The first midterm exam will be a take-home exam which will be made available in the Midterm 1 folder under the documents section of PWeb at 5pm on Friday, March 3rd and due at 11:59pm (uploaded to Gradescope) on Monday, March 6th.

Content. The exam will cover Chapter 7 (7.1 - 7.6) and parts of Chapter 8 (8.1 - 8.3, 8.6) of DeGroot and Schervish's *Probability and Statistics*. There will be some questions that ask you to use R.

Format. The exam is intended to take 2.5 hours to complete, although you may take up to 4 hours to complete it. These 4 hours do not need to be consecutive. You should monitor your own time, and record on the test your estimate for the total amount of time you actively worked on the exam.

Your solutions to the exam should be neatly neatly written or typed. If you scan a handwritten assignment, be sure to review the legibility of your scan on Gradescope after you submit.

Resources. You may use any notes you've taken for this class, your work on any previous homework or daily assignments, lecture notes I've posted on the course website, the recorded lecture video from Monday 2/20 and DeGroot and Schervish's *Probability and Statistics* textbook, as well as Blitzstein's *Introduction to Probability* textbook.

For problems asking you to do analysis or perform computations using R, you may use either a local installation of R or the Grinnell R Studio server, and you may reference any of the R help files (available by typing ?functionname in the console).

You may not use any other resources other than those listed above. If you have questions about whether a resource can be used, you are welcome to message me.

Preparation. The best preparation you can do for the exam is to organize your notes and/or homework to make finding information and examples as quick and efficient as possible. Beyond that, you should attempt to accurately assess what topics you have mastered and which you need to practice more. A good starting point is to review the list of objectives on each daily assignment. Another way to prepare is to create your own study guide with summaries of the important concepts, along with example problems you've designed and solved. Exam problems will be comparable in difficulty to those exhibited in class and assigned for homework. Some exam questions may be similar to problems you have seen before, while others will require you to synthesize your knowledge in new ways.

On the exam, you may be asked to do the following:

- Rephrase a key definition and/or theorem in your own words.
- Determine whether a given statement is true or false.
- Interpret or explain a statistics concept in everyday language.
- Sketch the proof of an important result discussed in class.
- Perform calculations using relevant techniques from the course.
- Provide a short, rigorous proof of a novel statement or result.
- Create and analyze a statistical model for a particular phenomenon.
- Use R to simulate a random phenomenon.

For extra practice, several additional review problems are printed below. Solutions to these problems can be found on the exams page of the course website. While these questions are representative of the typical scope and difficulty of individual exam questions, this review is not comprehensive, nor does it necessarily represent the total amount of time available for the exam.

Practice Problems.

(1) The method of *randomized response* is sometimes used to conduct surveys on sensitive topics. A simple version of the method can be described as follows:

A random sample of n people are drawn from a large population. For each person in the sample, there is probability 1/2 that the person will be asked a standard question and probability 1/2 that the person will be asked a sensitive question. Furthermore, this selection of the standard or sensitive question is made independently from person to person. If a person is asked the standard question, then there is probability 1/2 that the person will give a positive response; however, the person is asked the sensitive question, then there is an unknown probability p that they will give a positive response. The statistician can observe only the total number X of positive responses that were given by the n persons in the sample, but cannot observe which of these persons were asked the sensitive question.

Determine the MLE of p based on the observation X.

Solution. By the Law of Total Probability, probability that a particular person gives a positive response is

$$P(\text{positive}) = P(\text{positive}|\text{sensitive})P(\text{sensitive}) + P(\text{positive}|\text{standard})P(\text{standard})$$
$$= \frac{1}{2} \cdot \frac{1}{2} + p \cdot \frac{1}{2} = \frac{1}{4} + \frac{1}{2}p$$

Therefore, the number of positive responses X in the sample has the $Bin(n,\theta)$ distribution, where $\theta = \frac{1}{4} + \frac{1}{2}p$. Since $0 , then the parameter space for <math>\theta$ is $\frac{1}{2} < \theta < \frac{3}{4}$. By a previous homework problem, the likelihood function for $X \sim Bin(n,\theta)$ has its largest value at $\theta = \frac{x}{n}$. This is our MLE, if $\frac{x}{n}$ is within the parameter space. If not, the MLE corresponds to one of the endpoints of the parameter space interval. In particular, the MLE $\hat{\theta}$ for θ is

$$\hat{\theta} = \begin{cases} \frac{X}{n}, & \text{if } \frac{1}{4} \le \frac{X}{n} \le \frac{3}{4}; \\ \frac{1}{4}, & \text{if } \frac{X}{n} < \frac{1}{4}, \\ \frac{3}{4}, & \text{if } \frac{X}{n} > \frac{3}{4} \end{cases}$$

Now, solving for p in the equation $\theta = \frac{1}{4} + \frac{1}{2}p$ gives

$$p = 2\left(\theta - \frac{1}{4}\right)$$

Therefore, by the invariance property of the MLE, the MLE \hat{p} for p is

$$\hat{p} = 2\left(\hat{\theta} - \frac{1}{4}\right).$$

(2) Suppose that an observation X is from from a distribution with pdf

$$f(x|\theta) = \begin{cases} \frac{1}{\theta}, & \text{for } 0 < x < \theta, \\ 0, & \text{otherwise.} \end{cases}$$

Also, suppose that the prior pdf of θ is

$$\xi(\theta) = \begin{cases} \theta e^{-\theta}, & \text{for } \theta > 0, \\ 0, & \text{otherwise.} \end{cases}$$

Determine the Bayes estimator of θ with respect to:

- (a) The mean squared error loss function;
- (b) the absolute error loss function.

Solution. We first compute the posterior distribution. Suppose X = x. Using Bayes rule,

$$\xi(\theta|x) \propto f(x|\theta)\xi(\theta) = e^{-\theta} \quad \theta > x$$

This looks much like an Expo(1) distribution, except with support (x, ∞) instead of $(0, \infty)$. In particular, if $Y \sim \text{Expo}(1)$, then $\theta | x = Y + x$.

(a) Using squared loss, the Bayes estimate is the mean of the posterior distribution: $\hat{\theta} = E[\theta|x]$. In this case, since $\theta | x = Y + x$, then by linearity of conditional expectation, the posterior mean \mathbf{is}

$$\hat{\theta}(x) = E[\theta|x] = E[Y|x] + E[X|x] = E[Y] + x = 1 + x$$

Replacing x with X, the Bayes Estimator is $E[\theta|X] = 1 + X$.

(b) Using absolute loss, the Bayes estimate is the median of the posterior distribution. Since median of Y + x is equal to the median of Y, plus x, and since the median of $Y \sim \text{Expo}(1)$ is $\ln 2$, then the posterior median is

$$\theta(x) = \ln 2 + x.$$

Replacing x with X, the Bayes Estimator is $\hat{\theta}(X) = \ln 2 + X$.

(3) Suppose that the random variable X has a binomial distribution with unknown value n and a known value of p for 0 . Determine the MLE of n based on the observation X. Hint:consider the ratio e/ 1

$$\frac{f(x|n+1,p)}{f(x|n,p)}.$$

Solution. The likelihood function is

$$f(x|n,p) = \binom{n}{x} p^x (1-p)^{n-x}$$

To find the maximum value of the likelihood function, we consider the ratio of successive values of n: (n+1)!

$$\frac{f(x|n+1,p)}{f(x|n,p)} = \frac{\binom{n+1}{x}p^x(1-p)^{n+1-x}}{\binom{n}{x}p^x(1-p)^{n-x}} = \frac{\frac{(n+1)!}{x!(n+1-x)!}}{\frac{n!}{x!(n-x)!}}(1-p) = \frac{n+1}{n+1-x}(1-p)$$

Note that this ratio is a decreasing function of n, and so the maximum of the likelihood function is attained at the smallest value of n for which the ratio is less than 1. Setting the ratio equal to 1 and solving for n, we obtain:

$$n = \frac{x}{p} - 1$$

Thus, the MLE for n is the smallest integer greater than $\frac{x}{p} - 1$. If $\frac{x}{p} - 1$ is itself an integer, than both $\frac{x}{p} - 1$ and $\frac{x}{p}$ are MLEs.

(4) Suppose X_1, \ldots, X_m form a random sample from the Normal distribution with mean μ_1 and variance σ^2 and that Y_1, \ldots, Y_m form an independent sample from the Normal distribution with mean μ_2 and variance σ^2 . Let

$$S_X^2 = \sum_{i=1}^m (X_i - \bar{X})^2 \qquad S_Y^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

- (a) For what values of (α, β) is S² = αS²_X + βS²_Y an **unbiased** estimator of σ²?
 (b) Determine the values of α and β for which αS²_X + βS²_Y will be an unbiased estimator with minimum variance.

Solution. (a) Computing the mean of the estimator,

$$E[\alpha S_X^2 + \beta S_Y^2] = \alpha E[S_X^2] + \beta E[S_Y^2] = \alpha (m-1)\sigma^2 + \beta (n-1)\sigma^2$$

since $S_X^2/\sigma^2 \sim \chi^2(m-1)$ and $S_Y^2/\sigma^2 \sim \chi^2(n-1)$. Hence, this estimator is unbiased exactly when

$$\alpha(m-1) + \beta(n-1) = 1$$

Or by solving for α , when

$$\alpha = \frac{1}{m-1} - \beta \frac{(n-1)}{m-1}$$

(b) Since S_X^2 and S_Y^2 are independent, then

$$\operatorname{Var}(\alpha S_X^2 + \beta S_Y^2) = \alpha^2 \operatorname{Var}(S_X^2) + \beta^2 \operatorname{Var}(S_Y^2)$$
$$= 2\alpha^2 (m-1)\sigma^4 + 2\beta^2 (n-1)\sigma^4$$
$$= 2\sigma^4 [(m-1)\alpha^2 + (n-1)\beta^2]$$

Making the substitution from part (a),

$$\operatorname{Var}(\alpha S_X^2 + \beta S_Y^2) = 2\sigma^4 \left[(m-1) \left(\frac{1}{m-1} - \beta \frac{(n-1)}{m-1} \right)^2 + (n-1)\beta^2 \right]$$

Differentiating with respect to β and setting the result equal to 0, we get:

$$\beta = \frac{1}{n+m-2}$$

And substituting this into the equation for α :

$$\alpha = \frac{1}{n+m-2}$$

(5) Let X_1, \ldots, X_n be iid random variables with density function

$$f(x|\sigma) = \frac{1}{2\sigma} e^{-|x|/\sigma}$$

where σ is an unknown parameter with $\sigma > 0$.

- (a) Find a formula for the MLE of σ .
- (b) Find a formula for the Method of Moments estimator of σ .
- (c) Use R to simulate the sampling distribution for the MLE estimator and for the Method of Moments estimator when $\sigma = 2$ and n = 10. Use your simulation to estimate the mean and variance of each estimator.

Solution. (a) The likelihood function is

$$f(\mathbf{x}|\sigma) = \frac{1}{(2\sigma)^n} e^{-\frac{1}{\sigma}\sum |x_i|}$$

and the log-likelihood function is

$$\log f(\mathbf{x}|\sigma) = -n\log(2\sigma) - \frac{1}{\sigma}\sum |x_i|$$

Differentiating with respect to σ :

$$\frac{\partial}{\partial \sigma} \log f(\mathbf{x}|\sigma) = \frac{-n}{\sigma} + \frac{1}{\sigma^2} \sum |x_i|$$

Which is 0 when

$$\sigma = \frac{1}{n} \sum |x_i|$$

Taking the second derivative:

$$\frac{\partial^2}{\partial \sigma^2} \log f(\mathbf{x}|\sigma) = \frac{n}{\sigma^2} - \frac{2}{\sigma^3} \sum |x_i|$$

which is negative at $\sigma = \frac{1}{n} \sum |x_i|$. Hence, we have found a local maximum and so the MLE for σ is

$$\hat{\sigma} = \frac{1}{n} \sum |x_i|$$

(b) First, note that the mean of X_i is 0, and so we need to find the second moment of X_i :

$$E[X_i^2] = \int_{-\infty}^{\infty} x^2 \frac{1}{2\sigma} e^{-|x|/\sigma} \, dx = 2 \int_0^{\infty} x^2 \frac{1}{2\sigma} e^{-x/\sigma} \, dx$$

where the last equality follows by symmetry. But we can recognize this final expression as the second moment of an $\text{Expo}(1/\sigma)$ variable, and so

$$E[X_i^2] = 2\sigma^2$$

Solving for σ , and replacing $E[X_i^2]$ with the second sample moment, our method of moments estimator is

$$\hat{\sigma} = \sqrt{\frac{1}{2n} \sum X_i^2}$$

(c) First, we note that the PDF of X_i corresponds to the PDF of the variable $J \cdot Y$, where $J \sim \text{DUnif}\{-1, 1\}$ and $Y \sim \text{Expo}(1/\sigma)$ with J and Y independent.

We can use the following code to simulate 10000 experiments when σ = 2 and n = 10. n <- 10 sigma <- 2 trials <- 10000 mle <- rep(0,trials) mom <- rep(0, trials) for (i in 1:trials){ J = sample(c(-1,1), size = n, replace = T)

```
Y = rexp(n, rate = 1/sigma)
X = J*Y
mle[i] = mean(abs(X))
mom[i] = sqrt(1/(2*n)*sum(X^2))
}
The mean and variance of the MLE is
mean(mle)
# 2.005809
var(mle)
# 0.4063677
while the mean and variance of the MoM is
mean(mom)
# 1.900474
var(mom)
# 0.4199017
```

(6) The *Pareto* distribution with shape θ and minimum value 1 has PDF:

$$f(x) = \frac{\theta}{x^{\theta+1}} \qquad x > 1$$

Show that the family of Gamma distributions $\text{Gamma}(\alpha, \beta)$ is a conjugate family of prior distributions for θ , when samples are taken from a Pareto distribution with shape θ and minimum value 1.

Solution. Suppose X_1, \ldots, X_n are samples from the Pareto distribution with shape θ and minimum value 1, and let $\theta \sim \text{Gamma}(\alpha, \beta)$. Then the posterior distribution is

$$\xi(\theta|\mathbf{x}) \propto f(\mathbf{x}|\theta)\xi(\theta) = f(x_1|\theta) \cdots f(x_n|\theta) \frac{\beta^{\alpha}}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}$$
$$= \frac{\theta^n}{x_1^{\theta+1} \cdots x_n^{\theta+1}} \frac{\beta^{\alpha}}{\Gamma(\alpha)} \theta^{\alpha-1} e^{-\beta\theta}$$
$$\propto (x_1 \cdots x_n)^{-\theta} \theta^{\alpha+n-1} e^{-\beta\theta}$$
$$= e^{-\theta \log(x_1 \cdots x_n)} \theta^{\alpha+n-1} e^{-\beta\theta}$$
$$= \theta^{\alpha+n-1} e^{-(\beta+\log(x_1 \cdots x_n))\theta}$$

which we recognize as proportional to a $\text{Gamma}(\alpha + n, \beta + \log(x_1 \cdots x_n))$ distribution. Hence, the family of Gamma distributions is indeed a conjugate prior family of distributions for θ .