

Lab 2: Probability Theory

W203: Statistics for Data Science

Nishant Velagapudi

1. Meanwhile, at the Unfair Coin Factory...

You are given a bucket that contains 100 coins. 99 of these are fair coins, but one of them is a trick coin that always comes up heads. You select one coin from this bucket at random. Let T be the event that you select the trick coin. This means that $P(T) = 0.01$.

- a. To see if the coin you have is the trick coin, you flip it k times. Let H_k be the event that the coin comes up heads all k times. If you see this occur, what is the conditional probability that you have the trick coin? In other words, what is $P(T|H_k)$.

We know that:

$$P(T|H_k) = \frac{P(T \cap H_k)}{P(H_k)}$$

We also know that:

$$P(H_k) = \frac{1}{100} + \frac{99}{100} * \frac{1}{2^k}$$

Since the probability of the trick coin is $1/100$ (and any time we receive the trick coin, we are guaranteed to observe any amount of heads in a row):

$$P(T \cap H_k) = 1$$

Therefore:

$$P(T|H_k) = \frac{\frac{1}{100}}{\frac{1}{100} + \frac{99}{100} * \frac{1}{2^k}}$$

Which simplifies to:

$$P(T|H_k) = \frac{2^k}{2^k + 99}$$

- b. How many heads in a row would you need to observe in order for the conditional probability that you have the trick coin to be higher than 99%?

We are solving for k , when $P(T|H_k) = .99$:

$$.99 = \frac{2^k}{2^k + 99}$$

Simplifying yields:

$$98.01 = .01 * 2^k$$

We can then take the logarithm of each side:

$$\log(9801)/\log(2) = k, k = 13.258$$

Thus, we need at least 14 heads in a row (rounding k up to nearest int) to have 99% confidence that we have the trick coin. # 2. Wise Investments

You invest in two startup companies focused on data science. Thanks to your growing expertise in this area, each company will reach unicorn status (valued at \$1 billion) with probability $3/4$, independent of the other company. Let random variable X be the total number of companies that reach unicorn status. X can take on the values 0, 1, and 2. Note: X is what we call a binomial random variable with parameters $n = 2$ and $p = 3/4$.

- a. Give a complete expression for the probability mass function of X .

$$F(X) = \begin{cases} 0, & p = (.25) * (.25) \\ 1, & p = (.75) * (.25) + (.25) * (.75) \\ 2, & p = (.75) * (.75) \end{cases}$$

- b. Give a complete expression for the cumulative probability function of X .

$$E(X) = \begin{cases} 0, & p = 1/16 \\ 1, & p = 1/16 + 3/16 + 3/16 \\ 2, & p = 1/16 + 2 * (3/16) + 9/16 \end{cases}$$

- c. Compute $E(X)$. We can calculate this by taking the sum of $x * p(x)$ for all values of X .

$$E(X) = 0 * p(X = 0) + 1 * p(X = 1) + 2 * p(X = 2)$$

$$E(X) = 0 + 1 * 2 * (3/16) + 2 * (9/16), E(X) = 1.5$$

- d. Compute $var(X)$. We can calculate this by taking the sum of $(p(x) * (x - E(X)))^2$ for all values of x .

$$Var(X) = 1/16 * (0 - 1.5)^2 + 6/16 * (1 - 1.5)^2 + 9/16 * (2 - 1.5)^2$$

$$var(X) = 24/64$$

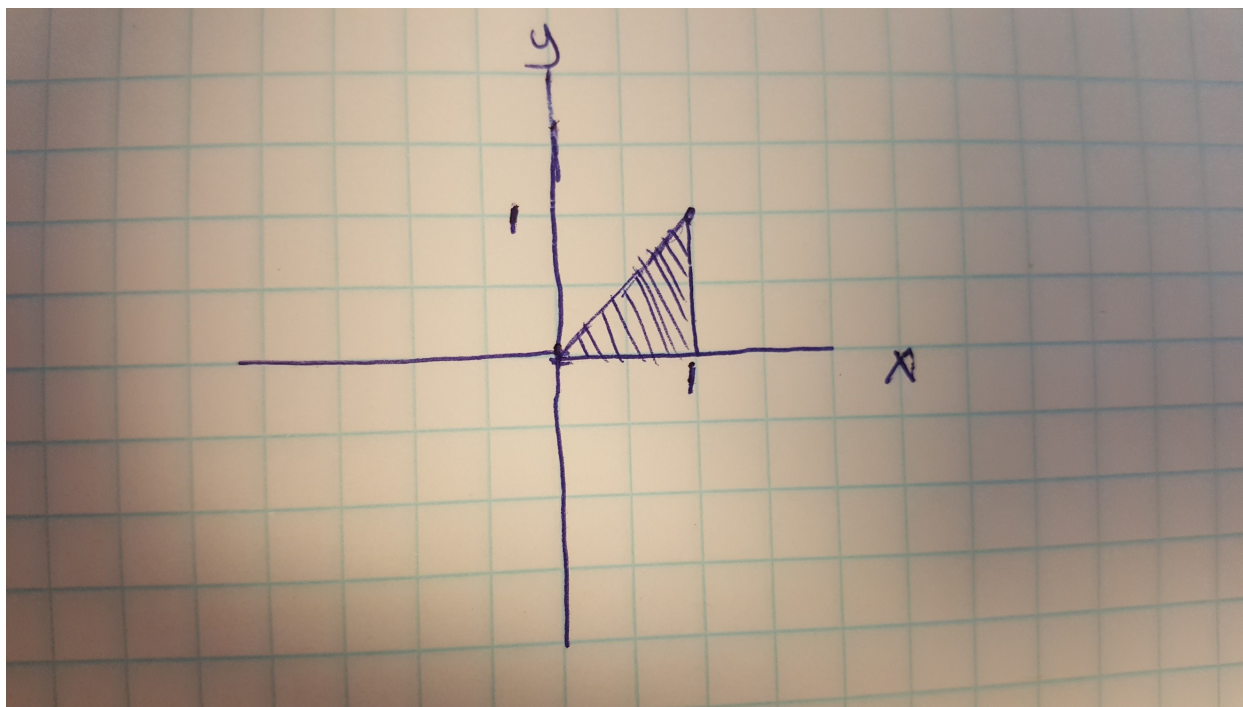
3. Relating Min and Max

Continuous random variables X and Y have a joint distribution with probability density function,

$$f(x, y) = \begin{cases} 2, & 0 < y < x < 1 \\ 0, & otherwise. \end{cases}$$

You may wonder where you would find such a distribution. In fact, if A_1 and A_2 are independent random variables uniformly distributed on $[0, 1]$, and you define $X = \max(A_1, A_2)$, $Y = \min(A_1, A_2)$, then X and Y will have exactly the joint distribution defined above.

- a. Draw a graph of the region for which X and Y have positive probability density.



- b. Derive the marginal probability density function of X , $f_X(x)$. Marginal Probability of X :

$$\int_{-\inf}^{\inf} f(x, y) dy$$

$$\int_0^x 2 dy = 2x$$

- c. Derive the unconditional expectation of X .

$$E(X) = \int_0^1 x * f_x(x) dx$$

$$E(x) = \int_0^1 2x^2 dx, E(x) = 2/3$$

- d. Derive the conditional probability density function of Y , conditional on X , $f_{Y|X}(y|x)$

$$f_{Y|X}(y|x) = \frac{f(x, y)}{f_x(x)}$$

Within the range $0 < y < x < 1$, $f(x, y)$ simplifies to 2. Therefore, for the range $0 < y < x < 1$:

$$f_{Y|X}(y|x) = \frac{2}{2x}, f_{Y|X}(y|x) = 1/x$$

- e. Derive the conditional expectation of Y , conditional on X , $E(Y|X)$.

$$E(Y|X = x) = \int_{-\inf}^i nfy * f(y|x) dy$$

$$E(Y|X = x) = \int_0^x y/x dy, E(Y|X = x) = x/2$$

- f. Derive $E(XY)$. Hint: if you take an expectation conditional on X , X is just a constant inside the expectation. This means that $E(XY|X) = XE(Y|X)$.

$$E(X * Y) = \int_0^1 \int_0^x 2 * x * y \, dy \, dx$$

$$E(X * Y) = \int_0^1 x^3 \, dx, E(X * Y) = 1/4$$

- g. Using the previous parts, derive $cov(X, Y)$

$$cov(X, Y) = E(X * Y) - E(X) * E(Y)$$

$$cov(X, Y) = 1/4 - (2/3) * (2/3) * (1/2)$$

$$cov(X, Y) = 1/36$$

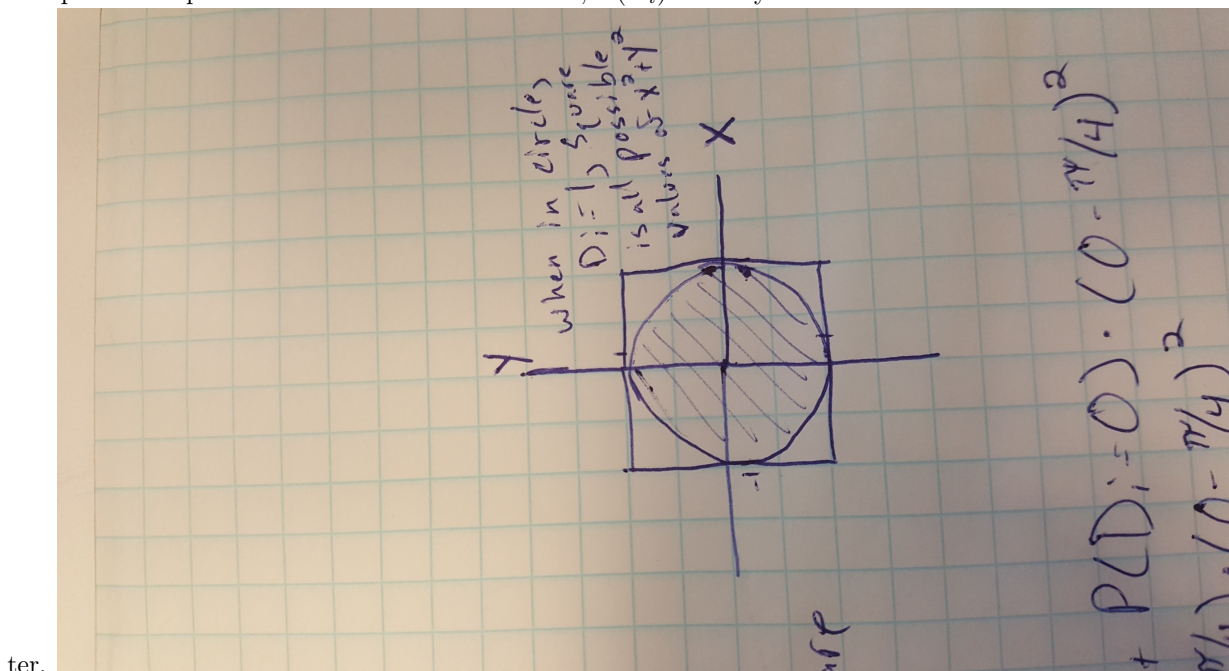
4. Circles, Random Samples, and the Central Limit Theorem

Let X_1, X_2, \dots, X_n and Y_1, Y_2, \dots, Y_n be independent random samples from a uniform distribution on $[-1, 1]$. Let D_i be a random variable that indicates if (X_i, Y_i) falls within the unit circle centered at the origin. We can define D_i as follows:

$$D_i = \begin{cases} 1, & X_i^2 + Y_i^2 < 1 \\ 0, & \text{otherwise} \end{cases}$$

Each D_i is a Bernoulli variable. Furthermore, all D_i are independent and identically distributed.

- a. Compute the expectation of each indicator variable, $E(D_i)$. Hint: your answer should involve a Greek letter.



The square in the above image represents the possible values of (x,y) - both are sampled from uniform distributions, so the probability of each value of (x,y) is equal. The circle concentric to the square of possible values represents all values of (x,y) for which D_i is 1. Thus, the probability of D_i equalling one is the ratio of the areas of the circle and the square. Because D_i is a binary variable with values 1 and 0, the probability of the variable having a value of one is equal to the expectation of the variable itself.

Area of square:

$$(1 - (-1))^2 = 4$$

Area of circle:

$$\begin{aligned} Area &= \pi * r^2, Area = \pi \\ E(D_i) &= A_{circle}/A_{square} = \pi/4 \end{aligned}$$

b. Compute the standard deviation of each D_i .

In this case, the variance of

$$D_i$$

is the sum of the products of each possible value of

$$D_i$$

minus the expectation of the variable squared and the corresponding probability of that possible value.

$$var(D_i) = P(D_i = 1) * (1 - E(D_i))^2 + P(D_i = 0) * (0 - E(D_i))^2$$

This derives to:

$$var(D_i) = (\pi/4) * (1 - (\pi/4))^2 + (1 - \pi/4) * (0 - (\pi/4))^2, var(D_i) = .1685$$

The standard deviation is just equal to the square root of the variance, thus:

$$\sigma D_i = \sqrt{var(D_i)}, \sigma D_i = .411$$

c. Let \bar{D} be the sample average of the D_i . Compute the standard error of \bar{D} . This should be a function of sample size n .

$$SE(\bar{D}) = \sigma/\sqrt{n}$$

d. Now let $n=100$. Using the Central Limit Theorem, compute the probability that \bar{D} is larger than $3/4$. Make sure you explain how the Central Limit Theorem helps you get your answer.

To convert a value from a normal distribution into a value from the Z distribution, we would use the following equation:

$$Z(x) = \frac{x - \bar{x}}{\sigma_x}$$

In our case, we are computing whether or not the sample mean will be greater than a certain value. Thus, we will modify this equation to be:

$$Z(\bar{D}) = \frac{\bar{D} - E(D)}{SE(\bar{D})}$$

Using the result from c, the standard error term is .0411 if $n=100$. Therefore:

$$Z(\bar{D}) = \frac{3/4 - \pi/4}{.0411}, Z(\bar{D}) = -.861$$

Looking this value up in a Z-table gives a value of .1949. We want to know the likelihood of the sample average being higher than $3/4$ - which means we want to the area under the z curve "to the right" of this value: thus, our probability is $1 - .1949$, or .8051.

- e. Now let $n = 100$. Use R to simulate a draw for X_1, X_2, \dots, X_n and Y_1, Y_2, \dots, Y_n . Calculate the resulting values for D_1, D_2, \dots, D_n . Create a plot to visualize your draws, with X on one axis and Y on the other. We suggest using a command like the following to assign a different color to each point, based on whether it falls inside the unit circle or outside it. Note that we pass $d + 1$ instead of d into the color argument because 0 corresponds to the color white.

```
x<-runif(100,min=-1,max=1)
y<-runif(100,min=-1,max=1)
D <- x^2 + y^2
plot(x,y, col=D+1, asp=1)
```

- f. What value do you get for the sample average, \bar{D} ? How does it compare to your answer for part a?

```
length(D[D<1])/100
```

Our experimentally observed value of .77 is similar to the theoretically derived value of $\pi/4$ ($\sim .785$)

- g. Now use R to replicate the previous experiment 10,000 times, generating a sample average of the D_i each time. Plot a histogram of the sample averages.

```
runs = 10000
Averages = array(runs)
for(i in 1:runs){
  x<-runif(100,min=-1,max=1)
  y<-runif(100,min=-1,max=1)
  D <- x^2 + y^2
  Averages[i] = length(D[D<1])/100
}
hist(Averages)
```

- h. Compute the standard deviation of your sample averages to see if it's close to the value you expect from part c.

```
sd(Averages)
```

The observed value of .041 is close to the theoretically derived value of the standard error, which was .0411

- i. Compute the fraction of your sample averages that are larger than $3/4$ to see if it's close to the value you expect from part d.

```
length(Averages[Averages > .75])
```

We anticipated that 80.5% of sample averages would be larger than $3/4$ and instead observe that about 77.4% of averages are large than $3/4$. The two numbers are relatively close, but the gap is a bit higher than I would have expected.