

# Practice Final Exam Solutions

1. In RNA sequencing experiments, the counts of expressed genes can be modeled as a Poisson distribution with pmf  $f(x) = \Pr(X = x) = \lambda^x e^{-\lambda} / x!$ . On the other hand, non-expressed genes have zero counts. This creates an issue in modeling the data, as genes with zero counts could be expressed and the zero is due to random sample or they could just not be expressed. To accommodate this, we can model all of the counts as a zero inflated Poisson distribution. To do this, we assume that for each gene there is some probability  $p$  of being expressed. Conditional upon being expressed, the counts follow a Poisson distribution with parameter  $\lambda$ . If the gene is not expressed, with probability  $1 - p$ , then the count is zero. Suppose we have 100 we investigate with counts  $x_1, \dots, x_{100}$ .

- a. If  $\lambda = 2$  and  $p = 0.25$ , calculate the probability that  $X = 0$ .

The pmf of  $X$  is  $f(x) = p1(x = 0) + (1 - p)\lambda^x e^{-\lambda} / x!$ . Therefore  $\Pr(X = 0) = p + (1 - p)e^{-\lambda} = 0.25 + 0.75 * e^{-2}$

- b. Write out the log-likelihood function.

The likelihood is equal to

$$L(\lambda, p | x_1, \dots, x_{100}) = \prod_{i=1}^{100} (p1(x_i = 0) + (1 - p)\lambda^{x_i} e^{-\lambda} / x_i!).$$

Therefore the log likelihood is equal to

$$\log L(\lambda, p | x_1, \dots, x_{100}) = \sum_{i=1}^{100} \log (p1(x_i = 0) + (1 - p)\lambda^{x_i} e^{-\lambda} / x_i!).$$

There is not much we can do to simplify this because we can not break up the sum inside the logarithm.

- c. Set up the maximum likelihood equations. Do not solve.

There are two parameters, so there will be two maximum likelihood equations.

$$\frac{\partial}{\partial p} \log L = \sum_{i=1}^{100} \frac{1}{p1(x_i = 0) + (1 - p)\lambda^{x_i} e^{-\lambda} / x_i!} (1(x_i = 0) - \lambda^{x_i} e^{-\lambda} / x_i!) = 0$$

$$\frac{\partial}{\partial \lambda} \log L = \sum_{i=1}^{100} \frac{1}{p1(x_i = 0) + (1 - p)\lambda^{x_i} e^{-\lambda} / x_i!} (x_i \lambda^{x_i-1} e^{-\lambda} / x_i! - e^{-\lambda} \lambda^{x_i} / x_i!) = 0.$$

- d. Calculate the expectation and variance of a zero-inflated Poisson distribution with parameters  $p$  and  $\lambda$ .

The expectation is equal to

$$E(X) = \sum_{x=0}^{\infty} x f(x) = \sum_{x=0}^{\infty} x (p1(x_i = 0) + (1 - p)\lambda^{x_i} e^{-\lambda} / x_i!)$$

. Note that the first part of the sum above is always zero, either  $x = 0$  or  $x \neq 0$  and the indicator function is zero. Therefore

$$E(X) = \sum_{x=0}^{\infty} x (1 - p) \lambda^x e^{-\lambda} / x! = (1 - p) \sum_{x=0}^{\infty} x \lambda^x e^{-\lambda} / x!.$$

Note that the last sum is the expected value of a Poisson random variable, equal to  $\lambda$  and therefore  $E(X) = (1 - p)\lambda$ .

To compute the second moment we do the same thing.

$$E(X^2) = \sum_{x=0}^{\infty} x^2 f(x) = \sum_{x=0}^{\infty} x^2 (p1(x_i = 0) + (1-p)\lambda^{x_i} e^{-\lambda}/x_i!) = (1-p) \sum_{x=0}^{\infty} x^2 \lambda^x e^{-\lambda}/x!.$$

The portion  $\sum_{x=0}^{\infty} x^2 \lambda^x e^{-\lambda}/x!$  is equal to the second moment of a Poisson random variable, which is equal to  $\lambda + \lambda^2$ , therefore  $E(X) = (1-p)(\lambda + \lambda^2)$ .

e. Set up the Method of Moments equations. Do not solve.

Two parameters means two equations.

$$E(X) = (1-p)\lambda = \bar{x}$$

$$E(X^2) = (1-p)(\lambda + \lambda^2) = \bar{x^2}$$

with the convention that  $\bar{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2$ , the average of the squared observations.

2. A survey was given to a random sample of college students to assess the support of gay marriage. 242 females and 207 males were surveyed. 176 females expressed support while 134 males expressed.

- a. Construct 95% confidence intervals for the percentage of support for gay marriage among college females and college males, separately.

For proportions, 95% confidence intervals are given by  $\hat{p} \pm z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p})/\sqrt{n}}$ .

$\hat{p}_f = \frac{176}{242}$  and so a 95% confidence interval will be given by  $(0.7272727 - 1.96 \cdot 0.445/15.56, 0.7272727 + 1.96 \cdot 0.445/15.56) = (0.671, 0.783)$

$\hat{p}_m = \frac{134}{207}$  and so a 95% confidence interval will be given by  $(0.647 - 1.96 \cdot 0.478/14.39, 0.647 + 1.96 \cdot 0.478/14.39) = (0.582, 0.712)$ .

- b. Construct a 95% confidence interval for the difference between the two proportions. Use the fact that the pooled standard error is  $\sqrt{\hat{p}(1-\hat{p})(\frac{1}{n_1} + \frac{1}{n_2})}$  with  $\hat{p}$  equal to the pooled proportion.

The 95% confidence interval for two proportions is given by  $(\hat{p}_f - \hat{p}_m) \pm z_{\alpha/2} s.e.m$  with  $s.e.m$  given above. Therefore a 95% confidence interval is given by

$$\left(\frac{176}{242} - \frac{134}{207}\right) \pm 1.96 * \sqrt{\frac{310}{449} \left(1 - \frac{310}{449}\right) \cdot \left(\frac{1}{242} + \frac{1}{207}\right)} = (-0.00585, 0.1657).$$

- c. If we were to do a hypothesis test to test for a difference between the two proportions, what would be the null and alternative hypotheses?

null:  $H_0 : p_f = p_m$  or  $p_f - p_m = 0$  alternative:  $H_1 : p_f \neq p_m$  or  $p_f - p_m \neq 0$

- d. Is there a statistically significant difference between the two proportions at the 95% confidence level?

The above 95% confidence interval for the two proportions contains 0 so we cannot reject the null hypothesis at the 95% confidence level.

3. Suppose that  $Y, X_1, X_2$  and  $\epsilon$  are all Gaussian (aka Normal) random variables and that  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$ .

a. If  $X_1, X_2, X_3$ , and  $\epsilon$  are all independent standard normal random variables (mean 0 and variance 1) and  $\beta_0 = 10, \beta_1 = 2, \beta_2 = 5$ , and  $\beta_3 = 1$ , what is the variance of  $Y$ ?

$\text{Var}(Y) = \text{Var}(10 + 2X_1 + 5X_2 + X_3 + \epsilon)$ . Since  $X_1, X_2, X_3$ , and  $\epsilon$  are all independent, we can use the fact that the variance of the sum of independent random variables is the sum of the variance.

$$\text{Var}(Y) = 4\text{Var}(X_1) + 25\text{Var}(X_2) + \text{Var}(X_3) + \text{Var}(\epsilon) = 4 + 25 + 1 + 1 = 31.$$

b. If  $X_1, X_2$ , and  $\epsilon$  are all independent standard normal random variables (mean 0 and variance 1) and  $\beta_0 = 10, \beta_1 = 2$ , and  $\beta_2 = 5$ , what is the correlation of  $X_1$  and  $Y$ ?

The correlation of  $X_1$  and  $Y$  equals the covariance of  $X_1$  and  $Y$  divided by the product of the respective standard deviations. The covariance of  $X_1$  and  $Y$  equals

$$\sigma_{X_1, Y} = E((X_1 - EX_1) \cdot (Y - E(Y))) = E((X_1 - EX_1) \cdot (10 + 2X_1 + 5X_2 + X_3 + \epsilon - E(10 + 2X_1 + 5X_2 + X_3 + \epsilon))).$$

Note that when we multiply things out, anything involving a product of  $X_1$  and anything that's not  $X_1$  will disappear due to independence.

$$\begin{aligned} &= E((X_1 - EX_1) \cdot 2(X_1 - EX_1) + (X_1 - EX_1) \cdot 5(X_2 - EX_2) + (X_1 - EX_1) \cdot (X_3 - EX_3) + (X_1 - EX_1) \cdot (\epsilon - E(\epsilon))) \\ &= E((X_1 - EX_1) \cdot 2(X_1 - EX_1)) = 2\text{Var}(X_1) = 2. \end{aligned}$$

Therefore the correlation is equal to  $2/(1 \cdot \sqrt{31}) = 0.3592$ .

c. We have 20 total observations of  $Y, X_1, X_2$ , and  $X_3$ . We fit linear models using  $X_1$  alone,  $X_1$  and  $X_2$ , and all three  $X$ 's. State which of the three models will be chosen according to each of the below criteria.

criterion	$X_1$	$X_1, X_2$	$X_1, X_2, X_3$
BIC	<b>125.5</b>	128.2	128.1
AIC	<b>122.5</b>	124.2	123.1
adj $R^2$	-0.003	-0.045	<b>0.048</b>
5-fold CV	135.9	162.3	<b>104.2</b>

The selections are indicated in bold. We want the lowest BIC, AIC, and cross validation error while we want the highest adjusted  $R^2$ .

d. Below is a table for fitting the model  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ . What is the interpretation of the coefficient for  $X_2$ ?

Table 2: Fitting linear model:  $Y \sim X_1 + X_2$

	Estimate	Std. Error	t value	Pr(> t )
$X_1$	1.969	0.4229	4.657	0.000226
$X_2$	4.747	0.2657	17.86	1.878e-12
<b>(Intercept)</b>	10.11	0.2968	34.07	4.337e-17

The interpretation for the coefficient for  $X_2$  is that a unit increase in  $X_2$  will result in a 4.747 increase in  $Y$  on average, holding  $X_1$  constant.

4. In a manufactured product you want to be certain that the delivered product is the weight advertised. The goal is 100 grams.

a. To test the process you weigh 5 products and get weights of 106, 104, 99, 95, and 97. Calculate the sample mean and standard deviation.

$$\bar{x} = (106 + 104 + 99 + 95 + 97)/5 = 100.2$$

$$s = \sqrt{((106 - 100.2)^2 + (104 - 100.2)^2 + (99 - 100.2)^2 + (95 - 100.2)^2 + (97 - 100.2)^2)/5} = 4.167$$

b. Obviously 5 samples is not sufficient so you take 20 more samples (for 25 total) and get a sample mean of 101.84 and sample standard deviation of 3.3. State the null and alternative hypotheses and then compute a  $p$ -value under the null distribution.

We have said nothing on whether the manufacturer cares about too little or too much material. Therefore the null is two sided.

$$H_0 : \mu = 100, H_1 \mu \neq 100.$$

The test statistic is given by  $(\bar{x} - \mu_0)/(s/\sqrt{n}) = (101.84 - 100)/(3.3/\sqrt{25}) = 2.7879$ . The corresponding  $p$ -value is given by  $\Pr(|Z| > 2.7879) = 2\Pr(Z > 2.7879) = 0.0053$ .

c. If the average weight is in fact 101 with variance equal to 4, what is the power with 25 observations compared to null distribution of  $N(100, 4)$  at the 95% confidence level?

Under the null distribution  $t = \frac{\bar{x}-100}{2/\sqrt{n}} \sim N(0, 1)$  and under the alternative  $t \sim N(\frac{101-100}{2/\sqrt{25}} = 2.5, 1)$ . Under the alternative, the probability that  $t > z_{\alpha/2} = 1.96$  is equal to  $\Pr(t > z_{\alpha/2}) = \Pr(z > 1.96 - 2.5) = 0.705$ .

5. Suppose that we want to estimate the size of the wolf population in Yellowstone. One method to do this is to go out, capture wolves, tag them, and release them. You then go out later and capture more wolves, noting how many are recaptured (this is known as mark-recapture). Suppose that there are in fact 400 wolves in Yellowstone. In the first round you capture and tag 50 wolves.

a. Suppose you capture another 50 wolves in the second round. What is the expected number of wolves you will capture?

The number of wolves captured in the second round is a hypergeometric random variable with  $N = 400$ ,  $K = 50$ , and  $n = 50$ . Therefore the expected number of captured wolves is  $50 \cdot 50/400 = 6.25$ .

b. Suppose you capture another 50 wolves in the second round. What is the probability you capture 10 tagged wolves in the second round of captures?

Again, this is a hypergeometric random variable and we can apply the formula, or we can reason it out. Total number of ways to capture 50 wolves  $= \binom{400}{50}$ . # ways with 10 tagged wolves  $= \binom{50}{10} \binom{350}{40}$ . Therefore the probability is equal to

$$\frac{\binom{50}{10} \binom{350}{40}}{\binom{400}{50}}.$$

c. The Lincoln-Peterson estimator for the total number of captured individuals is  $\frac{Kn}{k}$  where  $K$  is the total number of capture in the second round,  $n$  is the number of tagged wolves in the first round, and  $k$  is the number of tagged wolves seen in the second round. Write out the formula to calculate the expectation of the Lincoln-Peterson estimator when  $N = 400$ ,  $n = 50$ , and  $K = 50$  (the only random part is  $k$ ).

The expected value of  $\frac{Kn}{k}$  is equal to  $\sum_k \frac{Kn}{k} \Pr(k)$ . Recall that  $k$  is hypergeometric with  $\Pr(k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}$ . Therefore the expectation is equal to

$$\sum_{k=0}^K \frac{Kn}{k} \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}} = \sum_{k=0}^{50} \frac{50 \cdot 50}{k} \frac{\binom{50}{k} \binom{350}{n-k}}{\binom{400}{50}}$$

6. In basketball, baseball, and hockey playoffs, teams play each other for the best of 7 or 5 series.

- a. If team A and team B play each other and team A has a 60% chance of winning each game (whether home or away), what is the probability that team A will win the best of 7 series?

The number of times A wins is a binomial random variable with  $p = 0.6$  and  $n = 7$ . The probability A wins the best of 7 is the probability A wins 4 or more games, which equals

$$\binom{7}{4}0.6^40.4^3 + \binom{7}{5}0.6^50.4^2 + \binom{7}{6}0.6^60.4^1 + \binom{7}{7}0.6^70.4^0 = 0.71$$

- b. Suppose that team A has a 60% chance of winning home games and a 50% chance of winning away games and that games alternate home and away (starting with a home game for A). What is the probability that team A wins a best of 7 series?

A gets 4 home games and 3 away games. It has to win at least 4. It can win 4 home games with probability  $0.6^4$  and it won't matter what happens in the away games. It can win 3 home games with probability  $\binom{4}{3}0.6^30.4^1$  and then it has to win at least 1 away game, with probability  $1 - \binom{3}{0}0.5^00.5^3$ . It can win 2 home games with probability  $\binom{4}{2}0.6^20.4^2$  and then it needs to win at least two home games with probability  $1 - \binom{3}{0}0.5^00.5^3 - \binom{3}{1}0.5^10.5^2$ . Finally it can win 1 home games with probability  $\binom{4}{1}0.6^10.4^3$  and win at least 2 away games with probability  $\binom{3}{2}0.5^20.5^1 + \binom{3}{3}0.5^30.5^0$ . This gives total probability 0.6816.

- c. Does the order that the teams play matter (as long as team A plays 4 home games and 3 away games)?

No.



7. Your friend wants to sell his car, but they tell you they won't accept less than \$10,000. You think the car is worth more like \$8,000. Suppose offers will follow an exponential distribution with mean \$8,000.

a. What is the probability the first offer will be greater than \$10,000?

$$\Pr(X > 10000) = e^{-10000/8000} = e^{-1.25} = 0.287$$

b. What is the probability your friend will have to take more than 5 offers before getting one that is \$10,000 or larger?

This is the probability that none of the first 5 offers is \$10,000 or larger. The probability that each offer is less than \$10,000 is  $1 - e^{-1.25}$ . Therefore the probability that all 5 offers are less than \$10,000 is  $(1 - e^{-1.25})^5 = 0.185$ .

c. What is the expected number of offers they would have to take for an offer larger than \$10,000?

The number of offers required is the number of offers before the one larger than \$10,000 plus the desired offer. The number of offers before one larger than \$10,000 is a geometric random variable with  $p = e^{-1.25}$ . The expected time your friend would have to wait is  $\frac{1-p}{p} + 1 = \frac{1-e^{-1.25}}{e^{-1.25}} + 1 = 3.49$ .

d. Suppose they change their strategy to taking the largest of the first 5 offers. What would the expected sale price be?

First we need to figure out the distribution of the maximum of 5 geometric random variables is. Note that if the maximum is less than  $x$ , then all of the 5 random variables are less than  $x$ , so  $\Pr(\max(x_1, x_2, x_3, x_4, x_5) \leq x) = \Pr(x_1 \leq x, x_2 \leq x, \dots, x_5 \leq x) = \Pr(X \leq x)^5 = (1 - e^{-x/8000})^5$ . This is the cdf of the maximum. The pdf is given by the derivative in  $x$ ,  $f(x) = \frac{\partial}{\partial x}(1 - e^{-x/8000})^5 = 5(1 - e^{-x/8000})^4 \cdot e^{-x/8000} \cdot \frac{1}{8000}$ . The expectation is given by

$$\int_0^\infty xf(x)dx = \int_0^\infty x5(1 - e^{-x/8000})^4 \cdot e^{-x/8000} \cdot \frac{1}{8000}dx$$

and I'm not continuing further and neither should you.

8. Suppose you're on the game show *Let's Make a Deal*. The host reveals 3 doors. One contains a car and the other two contain goats, but you want the car. You pick one of the 3 doors. The host, Monty Hall, says, "Let's make a deal." He goes to the nearest (lowest in numerical order) door that contains a goat (for example, if you selected Door 1 and the car was indeed behind Door 1, then the host would always open Door 2, never Door 3). What are the probabilities that you will win the car if you stay with your original choice vs if you switch?

If you pick door 1, and the door contains the car then Monty will always show door 2. If door 2 contains the car then Monty will show door 3 and if door 3 contains the car then Monty will show door 2. If you pick door 2, and the door contains the car then Monty will always show door 3. If door 1 contains the car then Monty will show door 3 and if door 3 contains the car then Monty will show door 1. Similarly if you pick door 3, then Monty will show door 1 if the car is in door 2 or door 3 and will show door 2 if the car is in door 1.

If we stay with our original choice, then our chance of choosing right is  $1/3$ . If we switch then if Monty shows the greater door number we will always, and this occurs with probability  $1/3$ . If we switch when Monty shows the smaller door number then we will win with probability  $1/2$ . This occurs  $2/3$  of the time. Therefore our total probability of winning is  $1/3 + 2/3 \cdot 1/2 = 2/3$ .

9. A popular model for long tailed distributions (meaning they are likely to take on large values) is power law type distributions. This has been used to model word frequency and city population. If  $X$  is power law distributions the  $X$  has cdf that is proportional to  $x^{-\alpha}$  for  $x > 1$  and some  $\alpha > 0$ .

a. If  $\alpha = 3$  then  $f(x) = cx^{-3}$ . Calculate the normalizing constant  $c$ .

A pdf is valid if it integrates to 1. Therefore  $1 = c \int_1^\infty x^{-3} = c(-\frac{1}{2}x^{-2}|_1^\infty) = c\frac{1}{2}$  and  $c = 2$ .

b. If  $\alpha = 3$ , calculate the expectation of  $X$ .

The expectation of  $X$  is equal to  $\int_1^\infty xf(x)dx = \int_1^\infty x2x^{-3} = 2 \int_1^\infty x^{-2}dx = 2(-x^{-1})|_1^\infty = 2$ .

c. Set up the maximum likelihood equation.

The pdf is equal to  $cx^{-\alpha}$  for some constant  $c$  that depends on  $\alpha$ . Therefore the log likelihood function is equal to

$$\log L(\alpha|x_1, \dots, x_n) = \sum_{i=1}^n -\alpha \log(x_i) + \log c.$$

We take the derivative in terms of  $\alpha$  and obtain

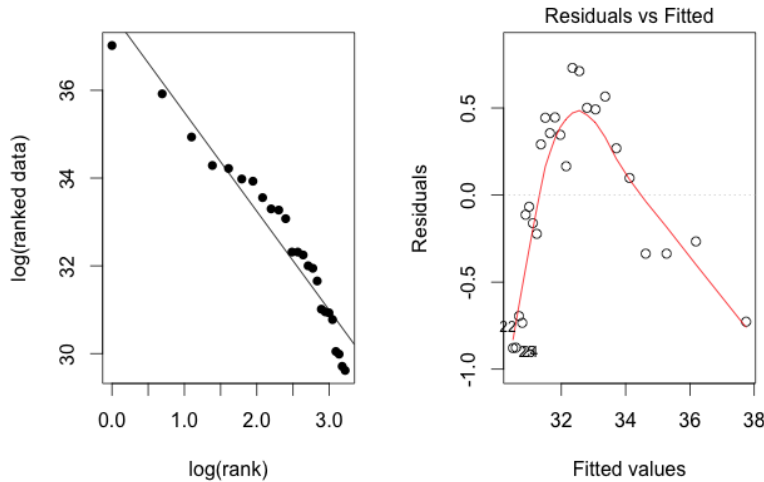
$$\frac{\partial}{\partial \alpha} \log L(\alpha|x_1, \dots, x_n) = \sum_{i=1}^n -\log x_i + \frac{\partial}{\partial \alpha} \log c = 0,$$

recalling that  $c$  depends on  $\alpha$ . If we want to calculate  $c$ , then  $c = (\int_1^\infty x^{-\alpha}dx)^{-1} = \frac{1}{-\alpha+1}x^{-\alpha+1}|_1^\infty = \frac{1}{\alpha-1}$ . This means that  $\alpha$  needs to be bigger than 1 for this to make sense. The complete MLE equation is

$$\frac{\partial}{\partial \alpha} \log L(\alpha|x_1, \dots, x_n) = \sum_{i=1}^n -\log x_i - \frac{1}{\alpha-1} = 0.$$

- d. A common method to fitting power law data is to plot the log of the ranks versus the log of the ranked data. Below is a figure of a power law linear fit. Does this fit look homoscedastic?

No, the error does not appear random. It's positive early, negative in the middle, and then positive later. There's a clear dependency, as shown by the right hand plot.



10. A random sample of 100 voters in a community produced 59 voters in favor of measure A. You would like to test to see if measure A passes.

a. State the null and alternative hypotheses.

We want to show measure A will pass. Therefore this is a one-sided test. Let  $p$  be the proportion of the population who support measure A. The null and alternatives are  $H_0 : p \leq 0.5$   $H_1 : p > 0.5$

b. Construct a 95% confidence interval for the support of measure A.

A 95% confidence interval is given by

$$\hat{p} \pm z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p})}/\sqrt{n} = \frac{59}{100} \pm 1.96 \sqrt{\frac{59}{100} \frac{41}{100}}/\sqrt{100} = (0.494, 0.686).$$

c. Compute the  $p$ -value. Will measure A win with 90% confidence?

The test statistic is given by

$$t = \frac{\hat{p} - p_0}{\sqrt{p_0(1-p_0)}/\sqrt{n}} = \frac{\frac{59}{100} - 0.5}{\sqrt{0.5 * 0.5}/\sqrt{100}} = 1.8.$$

For a one sided test we compute the  $p$ -value as  $Pr(Z > t)$ , where  $Z \sim N(0, 1)$ .  $Pr(Z > 1.8) = 0.0359$ . Since  $p < 0.1$ , we can say that the measure has greater than 50% support at the 90% confidence level and it should pass (with 90% confidence).

d. Suppose the true support for measure A is 58%. What is the minimum sample size needed for a power of 90% at a 95% confidence level?

The equation for the minimum sample size is  $\left( \frac{z_{1-\alpha} + z_{1-\beta}}{|p_1 - p_0|/\sigma} \right)^2$ . In this case  $z_{0.95} = 1.645$ ,  $z_{0.9} = 1.28$ ,  $|p_1 - p_0| = 0.58 - 0.5 = 0.08$ , and  $\sigma = \sqrt{0.5 \cdot 0.5} = 0.5$ . Therefore the minimum sample size is

$$\left( \frac{1.645 + 1.28}{0.08/0.5} \right)^2 = 334.2$$

and we would need  $n = 335$  at minimum (remember to round up).

12. After the 1936 presidential polling fiasco, more rigorous polling methods were developed. The best known was the Gallup poll, who used a method known as quota sampling. Quota sampling is nothing more than a systematic effort to force the sample to fit a certain national profile by using quotas: The sample should have so many women, so many men, so many blacks, so many whites, so many under 40, so many over 40 etc. The numbers in each category are taken to represent the same proportions in the sample as are in the electorate at large. If we assume that every important characteristic of the population is taken into account when setting up the quotas, it is reasonable to expect that quota sampling will produce a good cross-section of the population and therefore lead to accurate predictions. For the 1948 election between Thomas Dewey and Harry Truman, Gallup conducted a poll with a sample size of about 3250. Each individual in the sample was interviewed in person by a professional interviewer to minimize nonresponse bias, and each interviewer was given a very detailed set of quotas to meet. For example, an interviewer could have been given the following quotas: seven white males under 40 living in a rural area, five black males under 40 living in a rural area, six black females under 40 living in a rural area, etc. Other than meeting these quotas the ultimate choice of who was interviewed was left to each interviewer. Based on the results of this poll, Gallup predicted a victory for Dewey, the Republican candidate. The predicted breakdown of the vote was 50% for Dewey, 44% for Truman, and 6% for third-party candidates Strom Thurmond and Henry Wallace. The actual results of the election turned out to be almost exactly reversed: 50% for Truman, 45% for Dewey, and 5% for third-party candidates.

Can you identify a major problem or confounding factor with quota sampling as described above?

See [<https://www.math.upenn.edu/~deturck/m170/wk4/lecture/case2.html>]

The following are true / false questions.

13. When there are more predictor variables ( $x$ s) overfitting is less likely in linear models.

FALSE

14. Increasing the regularization parameter  $\lambda$  in lasso regression leads to sparser regression.

TRUE. When  $\lambda = 0$  the lasso is equivalent to regular linear regression and when  $\lambda = \infty$  the lasso will fit a constant function.

15. Increasing the regularization parameter  $\lambda$  in ridge regression leads to sparser regression.

FALSE. Ridge regression shifts parameters towards zero but not all the way.

16. For a fixed size of the training and test set, increasing the complexity of the model always leads to reduction of the test error.

FALSE. We don't know. If you add spurious variables, the test error will decrease.

17. A linear regression analysis allows you to make predictions for a variable using information about another variable and the relationship between them.

TRUE

18. The standard error of the estimate is a measure of how far off are our predictions for  $X$ .

TRUE

19. In any regression analysis, the total variance in  $Y$  can be broken down into two parts: the variance attributable to variation in  $X$  and that which is left over and unexplained.

TRUE

20. If you wish to use many predictors to predict a single criterion, you must use each predictor in its own regression analysis.

FALSE. That's what multiple regression is for.