a 6 and toss them. Continue doing this until all dice show a 6. Let X denote the total number of tosses made in this experiment.

- How many tosses do you expect to make if $n = 1$?
- Compute the cdf of X for general n . (Hint: Think about what happens with each individual die.)
- Suppose that Y is a discrete random variable whose support is contained in $\{0, 1, 2, \dots\}$. Show that

$$\sum_{y=0}^{\infty} [1 - F_Y(y)] = EY,$$

where $F_Y(\cdot)$ is the cdf of Y . (Hint: A double sum will help.)

- Use the result from (c) to find the expected value of X when $n = 2$.

7. Suppose that $Y|Z = z \sim N(\mu, 1/z)$ and that $Z \sim \text{Gamma}(\alpha/2, \beta)$ where $\alpha, \beta > 0$.

- (a) Find the marginal pdf of Y .
- (b) Show that, in general, for any two random variables S and T

$$\text{Var}(S) = E[\text{Var}(S|T)] + \text{Var}[E(S|T)] ,$$

provided the expectations exist.

- (c) Use the result in part (b) to find the marginal variance of Y .
- (d) Suppose that $T \sim t_\alpha$. Find a function g such that if $W = g(\mu, \alpha, \beta, T)$, then W and Y have the same distribution. Note that μ, α and β are constant.

8. Let X_1, \dots, X_n be iid random variables from a distribution with a finite second moment. Let $EX_1 = \mu$, $\text{Var}X_1 = \sigma^2$, $\bar{X} = n^{-1} \sum_{i=1}^n X_i$ and $S^2 = (n-1)^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$. Show that

- (a) $E(\bar{X}) = \mu$
- (b) $\text{Var}(\bar{X}) = \frac{\sigma^2}{n}$
- (c) $S^2 = \frac{1}{n-1} [\sum_{i=1}^n X_i^2 - n\bar{X}^2]$
- (d) $E(S^2) = \sigma^2$

For the remainder of this problem, assume X_1, \dots, X_n are iid $\text{Poisson}(\lambda)$ and consider estimating λ with

$$\delta_w(X_1, \dots, X_n) = w\bar{X} + (1-w)S^2 ,$$

where $w \in [0, 1]$.

- (e) Show that δ_w is an unbiased estimator of λ .
- (f) Find the Cramér-Rao lower bound for the variance of an unbiased estimator of λ .
- (g) How many values of $w \in [0, 1]$ yield estimators that attain the Cramér-Rao lower bound? (Warning: A correct answer with nothing to back it up is worth 0 points.)
- (h) Find

$$E[\delta_w(X_1, \dots, X_n) | \bar{X}] .$$

9. Suppose that $X \in \mathcal{X}$ is a random vector with pdf $f(x|\theta)$, $\theta \in \Theta$. Consider testing $H_0 : \theta \in \Theta_0$ against $H_1 : \theta \in \Theta_0^c$, where Θ_0^c denotes the complement of Θ_0 ; that is, $\Theta_0^c = \Theta \setminus \Theta_0$.

(a) Fix $c \in (0, 1)$. A LRT has rejection region

$$R_1 = \{x \in \mathcal{X} : \lambda(x) < c\},$$

where $\lambda(x)$ is the LRT statistic. Write down the definition of $\lambda(x)$.

Now suppose that $\Theta_0 = \bigcap_{\gamma \in \Gamma} \Theta_\gamma$ where $\Theta_\gamma \subset \Theta$ and Γ is an index set. Let $\lambda_\gamma(x)$ denote the LRT statistic for testing $H_{0\gamma} : \theta \in \Theta_\gamma$ against $H_{1\gamma} : \theta \in \Theta_\gamma^c$.

(b) Explain why the set

$$R_2 = \{x \in \mathcal{X} : \lambda_\gamma(x) < c \text{ for some } \gamma \in \Gamma\}$$

might be a reasonable rejection region for testing $H_0 : \theta \in \Theta_0$ against $H_1 : \theta \in \Theta_0^c$.

(c) Define $T(x) = \inf_{\gamma \in \Gamma} \lambda_\gamma(x)$ and note that

$$R_2 = \{x \in \mathcal{X} : T(x) < c\}.$$

Show that $T(x) \geq \lambda(x)$ for all $x \in \mathcal{X}$

(d) Let $\beta_1(\theta)$ and $\beta_2(\theta)$ denote the power functions corresponding to R_1 and R_2 , respectively. Write down the definition of $\beta_1(\theta)$ and show that $\beta_1(\theta) \geq \beta_2(\theta)$ for all $\theta \in \Theta$.

(e) Fix $\alpha \in (0, 1)$. Suppose that the test with rejection region R_1 is a *size* α test. Does this imply that the test with rejection region R_2 is *level* α ? Explain.

10. Suppose that the random variables Y_1, \dots, Y_n satisfy

$$Y_i = x_i \beta + \varepsilon_i, \quad i = 1, \dots, n$$

where x_1, \dots, x_n are fixed constants, $\varepsilon_1, \dots, \varepsilon_n$ are iid $N(0, \sigma^2)$ and σ^2 is known.

(a) Find the ML estimator of β , call it $\hat{\beta} = \hat{\beta}(Y)$.

(b) Find the distribution of $\hat{\beta}$.

(c) Find the distribution of the alternative estimator of β given by

$$\tilde{\beta} = \tilde{\beta}(Y) = \frac{\sum_{i=1}^n Y_i}{\sum_{i=1}^n x_i}.$$

(d) Find the posterior distribution of β under a normal prior with mean 0 and variance $\tau^2 / (\sum_{i=1}^n x_i^2)$.

(e) Show that the posterior expectation of β , call it $\beta_B = \beta_B(Y)$, can be written as a simple function of $\hat{\beta}$.

(f) Compare these three estimators using MSE. Does any one of the three dominate the others? Can any one of the three be ruled out?