

# CSE 546 HW1

DAVID FLEMING

## CONTENTS

Question 1	1
1.1	1
1.2	2
1.3	2
Question 2	3
2.1	3
2.2	3
2.3	3
Question 3	3
Question 4	3
4.1: Too Small a $\lambda$	3
4.2: Too Large a $\lambda$	4

## QUESTION 1

**1.1.** The expectation value for  $\max(X_1, X_2)$  is given by

$$(1) \quad E[X] = \int \max(X_1, X_2) f(x) dx_1 dx_2$$

where  $f(x) = 1$  is the PDF for the uniform random variates. We consider two regions, one below the line defined by  $X = X_1 = X_2$  and one above said line for our integrations. This gives us the following integral:

$$\begin{aligned} E[X] &= \int_0^1 \int_{x_2}^1 x_1 dx_1 dx_2 + \int_0^1 \int_{x_1}^1 x_2 dx_2 dx_1 \\ (2) \quad &= \int_0^1 \left( \frac{1}{2} - \frac{1}{2} x_2^2 \right) dx_2 + \int_0^1 \left( \frac{1}{2} - \frac{1}{2} x_1^2 \right) dx_1 \\ &= \left( \frac{1}{2} x_2 - \frac{1}{6} x_2^3 \right) \Big|_0^1 + \left( \frac{1}{2} x_1 - \frac{1}{6} x_1^3 \right) \Big|_0^1 \\ &= \frac{2}{3} \end{aligned}$$

**1.2.** Since

$$(3) \quad \text{Var}[X] = E[X^2] - E[X]^2$$

and we now know  $E[X] = 2/3$ , we solve the integral presented above but with the integrand squared as follows

$$(4) \quad \begin{aligned} E[X^2] &= \int_0^1 \int_{x_2}^1 x_1^2 dx_1 dx_2 + \int_0^1 \int_{x_1}^1 x_2^2 dx_2 dx_1 \\ &= \int_0^1 \left(\frac{1}{3} - \frac{1}{3}x_2^3\right) dx_2 + \int_0^1 \left(\frac{1}{3} - \frac{1}{3}x_1^3\right) dx_1 \\ &= \left(\frac{1}{3}x_2 - \frac{1}{12}x_2^4\right)\Big|_0^1 + \left(\frac{1}{3}x_1 - \frac{1}{12}x_1^4\right)\Big|_0^1 \\ &= \frac{1}{2} \end{aligned}$$

which when combined with the definition for  $\text{Var}[X]$  gives

$$(5) \quad \text{Var}[X] = \frac{1}{2} - \left(\frac{2}{3}\right)^2 = \frac{1}{18}$$

**1.3.** Since

$$(6) \quad \text{Cov}[X, X_1] = E[XX_1] - E[X]E[X_1]$$

and we know  $E[X]$  from 1.1 and  $E[X_1] = 1/2$  is the trivial result for the uniform distribution from  $[0, 1]$ , we need only calculate the first term as follows:

$$(7) \quad \begin{aligned} E[XX_1] &= \int_0^1 \int_{x_2}^1 x_1^2 dx_1 dx_2 + \int_0^1 \int_{x_1}^1 x_2 x_1 dx_2 dx_1 \\ &= \int_0^1 \left(\frac{1}{3} - \frac{1}{3}x_2^3\right) dx_2 + \int_0^1 x_1 \left(\frac{1}{2} - \frac{1}{2}x_1^2\right) dx_1 \\ &= \left(\frac{1}{3}x_2 - \frac{1}{12}x_2^4\right)\Big|_0^1 + \left(\frac{1}{4}x_1^2 - \frac{1}{8}x_1^4\right)\Big|_0^1 \\ &= \frac{3}{8} \end{aligned}$$

giving us

$$(8) \quad \text{Cov}[X, X_1] = E[XX_1] - E[X]E[X_1] = \frac{3}{8} - \frac{2}{3} \cdot \frac{1}{2} = \frac{1}{24}$$

Note: For this question, I collaborated with Matt Wilde.

## QUESTION 2

**2.1.** For the log-likelihood of  $G$  given  $\lambda$  and i.i.d. samples, we have

$$\begin{aligned}
 LL &= \log(P(G|\theta)) \\
 &= \prod_i^n P(G_i|\theta) \\
 (9) \quad &= \log\left(\frac{\lambda^{\sum_i^n k_i} e^{-\lambda n}}{\prod_i^n k_i!}\right) \\
 &= \sum_i^n k_i \log \lambda - \lambda n - \log \prod_i^n k_i!
 \end{aligned}$$

**2.2.** To compute the MLE for  $\lambda$  in general,

$$(10) \quad \frac{\partial}{\partial \lambda} [LL] = 0$$

which yields

$$\begin{aligned}
 0 &= \frac{\partial}{\partial \lambda} (\sum_i^n k_i \log \lambda - \lambda n - \log \prod_i^n k_i!) \\
 (11) \quad &= \frac{\sum_i^n k_i}{\lambda} - n \\
 \hat{\lambda}_{MLE} &= \frac{\sum_i^n k_i}{n}
 \end{aligned}$$

**2.3.** For the observed set  $G$ , I use Eqn. 11 to compute  $\lambda_{MLE}$  as

$$(12) \quad \hat{\lambda}_{MLE} = \frac{4 + 1 + 3 + 5 + 5 + 1 + 3 + 8}{8} = 3.75$$

Note: For this question, I collaborated with Matt Wilde.

## QUESTION 3

TODO

## QUESTION 4

Note: Both both subquestions 1 and 2, I assume that the too small or too large  $\lambda$ s bias the model to a too complex or too simple model relative to the optimum model complexity, respectively.

**4.1: Too Small a  $\lambda$ .**

*4.1.a.* For both LASSO and Ridge Regression (RR), the error on the training set would decrease as the penalty term for both regression techniques would be effectively negligible. With small  $\lambda$ , the regularization penalty is also small causing both regression techniques to tend towards least squares regression (LSR) and allow for more complex models which overfit the training data and hence leading to smaller training set error. For example, too small of a  $\lambda$  could push both LASSO and RR to favor a high-order polynomial model when the underlying data is linear.

*4.1.b.* For both LASSO and RR, too small a  $\lambda$  leads to overfitting on the training set. With a model overfit on the training set, it will do a poor job of generalizing to new data and hence will poorly fit the testing set leading to larger testing error.

*4.1.c.* For both LASSO and RR, too small a  $\lambda$  yields too small of a complexity penalty causing both algorithms to tend towards the LSR solution. In this case,  $\hat{\omega}$  could get large via overfitting. RR, however, will likely predict larger  $\hat{\omega}$  than LASSO as RR's  $l_2$  norm penalty primarily seeks to control the magnitude of the weight vector, hence biasing towards larger values in this case as the penalty decreases, while LASSO seeks a sparse solution.

*4.1.d.* With LASSO, too small of a  $\lambda$  would yield more non-zero parameters since a smaller regularization penalty will prevent many of the elements of  $\hat{\omega}$  from being set to 0 under LASSO's sharp  $l_1$  norm penalty which tries to make  $\hat{\omega}$  sparse. For RR, too small of a  $\lambda$  would have little effect on the number of non-zero parameters as RR's regularization penalty deals with constraining the magnitude of  $\hat{\omega}$  as opposed to forcing it to be sparse.

## 4.2: Too Large a $\lambda$ .

*4.2.a.* For both regression techniques, too large of a  $\lambda$  will yield larger training set error as the regularization penalty will select for simpler models that could poorly fit the data. For example, too large of a  $\lambda$  could push both LASSO and RR to favor a linear model when the underlying data is quadratic.

*4.2.b.* For both LASSO and RR, too large of a  $\lambda$  will yield larger testing set error since the models again will be biased towards lower-than-optimal model complexity and hence will likely yield a poor fit to future data similar to why the error on the training set is also larger in this case.

*4.2.c.* For both LASSO and RR, too large of a  $\lambda$  will yield smaller  $\hat{\omega}$  as the larger regularization penalty restricts the weight vector. RR, however, will likely predict smaller  $\hat{\omega}$  than LASSO as RR's  $l_2$  norm penalty primarily seeks to control the magnitude of the weight vector, hence biasing towards smaller values in this case, while LASSO seeks a sparse solution.

4.2.d. For LASSO, a larger  $\lambda$  will yield a sparser solution via its  $l_1$  norm regularization penalty and hence will have fewer nonzero elements in  $\hat{\omega}$  than RR. In the large  $\lambda$  limit, RR will also yield fewer nonzero elements in  $\hat{\omega}$  only because its large regularization penalty makes coefficients in its weight vector small and hence potentially pushing some towards 0.