

Section #4 Solutions

Based on the work of many prior CS109 staffs

- 1. 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:

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

- 2. Fairness in AI.** In their 2018 paper “Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification,” Joy Buolamwini and Timnit Gebru showed that commercial gender classifiers performed significantly worse on face images of darker-skinned females (error rates up to 34.7%) than lighter-skinned males (maximum error rate of 0.8%.) This disparity may result from training the classifiers on unbalanced datasets. To evaluate the classifiers, the authors designed their own dataset by collecting photos of national parliamentarians from three African countries and three European ones.

The probability table below shows the joint distribution of the dataset between two random variables: the demographic (D) of the photo subject and their country (C).

Demographic	South Africa	Senegal	Rwanda	Sweden	Finland	Iceland
Darker Female	0.12	0.05	0.04	0.01	0	0
Darker Male	0.15	0.07	0.02	0.01	0	0
Lighter Female	0.02	0	0	0.12	0.06	0.02
Lighter Male	0.05	0	0	0.14	0.09	0.03

- What is the marginal probability distribution for demographic D ? Provide your result as a mapping from values that D can take to probabilities.
- What is the conditional probability of country given that the subject is a lighter female, $P(C|D = \text{Lighter Female})$? Provide your result as a mapping from values that C can take to probabilities. Is this mapping a probability distribution?
- What is the conditional probability that the subject is from Senegal given their demographic, $P(C = \text{Senegal}|D)$? Provide your answer as a mapping from values that D can take to probabilities. Is this mapping a probability distribution?
- What are the pitfalls in using this dataset for a purpose beyond what the authors intended?

- For each assignment to D , sum over all the values of C that are consistent with that assignment.

$$P(\text{Darker Female}) = 0.12 + 0.05 + 0.04 + 0.01 + 0 + 0 = 0.22$$

$$P(\text{Darker Male}) = 0.15 + 0.07 + 0.02 + 0.01 + 0 + 0 = 0.25$$

$$P(\text{Lighter Female}) = 0.02 + 0 + 0 + 0.12 + 0.06 + 0.02 = 0.22$$

$$P(\text{Lighter Male}) = 0.05 + 0 + 0 + 0.14 + 0.09 + 0.03 = 0.31$$

-

$$P(\text{South Africa}|\text{Lighter Female}) = \frac{P(\text{South Africa, Lighter Female})}{P(\text{Lighter Female})} = \frac{0.02}{0.22} \approx 0.09$$

Similarly,

$$P(\text{Senegal}|\text{Lighter Female}) = 0$$

$$P(\text{Rwanda}|\text{Lighter Female}) = 0$$

$$P(\text{Sweden}|\text{Lighter Female}) \approx 0.55$$

$$P(\text{Finland}|\text{Lighter Female}) \approx 0.27$$

$$P(\text{Iceland}|\text{Lighter Female}) \approx 0.09$$

This mapping is the conditional probability distribution $P(C|D = \text{Lighter Female})$. Its probabilities sum to 1.

-
-

$$P(\text{Senegal}|\text{Darker Female}) = \frac{P(\text{Senegal, Darker Female})}{P(\text{Darker Female})} = \frac{0.05}{0.22} \approx 0.23$$

Similarly,

$$P(\text{Senegal}|\text{Darker Male}) \approx 0.28$$

$$P(\text{Senegal}|\text{Lighter Female}) = 0$$

$$P(\text{Senegal}|\text{Lighter Male}) = 0$$

This mapping is not a probability distribution because the conditioning event changes. We can also see that the probabilities do not sum to 1.

- d. This dataset does not come close to representing the diversity of the world; it draws subjects from just six countries. Even within those countries, the dataset may be unbalanced with respect to socioeconomic and cultural groups because the subjects are all parliamentarians. For example, the dataset may underrepresent ethnic minorities within those countries.

- 3. Hat-Check Again??** Recall the hat-check problem from section 2: n people go to a party and drop off their hats to a hat-check person. When the party is over, a different hat-check person is on duty, and returns the n hats randomly back to each person. Let X be the random variable representing the number of people who get their own hat back. We showed last time that $E[X] = 1$ for any n . What is $Var(X)$? Hint: Be careful when taking the variance of a sum of random variables.

Similarly to last time, let $X_i \sim \text{Bernoulli}(p = 1/n)$ be the indicator variable for whether the i^{th} person gets their hat back, so that $X = \sum_{i=1}^n X_i$. Then,

$$Var(X) = Var\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n Var(X_i) + 2 \sum_{i < j} Cov(X_i, X_j)$$

The first term is simply $np(1-p) = n\left(\frac{1}{n}\right)\left(1 - \frac{1}{n}\right) = \frac{n-1}{n}$ since each individual variance is $p(1-p)$.

To compute $Cov(X_i, X_j)$ for $i \neq j$, we note that $Cov(X_i, X_j) = E[X_i X_j] - E[X_i]E[X_j]$. The random variable $X_i X_j$ is also Bernoulli, with parameter $1/n \cdot 1/(n-1)$ since both have to get their hat back. So

$$Cov(X_i, X_j) = E[X_i X_j] - E[X_i]E[X_j] = \frac{1}{n(n-1)} - \left(\frac{1}{n}\right)^2 = \frac{1}{n^2(n-1)}$$

Noting that there are $\binom{n}{2} = \frac{n(n-1)}{2}$ identical terms in the summation over $i < j$ and putting this all together gives

$$Var(X) = \frac{n-1}{n} + 2\binom{n}{2} \frac{1}{n^2(n-1)} = \frac{n-1}{n} + \frac{1}{n} = 1.$$

So $E[X] = Var(X) = 1$. What a coincidence!