

First Year Examination
Department of Statistics, University of Florida
May 9, 2008, 8:00 am - 12:00 noon

Instructions:

1. You have four hours to answer questions in this examination.
2. You must show your work to receive credit.
3. **Write only on one side of the paper, and start each question on a new page.**
4. Questions 1 through 5 are the “applied” questions and questions 6 through 10 are the “theory” questions. You must do exactly four of the applied questions and exactly four of the theory questions.
5. While the 10 questions are equally weighted, some questions are more difficult than others.
6. The parts within a given question are not necessarily equally weighted.
7. You are allowed to use a calculator.

The following abbreviations and terminology are used throughout:

- ANOVA = analysis of variance
- iid = independent and identically distributed
- LRT = likelihood ratio test
- mgf = moment generating function
- ML = maximum likelihood
- MSE = mean squared error
- pdf = probability density function
- $\mathbb{R}^+ = (0, \infty)$
- $N(\mu, \sigma^2)$ = normal distribution with mean μ and variance σ^2
- H_0 = Null hypothesis
- H_a = Alternative hypothesis
- \sim = “is distributed as”

You may use the following facts/formulas without proof:

Normal density: $X \sim N(\mu, \sigma^2)$ means

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}$$

where $\mu \in \mathbb{R}$ and $\sigma^2 > 0$.

Beta density: $X \sim \text{Beta}(\alpha, \beta)$ means

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

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

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 $\beta > 0$. Also, $E(X) = \alpha\beta$ and $\text{Var}(X) = \alpha\beta^2$. The mgf is given by $m_X(t) = (1 - \beta t)^{-\alpha}$ for $t < 1/\beta$.

Iterated Expectation Formula: $E(X) = E[E(X|Y)]$.

Iterated Variance Formula: $\text{Var}(X) = E[\text{Var}(X|Y)] + \text{Var}[E(X|Y)]$.

1. For each value of a variable X , some observations of a variable Y are as shown in the following table:

X	Y values			
0	9	5	7	7
10	8	6		
20	1	0	1	6

- Compute the ordinary least squares estimates of the intercept and slope parameters in the simple linear regression of Y on X . (Be sure to identify which is which.)
- Compute an ANOVA table (corrected for the mean) for the simple linear regression. Based on the model, test whether there is any relationship between Y and X (making the usual assumptions about the data and using $\alpha = 0.05$). (Remember to state the null and alternative hypotheses.)
- Compute an unbiased estimate of the expected value of a new observation Y at $X = 5$ that follows the same simple linear regression model. Compute a 95% symmetric two-sided confidence interval for this quantity.
- Test whether there is any evidence that a quadratic regression model would provide a better fit to the data ($\alpha = 0.05$). (Remember to state the null and alternative hypotheses.)

2. Consider the following model

$$\begin{aligned}
 y_{ijk} &= \mu + \alpha_i + \eta_{k(i)} + \beta_j + \alpha\beta_{ij} + \epsilon_{ijk} & i &= 1, \dots, a \\
 \eta_{k(i)} &\sim \text{iid } N(0, \sigma_\eta^2), \quad \epsilon_{ijk} \sim \text{iid } N(0, \sigma^2), \quad \text{all independent} & j &= 1, \dots, b \\
 \sum_i \alpha_i &= \sum_j \beta_j = \sum_i \alpha\beta_{ij} = \sum_j \alpha\beta_{ij} = 0 & k &= 1, \dots, n
 \end{aligned}$$

which is often used for a split-plot design with one whole-plot factor (A) and one split-plot factor (B): y_{ijk} is the response from the split plot that receives level j of B , in the k^{th} whole plot that receives level i of A .

- Find the expected value and variance of y_{ijk} .
- Find the *correlation* between the responses from two different split plots in the same whole plot.
- Derive the expected value and variance of $\bar{y}_{ij\bullet} - \bar{y}_{i'j\bullet}$ for $i \neq i'$, where $\bar{y}_{ij\bullet} = \frac{1}{n} \sum_{k=1}^n y_{ijk}$.
- Derive the expected value of

$$\sum_{i=1}^a \sum_{k=1}^n (\bar{y}_{i\bullet k} - \bar{y}_{i\bullet\bullet})^2,$$

where $\bar{y}_{i\bullet k} = \frac{1}{b} \sum_{j=1}^b y_{ijk}$ and $\bar{y}_{i\bullet\bullet} = \frac{1}{bn} \sum_{j=1}^b \sum_{k=1}^n y_{ijk}$.

3. Twenty athletes are randomly split into 4 groups, with 5 athletes in each (using a method that makes all such splits equally likely). Three groups are each assigned to drink one of three different brands of “energy drink”, and the fourth group is given none. An hour later, all athletes compete in a 2 mile running race, and their times Y (minutes) to complete the race are recorded. On the previous day, all athletes had run the same 2 mile race without consuming energy drinks, and their completion times X (minutes) were recorded. Analysis of the variance of Y in terms of the group and the covariate X (but *not* their interaction) produces the following sequential (Type I) sums of squares:

Source	Sequential SS	Source	Sequential SS
Group	8.13	Previous Day's Time (X)	17.63
Previous Day's Time (X)	16.16	Group	6.66

(Note the different order of the sources in the two tables.)

The *sum of the squared residuals* is 4.13.

- Name the type of *design* that this experiment uses.
- For this experiment, why is it unnecessary to test whether the treatment affects the covariate X ? Explain *briefly*.
- Briefly explain the purpose of including X as a covariate in this experiment.
- Write a linear model equation for the analysis using the covariate, clearly defining all symbols and stating any conditions they satisfy.
- Test whether energy drink consumption affects running time. (Your test *should* account for the effect of the covariate X). Clearly state your conclusion. Use $\alpha = 0.05$.
- Perform the same test, but this time *ignoring* the covariate (i.e. based on the simple one-way analysis of variance, as if X had not been recorded). Clearly state your conclusion. Use $\alpha = 0.05$.

4. A medical study seeks to understand the effects of a new prescription drug and a nutritional supplement for treating a rare disorder, based on a single response variable. The study uses four treatment conditions: (1) placebo, (2) drug only, (3) supplement only, and (4) both drug and supplement. Human patients enter the study at different times. Once 4 patients have entered the study, the 4 treatments are randomly assigned to these patients (each patient receiving a different treatment). This process continues for each new block of 4 consecutive patients, until a total of 24 patients have entered the study. Let

y_{ij} = response of patient from j^{th} block who is assigned treatment condition i

for $i = 1, \dots, 4$ and $j = 1, \dots, 6$. Suppose

$$\bar{y}_{1\bullet} = 20, \quad \bar{y}_{2\bullet} = 15, \quad \bar{y}_{3\bullet} = 20, \quad \bar{y}_{4\bullet} = 5,$$

$$\sum_{j=1}^6 (\bar{y}_{\bullet j} - \bar{y}_{\bullet\bullet})^2 = 79, \quad \sum_{i=1}^4 \sum_{j=1}^6 (y_{ij} - \bar{y}_{\bullet\bullet})^2 = 1531$$

where

$$\bar{y}_{i\bullet} = \frac{1}{6} \sum_{j=1}^6 y_{ij} \quad \bar{y}_{\bullet j} = \frac{1}{4} \sum_{i=1}^4 y_{ij} \quad \bar{y}_{\bullet\bullet} = \frac{1}{24} \sum_{i=1}^4 \sum_{j=1}^6 y_{ij}$$

and suppose that, within each treatment condition, y_{ij} is normally distributed, with the same variance for all responses, regardless of treatment condition.

- Write a linear model equation appropriate for analyzing the responses y_{ij} from this experiment, clearly specifying the conditions on all terms, if any.
- Write an ANOVA table appropriate for your model.
- Test whether there are any differences between the treatment conditions ($\alpha = 0.05$). (Remember to state the null and alternative hypotheses in terms of your model.)
- Test whether the effects of the drug and supplement on the responses y_{ij} could be additive, i.e. whether or not there appears to be interaction between the effects of the drug and supplement ($\alpha = 0.05$). (Remember to state the null and alternative hypotheses in terms of your model.)
- Can you assert that the response mean is significantly *smaller* when the drug is given *without* the supplement versus when only a placebo is given? Test at level $\alpha = 0.05$. (Remember to state the null and alternative hypotheses in terms of your model.)

5. Consider the linear model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon, \quad \epsilon \sim N(0, \sigma^2).$$

Suppose you have n independent observations from this model, with the variable values for the i^{th} observation denoted $(Y_i, X_{i1}, X_{i2}, X_{i3})$. If \mathbf{X} is the matrix whose i^{th} row is $[1 \ X_{i1} \ X_{i2} \ X_{i3}]$ and \mathbf{Y} is the column vector whose i^{th} element is Y_i , summary statistics are as follows:

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 25 & 0 & 0 & 0 \\ 0 & 4 & -4 & -3.6 \\ 0 & -4 & 104 & -6.4 \\ 0 & -3.6 & -6.4 & 5.24 \end{bmatrix} \quad (\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} 0.04 & 0 & 0 & 0 \\ 0 & 1.26 & 0.11 & 1 \\ 0 & 0.11 & 0.02 & 0.1 \\ 0 & 1 & 0.1 & 1 \end{bmatrix} \quad \mathbf{X}'\mathbf{Y} = \begin{bmatrix} 55 \\ 4.4 \\ -34.4 \\ 4.54 \end{bmatrix}$$

$$\mathbf{Y}'\mathbf{Y} = 182.52$$

- Find n and the average of the values Y_1, \dots, Y_n .
- Note the zeros in the matrix $\mathbf{X}'\mathbf{X}$. What do they tell you about the independent variables?
- Compute the least squares estimates of $\beta_0, \beta_1, \beta_2$, and β_3 .
- Compute the sum of the squared residuals and the coefficient of determination (R^2).
- Find the *partial* sums of squares for X_1, X_2 , and X_3 . (Hint: Recall the relationship between t and F statistics for testing whether $\beta_j = 0$?)
- Suppose *backward elimination* is applied (until all variables are dropped). Which would be the first independent variable to be dropped? Explain briefly.
- Suppose *forward selection* is applied (until all variables have been added). Which would be the first independent variable to be added? Explain briefly.
- Test $H_0 : \beta_1 = \beta_3$ versus the alternative $H_a : \beta_1 \neq \beta_3$. Use $\alpha = 0.05$.

6. Let X_1, \dots, X_n be iid $\text{Uniform}(0, \theta)$ where $\theta > 0$; that is, the common pdf is

$$f(x|\theta) = \theta^{-1} I_{(0,\theta)}(x).$$

- Is θ a location parameter, scale parameter, or neither of these? Carefully explain your answer.
- Show that $X_{(n)}$ is a complete sufficient statistic.
- Find the ML estimator of θ , call it $\hat{\theta}(X)$.
- Find the ML estimator of $1/\theta$.
- Find the pdf of $\hat{\theta}(X)/\theta$.
- Show that $n\left(1 - \frac{\hat{\theta}(X)}{\theta}\right)$ converges in distribution and find the limiting distribution.
- Find the method of moments estimator of θ .
- Find the best unbiased estimator of θ .

7. Suppose that Y_1, Y_2, \dots, Y_n are independent random variables such that $Y_i \sim \text{Gamma}(\alpha, \beta_i)$ with $\alpha > 1$. In this problem, we consider the following expectation

$$\mathbb{E} \left[\frac{1}{\sum_{i=1}^n Y_i} \right]. \quad (1)$$

- (a) Give an exact formula for (1) in the special case where the β_i s are all equal with common value β .

When the β_i s are not equal, exact calculation is impossible so we develop upper and lower bounds.

- (b) Use Jensen's inequality to derive a simple lower bound for (1).
 (c) The upper bound requires more work. Let x_1, \dots, x_n be positive numbers and define the arithmetic mean (m_A), geometric mean (m_G), and harmonic mean (m_H) of these numbers as follows

$$m_A = \frac{1}{n} \sum_{i=1}^n x_i, \quad m_G = [x_1 x_2 \cdots x_n]^{\frac{1}{n}}, \quad \text{and} \quad m_H = \frac{1}{\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}}.$$

Let X be a discrete random variable that is uniform on the set $\{x_1, x_2, \dots, x_n\}$. Use X together with Jensen's inequality to prove that $m_H \leq m_G \leq m_A$.

- (d) Use the inequality $m_H \leq m_A$ to get an upper bound on (1) that is similar in form to the result from part (a).

8. Let X and Y be jointly continuous random variables with joint pdf given by

$$f(x, y) = \begin{cases} y e^{-y(x+1)} & x > 0 \text{ and } y > 0 \\ 0 & \text{otherwise} \end{cases}$$

- (a) What is the probability that X exceeds Y ? (You do not have to evaluate the integrals.)
 (b) Prove or disprove the following statement: X and Y are independent.
 (c) Find the marginal density of X and use it to calculate

$$\mathbb{E} \left[\frac{1}{(X+1)^p} \right]$$

for $p \in \mathbb{R}$.

- (d) Find the conditional density of Y given $X = x$.
 (e) Show that, conditional on $X = x$, the probability that Y exceeds 2 is no more than $\frac{1}{x+1}$.
 (f) Find the conditional mean and variance of Y given $X = x$.
 (g) Use the results from (f) together with the iterated variance formula to calculate the (marginal) variance of Y .
 (h) Find $P(X > 1 \mid Y > \pi)$. (You do not have to evaluate the integrals.)
 (i) Find $P(X > 1 \mid Y = \pi)$. (You do not have to evaluate the integrals.)

9. Suppose that $X \sim \text{Bernoulli}(p)$ and that the prior density for p is $\text{Beta}(\alpha, \beta)$.

- Find the posterior density of p given x , call it $\pi(p|x)$.
- Find the Bayes estimator of p , call it $\tilde{p}(X)$, and show that it can be written as a weighted average of the ML estimator of p , which is $\hat{p}(X) = X$, and the prior mean of p .
- Find the MSE of $\tilde{p}(X)$ and the MSE of $\hat{p}(X)$.
- Suppose that $\alpha = \beta = 1$. Does either of these estimators dominate the other in terms of MSE? In other words, is one of the MSE functions uniformly below the other?

Now consider an alternative prior that is a mixture of k beta densities; that is,

$$\pi^*(p) = \sum_{i=1}^k r_i \frac{\Gamma(\alpha_i + \beta_i)}{\Gamma(\alpha_i)\Gamma(\beta_i)} p^{\alpha_i-1} (1-p)^{\beta_i-1} I_{(0,1)}(p),$$

where, for each $i \in \{1, \dots, k\}$, $\alpha_i, \beta_i \in \mathbb{R}^+$, $r_i \in (0, 1)$ and $\sum_{i=1}^k r_i = 1$.

- Find the exact posterior density of p given x , call it $\pi^*(p|x)$.
- Is $\pi^*(p)$ a conjugate prior? (A yes/no answer is not sufficient here.)

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

$$Y_i = \beta x_i + \varepsilon_i,$$

for $i = 1, \dots, n$ where x_1, \dots, x_n are known constants, β is an unknown regression parameter, and $\varepsilon_1, \dots, \varepsilon_n$ are iid $N(0, \sigma^2)$ with σ^2 known. In this question, we develop the LRT of $H_0 : \beta = 0$ against $H_a : \beta \neq 0$. We start with some preliminary results.

- Derive the mgf of $V \sim N(\mu, \tau^2)$.
- Suppose that W_1, \dots, W_m are independent random variables such that $W_i \sim N(\mu_i, \tau_i^2)$. Find the distribution of $\sum_{i=1}^m a_i W_i$ where a_1, \dots, a_m are known constants.

Now back to the original problem.

- Find the ML estimator of β , call it $\hat{\beta}(Y)$.
- Construct the LRT statistic for testing $H_0 : \beta = 0$ against $H_a : \beta \neq 0$.
- Show that the LRT statistic can be written in such a way that it involves the data, Y , only through $T = \hat{\beta}^2(Y)$.
- Find the distribution of $T = \hat{\beta}^2(Y)$ under H_0 .
- The general LRT theory tells us to reject H_0 when the LRT statistic is small. Give an equivalent rejection rule in terms of T .
- Suppose that $n = 100$, $\sum_{i=1}^{100} x_i^2 = 10$ and $\sigma^2 = 5$. Give the *exact* rejection region of the size 0.10 LRT in terms of T . (You may use the fact that $P(Z > 1.645) = 0.05$ when $Z \sim N(0, 1)$.)