CS 70Discrete Mathematics and Probability TheorySpring 2025RaoHW 12

1 Coupon Collector Variance

Note 19

It's that time of the year again—Safeway is offering its Monopoly Card promotion. Each time you visit Safeway, you are given one of *n* different Monopoly Cards with equal probability. You need to collect them all to redeem the grand prize.

Let *X* be the number of visits you have to make before you can redeem the grand prize. Show that $Var(X) = n^2 (\sum_{i=1}^n i^{-2}) - \mathbb{E}[X].$

Solution:

Note that this is the coupon collector's problem, but now we have to find the variance. Let X_i be the number of visits we need to make before we have collected the *i*th unique Monopoly card actually obtained, given that we have already collected i-1 unique Monopoly cards. Then $X = \sum_{i=1}^{n} X_i$ and each X_i is geometrically distributed with p = (n - i + 1)/n. Moreover, the indicators themselves are independent, since each time you collect a new card, you are starting from a clean slate.

$$\operatorname{Var}(X) = \sum_{i=1}^{n} \operatorname{Var}(X_i) \qquad (\text{as the } X_i \text{ are independent})$$
$$= \sum_{i=1}^{n} \frac{1 - (n - i + 1)/n}{[(n - i + 1)/n]^2} \qquad (\text{variance of a geometric r.v. is } (1 - p)/p^2)$$
$$= \sum_{j=1}^{n} \frac{1 - j/n}{(j/n)^2} \qquad (\text{by noticing that } n - i + 1 \text{ takes on all values from 1 to } n)$$
$$= \sum_{j=1}^{n} \frac{n(n - j)}{j^2}$$
$$= \sum_{j=1}^{n} \frac{n^2}{j^2} - \sum_{j=1}^{n} \frac{n}{j}$$
$$= n^2 \left(\sum_{j=1}^{n} \frac{1}{j^2}\right) - \mathbb{E}[X] \qquad (\text{using the coupon collector problem expected value}).$$

2 Diversify Your Hand

Note 15 You are dealt 5 cards from a standard 52 card deck. Let *X* be the number of distinct values in your hand. For instance, the hand (A, A, A, 2, 3) has 3 distinct values.

- (a) Calculate $\mathbb{E}[X]$. (Hint: Consider indicator variables X_i representing whether *i* appears in the hand.)
- (b) Calculate Var(X). The answer expression will be quite involved; you do not need to simplify anything.

Solution:

(a) Let X_i be the indicator of the *i*th value appearing in your hand. Then, $X = X_1 + X_2 + ... + X_{13}$. (Here we let 13 correspond to K, 12 correspond to Q, and 11 correspond to J.) By linearity of expectation, $\mathbb{E}[X] = \sum_{i=1}^{13} \mathbb{E}[X_i]$.

We can calculate $\mathbb{P}[X_i = 1]$ by taking the complement, $1 - \mathbb{P}[X_i = 0]$, or 1 minus the probability that the card does not appear in your hand. This is $1 - \frac{\binom{48}{5}}{\binom{52}{5}}$.

Then,
$$\mathbb{E}[X] = 13 \mathbb{P}[X_1 = 1] = 13 \left(1 - \frac{\binom{48}{5}}{\binom{52}{5}} \right)$$

(b) To calculate variance, since the indicators are not independent, we have to use the formula $\mathbb{E}[X^2] = \sum_{i=j} \mathbb{E}[X_i^2] + \sum_{i \neq j} \mathbb{E}[X_i X_j].$

First, we have

$$\sum_{i=j} \mathbb{E}[X_i^2] = \sum_{i=j} \mathbb{E}[X_i] = 13 \left(1 - \frac{\binom{48}{5}}{\binom{52}{5}} \right).$$

Next, we tackle $\sum_{i \neq j} \mathbb{E}[X_i X_j]$. Note that $\mathbb{E}[X_i X_j] = \mathbb{P}[X_i X_j = 1]$, as $X_i X_j$ is either 0 or 1.

To calculate $\mathbb{P}[X_iX_j = 1]$ (the probability we have both cards in our hand), we note that $\mathbb{P}[X_iX_j = 1] = 1 - \mathbb{P}[X_i = 0] - \mathbb{P}[X_j = 0] + \mathbb{P}[X_i = 0, X_j = 0]$. Then

$$\sum_{i \neq j} \mathbb{E}[X_i X_j] = 13 \cdot 12 \mathbb{P}[X_i X_j = 1]$$

= 13 \cdot 12(1 - \mathbb{P}[X_i = 0] - \mathbb{P}[X_j = 0] + \mathbb{P}[X_i = 0, X_j = 0])
= 156 \left(1 - 2\frac{48}{52} + \frac{44}{52} \right) \right)

Putting it all together, we have

$$\operatorname{Var}(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$
$$= 13\left(1 - \frac{\binom{48}{5}}{\binom{52}{5}}\right) + 156\left(1 - 2\frac{\binom{48}{5}}{\binom{52}{5}} + \frac{\binom{44}{5}}{\binom{52}{5}}\right) - \left(13\left(1 - \frac{\binom{48}{5}}{\binom{52}{5}}\right)\right)^2$$

3 Double-Check Your Intuition Again

- Note 16 (a) You roll a fair six-sided die and record the result *X*. You roll the die again and record the result *Y*.
 - (i) What is cov(X+Y, X-Y)?
 - (ii) Prove that X + Y and X Y are not independent.

For each of the problems below, if you think the answer is "yes" then provide a proof. If you think the answer is "no", then provide a counterexample.

- (b) If X is a random variable and Var(X) = 0, then must X be a constant?
- (c) If X is a random variable and c is a constant, then is Var(cX) = c Var(X)?
- (d) If *A* and *B* are random variables with nonzero standard deviations and Corr(A, B) = 0, then are *A* and *B* independent?
- (e) If X and Y are not necessarily independent random variables, but Corr(X, Y) = 0, and X and Y have nonzero standard deviations, then is Var(X + Y) = Var(X) + Var(Y)?
- (f) If X and Y are random variables then is $\mathbb{E}[\max(X, Y) \min(X, Y)] = \mathbb{E}[XY]$?
- (g) If X and Y are independent random variables with nonzero standard deviations, then is

 $\operatorname{Corr}(\max(X,Y),\min(X,Y)) = \operatorname{Corr}(X,Y)?$

Solution:

(a) (i) Using bilinearity of covariance, we have

$$cov(X+Y,X-Y) = cov(X,X) + cov(X,Y) - cov(Y,X) - cov(Y,Y)$$
$$= cov(X,X) - cov(Y,Y),$$
$$= 0$$

where we use that cov(X, Y) = cov(Y, X) to get the second equality.

- (ii) Observe that $\mathbb{P}[X+Y=7, X-Y=0] = 0$ because if X-Y=0, then the sum of our two dice rolls must be even. However, both $\mathbb{P}[X+Y=7]$ and $\mathbb{P}[X-Y=0]$ are nonzero, so $\mathbb{P}[X+Y=7, X-Y=0] \neq \mathbb{P}[X+Y=7] \cdot \mathbb{P}[X-Y=0]$.
- (b) Yes. If we write $\mu = \mathbb{E}[X]$, then $0 = \operatorname{Var}(X) = \mathbb{E}[(X \mu)^2]$ so $(X \mu)^2$ must be identically 0 since perfect squares are non-negative. Thus $X = \mu$.
- (c) No. We have $\operatorname{Var}(cX) = \mathbb{E}[(cX \mathbb{E}[cX])^2] = c^2 \mathbb{E}[(X \mathbb{E}[X])^2] = c^2 \operatorname{Var}(X)$ so if $\operatorname{Var}(X) \neq 0$ and $c \neq 0$ or $c \neq 1$ then $\operatorname{Var}(cX) \neq c \operatorname{Var}(X)$. This does prove that $\sigma(cX) = c\sigma(X)$ though.
- (d) No. Let A = X + Y and B = X Y from part (a). Since A and B are not constants then part (b) says they must have nonzero variances which means they also have nonzero standard

deviations. Part (a) says that their covariance is 0 which means they are uncorrelated, and that they are not independent.

Recall from lecture that the converse is true though.

- (e) Yes. If Corr(X, Y) = 0, then cov(X, Y) = 0. We have Var(X + Y) = cov(X + Y, X + Y) = Var(X) + Var(Y) + 2 cov(X, Y) = Var(X) + Var(Y).
- (f) Yes. For any values x, y we have $\max(x, y) \min(x, y) = xy$. Thus, $\mathbb{E}[\max(X, Y) \min(X, Y)] = \mathbb{E}[XY]$.
- (g) No. You may be tempted to think that because $(\max(x,y),\min(x,y))$ is either (x,y) or (y,x), then $Corr(\max(X,Y),\min(X,Y)) = Corr(X,Y)$ because Corr(X,Y) = Corr(Y,X). That reasoning is flawed because $(\max(X,Y),\min(X,Y))$ is not always equal to (X,Y) or always equal to (Y,X) and the inconsistency affects the correlation. It is possible for X and Y to be independent while $\max(X,Y)$ and $\min(X,Y)$ are not.

For a concrete example, suppose X is either 0 or 1 with probability 1/2 each and Y is independently drawn from the same distribution. Then Corr(X,Y) = 0 because X and Y are independent. Even though X never gives information about Y, if you know max(X,Y) = 0 then you know for sure min(X,Y) = 0.

More formally, max(X, Y) = 1 with probability 3/4 and 0 with probability 1/4, and min(X, Y) = 1 with probability 1/4 and 0 with probability 3/4. This means

$$\mathbb{E}[\max(X,Y)] = 1 \cdot \frac{3}{4} + 0 \cdot \frac{1}{4} = \frac{3}{4}$$

and

$$\mathbb{E}[\min(X,Y)] = 1 \cdot \frac{1}{4} + 0 \cdot \frac{3}{4} = \frac{1}{4}$$

Thus,

$$\operatorname{cov}(\max(X,Y),\min(X,Y)) = \mathbb{E}[\max(X,Y)\min(X,Y)] - \frac{3}{16}$$

= $\frac{1}{4} - \frac{3}{16} = \frac{1}{16} \neq 0$

We conclude that $\operatorname{Corr}(\max(X,Y),\min(X,Y)) \neq 0 = \operatorname{Corr}(X,Y)$.

4 Dice Games

- Note 20
- (a) Alice rolls a die until she gets a 1. Let X be the number of total rolls she makes (including the last one), and let Y be the number of rolls on which she gets an even number. Compute E[Y | X = x], and use it to calculate E[Y].
 - (b) Bob plays a game in which he starts off with one die. At each time step, he rolls all the dice he has. Then, for each die, if it comes up as an odd number, he puts that die back, and adds a number of dice equal to the number displayed to his collection. (For example, if he rolls a

one on the first time step, he puts that die back along with an extra die.) However, if it comes up as an even number, he removes that die from his collection.

Compute the expected number of dice Bob will have after *n* time steps. (Hint: compute the value of $\mathbb{E}[X_k \mid X_{k-1} = m]$ to derive a recursive expression for X_k , where X_i is the random variable representing the number of dice after *i* time steps.)

Solution:

(a) Let's compute $\mathbb{E}[Y \mid X = x]$. If Alice makes *x* total rolls, then before rolling a 1, she makes x - 1 rolls that are not a 1. Since these rolls are independent, *Y* follows a binomial distribution with n = x - 1 and p = 3/5, and $\mathbb{E}[Y \mid X = x] = \frac{3}{5}(x - 1)$.

Now, we'd like to compute $\mathbb{E}[Y]$. With total expectation, we have

$$\mathbb{E}[Y] = \sum_{x} \mathbb{E}[Y \mid X = x] \mathbb{P}[X = x]$$
$$= \sum_{x} \frac{3}{5} (x - 1) \mathbb{P}[X = x]$$
$$= \frac{3}{5} \sum_{x} x \cdot \mathbb{P}[X = x] - \frac{3}{5} \sum_{x} \mathbb{P}[X = x]$$
$$= \frac{3}{5} \mathbb{E}[X] - \frac{3}{5}$$

Since *X* follows a geometric distribution with p = 1/6, $\mathbb{E}[X] = 6$, and

$$\mathbb{E}[Y] = \frac{3}{5}\mathbb{E}[X] - \frac{3}{5} = \frac{3}{5} \cdot 6 - \frac{3}{5} = 3.$$

(b) Let X_k be a random variable representing the number of dice after k time steps. In particular, this means that $X_0 = 1$. To compute the number of dice at step k, we first condition on $X_{k-1} = m$. Each one of the m dice is expected to leave behind 2 in its place, since there's a $\frac{1}{2}$ probability that it leaves behind 0 dice, a $\frac{1}{6}$ probability for each of 2, 4, and 6 dice, corresponding to rolling a 1, 3, and 5 respectively.

Therefore, we have $\mathbb{E}[X_k | X_{k-1} = m] = 2m$, so with total expectation, we have

$$\mathbb{E}[X_k] = \sum_m \mathbb{E}[X_k \mid X_{k-1} = m] \mathbb{P}[X_{k-1} = m]$$
$$= \sum_m 2m \cdot \mathbb{P}[X_{k-1} = m]$$
$$= 2\sum_m m \cdot \mathbb{P}[X_{k-1} = m]$$
$$= 2\mathbb{E}[X_{k-1}]$$

This means that we expect to have $\mathbb{E}[X_n] = 2\mathbb{E}[X_{n-1}] = 2^2\mathbb{E}[X_{n-2}] = \cdots = 2^n\mathbb{E}[X_0] = 2^n$ dice.

5 LLSE and Graphs

- Note 20 Consider a graph with *n* vertices numbered 1 through *n*, where *n* is a positive integer ≥ 2 . For each pair of distinct vertices, we add an undirected edge between them independently with probability *p*. Let D_1 be the random variable representing the degree of vertex 1, and let D_2 be the random variable representing the degree of vertex 2.
 - (a) Compute $\mathbb{E}[D_1]$ and $\mathbb{E}[D_2]$.
 - (b) Compute $Var(D_1)$.
 - (c) Compute $cov(D_1, D_2)$.
 - (d) Using the information from the first three parts, what is $L(D_2 | D_1)$?

Solution:

Throughout this problem, let $X_{i,j}$ be an indicator random variable for whether the edge between vertex *i* and vertex *j* exists, for i, j = 1, ..., n. Note that $X_{i,j} = X_{j,i}$.

(a) Observing that $D_1, D_2 \sim \text{Binomial}(n-1, p)$, we obtain $\mathbb{E}[D_1] = \mathbb{E}[D_2] = (n-1)p$.

Anyway, it is good to review how we derived the expectation of the binomial distribution in the first place. By linearity of expectation,

$$\mathbb{E}[D_1] = \mathbb{E}\left[\sum_{i=2}^n X_{1,j}\right] = \sum_{i=2}^n \mathbb{E}[X_{1,j}] = (n-1)\mathbb{E}[X_{i,j}] = (n-1)p$$

By symmetry, $\mathbb{E}[D_2] = (n-1)p$ also.

(b) Since $D_1, D_2 \sim \text{Binomial}(n-1, p)$, then $\text{Var} D_1 = \text{Var} D_2 = (n-1)p(1-p)$.

Again, it is good to review how we calculated the variance of the binomial distribution. Solution 1: Write the variance of D_1 as a sum of covariances.

 $\operatorname{Var}(D_1) = \operatorname{cov}\left(\sum_{i=2}^n X_{1,i}, \sum_{i=2}^n X_{1,i}\right) = (n-1)\operatorname{Var}(X_{1,2}) + \left((n-1)^2 - (n-1)\right)\operatorname{cov}(X_{1,2}, X_{1,3})$

$$= (n-1)p(1-p) + 0 = (n-1)p(1-p).$$

We used the fact that $X_{1,i}$ and $X_{1,j}$ are independent if $i \neq j$, so their covariance is zero. Solution 2: Compute the variance directly.

$$Var(X) = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \mathbb{E}\left[\left(\sum_{i=2}^n X_{1,i}\right)^2\right] - (n-1)^2 p^2$$

= $(n-1)\mathbb{E}[X_{1,2}^2] + ((n-1)^2 - (n-1))\mathbb{E}[X_{1,2}X_{1,3}] - (n-1)^2 p^2$
= $(n-1)p + (n^2 - 3n + 2)p^2 - (n-1)^2 p^2$
= $(n-1)p + (n-1)(n-2)p^2 - (n-1)^2 p^2 = (n-1)p(1 + (n-2)p - (n-1)p)$
= $(n-1)p(1-p)$

CS 70, Spring 2025, HW 12

(c) We can write

$$\operatorname{cov}(D_1, D_2) = \operatorname{cov}\left(\sum_{i=2}^n X_{1,i}, \sum_{i=1, i\neq 2}^n X_{2,i}\right) = \sum_{i=2}^n \sum_{j=1, j\neq 2}^n \operatorname{cov}(X_{1,i}, X_{2,j})$$

Note that all pairs of $X_{1,i}, X_{2,j}$ are independent except for when i = 2 and j = 1, so all terms in the sum are zero except for $cov(X_{1,2}, X_{2,1})$, and our covariance is just equal to $cov(X_{1,2}, X_{2,1}) = Var(X_{1,2}) = p(1-p)$.

(d) Since

$$L(D_2 | D_1) = \mathbb{E}[D_2] + \frac{\operatorname{cov}(D_1, D_2)}{\operatorname{Var}(D_1)}(D_1 - \mathbb{E}[D_1]),$$

we plug in our values from the first three parts to get that

$$L(D_2 | D_1) = (n-1)p + \frac{p(1-p)}{(n-1)p(1-p)} (D_1 - (n-1)p)$$

= $(n-1)p + \frac{1}{n-1} (D_1 - (n-1)p) = \frac{1}{n-1} D_1 + (n-2)p.$

14

6 Balls in Bins Estimation

Note 20 We throw n > 0 balls into $m \ge 2$ bins. Let *X* and *Y* represent the number of balls that land in bin 1 and 2 respectively.

- (a) Calculate $\mathbb{E}[Y \mid X]$. [*Hint*: Your intuition may be more useful than formal calculations.]
- (b) What is L[Y | X] (where L[Y | X] is the best linear estimator of Y given X)? [*Hint*: Your justification should be no more than two or three sentences, no calculations necessary! Think carefully about the meaning of the conditional expectation.]
- (c) Unfortunately, your friend is not convinced by your answer to the previous part. Compute $\mathbb{E}[X]$ and $\mathbb{E}[Y]$.
- (d) Compute Var(X).
- (e) Compute cov(X, Y).
- (f) Compute L[Y | X] using the formula. Ensure that your answer is the same as your answer to part (b).

Solution:

(a) $\mathbb{E}[Y | X = x] = (n - x)/(m - 1)$, because once we condition on x balls landing in bin 1, the remaining n - x balls are distributed uniformly among the other m - 1 bins. Therefore,

$$\mathbb{E}[Y \mid X] = \frac{n-X}{m-1}.$$

- (b) We showed that $\mathbb{E}[Y | X]$ is a linear function of *X*. Since $\mathbb{E}[Y | X]$ is the best *general* estimator of *Y* given *X*, it must also be the best *linear* estimator of *Y* given *X*, i.e. $\mathbb{E}[Y | X]$ and L[Y | X] coincide.
- (c) Let X_i be the indicator that the *i*th ball falls in bin 1. Then, $X = \sum_{i=1}^{n} X_i$, and by linearity of expectation, $\mathbb{E}[X] = \sum_{i=1}^{n} \mathbb{E}[X_i] = n/m$, since there are *n* indicators and each ball has a probability 1/m of landing in bin 1. By symmetry, $\mathbb{E}[Y] = n/m$ as well.
- (d) The number of balls that falls into the first bin is binomially distributed with parameters n and 1/m. Hence the variance is n(1/m)(1-1/m).
- (e) Let X_i be as before, and let Y_i be the indicator that the *i*th ball falls into bin 2.

$$\operatorname{cov}(X,Y) = \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{cov}(X_i,Y_j)$$

We can compute $cov(X_i, Y_i) = \mathbb{E}[X_iY_i] - \mathbb{E}[X_i]\mathbb{E}[Y_i] = 0 - (1/m)(1/m) = -1/m^2$ (note that $\mathbb{E}[X_iY_i] = 0$ because it is impossible for a ball to land in both bins 1 and 2). Also, we have $cov(X_i, Y_j) = 0$ because the indicator for the *i*th ball is independent of the indicator for the *j*th ball when $i \neq j$. Hence, $cov(X, Y) = n(-1/m^2) = -n/m^2$.

(f)

$$L[Y \mid X] = \mathbb{E}[Y] + \frac{\operatorname{cov}(X, Y)}{\operatorname{var}(X)} (X - \mathbb{E}[X])$$
$$= \frac{n}{m} + \frac{-n/m^2}{n(1/m)(1 - 1/m)} \left(X - \frac{n}{m}\right)$$
$$= \frac{n}{m} - \frac{1}{m-1} \left(X - \frac{n}{m}\right)$$
$$= \frac{mn - n - mX + n}{m(m-1)} = \frac{n - X}{m - 1}$$