## Final Exam

## Name:

- 1. (4 points) If you believe that only a small number of variables are important, which of the following regression approaches are appropriate ways to select the important variables:
  - i. Ridge regression
  - ii. Lasso regression
  - iii. p-values.
  - iv. Forward stepwise regression
- 2. (4 points) The  $R^2$  in the linear regression  $Y = \beta_0 + \beta_1 X + \epsilon$  from data  $(y_1, x_1), \dots, (y_n, x_n)$  represents which of the following:
  - i. the fraction of variance Y explained by  $\beta_0 + \beta_1 X$
  - ii. the expected prediction error for new data  $(y_{n+1}, x_{n+1}), \dots, (y_{2n}, x_{2n})$
  - iii. the square of the correlation between Y and  $\hat{\beta}_0 + \hat{\beta}_1 X$
- 3. (4 points) Identify and describe at least one problem or confounding factor in the following experimental design.

Arabadopsis thaliana is a flowering plant native to Europe. Strains in northern Europe tend to have an earlier flowering time than strains in southern Europe. To identify epigenetic factors that may contribute to the difference, researchers took samples from Tunisia and Norway. The samples were processed for bisulfite sequencing on site and then sent to a central sequencing center for processing and analysis.

4. (8 points) Suppose we have data on a response variable Y and three predictor variables  $X_1$  and  $X_2$ . We fit the models on all possible subsets. Circle the model that will be selected according the following criteria.

criterion	$X_1$	$X_2$	$X_1, X_2$
AIC	126.86	106.57	106.59
BIC	129.85	109.56	110.57
adj $R^2$	-0.023	0.63	0.64
10-fold CV	402.56	131.34	151.99

5. Using a subsample of the Structural Survey of Wages for Spain in 2006, the following model is estimated to explain wages (in tens of thousands of dollars) as a function of education (in years:

Table 2: Fitting linear model: wages  $\sim$  education

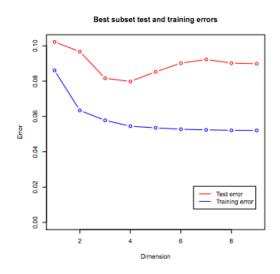
	Estimate	Std. Error	t value	Pr(> t )
$\begin{array}{c} \textbf{education} \\ \textbf{(Intercept)} \\ R^2 \end{array}$	0.0527 1.918 0.2445	0.0121 $0.7155$	5.762 5.268	1.33e-05 7.35e-03

a. (4 points) What is the interpretation of the education coefficient?

b. (4 points) Assuming the linear model is true, what is the expected wages for someone with 16 years of education?

c. (4 points) What is percentage of variance in wages explained by education?

6. Suppose we have 9 predictor variables and 1 response variable. We run best subset regression and obtain the following estimates of training and test error for the number of predictor variables ranging from 0 to 9 (labeled dimension).



- a. (4 points) What is the optimal number of predictors according to the training error?
- b. (4 points) What is the optimal number of predictors according to the test error?

7.	In species sampling experiments, the number of individuals captured from each species is thought to be Poisson distributed with pmf $f(x) = \Pr(X = x) = \lambda^x e^{-\lambda}/x!$ . Unfortunately, species with no captured individuals are not observed. In this case, the number of individuals captured for each observed species (those with at least one observed individual) follows a zero-truncated Poisson distribution with pmf $f(x) = \Pr(X = x   X > 0)$ . Suppose we observed 10 species with counts $x_1, \ldots, x_{10}$ .
	a. (4 points) If $\lambda = 2$ , what is the probability $X = 2$ ?
	b. (4 points) Write out the log-likelihood function.
	c. (4 points) Set up the maximum likelihood equations. Do not solve.
	d. (4 points) Calculate the expectation and variance of a zero-truncted Poisson distribution with parameter $\lambda$ .
	e. (4 points) Set up the Method of Moments equations. Do not solve.
	f. (3 points) Suppose we fit a zero-truncated Poisson and a zero truncated Negative Binomial distribution and obtain the following AIC values: 123.2 and 127.5953, respectively. Which model would be chosen by the AIC criterion?

per s and	second. To temperature specific		ers are getting s with an esti	at least the imated stand	advertised sp	peed, regulator	eeds of 1 gigabits
b.	(5 points) C	onstruct a $95\%$	confidence in	terval for th	ne average tra	ansmission spe	eed.
c.		an you reject th or why not you			9% confidenc	e level? Comp	at the  p-value o
d.	of $0.2 \text{ gb/s}$ .		proximate po	wer of the t	-	~ /	tandard deviation ution of $N(1,0.2)$

9.	To bid for government contracts, businesses submit bids and the government selects the lowest bid that satisfies the criteria. Suppose that bids are exponentially distributed with average bid equal to \$ 16 million.
	a. (4 points) What is the probability the first bid is below \$ 3 million?
	b. (5 points) If there are 10 bids, what is the expected value of the lowest bid?
	c. (5 points) If the government decides to take bids one by one until one is below \$ 3 million. How many bids will the government have to wait on average?
	d. (5 points) If you want to enter the auction (for 11 bids total, 1 being yours and 10 other bids) what bid would you have to submit to ensure a 90% probability of having the lowest bid?

```
Binomial(n, p):
pmf: f(x) = \binom{n}{r} p^x (1-p)^{n-x}
expectation: np
variance: np(1-p)
Poisson(\lambda):
pmf: f(x) = \lambda^x e^{-\lambda}/x!
expectation: \lambda
variance: \lambda
Geometric(p):
pmf: f(x) = (1 - p)^x p
expectation: \frac{1}{p} - 1
variance: \frac{1-p}{p^2}
Negative Binomial(r, p):

pmf: f(x) = {x+r-1 \choose 1-p} p^x (1-p)^r

expectation: \frac{pr}{(1-p)^2}

variance: \frac{pr}{(1-p)^2}
\begin{aligned} & \text{Hypergeometric}(N_1, N_2, n); \\ & \text{pmf: } f(x) = \frac{\binom{N_1}{N}\binom{N_2}{n-x}}{\binom{N_1+N_2}{N_1+N_2}} \\ & \text{expecation: } nN_1/(N_1+N_2) \\ & \text{variance: } n\frac{N_1}{N_1+N_2}\frac{N_2}{N_1+N_2}\frac{N_1+N_2-n}{N_1+N_2-1} \end{aligned}
Uniform(a, b):
pdf: f(x) = 1/(b-a)
cdf: F(x) = \begin{cases} 0 & x \le a \\ (x-a)/(b-a) & a < x < b \\ 1 & x > b \end{cases}
expectation: (a+b)/2
variance: (b-a)^2/12
Exponential(\lambda):
pdf: f(x) = \lambda e^{-\lambda x}, x > 0
cdf: F(x) = 1 - e^{-\lambda x}, x > 0
expectation: \lambda^{-1}
variance: \lambda^{-2}
Gamma(k, \theta):
pdf: f(x) = x^{k-1}e^{-x/\theta}/(\Gamma(k)\theta^k), x > 0
expectation: k\theta
variance: k\theta^2
Beta(\alpha, \beta):
pdf: f(x) = x^{\alpha-1}(1-x)^{\beta-1}/\mathcal{B}(\alpha,\beta), 0 < x < 1
expectation: \alpha/(\alpha+\beta)
variance: \alpha\beta/((\alpha+\beta)^2(\alpha+\beta+1))
```

 $\begin{array}{l} \operatorname{Normal}(\mu,\sigma^2) \colon \\ \operatorname{pdf:} \ f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/\sigma^2} \\ \operatorname{expectation:} \ \mu \end{array}$ 

variance:  $\sigma^2$ 

Standard normal table for  $z_{\alpha}$  (or  $z_{1-\alpha}$  depending on the notation)