

Section #3 Solutions

Based on the work of many CS109 staffs

- 1. Email Predictions:** Let's say that on average you get an email every 5 minutes. Assume that the time between email arrivals is exponentially distributed. What is the probability that you get no emails in the next 10 minutes?

Let X be the number of minutes until the next email. $X \sim \text{Exp}(\lambda = \frac{1}{5})$.

$$P(X > 10) = 1 - F_X(10) = 1 - (1 - e^{-\lambda 10}) = e^{-2} \approx 0.1353$$

Alternatively, let Y be the number of emails you get in the next 10 minutes. The average email rate is 2 emails per 10 minutes. $Y \sim \text{Poi}(\lambda = 2)$.

$$\begin{aligned} P(Y = 0) &= \frac{2^0 e^{-2}}{0!} \\ &= e^{-2} \approx 0.1353 \end{aligned}$$

- 2. Approximating Normal:** Your website has 100 users and each day each user independently has a 20% chance of logging into your website. Use a normal approximation to estimate the probability that more than 21 users log in.

The number of users that log in B is binomial: $B \sim \text{Bin}(n = 100, p = 0.2)$. It can be approximated with a normal that matches the mean and variance. Let C be the normal that approximates B . We have $E[B] = np = 20$ and $\text{Var}(B) = np(1 - p) = 16$, so $C \sim N(\mu = 20, \sigma^2 = 16)$. Note that because we are approximating a discrete value with a continuous random variable, we need to use the continuity correction:

$$\begin{aligned} P(B > 21) &\approx P(C > 21.5) \\ &= P\left(\frac{C - 20}{\sqrt{16}} > \frac{21.5 - 20}{\sqrt{16}}\right) \\ &= P(Z > 0.375) \\ &= 1 - P(Z < 0.375) \\ &= 1 - \phi(0.375) = 1 - 0.6462 = 0.3538 \end{aligned}$$

- 3. Continuous Random Variable:** Let X be a continuous random variable with the following probability density function:

$$f_X(x) = \begin{cases} c(e^{x-1} + e^{-x}) & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- a. Find the value of c that makes f_X a valid probability distribution.
 b. What is $P(X > 0.75)$?

a. We need $\int_{-\infty}^{\infty} f_X(x) dx = 1$.

$$\begin{aligned} \int_{-\infty}^{\infty} f_X(x) dx &= \int_0^1 c(e^{x-1} + e^{-x}) dx \\ 1 &= c \left[e^{x-1} - e^{-x} \right]_{x=0}^1 \\ 1 &= c(e^{1-1} - e^{-1} - (e^{0-1} - e^{-0})) \\ c &= \frac{1}{1 - e^{-1} - (e^{-1} - 1)} = \frac{1}{2 - \frac{2}{e}} \end{aligned}$$

b.

$$\begin{aligned} P(X > 0.75) &= \int_{0.75}^1 c(e^{x-1} + e^{-x}) dx \\ &= c \left[e^{x-1} - e^{-x} \right]_{x=0.75}^1 \\ &= c \left(e^{1-1} - e^{-1} - (e^{0.75-1} - e^{-0.75}) \right) \\ &= c \left(1 - e^{-1} - e^{-0.25} + e^{-0.75} \right) = \frac{1 - e^{-1} - e^{-0.25} + e^{-0.75}}{2 - \frac{2}{e}} \end{aligned}$$

4. Air Quality: Throughout the United States, the Environmental Protection Agency monitors levels of PM2.5, a type of dangerous air pollution. These PM2.5 measurements can be approximately modeled by a normal distribution.

- a. Let us model PM2.5 measurements with a normal distribution that has a mean of 8. If three-quarters of all measurements fall below 11.4, what is the standard deviation? Round to the nearest integer.
- b. PM2.5 values above 12 can pose some health risks, especially to sensitive populations. Using the standard deviation found above, what is the probability that a randomly selected PM2.5 measurement is over 12?
- c. What is the probability that a randomly selected PM2.5 measurement is between 7 and 8?

a. $\Phi\left(\frac{x-\mu}{\sigma}\right) = \Phi\left(\frac{11.4-8}{\sigma}\right) = 0.75 \implies \frac{3.4}{\sigma} \approx .68 \implies \sigma \approx 5$.

b. $P(q > 12) = 1 - P(q < 12) = 1 - \Phi\left(\frac{12-8}{5}\right) = 1 - \Phi(.8) = 1 - 0.7881 = 0.2119$.

c. $P(7 < h < 8) = P(h < 8) - P(h < 7) = \Phi\left(\frac{8-8}{5}\right) - \Phi\left(\frac{7-8}{5}\right)$

$$\begin{aligned}
 &= \Phi\left(\frac{8-8}{5}\right) - \Phi\left(\frac{-1}{5}\right) = \Phi\left(\frac{8-8}{5}\right) - (1 - \Phi\left(\frac{1}{5}\right)) \\
 &= \Phi(0) - (1 - \Phi(0.2)) = 0.5 - (1 - 0.5793) = 0.0793
 \end{aligned}$$

5. Elections: We would like to see how we could predict an election between two candidates in France (A and B), given data from 10 polls. For each of the 10 polls, we report below their sample size, how many people said they would vote for candidate A, and how many people said they would vote for candidate B. Not all polls are created equal, so for each poll we also report a value "weight" which represents how accurate we believe the poll was. The data for this problem can be found on the class website in polls.csv:

Poll	N samples	A votes	B votes	Weight
1	862	548	314	0.93
2	813	542	271	0.85
3	984	682	302	0.82
4	443	236	207	0.87
5	863	497	366	0.89
6	648	331	317	0.81
7	891	552	339	0.98
8	661	479	182	0.79
9	765	609	156	0.63
10	523	405	118	0.68
Totals:	7453	4881	2572	

- a. First, assume that each sample in each poll is an independent experiment of whether or not a random person in France would vote for candidate A (disregard weights).
 - Calculate the probability that a random person in France votes for candidate A.
 - Assume each person votes for candidate A with the probability you've calculated and otherwise votes for candidate B. If the population of France is 64,888,792, what is the probability that candidate A gets more than half of the votes?
- b. Nate Silver at fivethirtyeight pioneered an approach called the "Poll of Polls" to predict elections. For each candidate A or B, we have a random variable S_A or S_B which represents their strength on election night (like ELO scores). The probability that A wins is $P(S_A > S_B)$.
 - Identify the parameters for the random variables S_A and S_B . Both S_A and S_B are defined to be normal with the following parameters:

$$S_A \sim \mathcal{N}\left(\mu = \sum_i p_{A_i} \cdot \text{weight}_i, \sigma^2\right) \quad S_B \sim \mathcal{N}\left(\mu = \sum_i p_{B_i} \cdot \text{weight}_i, \sigma^2\right)$$

where p_{A_i} is the ratio of A votes to N samples in poll i , p_{B_i} is the ratio of B votes to N samples in poll i , weight_i is the weight of poll i , m_i is the N samples in poll i and:

$$\sigma = \frac{K}{\sqrt{\sum_i m_i}} \text{ s.t. } K = 350; \text{ thus } \sigma = 4.054.$$

- We will calculate $P(S_A > S_B)$ by simulating 100,000 fake elections. In each fake election, we draw a random sample for the strength of A from S_A and a random sample for the strength of B from S_B . If S_A is greater than S_B , candidate A wins. What do we expect to see if we simulate so many times? What do we actually see?
- c. Which model, the one from (a) or the model from (b) seems more appropriate? Why might that be the case? On election night candidate A wins. Was your prediction from part (b) "correct"?

a. $P(\text{random person votes for A}) = \frac{\text{votes for A}}{\text{total votes}} = \frac{4881}{7453} = 0.655$

Now, let X be the number of votes for candidate A. We assume that $X \sim \text{Bin}(64888792, 0.655)$.

- Since n is so large, we can approximate X using a normal $Y \sim N(np, np(1 - p))$.
- $\mu = np = 42502158.76$, Variance = $np(1 - p) = 14663244.77$ Std Dev = 3829.26
- Votes to win = $\frac{64888792}{2} = 32444396$
- $P(\text{A gets enough votes}) = P(X > 32444396) \approx P(Y > 32444396.5) = 1.00$

b. $S_A \sim N(5.324, 16.436)$

$S_B \sim N(2.926, 16.436)$

$P(S_A > S_B) \approx 0.66$

We can figure this out through simulation by drawing from S_A and S_B 100,000 times and seeing how often the S_A value is greater than the S_B value. Later in the quarter, when we learn the convolution of independent normals, you will be able to figure this out mathematically.

- c. Algorithm (a) makes very few assumptions, and simplicity can be useful, but it does assume that each voter is independent - which we definitely know isn't the case in real elections. Algorithm (b) allows us to model bias (using the weights we incorporated), and doesn't think of each voter as necessarily independent.