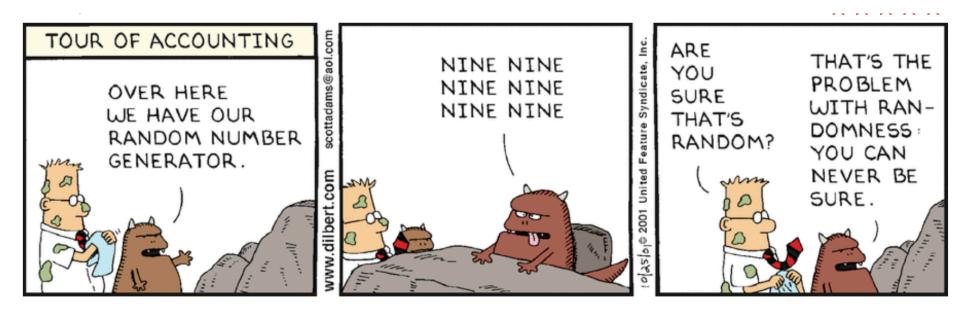
# Stat 88: Probability & Math. Stat in Data Science



#### Lecture 11: 2/9/2024

#### More examples of binomial and hypergeometric, and cdf

#### 3.5, 4.1

Shobhana Stoyanov

### Recap

- When we have two kinds of tickets in a box and we draw tickets at random from this box, each draw is called a *trial*
- We call the two kinds (binary) outcomes *Success*, and *Failure*
- Might be with replacement (like a coin toss) or without replacement (taking a simple random sample of Cal students and checking the number of people who are planning on going to cheer on the rugby team against BYU on Saturday.)
- Random variables (usually denoted by X, Y etc) are numbers that map the outcome space  $\Omega$  to real numbers, so they inherit a probability distribution.
- The probability distribution of a random variable X, is a description of the values taken by X, and the probabilities that X takes these values.
- The *probability mass function* of X, denoted by f(x), is a function that gives, for each value x taken by X, its chance P(X = x).

#### Recap

- Can think of these problems as *classifying* and *counting*.
- *n* independent trials each of which can result in one of two outcomes.
- We call these outcomes Success or Failure, and can represent the random experiment by drawing *n* tickets with replacement from a box with tickets marked 0 or 1, where the proportion of tickets marked 1 is equal to the probability of a success in a trial (*p*)
- Each draw can be represented by a **Bernoulli** random variable,  $X_k \sim Bernoulli(p)$
- If X is the number of successes in *n* trials, then X is the *sum of draws* from such a box as described above.
- We say that  $X \sim Bin(n, p)$  and  $P(X = k) = {n \choose k} \times p^k \times (1-p)^{n-k}, \ k = 0, 1, ... n$
- We might also draw *without* replacement, in which case, we say that X has the *hypergeometric*(N, G, n) distribution, and

$$P(X=g) = \frac{\binom{G}{g}\binom{N-G}{n-g}}{\binom{N}{n}}$$

• *N* is the number of tickets in the box, *G* is the number of successes possible, *n* is the number of draws (fixed beforehand).

Classify the following as binomial, hypergeometric, or neither

- 1.Number of heads in 12 tosses of a fair coin.
- 2.Number of tosses until we see two heads.
- 3.Number of queens in a five card hand
- 4.Number of Democrats in a simple random sample of 500 adult voters drawn from the SF Bay Area.

Go to pollev.com/shobhana to answer

#### Example

- A large supermarket chain in Florida occasionally selects employees to receive management training. A group of women there claimed that female employees were passed over for this training in favor of their male colleagues. The company denied this claim. (A similar complaint of gender bias was made about promotions and pay for the 1.6 million women who work or who have worked for Wal-Mart. The Supreme Court heard the case in 2011 and ruled in favor of Wal-Mart.)
- Suppose that the large employee pool of the Florida chain (more than a 1000 people) that can be tapped for management training is half male and half female.
  Since this program began, none of the 10 employees chosen have been female.
  What would be the probability of 0 out of 10 selections being female, if there truly was no gender bias?
- Method 1: pretend we are sampling with replacement, use Binomial dsn.

Are we really sampling with replacement?

## Problem solving techniques

- See if problem can be broken into smaller problems
- See which distribution applies to the situation
- Identify the parameters
- Use the addition and multiplication rules carefully

An advisor at a university provides guidance to **10** students. Each student has to meet with her **once a month** during the school year which consists of **nine** months.

Each month the advisor schedules one day of meetings. **Each** student has to sign up for one meeting that day. Students have the choice of meeting her in the **morning or in the afternoon**.

Assume that every month each student, independently of other students and other months, chooses to meet in the afternoon with probability 0.75.

What is the chance that she has **both** morning and afternoon meetings in **all** of the months except one?

## Advisors and their students

- Need to figure out a random variable. First fix **one** month, any month.
- Figure out the chance in that month, *all* the students choose the afternoon OR *all* the students choose the morning: this would mean that the meetings happen *only* in the morning OR *only* in the afternoon.
- We need the chance of the complement of this event.
- What is the random variable?

## Randomized Controlled Experiments

Two randomized controlled experiments are being run independently of each other. In each experiment, a simple random sample of **half** the participants will be assigned to the treatment group and the other half to control. Expt 1 has 100 participants of whom 20 are men. Expt 2 has 90 participants of whom 30 are men.

What is the chance that the treatment and control groups in Experiment 1 contain the same number of men?

#### Problems, continued

What is the chance that the treatment groups in the two experiments have the **same** number of men?

- Notice this is a bit tricky. There are many disjoint cases (each of the treatment groups has 1 man, or 2 men or 3 men etc. What is the max?
- We will have to split the chance into the chance of each of the cases and add them.

## Did the treatment have an effect?

- RCE with 100 participants, 60 in Treatment, 40 in Control
- T: 50 recover, out of 60 (83%), C: 30 recover out of 40 (75%)
- Suppose treatment had no effect, and these 80 just happened to recover. What is the chance they would have recovered no matter what and 50 were assigned to the treatment group by chance?

### Hypergeometric but don't know N

• A state has several million households, half of which have annual incomes over 50,000 dollars. In a simple random sample of 400 households taken from the state, what is the chance that more than 215 have incomes over 50,000 dollars?

How should we do this? n = 400, k = 215, G = N/2, N = ???

4.1: Back to random variables and their distributions

- X, f(x) = P(X = x)
- Consider X = number of H in 3 tosses, then  $X \sim Bin(3, \frac{1}{2})$
- We can also define a new function *F*, called the *cumulative distribution function*, that, for each real number x, tells us how much mass has been accumulated by the time X reaches x.

$$F(x) = P(X \le x) = \sum_{k \le x} {\binom{3}{k}} p^k (1-p)^{n-k}$$

x	0	1	2	3
f(x) = P(X=x)	1/8	3/8	3/8	1/8
$F(x) = P(X \le x)$				

## $F(x) \longrightarrow f(x)?$

• How to recover the pmf from the cdf? Draw the graph of F(x):

• What are the properties of F(x)? What is its domain? Range?

#### Exercise 4.5.2

• A random variable *W* has the distribution shown in the table below. Sketch a graph of the cdf of *W*.

W	-2	-1	0	1	3
f(w)	0.1	0.3	0.25	0.2	0.15
F(w)					