

DISTRIBUTIONS, VARIANCE, INEQUALITIES, CONFIDENCE INTERVALS 8

COMPUTER SCIENCE MENTORS 70

October 31 to November 4, 2016

1 Distributions

1.1 Introduction

Geometric Distribution: $\text{Geom}(p)$ Number of trials required to obtain the first success. Each trial has probability of success equal to p . The probability of the first success happening at trial k is:

$$P[X = k] = (1 - p)^{k-1} * p, k > 0$$

The expectation of a geometric distribution is:

$$E(X) = \frac{1}{p}$$

The variance of a geometric distribution is:

$$\text{Var}(X) = \frac{1 - p}{p^2}$$

Solution: Derivation of $E(X)$: The clever way to find the expectation of the geometric distribution uses a method known as the renewal method. $E(X)$ is the expected number of trials until the first success. Suppose we carry out the first trial, and one of two outcomes occurs. With probability p , we obtain a success and we are done (it

only took 1 trial until success). With probability $1 - p$, we obtain a failure, and we are right back where we started. In the latter case, how many trials do we expect until our first success? The answer is $1 + E(X)$: we have already used one trial, and we expect $E(X)$ more since nothing has changed from our original situation (the geometric distribution is memoryless). Hence $E(X) = p * 1 + (1 - p) * (1 + E(X))$

Binomial Distribution: $\text{Bin}(n, p)$ Number of successes when we do n independent trials. Each trial has a probability p of success. The probability of having k successes:

$$P[X = k] = \binom{n}{k} * p^k * (1 - p)^{n-k}$$

The expectation of a binomial distribution is:

$$E(X) = np$$

The variance of a binomial distribution is:

$$\text{Var}(X) = np(1 - p)$$

Solution: Can walk through the derivation of $E(X)$: We would have to compute this sum:

$$E(X) = \sum_k k * P[X = k] = \sum_{k=0}^n k * \binom{n}{k} * p^k * (1 - p)^{n-k}$$

Instead of doing that just use Bernoulli variables:

$$X = X_1 + \dots + X_n$$

And now use linearity of expectation:

$$E(X) = E(X_1 + \dots + X_n) = E(X_1) + \dots + E(X_n)$$

Since the probability of a success happening at each step is p , and there are n steps, we are just summing p n times.

Poisson Distribution: $\text{Pois}(\lambda)$ This is an approximation to the binomial distribution. Let the number of trials approach infinity, let the probability of success approach 0, such that $E(X) = np = \lambda$. This is an accepted model for rare events. The probability of having k successes:

$$P[X = k] = \frac{e^{-\lambda} * \lambda^k}{k!}$$

The expectation of a poisson distribution is:

$$E(X) = \lambda$$

The variance of a poisson distribution is:

$$\text{Var}(X) = \lambda$$

Solution: Can walk through the derivation of $P(X)$:

$$\begin{aligned} P[X = k] &= \binom{n}{k} * p^k * (1 - p)^{n-k} \\ &= \frac{n!}{k! * (n - k)!} * p^k * (1 - p)^{n-k} \\ &\approx \frac{n^k * p^k}{k!} * \left(1 - \frac{\lambda}{n}\right) \\ &\approx \frac{\lambda^k * e^{-\lambda}}{k!} \end{aligned}$$

$$\begin{aligned} E(X) &= \sum_{k=0}^{\infty} k * \frac{e^{-\lambda} * \lambda^k}{k!} \\ &= \sum_{k=1}^{\infty} k * \frac{e^{-\lambda} * \lambda^k}{k!} \\ &= e^{-\lambda} * \lambda * \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} \\ &= e^{-\lambda} * \lambda * \sum_{k=1}^{\infty} \frac{\lambda^k}{k!} \\ &= e^{-\lambda} * \lambda * e^{\lambda} \\ &= \lambda \end{aligned}$$

1.2 Questions

1. You are Eve, and as usual, you are trying to break RSA. You are trying to guess the factorization of N , from Bobs public key. You know that N is approximately 1,000,000,000,000. To find the primes p and q , you decide to try random numbers from 2 to 1,000,000 $\approx \sqrt{N}$, and see if they divide N .

To do this, you roll a 999,999-sided die to choose the number, and see if it divides N using your calculator, which takes five seconds. Of course, there will be one number in this range that does divide N namely, the smaller of p and q .

- (a) What kind of distribution would you use to model this?

Solution: Geometric probability of success each time is $p = \frac{1}{999,999}$

- (b) What is the expected amount of time until you guess the correct answer, if it takes five seconds per guess (you only have a calculator)? Answer in days.

Solution:

$$E(x) = \frac{1}{p} = 999,999 \text{ tries}$$

$$(999,999 * 5 \text{ sec}) * \frac{1 \text{ min}}{60 \text{ sec}} * \frac{1 \text{ hr}}{60 \text{ min}} * \frac{1 \text{ day}}{24 \text{ hr}} \approx 57.9 \text{ days}$$

2. Now you are trying to guess the 6-digit factorization digit by digit. Lets assume that when you finish putting these digits together, you can figure out how many digits you got right. Use zeros for blank spaces. For example, to guess 25, you would put 000025

(a) What kind of distribution would you use to model this?

Solution: Binomial, since this is multiple independent trials that can either succeed or fail.

(b) What is the probability that you get exactly 4 digits right?

Solution: $\binom{6}{4} * \frac{1}{10}^4 * \frac{9}{10}^2$

(c) What is the probability that you get less than 3 correct?

Solution: $\binom{6}{2} * \frac{1}{10}^2 * \frac{9}{10}^4 + \binom{6}{1} * \frac{1}{10} * \frac{9}{10}^5$

3. You are Alice, and you have a high-quality RSA-based security system. However, Eve is often successful at hacking your system. You know that the number of security breaches averages 3 a day, but varies greatly.

(a) What kind of distribution would you use to model this?

Solution: Poisson! That's what we use to model the probably frequencies of rare events.

(b) What is the probability you experience exactly seven attacks tomorrow? At least seven (no need to simplify your answer)?

Solution:

$$P[X = 7] = \frac{\lambda^7}{7!} * e^{-\lambda} = \frac{3^7}{7!} * e^{-3} \approx 0.0216$$

$$P[X \geq 7] = \sum_{i=7}^{\infty} P[X = i] = \sum_{i=7}^{\infty} \frac{3^i}{i!} * e^{-3}$$

(c) What is the probability that, on some day in April, you experience exactly six attacks?

Solution:

$$P[X = 6] = \frac{3^6}{6!} * e^{-3} \approx 0.0504$$
$$1 - (1 - 0.0504)^{30} \approx 0.788 = 78.8\%$$

2 Variance

2.1 Introduction

For a random variable X with expectation $E(X) = \mu$, the variance of X is:

$$\text{Var}(X) = E((X - \mu)^2)$$

The square root of $\text{Var}(X)$ is called the standard deviation of X

Theorem: For a random variable X with expectation $E(X) = \mu$ and a constant c ,

$$\text{Var}(X) = E(X^2) - \mu^2$$

$$\text{Var}(cX) = c^2 * \text{Var}(X)$$

Theorem: For a random variable X , expectation $E(X) =$

$$\sum_a a * Pr[X = A]$$

2.2 Questions

1. Let's consider the classic problems of flipping coins and rolling dice. Let X be a random variable for the number of coins that land on heads and Y be the value of the die roll.
 - (a) What is the expected value of X after flipping 3 coins? What is the variance of X ?

Solution:

$$E(X) = 0 * \frac{1}{8} + 1 * \frac{3}{8} + 2 * \frac{3}{8} + 3 * \frac{1}{8} = \frac{3}{2}$$

$$E(X^2) = 0^2 * \frac{1}{8} + 1^2 * \frac{3}{8} + 2^2 * \frac{3}{8} + 3^2 * \frac{1}{8} = \frac{24}{8} = 3$$

$$E(X)^2 = \frac{9}{4} \text{Var}(X) = 3 - \frac{9}{4} = \frac{3}{4}$$

(b) Let Y be the sum of rolling a dice 1 time. What is the expected value of Y ?

Solution: $E(Y) = \frac{1}{6} * (1 + 2 + 3 + 4 + 5 + 6) = \frac{7}{2}$

(c) What is the variance of Y ?

Solution: $E(Y^2) = [\frac{1}{6}(1^2+2^2+3^2+4^2+5^2+6^2)] = \frac{91}{6}$ $\text{Var}(Y) = E(Y^2) - (E(Y))^2 = \frac{91}{6} - \frac{49}{4} = \frac{35}{12}$

2. Say you're playing a game with a coin and die, where you flip the coin 3 times and roll the die once. In this game, your score is given by the number of heads that show multiplied with the die result. What is the expected value of your score? What's the variance?

Solution: $E(XY) = E(X)E(Y) = \frac{21}{4}$ since X and Y are independent. $\text{Var}(XY) = E(X^2Y^2) - E(XY)^2 = E(X^2)E(Y^2) - E(X)^2E(Y)^2 = 3 * \frac{91}{6} - \frac{3^2}{2} * \frac{7^2}{2} = \frac{91}{2} - \frac{9}{4} * \frac{49}{4} = 17.9375 = \frac{287}{16}$

3. You are at a party with n people where you have prepared a red solo cup labeled with their name. Before handing red cups to your friends, you pick up each cup and put a sticker on it with probability $\frac{1}{2}$ (independently of the other cups). Then you hand back the cups according to a uniformly random permutation. Let X be the number of people who get their own cup back AND it has a sticker on it.

(a) Compute the expectation $E(X)$.

Solution: Define $X_i = 1$ if the i -th person gets their own cup back and it has a sticker on it 0 otherwise Hence $E(X) = E(\sum_{i=1}^n X_i) = \sum_{i=1}^n E(X_i)$ $E(X_i) = P[X_i = 1] = \frac{1}{2n}$ since the i -th student will get his/her cup with probability $\frac{1}{n}$ and has a sticker on it with probability $\frac{1}{2}$ and stickers are put

independently. Hence $E(X) = n \cdot \frac{1}{2n} = \frac{1}{2}$.

(b) Compute the variance $\text{Var}(X)$

Solution: To calculate $\text{Var}(X)$, we need to know $E(X^2)$

$$E(X^2) = E(X_1 + X_2 + \dots + X_n)^2 = E\left(\sum_{i,j} (X_i * X_j)\right) = \sum_{i,j} (E(X_i * X_j))$$

(by linearity of expectation)

Then we consider two cases, either $i = j$ or $i \neq j$. Hence

$$\sum_{i,j} E(X_i * X_j) = \sum_i E(X_i^2) + \sum_{i \neq j} E(X_i * X_j)$$

$E(X_i^2) = \frac{1}{2n}$ for all i . To find $E(X_i * X_j)$, we need to calculate $P[X_i X_j = 1]$. $P[X_i * X_j = 1] = P[X_i = 1]P[X_j = 1 | X_i = 1] = \frac{1}{2n} * \frac{1}{2*(n-1)}$ since if student i has received his/her own cup, student j has $n - 1$ choices left. Hence

$$E(X^2) = n * \frac{1}{2n} + n * (n - 1) * \frac{1}{2n} * \frac{1}{2 * (n - 1)} = \frac{3}{4}$$

$$\text{Var}(X) = E(X^2) - E(X)^2 = \frac{3}{4} - \frac{1}{4} = \frac{1}{2}.$$

4. a. Prove that for independent random variables X and Y , $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$.

Solution:

$$\begin{aligned}\text{Var}(X + Y) &= E((X + Y)^2) - E(X + Y)^2 \\ &= E(X^2) + E(Y^2) + 2 * E(XY) - (E(X) + E(Y))^2 \\ &= (E(X^2) - E(X)^2) + (E(Y^2) - E(Y)^2) + 2 * E(XY) - E(X) * E(Y) \\ &= \text{Var}(X) + \text{Var}(Y) + 2 * (E(XY) - E(X)E(Y))\end{aligned}$$

- b. Is the above result true for non-independent random variables? Prove or give a counterexample.

Solution: No! One simple counterexample is $X = Y$. Then

$$\begin{aligned}\text{Var}(X + Y) &= \text{Var}(2X) \\ &= E((2X)^2) - E(2X)^2 \\ &= E(4X^2) - (2 * E(X))^2 \\ &= 4E(X^2) - 4(E(X))^2 \\ &= 4(E(X^2) - E(X)^2) \\ &= 4 * \text{Var}(X)\end{aligned}$$

5. Consider the random variable $X = X_1 + \dots + X_n$, where X_i equals i with probability $\frac{1}{i}$ and 0 otherwise.

- (a) What is the variance of X ? (Assume that X_i and X_j are independent for $i \neq j$)

Solution:

$$\begin{aligned}\text{Var}(X) &= \text{Var}(X_1) + \dots + \text{Var}(X_n) \\ E(X_i^2) &= P[X_i = i] * i^2 + P[X_i = 0] * 0^2 = \frac{1}{i} * i^2 + 0 = i \\ (E(X_i))^2 &= (P[X_i = i] * i + P[X_i = 0] * 0)^2 = \left(\frac{1}{i} * i + 0\right)^2 = 1 \\ \text{Var}(X_i) &= E(X_i^2) - (E(X_i))^2 = i - 1\end{aligned}$$

Recall,

$$\sum_{k=1}^n k = \frac{n * (n + 1)}{2}$$

$$\begin{aligned}
\text{Var}(X) &= \sum_i \text{Var}(X_i) \\
&= \sum_i i - 1 \\
&= -n + \sum_i i \\
&= -n + \frac{n * (n + 1)}{2} \\
&= \frac{n * (n + 1) - 2n}{2} \\
&= \frac{n^2 + n - 2n}{2} \\
&= \frac{n^2 - n}{2} \\
&= \frac{n * (n - 1)}{2}
\end{aligned}$$

(b) For what value of n does $E(X) = \text{Var}(X)$?

Solution:

$$\begin{aligned}
E(X_i) &= P[X_i = i] * i + 0 \\
&= \frac{1}{i} * i = 1
\end{aligned}$$

$$\begin{aligned}
E(X) &= E(X_1 + \dots + X_n) \\
&= E(X_1) + \dots + E(X_n) \\
&= n
\end{aligned}$$

$$\begin{aligned}
n &= \frac{n * (n - 1)}{2} \rightarrow 1 = \frac{n - 1}{2} \\
&\rightarrow 2 = n - 1 \\
&\rightarrow 3 = n
\end{aligned}$$

(c) For what value of n does $E(X) = SD(X) * \sqrt{2} + 100$?

Solution: $E(X) = n, SD(X) = \sqrt{\text{Var}(X)}$

$$n = \sqrt{\frac{n * (n - 1)}{2}} * \sqrt{2} + 100$$

$$n = \sqrt{n * (n - 1)} + 100$$

$$(n - 100)^2 = n * (n - 1)$$

$$n^2 - 200n + 10000 = n^2 - n$$

$$10000 = 199n$$

$$n = \frac{10000}{199}$$

6. An urn contains n balls numbered $1, 2, \dots, n$. We remove k balls at random (without replacement) and add up their numbers. Find the mean and variance of the total.

Solution: The required total is $T = \sum_{i=1}^k X_i$, where X_i is the number shown on the i th ball. Hence $E(T) = k * E(X_1) = \frac{1}{2} * k * (n + 1)$. Now calculate:

$$\begin{aligned} E\left(\left(\sum_{i=1}^k X_i\right)^2\right) &= kE(X_1^2) + k * (k - 1) * E(X_1 * X_2) \\ &= \frac{k}{n} \sum_1^n j^2 + \frac{k * (k - 1)}{n * (n - 1)} * 2 * \sum_{i>j} i * j \\ &= \frac{k}{n} \left(\frac{1}{3} * n * (n + 1) * (n + 2) - \frac{1}{2} * n * (n + 1) \right) \\ &\quad + \frac{k * (k - 1)}{n * (n - 1)} * \sum_{j=1}^n j * (n * (n + 1) - j * (j + 1)) \\ &= \frac{1}{6} * k * (n + 1) * (2n + 1) + \frac{1}{12} * k * (k - 1) * (3n + 2) * (n + 1) \end{aligned}$$

Hence,

$$\begin{aligned} \text{Var}(T) &= k(n + 1) \left(\frac{1}{6} k(n + 1)(2n + 1) + \frac{1}{12} k(k - 1)(3n + 2)(n + 1) - \frac{1}{4} k(n + 1) \right) \\ &= \frac{1}{12} (n + 1) k(n - k) \end{aligned}$$

3 Markov, Chebyshev

3.1 Introduction

Markov's Inequality

For a non-negative random variable X with expectation $E(X) = \mu$, and any $\alpha > 0$:

$$P[X \geq \alpha] \leq \frac{E(X)}{\alpha}$$

Solution: Proof of Markov's Inequality

$$\begin{aligned} E(X) &= \sum_a a * Pr[X = a] \\ &\geq \sum_{a \geq \alpha} a \geq \alpha a * Pr[X = a] \\ &\geq \alpha \sum_{a \geq \alpha} Pr[X = a] \\ &= \alpha Pr[X \geq \alpha] \end{aligned}$$

Chebyshev's Inequality

For a random variable X with expectation $E(X) = \mu$, and any $\alpha > 0$:

$$P[|X - \mu| \geq \alpha] \leq \frac{\text{Var}(X)}{\alpha^2}$$

3.2 Questions

1. Use Markov's to prove Chebyshev's Inequality:

Solution: Define the random variable $Y = (X - \mu)^2$. Note that $E(Y) = E((X - \mu)^2) = \text{Var}(X)$. Also, notice that the event that we are interested in, $|X - \mu| \geq \alpha$ is exactly the same as the event $Y \geq \alpha^2$. Therefore, $P[|X - \mu| \geq \alpha] = P[Y \geq \alpha^2]$. Moreover, Y is non-negative, so we can apply Markov's inequality to it to get:

$$P[Y \geq \alpha^2] \leq \frac{E(Y)}{\alpha^2} = \frac{\text{Var}(X)}{\alpha^2}$$

2. Squirrel Standard Deviation

As we all know, Berkeley squirrels are extremely fat and cute. The average squirrel is 40% body fat. The standard deviation of body fat is 5%. Provide an upper bound on the probability that a randomly trapped squirrel is either too skinny or too fat? A skinny squirrel has less than 27.5% body fat, and a fat squirrel has more than 52.5% body fat?

Solution: We use Chebyshev's inequality. We are looking for the probability we fall within 2.5 standard deviations of the mean. By Chebyshev's inequality, the probability we are within this range is $\frac{1}{2.5^2}$, or $\frac{4}{25} = 0.16$. If we were to use Markov's inequality, we would get probabilities over 1, which yields a non-helpful value.

3. Bound It

A random variable X is always strictly larger than -100. You know that $E(X) = 60$. Give the best upper bound you can on $P[X \geq 20]$.

Solution: Notice that we do not have the variance of X , so Chebyshev's bound is not applicable here. Since X is also not a sum of other random variables, other bounds or approximations (Chernoff, Hoeffding's inequality. Don't worry about them if they are not covered.) are not available. This leaves us with just Markov's Inequality. But Markov Bound only applies on a non-negative random variable, whereas X can take on negative values.

This suggests that we want to shift X somehow, so that we can apply Markov's Inequality on it. Define a random variable $Y = X + 100$, which means Y is strictly larger than 0, since X is always strictly larger than -100. Then, $E(Y) = E(X + 100) = E(X) + 100 = 60 + 100 = 160$. Finally, the upper bound on X that we want can be calculated via Y , and we can now apply Markov's Inequality on Y since Y is strictly positive.

$$P[X \geq 20] = P[Y \geq 120] \leq \frac{E(Y)}{120} = \frac{160}{120} = \frac{4}{3}$$

Hence, the best upper bound on $P[X \geq 20]$ is $\frac{4}{3}$.

4. Give a distribution for a random variable where the expectation is 1,000,000 and the probability that the random variable is zero is 99%.

Solution: X is 100,000,000 with probability 0.01, and 0 otherwise.

5. Consider a random variable Y with expectation μ whose maximum value is $\frac{3\mu}{2}$, prove that the probability that Y is 0 is at most $\frac{1}{3}$.

Solution:

$$\begin{aligned}
\mu &= \sum_a a * P[Y = a] \\
&= \sum_{a \neq 0} a * P[Y = a] \\
&\leq \sum_{a \neq 0} \frac{3\mu}{2} * P[Y = a] \\
&= \frac{3\mu}{2} * \sum_{a \neq 0} P[Y = a] \\
&= \frac{3\mu}{2} * (1 - P[Y = 0])
\end{aligned}$$

This implies that $P[Y = 0] \leq \frac{1}{3}$

6. Let X be the sum of 20 i.i.d. Poisson random variables X_1, \dots, X_{20} with $E(X_i) = 1$. Find an upper bound of $P[X \geq 26]$ using,

(a) Markov's inequality:

Solution:

$$\begin{aligned}
P[X \geq a] &\leq \frac{E(X)}{a} \text{ for all } a > 0 \\
P[X \geq 26] &\leq \frac{20}{26} \\
&\approx 0.769
\end{aligned}$$

(b) Chebyshev's inequality:

Solution:

$$\begin{aligned}
P[|X - E(X)| \geq c] &\leq \frac{\sigma_X^2}{c^2} \\
P[|X - 20| \geq 6] &\leq \frac{20}{36} \\
&\approx 0.5556
\end{aligned}$$

4 Confidence Intervals

4.1 Questions

1. Define i. i. d. variables $A_k \sim \text{Bern}(p)$ where $k \in [1, n]$. Assume we can declare that $P[|\frac{1}{n} \sum_k A_k - p| > 0.25] = 0.01$.

(a) Please give a 99% confidence interval for p if given A_k .

Solution: $[\frac{1}{n} \sum_i A_k - 0.25, \frac{1}{n} \sum_i A_k + 0.25]$

- (b) We know that the variables X_i , for i from 1 to n , are i.i.d. random variables and have variance. We also have a value (an observation) of $A_n = \frac{X_1 + \dots + X_n}{n}$. We want to guess the mean, μ , of each X_i .

Prove that we have 95% confidence μ lies in the interval $[A_n - 4.5 \frac{\sigma}{\sqrt{n}}, A_n + 4.5 \frac{\sigma}{\sqrt{n}}]$

That is, $P[\mu \in [A_n - 4.5 \frac{\sigma}{\sqrt{n}}, A_n + 4.5 \frac{\sigma}{\sqrt{n}}]] \geq 95\%$

Solution: To do this, we use Chebyshev's. Because $E[A_n] = \mu$ (A_n is the average of the X_i s), we bound the probability that $|A_n - \mu|$ is more than the interval size at 5%:

$$P[|A_n - \mu| \geq 4.5 \frac{\sigma}{\sqrt{n}}] \leq \frac{\text{Var}(A_n)}{(4.5 \frac{\sigma}{\sqrt{n}})^2} \approx \frac{\frac{\sigma^2}{n}}{\frac{20\sigma^2}{n}} = \frac{1}{20} = 5\%$$

.

Thus, the probability that μ is in the interval is 95

- (c) Give the 99% confidence interval for μ :

Solution: Solution is similar to that of the 95% confidence interval.

$[A_n - 10 \frac{\sigma}{\sqrt{n}}, A_n + 10 \frac{\sigma}{\sqrt{n}}]$, because $P[|A_n - \mu| \geq 10 \frac{\sigma}{\sqrt{n}}] \leq \frac{\text{Var}(A_n)}{(10 \frac{\sigma}{\sqrt{n}})^2} \approx \frac{\frac{\sigma^2}{n}}{\frac{100\sigma^2}{n}} = \frac{1}{100} = 1\%$.

2. We have a die whose 6 faces are values of consecutive integers, but we don't know where it starts (it is shifted over by some value k ; for example, if $k = 6$, the die faces would take on the values 7, 8, 9, 10, 11, 12). If we observe that the average of the n samples (n is large enough) is 15.5, develop a 99% confidence interval for the value of k .

Solution: PUT SOLUTION HERE