

Stat 88: Probability & Math. Stat in Data Science



<https://xkcd.com/1016/>

Lecture 13: 2/14/2024
Examples, waiting times
4.2
Shobhana Stoyanov

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? $n = 400$, $k = 215$, $G = N/2$, $N = ???$

4.2: Waiting times

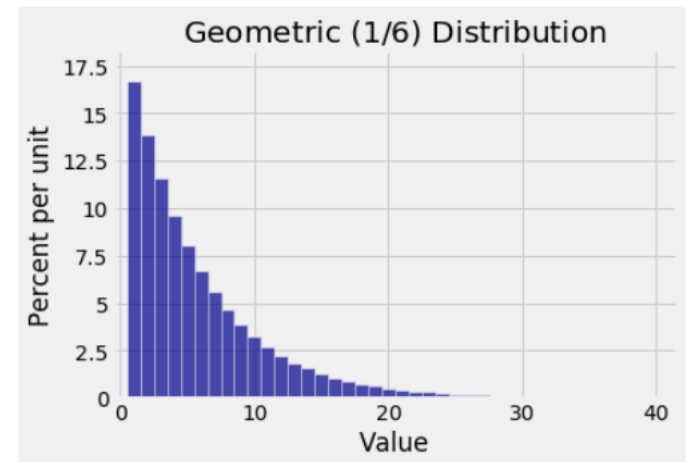
- Say Ali keeps playing roulette, and betting on red each time. The waiting time of a red win is the number of spins until they see a red (so the number of spins until and including the time the ball lands on a red pocket).

What is the probability that Ali will wait for 4 spins before their first win? (That is, the first time the ball lands in red is the 4th spin or trial)

- Say we have a sequence of **independent** trials (roulette spins, coin tosses, die rolls etc) each of which has outcomes of success or failure, and $P(S) = p$ on each trial.
- Let T_1 be the number of trials up to and including the first success. Then T_1 is the **waiting time until the first success**.
- What are the values T_1 takes? What is its pmf $f(x)$?

Geometric distribution

- Say T_1 has the **geometric distribution**, denoted $T_1 \sim \text{Geom}(p)$ on $\{1, 2, 3, \dots\}$
- $f(k) = P(T_1 = k) =$
- Check that it sums to 1. What is the cdf for this distribution? Can you think of an easy way to write down the cdf?



Waiting time until r^{th} success

- Say we roll a 8 sided die.
- What is the chance that the first time we roll an eight is on the 11th try?
 $= P(\text{FFFFFFFFFS}) =$
- What is the chance that it takes us 15 times until the 4th time we roll eight? (That is, the waiting time until the 4th time we roll an eight is 15)

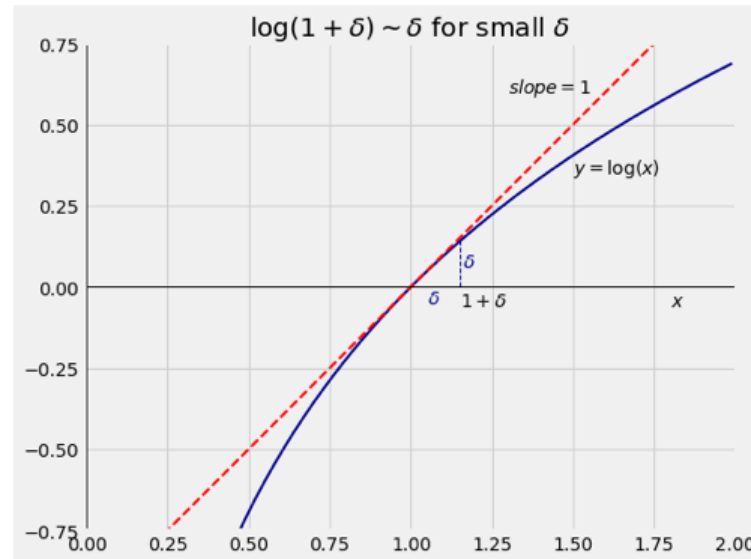
$$= P(\text{-----}S)$$

- What is the chance that we need **more** than 15 rolls to roll an eight 4 times?
- Notice that the **right-tail** probability of T_4 is a left hand (cdf) of the Binomial distribution for $(15, 1/8)$, and where $k=3$.

- In general, $P(T_r = k) =$

- And $P(T_r > k) =$

4.3 Exponential Approximations



Very useful approximation: $\log(1 + \delta) \approx \delta$, for δ close to 0

How to use this approximation

- Approximate the value of $x = \left(1 - \frac{3}{100}\right)^{100}$

- $x = \left(1 - \frac{2}{1000}\right)^{5000}$

- $x = (1 - p)^n$, for large n and small p

Example

- A book chapter $n = 100,000$ words and the chance that a word in the chapter has a typo (independently of all other words) is very small :
 $p = 1/1,000,000 = 10^{-6}$.

Give an approximation of the chance the chapter *doesn't* have a typo.
(Note that a typo is a *rare event*)

Bootstraps and probabilities

- Bootstrap sample: sample of size n drawn with replacement from original sample of n individuals
- Suppose one particular individual in the original sample is called Ali. What is the probability that Ali is chosen *at least once* in the bootstrap sample? (Use the complement.)