

# M50 Homework 1

Alex Craig

## Exercise 1.

### Part A.

Suppose that  $Y \sim \text{Bernoulli}(q)$  and let  $Z = \frac{1}{1+Y} + Y$ . What is the probability function of  $Z$ ?

### Solution

We can start with the probability mass function of  $Y$ , defined as the probability mass function of any Bernoulli random variable with parameter  $q$ :

$$f_Y(y) = \Pr(Y = y) = \begin{cases} q & \text{if } y = 1 \\ 1 - q & \text{if } y = 0 \\ 0 & \text{o.w.} \end{cases}$$

We can then use this to find the probability mass function of  $Z$  by defining the two cases of when  $Y = 1$  and  $Y = 0$ . When  $Y = 1$ , we now that  $Z = \frac{1}{1+1} + 1 = \frac{3}{2}$ . When  $Y = 0$ , we know that  $Z = \frac{1}{1+0} + 0 = 1$ . Thus, we can define the probability mass function of  $Z$  as:

$$f_Z(z) = \Pr(Z = z) = \begin{cases} q & \text{if } z = \frac{3}{2} \\ 1 - q & \text{if } z = 1 \\ 0 & \text{o.w.} \end{cases}$$

### Part B.

Suppose a coin is flipped. If the coin is heads, we write down 0. If the coin is tails, we roll a dice and write down the number. Write down the probability distribution for  $Y$ , the number that we write down.

### Solution

Let us start by defining random variables  $X$  and  $Z$ . Let  $X$  be the random variable that represents the result of the coin flip (let 0 denote heads, and 1 denote tails), and let  $Z$  be the random variable that represents the dice roll. We know that  $X \sim \text{bernoulli}(\frac{1}{2})$ , and we know that  $Z$  is a discrete uniform random variable that can take on values in the set  $\{1, 2, 3, 4, 5, 6\}$ , each with probability  $\frac{1}{6}$  given that we flip tails on the initial coin flip. We can therefore define the probability mass function of  $X$  and  $Z$  as:

$$f_X(x) = \Pr(X = x) = \begin{cases} \frac{1}{2} & \text{if } x = 0 \\ \frac{1}{2} & \text{if } x = 1 \\ 0 & \text{o.w.} \end{cases}$$

$$f_Z(z) = \Pr(Z = z) = \begin{cases} \frac{1}{6} & \text{if } z \in \{1, 2, 3, 4, 5, 6\} \wedge X = 1 \\ 0 & \text{o.w.} \end{cases}$$

Let us now define our sample space  $S$  for the outcomes of  $Y$ :

$$S = \{0, 1, 2, 3, 4, 5, 6\}$$

$Y = 0$  in the case where we flip heads on the initial coin flip. Therefore:

$$f_Y(0) = \Pr(Y = 0) = \Pr(X = 0) = \frac{1}{2}$$

For all other values of  $Y$  in the sample space, we know that we must have flipped tails on the initial coin flip, and then rolled that value on the dice. We can therefore define the probability of  $Y$  being these values as the probability of the intersection of flipping tails on the initial coin flip and flipping that value on the dice roll:

$$\begin{aligned}
 f_Y(z) &= Pr(Y = z \cap X = 1) = Pr(Z = z \cap X = 1) = Pr(Z = z \mid X = 1) \times Pr(X = 1) \\
 &= \frac{1}{6} \times \frac{1}{2} = \frac{1}{12} : z \in \{1, 2, 3, 4, 5, 6\}
 \end{aligned}$$

We can therefore define the probability mass function of  $Y$  as:

$$f_Y(y) = Pr(Y = y) = \begin{cases} \frac{1}{2} & \text{if } y = 0 \\ \frac{1}{12} & \text{if } y \in \{1, 2, 3, 4, 5, 6\} \\ 0 & \text{o.w.} \end{cases}$$

### Part C.

For the previous problem, conditioned on the dice rolling a 4, what is the probability we write down 0? Conditioned on the coin being tails, what is the probability the dice rolls a 3?

### Solution

Given that the dice rolled a 4, we know that  $Y = 4$  because in all the cases where we roll a dice, we write down the value that we roll. Therefore, the probability that we write down 0 is 0.

Given that the coin is tails, we know that  $X = 1$ . We may still define the probability that the dice rolls a 3 as the probability of the intersection of the dice rolling a 3 and the coin being tails, but in this case  $Pr(X = 1) = 1$ :

$$Pr(Z = 3 \cap X = 1) = Pr(Z = 3 \mid X = 1) \times Pr(X = 1) = \frac{1}{6} \times 1 = \frac{1}{6}$$

### Part D.

Consider the geometric distribution discussed in lecture. What are 3 examples of variables in the real world for which this might be a good model and what are some limitations of these models.

### Solution

#### Examples:

1. Shooting basketball free throws until you make one.
2. Playing lottery scratch tickets until you win.
3. Flying on an airplane until you get into a fatal crash.

#### Limitations:

1. The probability of success must be constant for each trial, so changes in the probability of success over time are not accounted for.
2. The trials must be independent, so the outcome of one trial cannot affect the outcome of another trial.

### Exercise 2.

Consider the following code:

```

for i in range(5):
    for j in range(i + 1):
        print(i, end = " ")
    print("")

```

```

0
11
222
3333
44444

```

Modify the code so that it prints out:

```
0
01
012
0123
01234
012345
```

### Solution

```
for i in range(6):
    for j in range(i + 1):
        print(j, end = "")
    print("")
```

```
0
01
012
0123
01234
012345
```

### Exercise 3.

(Washington post data): Below I load some data on homicide victims in US from the washington post. Don't worry about how I process it, all you need to work with is the DataFrame "data" on the very last line.

```
import pandas as pd

data = pd.read_csv("https://raw.githubusercontent.com/washingtonpost/data-homicides/master/homicide-data.csv",

data["victim_age"] = pd.to_numeric(data["victim_age"], errors="coerce")
```

### Part A.

For each age  $a = 1, \dots, 100$  determine the number of victims  $n(a)$  with an age  $< a$ . You can ignore the effects of those entries with missing ages. Make a plot of  $n(a)$  vs.  $a$ .

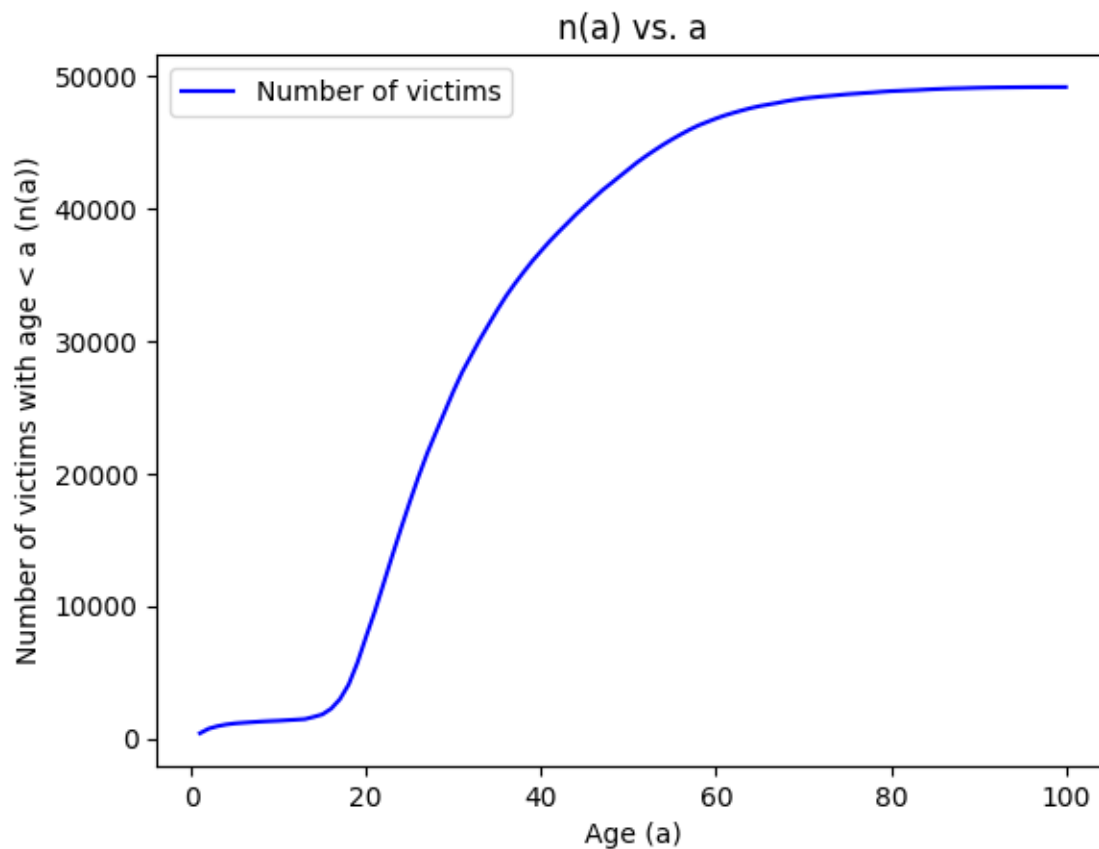
### Solution

```
import matplotlib.pyplot as plt

# Function to calculate number of victims for a given age range and data
def calculate_victims(df):
    n_values = []
    for a in range(1, 101):
        n = df[df["victim_age"] < a].shape[0]
        n_values.append(n)
    return n_values

# Create an array to store number of victims for each age a
n_values = calculate_victims(data)

# Plotting the results
plt.plot(range(1, 101), n_values, label='Number of victims', color='blue')
plt.xlabel('Age (a)')
plt.ylabel('Number of victims with age < a (n(a))')
plt.title('n(a) vs. a')
plt.legend()
plt.grid(True)
plt.show()
```



### Part B.

What do you notice about the plot. Does it have the expected behavior?

### Solution

There are very few homicide victims under the age of 18, but the number of homicide victims increases most rapidly between the ages of 20 and 40. This makes sense because this is about the age range where people are most likely to be involved in gang violence or other violent crime, which is a common cause of homicide.

### Part C.

Break the data up into white and non-white victims. Then, for each group make the plot from part (a). Comment on what you find.

### Solution

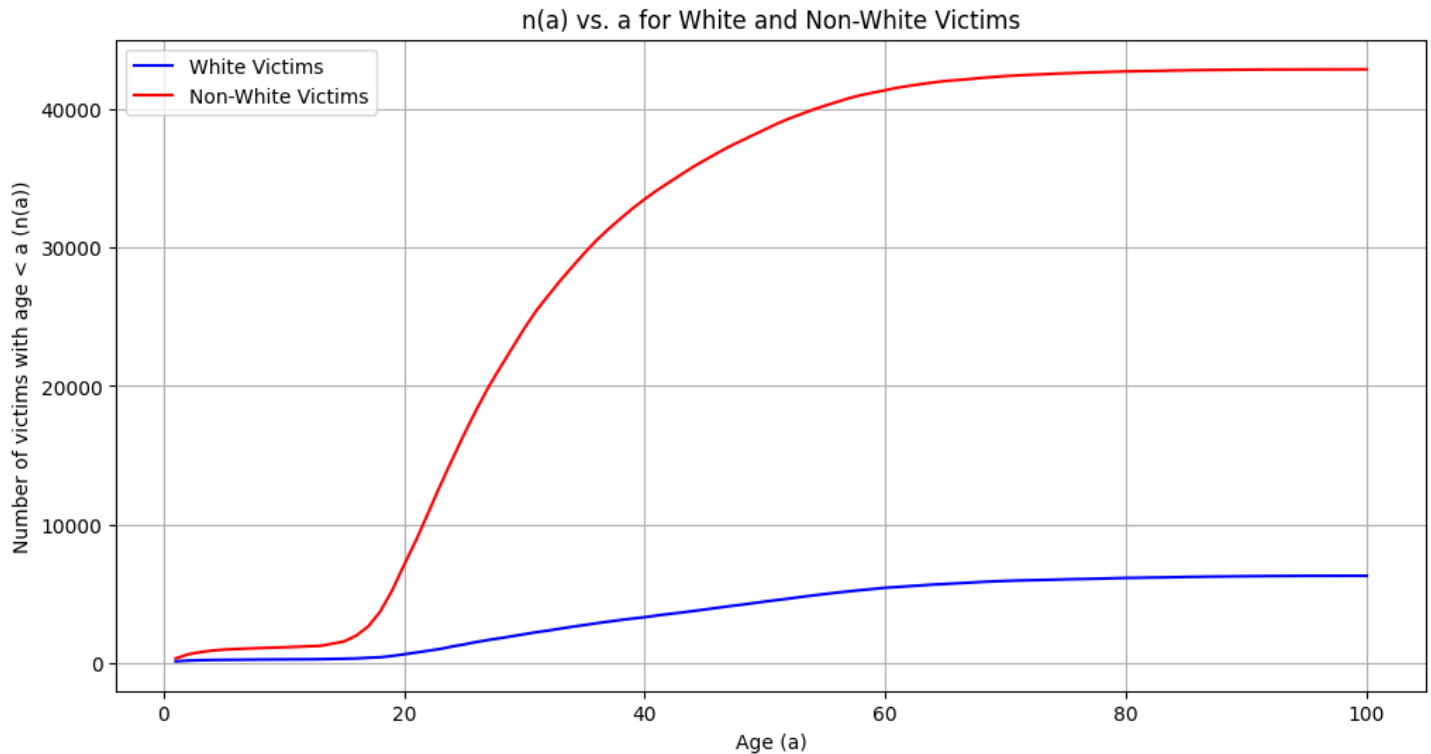
```
# Filter data for white victims and non-white victims
white_victims = data[data["victim_race"] == "White"]
non_white_victims = data[data["victim_race"] != "White"]

# Calculate for white and non-white victims
white_n_values = calculate_victims(white_victims)
non_white_n_values = calculate_victims(non_white_victims)

# Plotting the results
plt.figure(figsize=(12,6))

plt.plot(range(1, 101), white_n_values, label='White Victims', color='blue')
plt.plot(range(1, 101), non_white_n_values, label='Non-White Victims', color='red')
plt.xlabel('Age (a)')
plt.ylabel('Number of victims with age <= a (n(a))')
plt.title('n(a) vs. a for White and Non-White Victims')
```

```
plt.grid(True)
plt.legend()
plt.show()
```



The results between White and Non-White victims are quite different. The Non-White victims look very similar to the graph we previously saw, where there is a steep increase in homicide victims between the ages of 20 and 40, and a less steep increase between 40 and 60. The White victims, however, see a fairly consistent increase in homicide victims between the ages of 20 and 60, and then very few more homicide victims after that.

#### Exercise 4.

(Getting a sequence of wins): Let  $J$  denote a random variable representing the number of times a fair coin is flipped before two heads appear in a row. As we saw in class, the following code generates simulations of  $J$ :

```
import numpy as np

def flip_until_two():
    num_heads = 0
    total_flips = 0
    while num_heads < 2:
        y = np.random.choice([0,1])
        if y == 0:
            num_heads = 0
        else:
            num_heads = num_heads + 1
            total_flips = total_flips + 1
    return total_flips
```

#### Part A.

By changing the code above, write a function `roll_until_n()` that rolls a dice until we get  $n$  ones in a row. You should change the variable names accordingly. We will call this random variable `Roll_Count(n)`.

#### Solution

```
def roll_until_n(n: int):
    num_ones = 0
```

```

total_rolls = 0
while num_ones < n:
    y = np.random.choice([1,2,3,4,5,6])
    if y != 1:
        num_ones = 0
    else:
        num_ones = num_ones + 1
    total_rolls = total_rolls + 1
return total_rolls

```

```

Roll_Count = roll_until_n(2)
print(Roll_Count)

```

78

## Part B.

Make a DataFrame where each column represents a value of  $n$  from 1 to 6 and each row is a simulation from the model  $\text{Roll\_Count}(n)$ . There should be 100 rows.

## Solution

```

# Simulation
num_simulations = 100
n_values = [1,2,3,4,5,6] # for n from 1 to 6

# Create an empty DataFrame with columns from 1 to 6
df = pd.DataFrame(columns=n_values)

# For 100 simulations
for i in range(num_simulations):
    # Calculate the number of rolls for each n
    row_data = [roll_until_n(n) for n in n_values]
    # Add the row to the DataFrame
    df.loc[i] = row_data

print(df)

```

	1	2	3	4	5	6
0	25	15	42	1364	20368	74791
1	10	56	47	538	2711	88662
2	19	6	262	1010	5477	3444
3	7	27	493	2949	5834	47639
4	4	42	462	315	8139	27870
..	..	...	...	...	...	...
95	3	12	171	573	16493	5616
96	8	119	173	940	11946	58392
97	1	67	19	386	14603	1559
98	9	17	645	1439	307	205269
99	7	196	547	195	15179	51001

[100 rows x 6 columns]

## Part C.

Make a plot of the maximum and minimum values of  $J$  as a function of  $n$  on the same plot. You might notice one of these increases much faster than the other.

## Solution

```

# Extracting max and min values for each n
max_vals = df.max()

```

```
min_vals = df.min()
```

```
# Plotting
```

```
plt.figure(figsize=(10,6))
```

```
plt.plot(n_values, max_vals, '-o', label='Maximum')
```

```
plt.plot(n_values, min_vals, '-o', label='Minimum')
```

```
plt.xlabel('n')
```

```
plt.ylabel('Roll_Count')
```

```
plt.title('Maximum and Minimum values of Roll_Count as a function of n')
```

```
plt.legend()
```

```
plt.grid(True)
```

```
plt.show()
```

