

Section #3 Solutions

Problems by Chris

1 Conditional Flu (Optional)

If a person has the flu, the distribution of their temperature is Gaussian with mean 101 and variance 1. If a person does not have the flu, the distribution of their temperature is Gaussian with mean 98 and variance 1. All you know about a person is that they have a temperature of 100. What is the probability they have the flu? Historically, 20% of people you analyze have had the flu.

This is an inference problem, as it involves doing Bayes' Theorem with random variables. We are going to define two random variables:

F is an indicator variable which is 1 if the person has the flu.

X is the distribution of the person's temperature.

The question asks: what is $P(F = 1|X = 100)$?

The problem tells us that $F \sim \text{Bern}(p = 0.2)$ and that X is distributed as following, conditioned on F :

$$X|F = 1 \sim N(\mu = 101, \sigma^2 = 1)$$

$$X|F = 0 \sim N(\mu = 98, \sigma^2 = 1)$$

We can solve this using the inference version of Bayes, which allows for a mixture of discrete and continuous random variables.

$$P(F = 1|X = 100) = \frac{f(X = 100|F = 1)P(F = 1)}{f(X = 100|F = 1)P(F = 1) + f(X = 100|F = 0)P(F = 0)}$$

The next step is to substitute the PDF of the Normal distribution:

$$\begin{aligned} P(F = 1|X = 100) &= \frac{\frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}(100-101)^2} \cdot 0.2}{\frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}(100-101)^2} \cdot 0.2 + \frac{1}{\sqrt{2\pi}}e^{-\frac{1}{2}(100-98)^2} \cdot 0.8} \\ &= \frac{e^{-\frac{1}{2}} \cdot 0.2}{e^{-\frac{1}{2}} \cdot 0.2 + e^{-2} \cdot 0.8} \\ &\approx .528 \end{aligned}$$

2 Algorithmic Fairness

An AI model makes a binary prediction (G for guess) for whether a person will repay a loan. It is important to show that the model is fair with respect to a binary demographic (D for demographic). But what does fair mean? Let's analyze the historical predictions of the model and compare the predictions to the true outcome (T for truth). Consider the following joint probability table from the model's history:

	D = 0		D = 1	
	G = 0	G = 1	G = 0	G = 1
T = 0	0.21	0.32	0.01	0.01
T = 1	0.07	0.28	0.02	0.08

D : is the demographic of an individual (binary).

G : is the "repay" prediction made by the algorithm. 1 means predicted repay.

T : is the true "repay" result. 1 means did repay.

Recall that cell ($D = i, G = j, T = k$) is the probability $P(D = i, G = j, T = k)$.

- a. (4 points) What is $P(D = 1)$?

$$P(D = 1) = 0.01 + 0.01 + 0.02 + 0.08 = 0.12$$

- b. (4 points) What is $P(G = 1|D = 1)$?

$$P(G = 1|D = 1) = (0.01 + 0.08) / 0.12 = 0.75$$

- c. (6 points) Fairness definition 1: Parity

An algorithm satisfies "parity" if the probability that the algorithm makes a positive prediction ($G = 1$) is the same regardless of the demographic variable. Does this algorithm satisfy parity?

We want to see if $P(G = 1|D = 1) = P(G = 1|D = 0)$.

$$P(G = 1|D = 0) = (0.32 + 0.28) / (0.21 + 0.07 + 0.32 + 0.28) = 0.60/0.88 = 0.68.$$

Thus, we see that $P(G = 1|D = 1) > P(G = 1|D = 0)$ and the algorithm does not satisfy parity.

- d. (6 points) Fairness definition 2: Calibration (Optional)

An algorithm satisfies "calibration" if the probability that the algorithm is correct ($G = T$) is the same regardless of demographics. In terms of "correct" we need to separately check whether the algorithm is equally likely to be correct for both outcomes ($G = T = 1$ and $G = T = 0$). Does this algorithm satisfy calibration?

We essentially want to see if $P(G = 0, T = 0|D = 0) = P(G = 0, T = 0|D = 1)$ and if $P(G = 1, T = 1|D = 0) = P(G = 1, T = 1|D = 1)$.

First we check if $P(G = 0, T = 0|D = 0) = P(G = 0, T = 0|D = 1)$.

$$P(G = 0, T = 0|D = 0) = 0.21 / (0.21 + 0.07 + 0.32 + 0.28) = 0.239$$

$$P(G = 0, T = 0|D = 1) = 0.01 / (0.01 + 0.02 + 0.01 + 0.08) = 0.083$$

So we can see that the algorithm does not satisfy calibration.

e. (6 points) Fairness definition 3: Equality of odds (Optional)

An algorithm satisfies “equality of odds” if the probability that the algorithm predicts a positive outcome given that the true outcome is positive ($G = 1|T = 1$) is the same regardless of demographics. Does this algorithm satisfy equality of odds?

$$P(G = 1|T = 1, D = 0) = 0.28 / (0.28 + 0.07) = 0.8$$

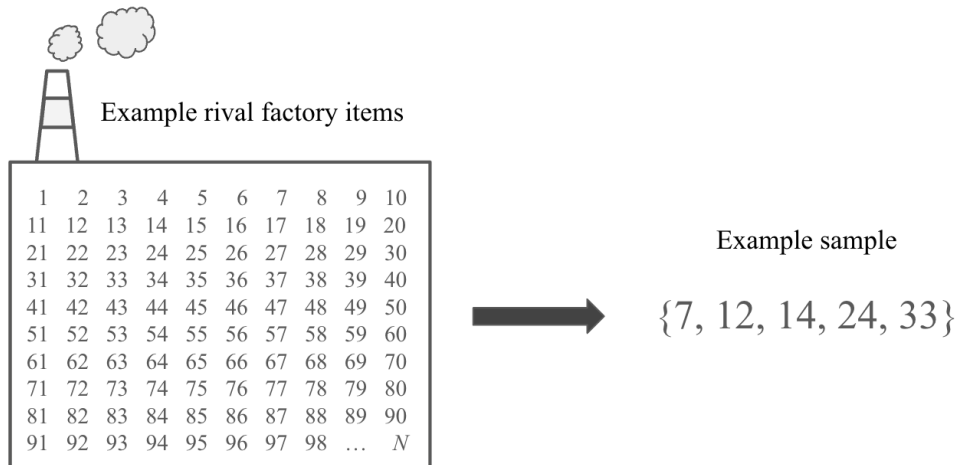
$$P(G = 1|T = 1, D = 1) = 0.08 / (0.08 + 0.02) = 0.8$$

We see that $P(G = 1|T = 1, D = 0) = P(G = 1|T = 1, D = 1)$ and thus, the algorithm does satisfy equality of the odds!

3 Tank Probability

A rival is producing items. We would like to estimate the number of items, N , that they have produced. We notice that each item has a unique serial number and we assume that when we acquire (sample) items each serial number on the item is a positive integer equally likely to be any number from the set $\{1, 2, \dots, N\}$.

For example, if you randomly acquired (sampled) 5 items produced at the factory, you might see the serial numbers $\{7, 12, 14, 24, 33\}$ which should give you a clue as to what N could be!



- a. (7 points) For part (a) only, assume $N = 100$. We sample 5 items. What is the probability that the largest serial number in our sample is 33? Let's solve this with equally likely outcome spaces!

In this problem, we are sampling without replacement, since each unique combination of five distinct integers between 1 and N is a potential outcome. Since all items are equally likely to be sampled, each set of five samples is equally likely. We can thus compute probabilities using the event space and sample space method.

If we consider the randomly sampled items as unordered, we have $|S| = \binom{100}{5}$, i.e. the number of unique ways to choose 5 objects from 100 objects. To count the event space (number of samples with 33 as the highest), use the following generative story:

- (1) Pick a set of four distinct numbers out of the range from 1 to 32,
- (2) Pick 33 as the fifth number.

Thus, $|E| = \binom{32}{4} \cdot 1$, and $P(\text{largest serial number is } 33 | N = 100) = \frac{\binom{32}{4}}{\binom{100}{5}}$.

You can similarly consider the randomly sampled items as ordered. $|S| = \frac{100!}{95!} = 100 \cdot 99 \cdot 98 \cdot 97 \cdot 96$, and $|E| = 5(1 \cdot 32 \cdot 31 \cdot 30 \cdot 29)$. (Note that we have to multiply by 5 to account for the fact that we could choose to sample 33 in any of the possible positions in our order.)

- b. (10 points) Your prior belief is that every value of N between 33 and 100 (inclusive) is equally likely. What is your updated probability mass function for N , given that you sampled 5 items and the largest serial number was 33?

This is a Bayesian inference problem. We are given the random variable N , the (unknown) total number of values or objects. Let L be the largest value in a sample of 5. We observe that $L = 33$. Setting up Bayes' Theorem and using Law of Total Probability in the denominator:

$$\begin{aligned} P(N = n|L = 33) &= \frac{P(L = 33|N = n) \cdot P(N = n)}{P(L = 33)} \\ &= \frac{P(L = 33|N = n) \cdot P(N = n)}{\sum_{i=33}^{100} P(L = 33|N = i) \cdot P(N = i)} \end{aligned}$$

Since there are $100 - 33 + 1 = 68$ values in the range $[33, 100]$, and we are told each of those possibilities is equally likely, we can use the prior

$$P(N = n) = \frac{1}{68} \text{ for } 33 \leq n \leq 100.$$

From part (a), we obtain the likelihood

$$P(L = 33|N = n) = \frac{\binom{32}{4}}{\binom{n}{5}} \text{ for } n \geq 33.$$

(since in part a, we assumed $N = 100$). Plugging these in:

$$P(N = n|L = 33) = \frac{\frac{\binom{32}{4}}{\binom{n}{5}} \cdot \frac{1}{68}}{\sum_{i=33}^{100} \frac{\binom{32}{4}}{\binom{i}{5}} \cdot \frac{1}{68}}$$

We only sum from $i = 33$ to 100 since the prior is 0 outside of this range.

- c. (3 points) Given that you sampled 5 items and the largest serial number was 33, what is the probability that $N < 50$?

Let $P(N = n|L = 33)$ be the PMF computed in part (b). Using the same prior from part (b), the outcomes that satisfy this event are those where N is between 33 and 49,

inclusive, so we sum over the PMF for this range of values of N :

$$\begin{aligned}
 P(N < 50|L = 33) &= \sum_{n=33}^{49} P(N = n|L = 33) \\
 &= \sum_{n=33}^{49} \frac{\binom{32}{4} \cdot \frac{1}{68}}{\sum_{i=33}^{100} \binom{32}{i} \cdot \frac{1}{68}}
 \end{aligned}$$

Historical Context: During World War 2, the Allies needed to know how many tanks Nazi Germany was producing. First they sent spies to Germany who estimated Germany produced 1,400 tanks per month. Separately, they noticed that the serial numbers on gear boxes on German tanks were unique and sequential. Using this observation, and a sample of gear box serial codes, mathematicians at the US Statistical Research Group used the math you derived to estimate the amount of tanks produced by Nazi Germany. They estimated expected production was 270 tanks per month. After the war, German records confirmed an actual production rate of 276 tanks per month—the probabilistic method was incredibly accurate!