

**1. Working with the Law of Large Numbers**

- (a) A fair coin is tossed and you win a prize if there are more than 60% heads. Which is better: 10 tosses or 100 tosses? Explain.
- (b) A fair coin is tossed and you win a prize if there are more than 40% heads. Which is better: 10 tosses or 100 tosses? Explain.
- (c) A coin is tossed and you win a prize if there are between 40% and 60% heads. Which is better: 10 tosses or 100 tosses? Explain.
- (d) A coin is tossed and you win a prize if there are exactly 50% heads. Which is better: 10 tosses or 100 tosses? Explain.

**Solution:**

- (a) 10 tosses. By LLN, the sample mean should have higher probability to be close to the population mean as n increases. Therefore the average proportion of coins that are heads should be closer to 0.50, and has a lower chance of being greater than 0.60 if there are 100 tosses compared with 10 tosses.
- (b) 100 tosses. Based on the first part, consider the inverse of the event “more than 60% heads” and the symmetry of heads and tails.
- (c) 100 tosses. Based on the first part, consider the union of the events “more than 60% heads” and “more than 60% tails” (“less than 40% heads”).
- (d) 10 tosses. Compare the probability of getting equal number of heads and tails between  $2n$  and  $2n + 2$  tosses.

$$\begin{aligned}
 \Pr[n \text{ heads in } 2n \text{ tosses}] &= \binom{2n}{n} / 2^{2n} \\
 \Pr[n + 1 \text{ heads in } 2n + 2 \text{ tosses}] &= \binom{2n + 2}{n + 1} / 2^{2n + 2} \\
 &= \frac{(2n + 2)!}{(n + 1)!(n + 1)!} \cdot \frac{1}{2^{2n + 2}} \\
 &= \frac{(2n + 2)(2n + 1)2n!}{(n + 1)(n + 1)n!n!} \cdot \frac{1}{2^{2n + 2}} \\
 &= \frac{2n + 2}{n + 1} \cdot \frac{2n + 1}{n + 1} \binom{2n}{n} \cdot \frac{1}{2^{2n + 2}} \\
 &< \left( \frac{2n + 2}{n + 1} \right)^2 \binom{2n}{n} \cdot \frac{1}{2^{2n + 2}} \\
 &= 4 \binom{2n}{n} \cdot \frac{1}{2^{2n + 2}} = \binom{2n}{n} / 2^{2n} = \Pr[n \text{ heads in } 2n \text{ tosses}]
 \end{aligned}$$

The larger  $n$  is, the less probability we'll get 50% heads. □

## 2. Playing Pollster

As an expert in probability, the staff members at the Daily Californian have recruited you to help them conduct a poll to determine the percentage  $p$  of Berkeley undergraduates that plan to participate in the student sit-in. They've specified that they want your estimate  $\hat{p}$  to have an error of at most  $\varepsilon$  with confidence  $1 - \delta$ . That is,

$$P(|\hat{p} - p| \leq \varepsilon) \geq 1 - \delta.$$

Assume that you've been given the bound

$$P(|\hat{p} - p| \geq \varepsilon) \leq \frac{1}{4n\varepsilon^2},$$

where  $n$  is the number of students in your poll.

- (a) Using the formula above, what is the smallest number of students  $n$  that you need to poll so that your poll has an error of at most  $\varepsilon$  with confidence  $1 - \delta$ ?
- (b) At Berkeley, there are about 26,000 undergraduates and about 10,000 graduate students. Suppose you only want to understand the frequency of sitting-in for the undergraduates. If you want to obtain an estimate with error of at most 5% with 98% confidence, how many undergraduate students would you need to poll? Does your answer change if you instead only want to understand the frequency of sitting-in for the graduate students?
- (c) It turns out you just don't have as much time for extracurricular activities as you thought you would this semester. The writers at the Daily Californian insist that your poll results are reported with at least 95% confidence, but you only have enough time to poll 500 students. Based on the bound above, what is the worst-case error with which you can report your results?

### Solution:

- (a) We know we need to have

$$P(|\hat{p} - p| \leq \varepsilon) \geq 1 - \delta.$$

Subtracting both sides from 1, it follows that we must have

$$P(|\hat{p} - p| > \varepsilon) \leq \delta.$$

Therefore if we choose  $n$  such that

$$\frac{1}{4n\varepsilon^2} \leq \delta,$$

we will have

$$P(|\hat{p} - p| \geq \varepsilon) \leq \delta,$$

and since  $P(|\hat{p} - p| > \varepsilon) \leq P(|\hat{p} - p| \geq \varepsilon)$ , this will meet the requirement that

$$P(|\hat{p} - p| > \varepsilon) \leq \delta.$$

Thus we must have that

$$\begin{aligned}\frac{1}{4n\varepsilon^2} &\leq \delta \\ \frac{1}{n} &\leq 4\varepsilon^2\delta \\ n &\geq \frac{1}{4\varepsilon^2\delta}.\end{aligned}$$

- (b) Plugging in to the bound you found above, you get that  $n \geq 5000$ . The answer is the same for graduate students; the size of the population does not affect the number of samples you need.
- (c) If you only have time to poll 500 people and want to report your results with 95% confidence, you must report that the error in your estimate is at most 10%. You can find this by plugging in  $\frac{1}{4 \cdot 500 \cdot \varepsilon^2} = .05$  and solving for  $\varepsilon$ .

### 3. Covariance

We have a bag of 5 red and 5 blue balls. We take two balls from the bag without replacement. Let  $X_1$  and  $X_2$  be indicator random variables for the first and second ball being red. What is  $Cov(X_1, X_2)$ ?

**Solution:**

We can use the formula  $Cov(X_1, X_2) = E(X_1X_2) - E(X_1)E(X_2)$ .

$$\begin{aligned}E(X_1) &= \frac{5}{10} \times 1 + \frac{5}{10} \times 0 = \frac{1}{2} \\ E(X_2) &= \frac{5}{10} \times 1 + \frac{5}{10} \times 0 = \frac{1}{2} \\ E(X_1X_2) &= \frac{5}{10} \cdot \frac{4}{9} \times 1 + \left(1 - \frac{5}{10} \cdot \frac{4}{9}\right) \times 0 = \frac{2}{9}\end{aligned}$$

Therefore,

$$E(X_1X_2) - E(X_1)E(X_2) = \frac{2}{9} - \frac{1}{2} \times \frac{1}{2} = \frac{-1}{36}$$

### 4. LLSE

We have two bags of balls. The fractions of red balls and blue balls in bag A are  $\frac{2}{3}$  and  $\frac{1}{3}$  respectively. The fractions of red balls and blue balls in bag B are  $\frac{1}{2}$  and  $\frac{1}{2}$  respectively. Someone gives you one of the bags (unmarked) uniformly at random. Then we draw 6 balls from the same bag with replacement. Let  $X_i$  be the indicator random variable that ball  $i$  is

red. Now, let us define  $X = \sum_{1 \leq i \leq 3} X_i$  and  $Y = \sum_{4 \leq i \leq 6} X_i$ . Find  $LLSE(Y|X)$ . [Hint: recall that  $LLSE(Y|X) = E(Y) + \frac{Cov(X,Y)}{Var(X)}(X - E(X))$ ]

**Solution:**

$$\begin{aligned} E(X) &= 3 \cdot E(X_1) \\ &= 3 \cdot P(X_1 = 1) \\ &= 3 \cdot \left( \frac{1}{2} \cdot \frac{2}{3} + \frac{1}{2} \cdot \frac{1}{2} \right) \\ &= \frac{7}{4} \end{aligned}$$

$$E(Y) = E(X) = \frac{7}{4}$$

$$\begin{aligned} Cov(X, Y) &= Cov\left( \sum_{1 \leq i \leq 3} X_i, \sum_{4 \leq j \leq 6} X_j \right) \\ &= 9 \cdot Cov(X_1, X_4) \\ &= 9 \cdot (E(X_1 X_4) - E(X_1) \cdot E(X_4)) \end{aligned}$$

$$\begin{aligned} E(X_1 X_4) - E(X_1)E(X_4) &= P(X_1 = 1, X_4 = 1) - P(X_1 = 1)^2 \\ &= \left[ \frac{1}{2} \cdot \left( \frac{2}{3} \right)^2 + \frac{1}{2} \cdot \left( \frac{1}{2} \right)^2 \right] - \left[ \frac{1}{2} \cdot \left( \frac{2}{3} \right) + \frac{1}{2} \cdot \left( \frac{1}{2} \right) \right]^2 \\ &= \frac{1}{144} \end{aligned}$$

$$\begin{aligned} Var(X) &= Cov\left( \sum_{1 \leq i \leq 3} X_i, \sum_{1 \leq j \leq 3} X_j \right) \\ &= 3 \cdot Var(X_1) + 6 \cdot Cov(X_1, X_2) \\ &= 3(E(X_1^2) - E(X_1)^2) + 6 \cdot \frac{1}{144} \\ &= 3\left( \frac{7}{12} - \left( \frac{7}{12} \right)^2 \right) + 6 \cdot \frac{1}{144} \\ &= \frac{111}{144} \end{aligned}$$

$$\text{So, } LLSE(Y|X) = \frac{7}{4} + \frac{9}{111} \left( X - \frac{7}{4} \right) = \frac{3}{37} X + \frac{119}{74}$$