

# Sp 2020 Stat 88 final solns

Adam Lucas

May 2020

(1) **(MIDTERM, symmetry, binomial, expectation with indicators)**

Halloween will be here before you know it and the children in your neighborhood will come trick-or-treating (that is, they will come to your door and demand candy). Suppose there are 20 children in your neighborhood and 30 houses (one of which is yours). Each child independently chooses 10 houses at random without replacement to visit.

- a What is the probability that a specific child will visit your house?
- b What is the probability that exactly 10 children visit your house?
- c What is the expected number of houses that exactly 10 children visit?

**Solution**

- a  $1/3$
- b  $\binom{20}{10}(1/3)^{10}(2/3)^{10}$ . This is a binomial problem with probability from the last part.
- c  $30 * (\binom{20}{10}(1/3)^{10}(2/3)^{10})$ . Here we have a sum of indicators where  $I_2 = 1$  if the second house (out of 30 houses) has exactly 10 children visit.

(2) **(continuing the last question) (MIDTERM, challenging probability)**

You have a bunch of candy in your house, ranging from very cheap to very fancy. Suppose that, if you ordered the candy in your house from cheapest to fanciest, the  $i$ th candy cost  $30i$  cents to buy (so the cheapest candy cost 30 cents, the 10th cheapest cost 300 cents, etc). You want to keep the fancy candy for yourself (the kids would not appreciate the fancy stuff anyway), so when a child comes to your door, you give them the cheapest remaining candy – thus, if 10 children come to your door, you will give out the 10 cheapest pieces of candy that you have.

- a If  $k$  children come to your door, what is the total cost of the candy denoted as  $T_k$  in terms of cents, that you will give out? (Hint: This uses a summation identity).
- b What is the probability that you will give away more than \$15 worth of candy? Hint: For what  $k$  is  $T_k > 1500$  and start from there.

**Solution**

- a Let  $T_k$  be the cost in dollars of giving away  $k$  pieces of candy. Then  $T_k = 30 \sum_{i=1}^k i = 30k(k+1)/2 = 15k(k+1)$ .
- b  $P(T_k > 1500) = P(15k(k+1) > 1500) = P(k \geq 10)$ . This is  $\sum_{i=10}^{30} \binom{20}{i}(1/3)^i(2/3)^{20-i}$  from part (b) of the previous problem.

(3) **(MIDTERM, binomial, indicators, properties of expectation and variance, conditional expectation, Poisson)**

- a Define  $X$  to be the sum of 10 indicator random variables, so  $X = \sum_{i=1}^{10} I_i$ . Mark all of the statements that *must* be true for  $X$  to be a binomial random variable.
  - i Each  $I_i$  has the same probability of taking the value 1.
  - ii The indicators are all independent of each other.
  - iii Each  $I_i$  should have probability  $1/2$  of taking the value 1.
  - iv None of the above.
- b We have two random variables  $X$  and  $Y$ . Which of the following statements must be true to apply the addition rule for expectations,  $E(X + Y) = E(X) + E(Y)$ ?
  - i  $X$  and  $Y$  must be independent.
  - ii  $X$  and  $Y$  must be identically distributed.
  - iii None of the above.

- c We have two random variables  $X$  and  $Y$ . Which of the following statements must be true to apply the addition rule for variance,  $Var(X + Y) = Var(X) + Var(Y)$ ?
- i  $X$  and  $Y$  must be independent.
  - ii  $X$  and  $Y$  must be identically distributed.
  - iii None of the above.
- d We have two random variables  $X$  and  $Y$ . Which of the following statements is true (there can be more than one):
- i  $E(Y|X)$  is a function of  $X$ .
  - ii  $E(Y|X)$  is a function of  $Y$ .
  - iii  $E(E(Y|X))$  is a non constant random variable.
  - iv  $E(E(E(Y|X))) = E(Y)$ .
- e You flip a fair coin  $N$  times where  $N$  is a random variable,  $N \text{ Poisson}(5)$ . What is the expected number of heads you will get?
- i 5
  - ii  $5/2$
  - iii  $N/2$
  - iv None of the above

### Solution

- a i,ii  
b iii  
c i  
d i,iv (since  $E(Y) = E(E(Y|X))$  is a constant and the expectation of a constant is itself)  
e ii since if  $H|N$  is the number of heads in  $N$  trials,  $E(H|N) = N/2$  and  $E(H) = E(N/2) = 5/2$ .

- (4) (**MIDTERM, hypergeometric, expectation of geometric**) Your family is trying to decide sleeping arrangements in the house. There are 10 young people staying in the house who will be divided into two rooms. You will choose a set of 5 people to sleep in the first room and the rest will sleep in the second room. A sleeping arrangement is acceptable if it meets the following two criteria:
- Your two youngest cousins (say, 7 and 8 years old) always fight, so exactly one of them must be in the first room.
  - There are 4 people over 18 years old, and exactly 2 of them must be in the first room to supervise the others.

You are too lazy to work out the arrangement by hand, so you decide to do this randomly by choosing a set of 5 people at random to sleep in the first room, and repeat this until you get an acceptable arrangement.

- a For a single draw of 5 people, what is the probability that you will draw an acceptable sleeping arrangement?  
b What is the expected number of draws of 5 people until you get an acceptable arrangement? Note: Each drawing of 5 people is independent of every other drawing of 5 people.

### Solution

- a  $P(\text{acceptable}) = \frac{\binom{2}{1}\binom{4}{2}\binom{4}{2}}{\binom{10}{5}}$   
b Let  $X$  be the number of drawings of 5 people until you get an acceptable arrangement. Each drawing of 5 people is independent of every other drawing of 5 people.  $X$  is Geometric with probability  $p = \frac{\binom{2}{1}\binom{4}{2}\binom{4}{2}}{\binom{10}{5}}$  of success. Hence,  $E[X] = 1/p = \frac{\binom{10}{5}}{\binom{2}{1}\binom{4}{2}\binom{4}{2}} = 3.5$ .

- (5) (**SD of mean, Chebyshev, chap 6**)

You are using a telescope to measure the speed at which the planet Saturn crosses the night sky. To do this you draw two lines on your lens, and measure the time it takes for Saturn to cross between the two lines. However, your time measurement is noisy, so you will conduct this observation several times and average their results.

Let  $X_i$  represent the time measurement from the  $i$ th observation. Your measurements are well calibrated, so for each  $i$ ,  $E(X_i) = \mu_X$ , where  $\mu_X$  is the true time it takes Saturn to cross between the lines. Each measurement also has standard deviation  $SD(X_i) = 0.03$  seconds.

- a You will take  $n$  measurements,  $X_1, \dots, X_n$ , using the same procedure, and use the sample average  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  to estimate  $\mu_X$ . In terms of  $n$ , what is  $SD(\bar{X})$ ?
- b What is the smallest number of measurements you will need to take so that your estimate  $\bar{X}$  has at most a  $\frac{1}{25}$  probability of falling outside the interval  $\mu_X \pm 0.003$  seconds? (Hint: Chebyshev)

**Solution**

a  $SD(\bar{X}) = \frac{SD(X_1)}{\sqrt{n}} = \frac{0.03}{\sqrt{n}}.$

- b Applying Chebyshev,  $P(|\bar{X}| \geq \mu_X + 0.003) \leq \frac{1}{5^2}$ , where 0.003 is  $k = 5$  SD away from the mean. So,  $0.003 = 5(\frac{.03}{\sqrt{n}})$ . Solving for  $n$  we get  $\sqrt{n} = 50$ , so  $n = 2500$ .

(6) **(confidence interval, hypothesis test, 2 sample z test, chap 9)**

You are interested in knowing the preference at Cal between two statistical software packages R and Python. In a SRS of 250 students, 80% prefer R over Python, and the rest (20%) prefer Python over R.

- a Construct a 95% CI for the percentage of students who prefer R over Python. Is there a 95% chance that the true percentage lies in your interval?
- b Cal Data Science believes that 70% of Cal students prefer R over Python. Based on the data from your SRS, would you agree with Cal Data Science's belief?
- c Your high school friend, who studies Data Science at Stanford, finds from a SRS of 250 students, that only 75% of Stanford Data Science students prefer R over Python and the rest (25%) prefer Python over R. Do a hypothesis test (at a 5% level of significance) to decide whether this difference in R preferences is due to chance (make a one sided alternative). Please state your null hypothesis, alternative hypothesis, test statistics and p value.

**Solution Update**

**The old solution for this problem contains many errors, please refer to the solution below.**

- a  $80\% \pm 2SE$  where  $SE = \frac{\sqrt{(.8)(.2)}}{\sqrt{250}} * 100\% = 2.5\%$ . So the 95% CI is [75%, 85%]. No, the true percentage is a fixed constant and it either lies in our interval or not. We can say that there is a 95% chance that a random 95% CI, such as this one, contains the true percentage.
- b The Null is that  $p = 70\%$  and the alternative is that  $p \neq 70\%$ . Under the null,  $SE = \frac{\sqrt{(.7)(.3)}}{\sqrt{250}} * 100\% = 2.9\%$ . Let T be test statistics,  $P(|T - 70\%| > |80\% - 70\%|) = 2 * (1 - \Phi(\frac{80\% - 70\%}{2.9\%})) = 2 * (1 - \Phi(3.45)) \approx 0.00056 < 0.05$ . Therefore, we reject the null hypothesis. Alternatively, if you decide to do an one-sided test, i.e.  $p > 70\%$ , you will also reject the null.
- c This is a two sample Z test. Let  $p_c, p_s$  denote the true proportions for Cal and Stanford. The Null is that  $p_c - p_s = 0$ . The alternative is that  $p_c - p_s > 0$ . First, we calculate  $p = \frac{p_c * 250 + p_s * 250}{500} = 0.775$ , which is the sample proportion of the combined sample.  $SE_{diff} = \sqrt{(0.775) * (0.225)/250 + (0.775) * (0.225)/250} * 100\% = 3.7\%$ . Hence the T statistics is  $T = (80\% - 75\%)/3.7\% = 1.35$ , and for this **one-sided test**, p-value =  $1 - \Phi(T) = 1 - \Phi(1.35) = 1 - 0.911 = 0.089 > 0.05$  so we accept the null that the two schools have the same preference for R.

(7) **(independent exponential distribution, cdf, chap 10)** Let  $T_1$  and  $T_2$  be independent Exponential random variables with rates  $\lambda_1$  and  $\lambda_2$ .

- a Show that the CDF of  $T = \text{Max}(T_1, T_2)$  is  $F(t) = (1 - e^{-\lambda_1 t})(1 - e^{-\lambda_2 t})$ .
- b An electric circuit consists of 2 components in the following diagram.



The lifetimes of the components, measured in days, have independent exponential distributions with means given in the diagram. Let  $T$  be the lifetime of the circuit. Find the CDF of  $T$ . Hint: Is  $T$  a Max or a Min of the two components lifetimes?

**Solution**

**a**  $F(t) = P(T < t) = P(T_1, T_2) = P(T_1 < t, T_2 < t) = (1 - e^{-\lambda_1 t})(1 - e^{-\lambda_2 t})$ .

**b**  $T = \text{Min}(T_1, T_2)$  where  $T_1 \sim \text{Exp}(1/5)$  and  $T_2 \sim \text{Exp}(1/4)$ .

Then  $F(t) = 1 - P(T > t) = 1 - P(T_1 > t)P(T_2 > t) = 1 - e^{-t/5}e^{-t/4}$ .

- (8) (**unbiased estimators chap 11**) Let  $X_1, X_2, \dots, X_n$  be iid Uniform(0,  $2\theta$ ) for an unknown parameter  $2\theta$ . Let  $M = \text{Max}(X_1, X_2, \dots, X_n)$ ,  $N = \text{Min}(X_1, X_2, \dots, X_n)$ .

**a** What is the expected distance of  $N$  from zero? **Clarify:**  $E[X_1] = \theta$ ,  $E[N]$  is what the question meant to ask.

**b** Find the expectation of  $M$

**c** Find an unbiased estimator of  $2\theta$  in terms of  $M$ .

**Solution**

**a**  $2\theta/(n+1)$  since the total gap is  $2\theta$  and the total number of gaps is  $n+1$ .

**b**  $E(M) = n2\theta/(n+1)$  since  $E(M)$  is  $n$  times the average distance of  $X_1$  from 0.

**c**  $T = M(n+1)/n$  is an unbiased estimator of  $2\theta$  since  $2\theta = E(M)(n+1)/n$  from part (b).

- (9) (**correlation, chap 11**)

You run a study at Thanksgiving. Let  $X$  be how many pounds of turkey each person ate, and  $Y$  be how many hours past midnight they slept. You find that  $r(X, Y) = 0.6$ .

**a** Qualitatively, what sort of pattern does this correlation indicate?

**b** You rerun the study, this time defining  $X$  as grams of turkey consumed, and  $Y$  as minutes slept past midnight. Does  $r(X, Y)$  change?

**c** Now you define  $Z$  how much money each person saves by attending Black Friday sales, which start (in your town) at midnight. You estimate that every minute spent shopping (i.e. not sleeping) after midnight is worth 2 dollars in savings. What is  $r(X, Z)$ , or the correlation between turkey eaten and dollars saved?

**Solution.**

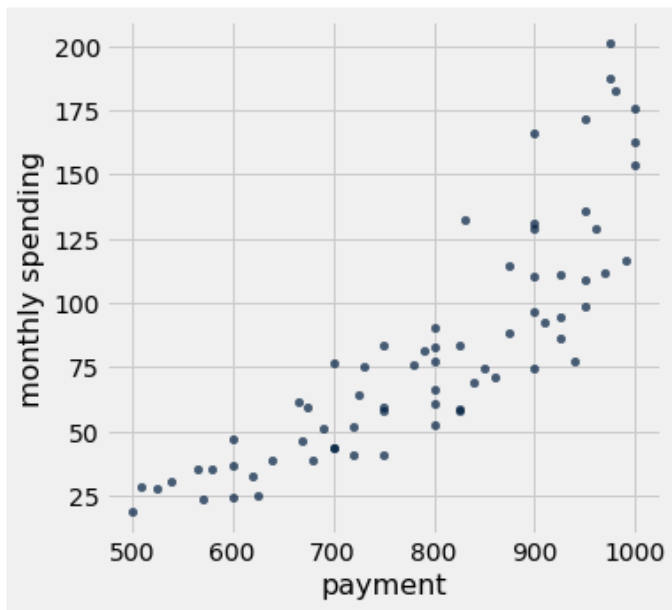
**a** In general, people who ate more turkey slept more past midnight; people who ate less turkey slept less past midnight.

**b** No because correlation doesn't change due to a change of scale  $ax + b$  where  $a > 0$ .

**c**  $r(X, Z) = r(X, -2Y) = -r(X, Y) = -0.6$ .

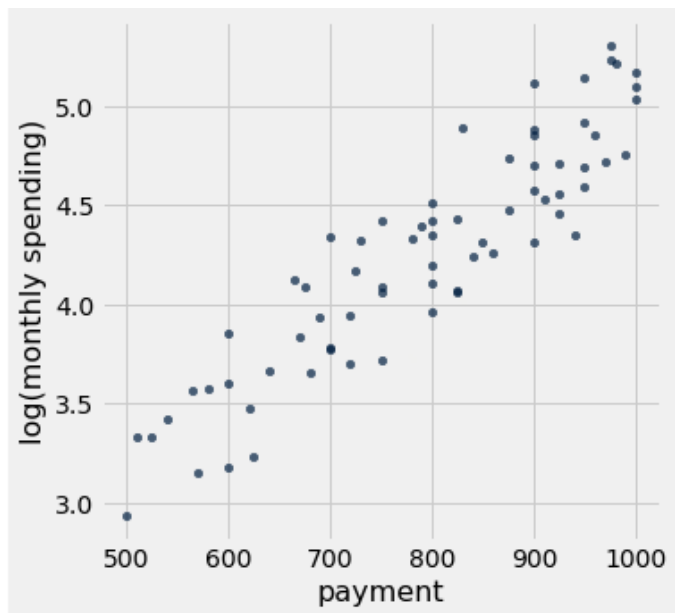
- (10) (**Simple linear regression, chap 12**)

As a statistical intern for the US Treasury you are studying the effect of the stimulus payment (in dollars) on monthly household spending (in dollars) during Covid-19. You decide to make a linear regression model of monthly household spending as a function of stimulus payment size. You create a scatter plot based on a random sample of 70 households.



- a Does it look like a regression line is a good fit for the data? Explain your answer in terms of the assumptions of the Simple Linear Regression Model.

You determine the regression line after taking the natural log (base e) of all the observations for the response variable. That is to say, for data points  $i = 1, 2, 3, \dots$ , the response is assumed to be  $\log(Y_i) = \beta_0 + \beta_1 x_i + \epsilon_i$ .



You now run a linear regression on this transformed data and print out a summary of the output, but some of the entries are removed for national security reasons.

	coef	std err	t	P> t	[0.025	0.975]
const	1.3131	0.160	8.182	0.000	0.993	1.633
payment		0.000	18.612			

You also know,

$$\mu_x = 792.8 \text{ and } \sigma_x = 140.1$$

$$\sigma_{\log(Y)} = 0.6 \text{ and } r = .91$$

where  $x = \text{payment}$  and  $\log(Y) = \log(\text{monthly spending})$ .

- b** Find the equation of the regression line.
- c** What is the predicted monthly spending, in dollars, for a household receiving a stimulus payment of \$900?
- d** Find the average  $\log(\text{monthly spending})$ ,  $\mu_{\log Y}$ .
- e** Conduct a hypothesis test  $H_0 : \beta_1 = 0$  versus alternative  $H_A : \beta_1 \neq 0$ . What can you conclude about the effect of payment size and household spending?

**Solution. Update: Some numerical values for this question were inaccurate.**

- a** No, our linear model is  $Y_i = \beta_0 + \beta_1 + \epsilon_i$  where  $\epsilon \sim N(0, \sigma^2)$  where  $\sigma$  doesn't change with  $x$ . If we were to make a linear Tyche line the errors in the scatterplot wouldn't have a constant variance (i.e. the variance gets bigger for bigger  $x$ ).
- b** We have  $\beta_1 = r\sigma_{\log(Y)}/\sigma_x$ , so  $\beta_1 = 0.0039$ . Hence the regression line is  $\widehat{\log(Y)} = 1.31 + 0.0039x$ .
- c** We have  $\log(Y) = 1.31 + 0.0039(900) = 4.82$  in log dollars. It follows that we predict a household monthly spending of  $e^{4.82} = 123.97$  dollars.
- d** The point of averages lies on the regression line so  $\mu_{\log(Y)} = 1.31 + .0039\mu_X = 4.40$  since we are given  $\mu_x = 792.8$ .
- e** The 95% CI of  $\beta_1$  is  $\hat{\beta}_1 \pm 2SE = .0039 \pm 0$  which doesn't contain zero so we reject the null. Alternatively, the p-value is  $t=18.612$  is very close to zero which is less than 0.05 so we reject the null. We conclude that payment and  $\log(\text{spending})$  are dependent and hence payment and spending are dependent.