

# CSE 544, Spring 2021, Probability and Statistics for Data Science

## Assignment 1: Probability Theory review

Due: 2/23, 1:15pm, via Blackboard

(8 questions, 70 points total)

I/We understand and agree to the following:

- (a) Academic dishonesty will result in an 'F' grade and referral to the Academic Judiciary.
- (b) Late submission, beyond the 'due' date/time, will result in a score of 0 on this assignment.

(write down the name of all collaborating students on the line below)

---

### 1. Nerdy NBA

(Total 15 points)

In the 2020 NBA Western Conference semi-finals, the seemingly invincible Los Angeles Clippers (LAC) played the relatively inexperienced team of the Denver Nuggets (DEN) in a best-of-7 series where the first team to win 4 games wins the series. Assume that the outcome of each game is independent.

- (a) Assuming that either team has a win probability of 0.5, what is the probability that after the first 4 games, LAC would be up 3-1? Clearly show all your steps. (2 points)
- (b) LAC-DEN were, in fact, 3-1 at the end of 4 games, all but sealing the fate of the Denver Nuggets. Assuming the either team has a 0.5 probability of winning each game, draw the decision tree for the subsequent games starting from 3-1; note that if a team ends up winning 4 games total, subsequent games will not be held. (3 points)
- (c) Using the decision tree from (b) (so starting from 3-1), compute the probability of DEN winning the series 4-3. DEN in fact did win the series 4-3, ending the Championship hopes of LAC. (1 point)
- (d) Repeat part (b), but now with the assumption that the home team has a 0.75 probability of winning the game. Game 5, 6, and 7 were to be held in LAC, DEN, LAC, respectively. (due to the pandemic, all games were held in a bubble in Orlando, but ignore that for this question.) (3 points)
- (e) Repeat part (c), but using the decision tree of part (d) (1 point)
- (f) The frequentist interpretation of probability based on a large number  $N$  of repetitions of an experiment is  $P(A) \approx \frac{N_A}{N}$ , where  $N_A$  is the number of times  $A$  occurs and  $N$  is the total number of times the experiment is repeated. Similarly, the conditional probability  $P(B | A) \approx \frac{N_{BA}}{N_A}$ , where  $N_{BA}$  is the number of times  $B \cap A$  occurs and  $N_A$  is the number of times  $A$  occurs. Use simulations (coded in Python) to verify the results of part (a), (c) and (e). For instance, one can let  $A$  be the event that LAC is up 3-1 after 4 games and  $B$  be the event that DEN wins the matchup 4-3. Then we can simulate "series", a sequence of games until one of the teams wins 4 games,  $N$  times. Among those  $N$  repetitions of series, one can compute  $N_A$ , the number of times LAC is up 3-1 after 4 games,  $N_{BA}$ , number of times LAC was up 3-1 after 4 games and DEN eventually won 4-3. Finally, we can approximate  $P(A) \approx \frac{N_A}{N}$  and  $P(B|A) \approx \frac{N_{BA}}{N_A}$ . To verify part (e), assume that games 1 and 2 were played in LAC and games 3 and 4 were in DEN. Try  $N = 10^n$  for  $n = 3, 4, 5, 6, 7$ . What do you observe as  $N$  increases?

Hint: In Python, `numpy.random.binomial(1, p)` can simulate a Bernoulli trial with probability  $p$ .

For this programming assignment, you should **submit a Python script** named `nba.py` as part of your zipped submission on Blackboard. The script should have a variable named `N`, the number of times the experiment is repeated, at the very beginning of the program so that TAs can try out different values for `N`. The program should print the results of part (a), (c) and (e) as follows:

For `N = ...`, the simulated value for part (a) is ...

For `N = ...`, the simulated value for part (c) is ...

For `N = ...`, the simulated value for part (e) is ...

You should **also report the answers in your digital assignment submission.**

(5 points)

## 2. Free yourself

(Total 10 points)

In the near future, you realize that you have spent far too much money on buying and hoarding phones and decide to rid yourself of all your hoarded iPhones. Turns out you have  $n$  iPhones, with each iPhone belonging to a unique generation from iPhone 1 to iPhone  $n$ . So, to cleanse your digital life, you play a risky game. In step 1, you randomly pick an iPhone from this pile; if the selected iPhone is iPhone 1, you keep it, else you discard it. In step 2, you again randomly pick an iPhone from the remaining  $(n-1)$  iPhones and if the selected iPhone is iPhone 2, you keep it, else you discard it. You repeat this immensely satisfying exercise  $n$  times. We would like to find out, at the end of this exercise, what is the probability that you have at least one undiscarded iPhone? Solve this problem using the principle of inclusion-exclusion (PIE). For  $n$  events  $E_1, E_2, \dots, E_n$ , PIE says

$$\Pr\left(\bigcup_{i=1}^n E_i\right) = \sum_i \Pr(E_i) - \sum_{i < j} \Pr(E_i \cap E_j) + \sum_{i < j < k} \Pr(E_i \cap E_j \cap E_k) - \dots + (-1)^{n+1} \Pr(E_1 \cap \dots \cap E_n).$$

Choose the events  $E_1, E_2, \dots, E_n$  carefully so that you can obtain the required probability.

### 3. The One Ring?

(Total 10 points)

Bilbo Baggins of the Shire has a ring. It is known that there are only 10,000 rings in Middle Earth. Gandalf the Wizard, however, fears that Bilbo's ring may, in fact, be the One Ring!

- (a) If the ring is the One Ring, there is a 95% chance that the owner will have an above-average lifespan. If the ring is not the One Ring, there is a 75% chance that the owner will not have an above-average lifespan. What is the probability that, given Bilbo is pushing 111 years (above-average for Hobbits), his ring is, in fact, the One Ring? (3 points)
- (b) To be absolutely sure, Gandalf administers another test and throws the ring into a fireplace. If it is the One Ring, writing will appear on it with probability 0.9; if it is not the One Ring, writing may still appear on it with probability 0.05. Given that writing appears on it, and that Bilbo has an above-average lifespan, what is the probability that this is the One Ring? Assume that the tests are independent conditioned on the ring being the One Ring, and the tests are independent conditioned on the ring not being the One Ring. Do not assume that the tests are independent. (7 points)

**4. Alternative expression for expectation****(Total 5 points)**

Let  $X$  be a non-negative, integer-valued RV. Prove that:

$$E[X] = \sum_{x=0}^{\infty} \Pr[X > x]$$

(Hint: One approach is to consider double summations and carefully switch the summations)

**5. Practice with discrete distributions****(Total 10 points)**

- (a) For the Indicator RV introduced in class (with event  $E$ ), compute  $E[I_E]$  in terms of  $\Pr(E)$ . (1 point)
- (b) For part (a), compute  $\text{Var}(I_E)$  in terms of  $\Pr(E)$ . (2 points)
- (c) Let  $X \sim \text{Geometric}(p)$ , with  $p < 1$ . Using the definition of a Geometric RV as in class, compute  $E[X]$ .  
You may assume that  $\sum_{i=0}^{\infty} x^i = \frac{1}{1-x}$ , for  $x < 1$ . If you use any other result, prove it. (4 points)
- (d) For part (c), compute  $\text{Var}(X)$ . (3 points)

## 6. Poisson distribution

(Total 5 points)

The Poisson distribution,  $X \sim \text{Poisson}(\lambda)$ , is a discrete distribution with p.m.f. given by:

$$p_X(i) = \frac{e^{-\lambda} \lambda^i}{i!}, i \geq 0$$

- (a) Ensure that the p.m.f. adds up to 1 (2 points)  
(Hint: You will need to use the infinite series expansion of an Exponential)
- (b) Find  $E[X]$  (3 points)

## 7. Pareto distribution

(Total 10 points)

The Pareto distribution,  $X \sim \text{Pareto}(\alpha)$ ,  $1 < \alpha < 2$ , is a continuous distribution with p.d.f. given by:

$$f_X(x) = \alpha x^{-\alpha-1}, x \geq 1$$

- (a) Ensure that the p.d.f. integrates to 1 (2 points)
- (b) Find  $E[X]$  (3 points)
- (c) Find  $\text{Var}[X]$  (5 points)



**8. One distribution to rule them all****(Total 5 points)**

Let  $F$  be a CDF such that  $F$  is strictly increasing on the support of the distribution. Since  $F$  is continuous and strictly increasing, its inverse,  $F^{-1}$  exists.

- (a) Let  $U \sim \text{Uniform}(0, 1)$  and  $X$  be a RV such that  $X = F^{-1}(U)$ . Prove that the CDF of  $X$  is  $F$ . (3 points)
- (b) Let  $Y$  be a random variable with CDF  $F$ . Prove that  $F(Y) \sim \text{Uniform}(0, 1)$ . (2 points)