

Qualifying Exam August 2018 — Statistical Methods

Instructions:

- NOTE: You are not allowed to use internet or email except for downloading the data and emailing the final report (optional). To use these for any other purpose, ask the proctor.
- Go to <http://www.utdallas.edu/~swati.biswas/QE> and download two datasets Snail2.csv and rat.csv. Let the proctor know if you have any problems with this step.
- You can use any software of your choice. You can use the lab machines or your own laptop.
- Your report should clearly explain the steps, results, conclusions, and justification for the conclusions. Also, include your codes (with brief comments explaining each step) and outputs (ONLY relevant parts; highlighted wherever possible). Do not attach the parts of the output that were not used in answering questions.
- Submit a report (written or typed), hard copy or by **email to Angie.Bustamante@utdallas.edu**. If you choose to email, then attach only one single PDF file with the whole report. DO NOT email separate files for codes or outputs.
- **Write your QE ID number (given to you by Angie) on all answer sheets. DO NOT put your name, UTD ID, or any other identifying information on any of your answer sheets. DO NOT email your exam to anyone other than Angie.**

1. Consider the Snail2.csv data. We would like to understand how Length is related to the other variables in this data set. As part of this investigation, we would like to determine if ShellType interacts with any of other predictor variables, and we would like to remove unimportant predictor variables from the model. [25 points]
 - (a) Fit a model to predict *Length* based on the other variables. Check assumptions and perform any transformations needed to obtain a model that is reasonable with respect to the standard assumptions for linear models. [8 points]
 - (b) Reduce your model by removing any unimportant variables if such variables exist. Interpret the reduced model including coefficients and r-squared. Perform a statistical test that compares the full model to the reduced model. Clearly state the hypotheses associated with this test and interpret the results. [5 points]
 - (c) Use your final model to obtain 95% confidence and 95% prediction intervals for the Length of a Type1 snail when the other variables equal the means within Type1 snails. Repeat for Type2 snails. [5 points]
 - (d) Construct a model to predict ShellType based on the other variables. Reduce this model by removing any unimportant predictors and interpret the resulting model. [7 points]

2. A study investigated agreement between the conclusions (positive or negative) of two types of studies: (1) a meta analysis (a synthesis of evidence in the literature) and (2) a large randomized clinical trial for 40 health-related outcomes. The cell counts for the data are as follows:

Result of Meta Analysis	Result of Clinical Trial	
	Positive	Negative
Positive	13	6
Negative	7	14

This type of study is referred to as an inter-rater agreement study. It measures the amount of agreement between the ratings given by two raters. In the above data, “meta analysis” and “clinical trial” may be viewed as the two raters and positive and negative conclusions as the two possible ratings. In general, such a contingency table with cell counts and their totals can be represented as the following:

Rater 1 rating	Rater 2 rating		Total
	Positive	Negative	
Positive	n_{++}	n_{+-}	n_{+}
Negative	n_{-+}	n_{--}	n_{-}
Total	$n_{\cdot+}$	$n_{\cdot-}$	$n_{\cdot\cdot} = n$

The vector of counts $(n_{++}, n_{+-}, n_{-+}, n_{--})$ is assumed to follow a multinomial distribution with $(p_{++}, p_{+-}, p_{-+}, p_{--})$ as the vector of corresponding cell probabilities. The parameter of interest in such inter-rater agreement problems is the well-known *Cohen's kappa coefficient*, defined as

$$\kappa = \frac{2(p_{++}p_{--} - p_{+-}p_{-+})}{p_{+-} + p_{-+} + 2(p_{++}p_{--} - p_{+-}p_{-+})}.$$

Its maximum value is one, which indicates perfect agreement between the raters. Its value equals zero when there is no agreement beyond what is expected by chance alone. The measure is estimated by replacing the probabilities p_{ij} by their maximum likelihood estimates, $\hat{p}_{ij} = n_{ij}/n$. It can be written as

$$\hat{\kappa} = \frac{2(n_{++}n_{--} - n_{+-}n_{-+})}{n(n_{+-} + n_{-+}) + 2(n_{++}n_{--} - n_{+-}n_{-+})}.$$

- Compute $\hat{\kappa}$ for the above data and interpret the result. [4 points]
- Use bootstrap with 1,000 replications to compute standard error of $\hat{\kappa}$. Also compute 95% two-sided confidence intervals for κ using both basic bootstrap and percentile bootstrap methods. [27 points]

(c) What do you conclude about the extent of agreement between meta analyses and clinical trials? [4 points]

3. “Conditioned Suppression Data”. In this experiment researchers investigated the degree to which three treatments, paired with two different phases (“shock”/“No Shock”), would suppress the rate of an ongoing bar-pressing response in rats. 24 rats were randomly assigned to the three treatment groups with 8 rats in each group. For each rat in a group, the response were measured at four cycles (i.e., four different time points), in each of which two phases were administered in a random order.

#	Variable	Description
1	rat_id	ID of rats
2	group	Treatment groups: 1, 2, 3
3	cycle	Time points: 1,2,3,4
4	phase	Phases: 1 (“Shock”), 2 (“No Shock”)
6	response	bar-pressing response rate of rats

(a) First, we examine the effects of treatments (group) and phase on the response in the first cycle (i.e., cycle=1). We assume that there are no interactions between treatments and subjects (rat_id), but we want to examine the interaction effects between group and phase. [25 points]

- i. Write the mathematical formula of an appropriate ANOVA model that incorporates the above specifications. What is this experiment design? [5 points]
- ii. Fit the above model. Test whether the interactions effects exist. Test for the two main effects if appropriate. Use $\alpha = 0.05$ in all tests. Clearly state your conclusions. [10 points]
- iii. Obtain confidence intervals for all pairwise comparisons for the main effect(s) that you deem appropriate based on the above tests. Control the family-wise confidence level at 95%. Interpret your result and state your conclusions. [10 points]

(b) Next, we study the effects of all three factors: group, cycle and phase. [15 points]

- i. Write the mathematical formula of an appropriate ANOVA model. [5 points]
- ii. Fit the above model. Are there any significant interaction terms using $\alpha = 0.05$? [5 points]
- iii. Test for the main effects if appropriate. Use $\alpha = 0.05$ in all tests. Clearly state your conclusions. [5 points]

STATISTICS Ph.D. QUALIFYING EXAM
STATISTICAL INFERENCE
August 2018

General Instructions: Write your QE ID number (given to you by Angie) on all answer sheets. Do NOT put your name, UTD ID, or any other identifying information on any of your answer sheets. Show all work/proofs/references. Justify all arguments. Simplify answers as much as possible. Please write neatly so that it is easy to read your solution. Total points = 100.

1. (20 points) Suppose that X is a continuous random variable with probability density function $f(x)$ and cumulative distribution function $F(x)$. Let $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ denote the order statistics for a random sample X_1, X_2, \dots, X_n from this distribution.
 - (a) (5 points) Show that the distribution of $F[X_{(k)}]$ ($k = 1, 2, \dots, n$) is Beta with parameters k and $n - k + 1$.
 - (b) (7 points) Show that the distribution of $F[X_{(j)}] - F[X_{(i)}]$ is the same as the distribution of $F[X_{(j-i)}]$, where $1 \leq i < j \leq n$.
 - (c) (5 points) A random interval (L, U) , $-\infty < L < U < \infty$, is called a two-sided β -expectation tolerance interval for the distribution of X if

$$E\left(\int_L^U f(x)dx\right) = \beta.$$

Explain how you will choose an integer $r \in \{1, \dots, n\}$ such that $(X_{(r)}, X_{(n-r+1)})$ is a two-sided β -expectation tolerance interval.

- (d) (3 points) Is the tolerance interval in (c) distribution-free? Explain.
2. (20 points) Consider a population represented by a random variable $X \sim N(\mu, \sigma^2)$, with both parameters unknown. Let X_1, \dots, X_n denote a random sample from this population. For a specified constant x_0 , define the probability

$$p = Pr(X \leq x_0) = \Phi((x_0 - \mu)/\sigma),$$

where $\Phi(\cdot)$ is the cumulative distribution function of a $N(0, 1)$ distribution. Derive a closed-form level α test for the hypotheses

$$H_0 : p \geq p_0 \text{ versus } H_1 : p < p_0,$$

where $0 < p_0 < 1$ is a specified probability. (You may use the fact that if $Z \sim N(0, 1)$, $Y \sim \chi_k^2$, and the two are independent, then $T_k(\Delta) = (Z + \Delta)/\sqrt{Y/k}$ follows a noncentral t_k -distribution with noncentrality parameter Δ .)

3. (25 points) Suppose that X is a single observation from a $N(\theta, \theta)$ distribution where $\theta > 0$ is unknown. Use the following steps to construct the uniformly minimum variance unbiased estimator (UMVUE) for θ based on X .

(a) (5 points) Show that $T = |X|$ is complete and sufficient for θ .

(b) (5 points) Show that the density of T is

$$g(t) = \begin{cases} (2\pi\theta)^{-1/2} \exp\{-\frac{1}{2\theta}(t - \theta)^2\} \{1 + \exp(-2t)\}, & \text{if } t > 0; \\ 0, & \text{otherwise.} \end{cases}$$

(c) (5 points) Show that for $t > 0$,

$$Pr(X = t|T = t) = \exp(t)/\{\exp(-t) + \exp(t)\}.$$

(d) (5 points) Show that $E(X|T = t) = t \tanh(t)$, $t > 0$, where $\tanh(t) = \frac{\exp(t) - \exp(-t)}{\exp(t) + \exp(-t)}$.

(e) (5 points) Argue that $|X| \tanh(|X|)$ is the UMVUE for θ .

4. (20 points) Consider the same population distribution as in the previous problem, i.e., $X \sim N(\theta, \theta)$ where $\theta > 0$ is unknown. Find a pivotal quantity and use it to construct a $100(1 - \alpha)\%$ confidence interval for θ . (Note: If you obtain a confidence set for θ that may not always be an interval, be sure to explain when the set is an interval and when it is not.)

5. (15 points) Let X be an observation from a continuous distribution with density $f(x|\theta)$, where $-\infty < \theta < \infty$ is an unknown parameter. Let $\pi(\theta)$ be the prior distribution for θ . Consider estimation of θ under the LINEX (LINear-EXponential) loss function, given by

$$L(\theta, a) = \exp\{c(a - \theta)\} - c(a - \theta) - 1,$$

where c is a specified positive constant.

(a) (6 points) Show that the Bayes estimator of θ is given by

$$\delta^\pi(X) = \frac{-1}{c} \log[E(\exp\{-c\theta\}|X)].$$

- (b) (9 points) Let X_1, \dots, X_n be a random sample from a $N(\theta, \sigma^2)$ distribution, where σ^2 is known, and suppose that θ has the noninformative prior $\pi(\theta) = 1$. Use the result in part (a) to show that the Bayes estimator of θ under the LINEX loss is given by $\delta^\pi(\bar{X}) = \bar{X} - (c\sigma^2/(2n))$.

August 2018 Qualifying Exam in Probability Theory

- This is a closed-book test.
- There are 5 questions; all have multiple parts.
- Answer each question as fully as possible.
- Show and justify all steps of your solutions.
- Refer clearly to any known results that you are using, **stating such results precisely.**
- Show how the assumptions of a result you are using are satisfied in your application of the result.
- Indicate how the assumptions given in the question are used in the solution.
- Write your solutions on the blank sheets of paper that are provided.
- Write your QE ID number (given to you by Angie) on all answer sheets. DO NOT put your name, UTD ID, or any other identifying information on any of your answer sheets.
- On each sheet, identify which question and part is being answered.
- Begin each question on a new sheet.
- *When finished, arrange your sheets in order, number each sheet, and be sure that your QE ID number (given by Angie) is on each sheet.*
- Total points = 100.

PROBLEM 1 (20 points)

Let $X_n, n \geq 1$, denote a sequence of independent random variables with $E(X_n) = \mu$. Consider the sequence of random variables

$$\hat{v}_n = \frac{1}{n(n-1)} \sum_{i \neq j} X_i X_j,$$

which is an unbiased estimator of μ^2 . Does

(a) $\hat{v}_n \xrightarrow{P} \mu^2$?

(b) $\hat{v}_n \xrightarrow{\text{a.s.}} \mu^2$?

(c) $\hat{v}_n \rightarrow \mu^2$ in mean square?

(d) Does the estimator \hat{v}_2 follow a normal distribution if $n \rightarrow \infty$? Verify your result.

Justify and prove each item separately.

PROBLEM 2 (20 points) A gambler wins or loses one dollar in each round of betting, with equal chances and independently of the past events. She starts betting with the firm determination that she will stop gambling when either she won a dollars or she lost b dollars.

- (a) Compute the probability that she will be winning when she stops playing further.
- (b) Compute the expected number of her betting rounds before she will stop playing further.

PROBLEM 3 (30 points) Let $\xi_{n,k}, n = 1, 2, \dots; k = 1, 2, \dots,$ be independent and identically distributed random variables which take values in $\mathbb{N} = \{0, 1, 2, \dots\}$. Assume that $\mu = E(\xi_{n,k})$ and $\sigma^2 = \text{Var}(\xi_{n,k})$ exist. Consider the stochastic process

$$Z_0 = 1 \quad Z_{n+1} = \sum_{k=1}^{Z_n} \xi_{n+1,k}$$

with respect to the filtration $\mathcal{G}_n = \tau(Z_j, 0 \leq j \leq n)$.

(a) Define the process

$$M_n = \mu^{-n} Z_n, \quad n = 0, 1, 2, \dots,$$

and show that (M_n, \mathcal{G}_n) is a martingale.

(b) Show that $E(Z_{n+1}^2 | \mathcal{G}_n) = \mu^2 Z_n^2 + \sigma^2 Z_n$

(c) Using the result from (b) show that

$$N_n = \begin{cases} M_n^2 - \frac{\sigma^2}{\mu^{n+1}} \frac{\mu^n - 1}{\mu - 1} M_n & \text{if } \mu \neq 1 \\ M_n^2 - n\sigma^2 M_n & \text{if } \mu = 1 \end{cases}$$

is also a martingale with respect to \mathcal{G}_n .

(d) Using the result from (c) show that if $\mu > 1$ then $\sup_{0 \leq n < \infty} E(M_n^2) < \infty$, while if $\mu \leq 1$ then $\lim_{n \rightarrow \infty} E(M_n^2) = \infty$.

(e) Interpret the result and describe the meaning of the phenomenon obtained in (d) with 1 - 2 sentences.

PROBLEM 4 (10 points) Let X be a random variable defined on a probability space (Ω, \mathcal{F}, P) . Let $\mathcal{F}_1 \subset \mathcal{F}_2 \subset \mathcal{F}$ and assume that $\text{Var}(X) = \sigma^2 < \infty$. Show that

$$E[(X - E(X|\mathcal{F}_2))^2] \leq E[(X - E(X|\mathcal{F}_1))^2].$$

Explain the meaning of this result using 1 - 2 sentences.

PROBLEM 5 (20 points) The following problems are not related to each other.

(a) Let $X_n \sim \text{Exp}(\lambda)$, $\lambda > 0$. Show that

$$\limsup_n \frac{X_n}{\log n} = \frac{1}{\lambda}.$$

Remember that the p.d.f. is given by $f(x) = \lambda \exp(-\lambda x) \mathcal{I}_{x>0}$.

(b) Let A_n denote a sequence of events. Show that if $P(A_n) \rightarrow 0$ and $\sum_{n=1}^{\infty} P(A_n \cap A_{n+1}^c) < \infty$, then $P(A_n \text{ i.o.}) = 0$.

(c) Let X_n be independent and identically distributed with $E(|X_1|) < \infty$. Define

$$M_n = \frac{1}{n} \max(X_1, \dots, X_n).$$

Show that $M_n \xrightarrow{\text{a.s.}} 0$.

August 2018 Qualifying Exam in Linear Models

- This is a closed-book test.
- There are 3 questions; some have multiple parts.
- Answer each question as fully as possible.
- Show and justify all steps of your solutions.
- Refer clearly to any known results that you are using, **stating such results precisely.**
- Show how the assumptions of a result you are using are satisfied in your application of the result.
- Indicate how the assumptions given in the question are used in the solution.
- Write your solutions on the blank sheets of paper that are provided.
- Write your QE ID number (given to you by Angie) on all answer sheets. DO NOT put your name, UTD ID, or any other identifying information on any of your answer sheets.
- On each sheet, identify which question and part is being answered.
- Begin each question on a new sheet.
- *When finished, arrange your sheets in order, number each sheet, and be sure that your QE ID number (given by Angie) is on each sheet.*
- Total points = 100.

EXERCISE 1 (25 points) Consider a 2×1 vector $\mathbf{Y} = (Y_1, Y_2)^T$ that follows a $N_2(\mathbf{0}, \Sigma)$ distribution, where $\Sigma = (\sigma_{ij})$ is a 2×2 non-singular covariance matrix. Prove that the random variable

$$\mathbf{Y}^T \Sigma^{-1} \mathbf{Y} - \frac{Y_1^2}{\sigma_{11}}$$

follows a chi-square distribution with 1 degree of freedom.

EXERCISE 2 (30 points) Let

$$\begin{aligned} Y_1 &= \alpha_1 + \epsilon_1, \\ Y_2 &= 2\alpha_1 - \alpha_2 + \epsilon_2, \\ Y_3 &= \alpha_1 + 2\alpha_2 + \epsilon_3, \end{aligned}$$

where $\epsilon \sim N_3(\mathbf{0}, \sigma^2 \mathbf{I}_3)$. We would like to perform an F -test for testing $H_0 : \alpha_1 = \alpha_2$ against $H_1 : \alpha_1 \neq \alpha_2$.

1. (20 points) Obtain an explicit expression for the F -statistic. Simplify the expression as much as possible.
2. (5 points) What is the null distribution of the test statistic?
3. (5 points) Provide the rejection region for a level α F -test of the hypotheses.

EXERCISE 3 (45 points) Suppose that a $n \times 1$ vector of responses (or observations) \mathbf{Y} follows the mixed-effects model

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e},$$

where $\boldsymbol{\beta}$ is a $p \times 1$ vector of unknown coefficients, \mathbf{u} is a $q \times 1$ vector of unknown effects of random variables, \mathbf{e} is a $n \times 1$ vector of random errors, and \mathbf{X} and \mathbf{Z} are known design matrices of order $n \times p$ and $n \times q$, respectively, which relate the elements of $\boldsymbol{\beta}$ and \mathbf{u} to elements of \mathbf{Y} . The design matrices are assumed to have full rank. The elements of $\boldsymbol{\beta}$ are considered to be fixed effects while the elements of \mathbf{u} are considered to be random effects. This is why this model is called a mixed-effects model. We assume that

$$E \begin{pmatrix} \mathbf{u} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \quad \text{cov} \begin{pmatrix} \mathbf{u} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{G} & \mathbf{0} \\ \mathbf{0} & \mathbf{R} \end{pmatrix},$$

where the $q \times q$ matrix \mathbf{G} and the $n \times n$ matrix \mathbf{R} are non-singular covariance matrices. Assume that $\boldsymbol{\beta}$ is unknown but \mathbf{G} and \mathbf{R} are known. Define the $n \times n$ matrix $\mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}^T + \mathbf{R}$.

1. (10 points) Show that $E(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$, $\text{var}(\mathbf{Y}) = \mathbf{V}$, and $\text{cov}(\mathbf{u}, \mathbf{Y}) = \mathbf{G}\mathbf{Z}^T$.
2. (35 points) We would like to jointly estimate $\boldsymbol{\beta}$ and \mathbf{u} . For this, our strategy is to find the best linear unbiased predictor (BLUP) of the scalar quantity

$$\mathbf{s}^T \mathbf{X}\boldsymbol{\beta} + \mathbf{t}^T \mathbf{u},$$

where \mathbf{s} and \mathbf{t} are arbitrary known vectors of appropriate dimensions. This involves finding a predictor of the form $\mathbf{a}^T \mathbf{Y}$ that minimizes the mean squared prediction error

$$E\{(\mathbf{a}^T \mathbf{Y} - \mathbf{s}^T \mathbf{X}\boldsymbol{\beta} - \mathbf{t}^T \mathbf{u})^2\}$$

with respect to the vector \mathbf{a} , subject to the unbiasedness condition

$$E(\mathbf{a}^T \mathbf{Y}) = E(\mathbf{s}^T \mathbf{X}\boldsymbol{\beta} + \mathbf{t}^T \mathbf{u}).$$

- (a) (6 points) Show that the unbiasedness condition is equivalent to $\mathbf{X}^T \mathbf{a} = \mathbf{X}^T \mathbf{s}$. Deduce that, under this condition, the mean squared prediction error is the variance of the prediction error, given as

$$\text{var}(\mathbf{a}^T \mathbf{Y} - \mathbf{s}^T \mathbf{X}\boldsymbol{\beta} - \mathbf{t}^T \mathbf{u}) = \mathbf{a}^T \mathbf{V} \mathbf{a} + \mathbf{t}^T \mathbf{G} \mathbf{t} - 2\mathbf{a}^T \mathbf{Z} \mathbf{G} \mathbf{t}.$$

- (b) (6 points) Let $2\boldsymbol{\lambda}$ be a vector of Lagrange multipliers. Then, the expression to be minimized with respect to \mathbf{a} and $\boldsymbol{\lambda}$ is

$$\text{var}(\mathbf{a}^T \mathbf{Y} - \mathbf{s}^T \mathbf{X}\boldsymbol{\beta} - \mathbf{t}^T \mathbf{u}) + 2\boldsymbol{\lambda}^T (\mathbf{X}^T \mathbf{a} - \mathbf{X}^T \mathbf{s}).$$

Show that \mathbf{a} and $\boldsymbol{\lambda}$ that minimize this expression are solutions of

$$\begin{pmatrix} \mathbf{V} & \mathbf{X} \\ \mathbf{X}^T & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{a} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} \mathbf{Z} \mathbf{G} \mathbf{t} \\ \mathbf{X}^T \mathbf{s} \end{pmatrix}.$$

- (c) (18 points) Show that solving for \mathbf{a} leads to the desired BLUP of $\mathbf{s}^T \mathbf{X}\boldsymbol{\beta} + \mathbf{t}^T \mathbf{u}$ as

$$\mathbf{a}^T \mathbf{Y} = \mathbf{s}^T \mathbf{X} \tilde{\boldsymbol{\beta}} + \mathbf{t}^T \mathbf{G} \mathbf{Z}^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X} \tilde{\boldsymbol{\beta}}),$$

where

$$\tilde{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{V}^{-1} \mathbf{Y}.$$

[Hint: You may use the fact that if all inverses exist,

$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{A}_{11}^{-1} + \mathbf{B}_{12} \mathbf{B}_{22}^{-1} \mathbf{B}_{21} & -\mathbf{B}_{12} \mathbf{B}_{22}^{-1} \\ -\mathbf{B}_{22}^{-1} \mathbf{B}_{21} & \mathbf{B}_{22}^{-1} \end{pmatrix},$$

where $\mathbf{B}_{22} = \mathbf{A}_{22} - \mathbf{A}_{21} \mathbf{A}_{11}^{-1} \mathbf{A}_{12}$, $\mathbf{B}_{12} = \mathbf{A}_{11}^{-1} \mathbf{A}_{12}$, and $\mathbf{B}_{21} = \mathbf{A}_{21} \mathbf{A}_{11}^{-1}$.]

- (d) (5 points) Deduce that $\mathbf{G} \mathbf{Z}^T \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X} \tilde{\boldsymbol{\beta}})$ is the BLUP of \mathbf{u} .