# Lab 2: Probability Theory

### W203: Statistics for Data Science

## 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 $P(T|H_k)$.

    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%?

## 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 $X$.
    b. Give a complete expression for the cumulative probability function of $X$.
    c. Compute $E(X)$.
    d. Compute $var(X)$.

## 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 $X$, $f_X(x)$.
    c. Derive the unconditional expectation of $X$.
    d. Derive the conditional probability density function of $Y$, conditional on $X$, $f_{Y|X}(y|x)$
    e. Derive the conditional expectation of $Y$, conditional on $X$, $E(Y|X)$.
    f. Derive $E(XY)$. Hint: if you take an expectation conditional on $X$, $X$ is just a constant inside the expectation. This means that $E(XY|X) = XE(Y|X)$.
    g. Using the previous parts, derive $cov(X, Y)$

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

Let $X_1, X_2, ..., X_n$ and $Y_1, Y_2, ..., Y_n$ 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:

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

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.

b. Compute the standard deviation of each $D_i$.

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 $n$.

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.

e. Now let $n = 100$. Use R to simulate a draw for $X_1, X_2, ..., X_n$ and $Y_1, Y_2, ..., Y_n$. Calculate the resulting values for $D_1, D_2, ...D_n$. 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.

```
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?

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.

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

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