

MIT Introduction to Statistics 18.05 Problem Set 3

John Hancock

March 18, 2017

Contents

1	References and License	1
2	Independence	2
2.1	Pairwise and mutual independence	2
2.2	Venn diagram	2
2.3	How many kids	3
3	R simulation	3
3.1	Mean and standard deviation	4
3.2	Histogram of the data	4
3.3	Conclusions on the data	4
3.4	Effectiveness of the treatment	5
4	Dice	5
4.1	Standard deviation of X , Y , and Z	5
4.2	Graph pmf and cdf of Z	6
4.3	Game	7

1 References and License

We are answering questions in the material from MIT OpenCourseWare course 18.05, Introduction to Probability and Statistics.

In this document we are answering questions Orloff and Bloom ask in [2].

Please see the references section for detailed citation information.

The material for the course is licensed under the terms at <http://ocw.mit.edu/terms>.

We use documentation in [11], [7], [8] to write L^AT_EXsource code for this document.

2 Independence

In this section we answer a problem in [2] that involves rolling two six sided dice.

2.1 Pairwise and mutual independence

We define two events, A , and B to be pairwise independent if $P(A \cap B) = P(A)P(B)$.

For this problem Orloff and Bloom give us the definition of mutual independence for three events, A , B , and C . A , B , and C are mutually independent if

$$P(A \cap B \cap C) = P(A)P(B)P(C) \quad (1)$$

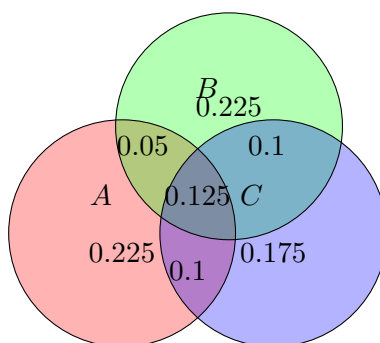
In this section, Orloff and Bloom give the following definitions for events A , B , and C :

- A is the event that we roll an odd number with the first die.
- B is the event that we roll an odd number with the second die.
- C is the event that the sum of the numbers we roll is odd.

A , B , and C are not mutually independent. Whatever the A , B , and C probabilities of A , B , and C are individually, the probability of $P(A \cap B \cap C)$ is 0 since the sum of two odd numbers is always an even number.

2.2 Venn diagram

Orloff and Bloom give the following Venn diagram:



And ask us whether or not the events in the Venn diagram above are mutually independent.

These events are not mutually independent because

$$P(A)P(B)P(C) = 0.225 \times 0.225 \times 0.175 = 0.008859375. \quad (2)$$

However, in the Venn diagram above, Orloff and Bloom give us that $P(A \cap B \cap C) = 0.125$

Therefore the events are not mutually independent.

2.3 How many kids

For this question we use the same assumptions about the probability of the gender that a child is born with that Orloff and Bloom use in example 9 of [5].

We define the following events:

- A is the event that the children in a family are both boys and girls.
- B is the event that at most one of the children is a girl.
- $C_{i,b}$ is the event that child number i is a boy.
- $C_{i,g}$ is the event that child number i is a girl.

Our goal is to construct a sample space such that A and B are independent. The definition of independent events is in [4].

We rely on the same assumption that Orloff and Bloom make in [5] regarding the probability of the genders of sequences of children.

Therefore we assume $P(C_{i,b}) = 0.5$, and $P(C_{i,g}) = 0.5$, independent of the event that any other child is a boy or a girl.

We write the following table to discover the number of children where A , and B will meet the definition of independent events.

We fill in one cell in the table below for each possible sequence of three children in the family being boys or girls.

$C_{1,b}C_{2,b}C_{3,b}$	$C_{1,b}C_{2,b}C_{3,g}$	$C_{1,b}C_{2,g}C_{3,b}$	$C_{1,b}C_{2,g}C_{3,g}$
$C_{1,g}C_{2,b}C_{3,b}$	$C_{1,g}C_{2,b}C_{3,g}$	$C_{1,g}C_{2,g}C_{3,b}$	$C_{1,g}C_{2,g}C_{3,g}$

In the table above there are 6 sequences that are in A , so $P(A) = \frac{6}{8}$.

Also, there are 4 sequences in B , so $P(B) = \frac{4}{8}$.

Moreover, there are 3 sequences where there is at most one girl, and the children are both boys and girls. Therefore $P(A \cap B) = \frac{3}{8}$.

A and B are independent since

$$P(A)P(B) = \left(\frac{6}{8}\right)\left(\frac{4}{8}\right) = \frac{24}{64} = \frac{3}{8}. \quad (3)$$

Therefore $P(A)P(B) = P(A \cap B)$, so A and B must be independent events.

We made these calculations assuming that there are 3 children, therefore the number of children we require in order for A , and B to be independent events is 3.

3 R simulation

Note: we needed to refer to [10] to recall the R function for standard deviation.

3.1 Mean and standard deviation

We use R functions to compute the mean and standard deviation of the data that Orloff and Bloom give us for this problem:

```
> mean(x)
[1] 2.554528
> sd(x)
[1] 2.07408
```

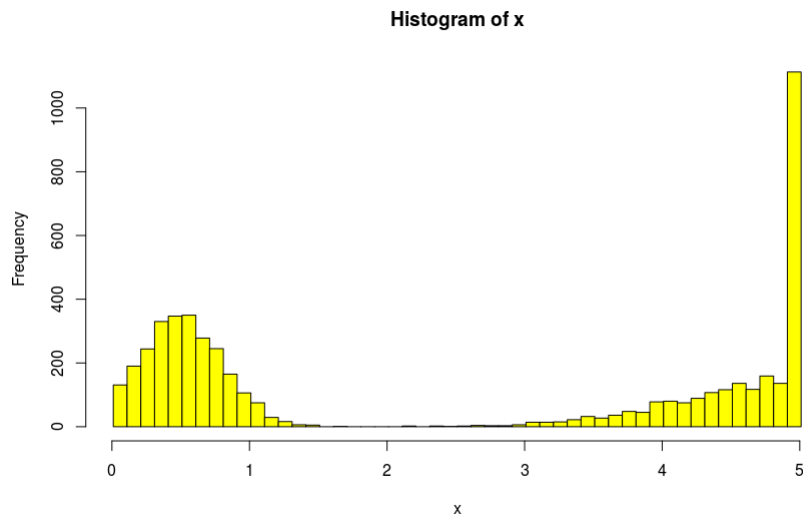
3.2 Histogram of the data

We rely on the examples in [3] in order to write the R code to produce the histogram Orloff and Bloom require.

The R code is

```
binwidth = .1
bins = seq(min(x), max(x)+binwidth, binwidth)
hist(x, breaks=bins, col='yellow', freq=TRUE)
```

The histogram this code produces is:



3.3 Conclusions on the data

The mean value of the data that we use R to compute above tells us that the mean number of years a patient lives after receiving the treatment is approximately 2.555 years. However, the standard deviation we compute using R tells us that the spread of values in the data is approximately $\frac{4}{5}$ of the mean number of years that a patient is expected to live after receiving the treatment. That means that we have many values in the data that are

larger and smaller than the mean value of the data. The histogram we plot bears this out; we see taller bins on the left and right hand sides of the histogram.

3.4 Effectiveness of the treatment

The tallest bin is on the right hand side of the histogram. To us, this implies that for some patients, the treatment is effective. We would recommend research to determine what thing or things the patients whose data lies in the rightmost bins have in common. Furthermore we would recommend that we conduct experiments to determine which of the things these patients have in common is or are the thing(s) that implies or imply the treatment is effective.

4 Dice

In this section we will deal with problems that Orloff and Bloom ask about the random variable X , that is equal to the value we roll with a fair 4-sided die, the random variable Y , that is equal to the value we roll with a fair 6 sided die, and the random variable Z , that is equal to the average of X and Y .

4.1 Standard deviation of X , Y , and Z

We use the definition of variance and standard deviation in [6] to calculate the standard deviations $\sigma(X)$, $\sigma(Y)$.

We use the exact same method to calculate the variance of a discrete random variable many times. For details on how to do the calculation see [1]. We calculate the variance of X , and Y , then take the square root of the variance to obtain the standard deviation.

Here are the results:

$$\sigma_X \approx 1.118. \quad (4)$$

$$\sigma_Y \approx 1.708. \quad (5)$$

In order to calculate the variance of Z , we can use properties of variance that Orloff and Bloom show in [6].

X and Y are independent random variables; the value we roll with a fair four sided die has no effect on the value we roll with a fair six sided die. Therefore the equation

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y). \quad (6)$$

We will also use the property of variance Orloff and Bloom show in [6] that

$$\text{Var}(X) = a^2 \text{Var}(X). \quad (7)$$

Therefore, in the scope of this problem,

$$\text{Var}\left(\frac{X}{2} + \frac{Y}{2}\right) = \text{Var}\left(\frac{X}{2} + \frac{Y}{2}\right) = \frac{1}{4} (\text{Var}(X) + \text{Var}(Y)). \quad (8)$$

We use the values we calculated for $\text{Var}(X)$, and $\text{Var}(Y)$ above,

$$\text{Var}(Z) \approx \frac{1}{4} (1.118 + 1.708) \approx 0.707 \quad (9)$$

Standard deviation is the square root of variance, so

$$\sigma_z \approx \sqrt{\frac{1}{4} (1.118^2 + 1.708^2)} \approx 1.021 \quad (10)$$

4.2 Graph pmf and cdf of Z

We use R to simulate rolling the dice. We find documentation in [9] helpful in writing this source code. Here is a listing of the R source code:

```
y=sample(c(1:6), replace = TRUE, 1000000)
x=sample(c(1:4), replace = TRUE, 1000000)
z=(x+y)/2
zTable = table(z)
zTable/1000000
```

The output of this program is:

	z								
	1	1.5	2	2.5	3	3.5			
4	4.5	5							
	0.041375	0.083538	0.125244	0.167014	0.166669	0.167029	0.124364	0.083153	0

We look at the numerical inverses of the values in the second row of the table above to get a clue about the probability mass of each possible value of Z . Here again we utilize R to compute the inverses:

```
freqs=c(0.041375, 0.083538, 0.125244, 0.167014, 0.166669, 0.167029, 0.124364, 0.083153, 0)
1/freqs
```

We execute the code above to get the result

```
24.169184 11.970600 7.984414 5.987522 5.999916 5.986984
8.040912 12.026024 24.030374
```

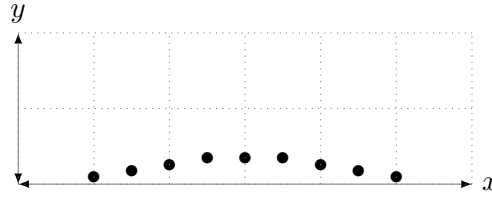
Now it becomes clear that we can approximate the frequencies of the values of Z as fractions with 24 in the denominator. We summarize this approximation in the table below, and make a guess as to the value of the pmf of Z :

Z	1	1.5	2	2.5	3	3.5	4	4.5	5
pmf $p(z)$?	$\frac{1}{24}$	$\frac{2}{24}$	$\frac{3}{24}$	$\frac{4}{24}$	$\frac{4}{24}$	$\frac{3}{24}$	$\frac{2}{24}$	$\frac{1}{24}$	$\frac{1}{24}$

The data from the simulation inspires us to make the following reasoning about the pmf of Z : There are 6×4 possible pairs of values we can roll with the two dice. To compute the value of Z we add the pair of values, and divide by 2. Some of the values of Z will occur more frequently because there are more values of X and Y divided by 2 that equal a particular value of Z . $1.5 = \frac{1+2}{2}$, and $1.5 = \frac{2+1}{2}$. The order of terms in the sum is important because the first number is the value we roll with the 4-sided die, and the second number is the value we roll with the 6-sided die. $2 = \frac{1+3}{2} = \frac{2+2}{2} = \frac{3+1}{2}$. $2.5 = \frac{1+4}{2} = \frac{2+3}{2} = \frac{3+2}{2} = \frac{4+1}{2}$. $3 = \frac{1+5}{2} = \frac{2+4}{2} = \frac{3+3}{2} = \frac{4+2}{2}$. $3.5 = \frac{1+6}{2} = \frac{2+5}{2} = \frac{3+4}{2} = \frac{4+3}{2}$. $4 = \frac{2+6}{2} = \frac{3+5}{2} = \frac{4+4}{2}$. $4.5 = \frac{3+6}{2} = \frac{4+5}{2}$. $5 = \frac{4+6}{2}$. Hence the pmf for a particular value of Z is the number of ways of summing a value between 1 and 4, and a value between 1 and 6, and dividing by 2 to equal Z .

Therefore the tentative pmf we write in the table above is the pmf for Z .

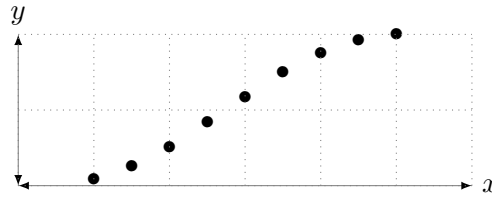
Here is a plot of the pmf:



Orloff and Bloom ask for a plot of the of the cumulative distribution function (CDF) of Z .

Orloff and Bloom define and explain how to calculate the CDF in [5].

Here is a plot of the cdf of Z :



4.3 Game

In this section we answer the question Orloff and Bloom pose regarding a game where we win $2X$ dollars if $X > Y$ and we lose 1 dollar otherwise. They ask us after playing the game 60 times, what is our expected total gain or loss.

In order to answer, we need to know the probability $P(X > Y)$.

We wrote the following R simulation to get an idea of what $P(X > Y)$ is.

```

winCount=0
loseCount=0
for ( i in 1:100000){
  x=sample(c(1:4), replace=TRUE, 1)
  y=sample(c(1:6), replace=TRUE, 1)
  if (x > y){
    winCount = winCount + 1
  } else{
    loseCount = loseCount + 1
  }
}
print(winCount/1000);

```

The output of this simulation is close to 0.25 every time we run it, so we should be able to find a reason why $P(X > Y) = 0.25$.

These are the possible combinations of values we can roll with a 6-sided die, and a 4-sided die. We denote a pair of values where we win with a W , and a pair of values where we lose with an L .

The first value in each pair is the value we roll with the 4-sided die, and the second value in each pair is the value we roll with the 6-sided die.

$1, 1L, 2, 1W, 3, 1W, 4, 1W$
 $1, 2L, 2, 2L, 3, 2W, 4, 2W$
 $1, 3L, 2, 3L, 3, 3L, 4, 3W$
 $1, 4L, 2, 4L, 3, 4L, 4, 4L$
 $1, 5L, 2, 5L, 3, 5L, 4, 5L$
 $1, 6L, 2, 6L, 3, 6L, 4, 6L$

Now, to calculate the probability that $X > Y$, we count the number of W 's above, and divide by the sum of the number of W 's and L 's.

There are 6 W 's and 24 L 's, so the probability that $X > Y$ is $\frac{6}{24} = \frac{1}{4}$.

We calculate the expected number of dollars we win or lose playing one round of the game, then multiply that number by 60 to compute our expected profit or loss playing 60 rounds of the game.

We define W as the event that we win. $P(W)$ is the probability of winning one round of the game, and $1 - P(W)$ is the probability of losing one round of the game. Then by what we show immediately above, $P(W) = \frac{1}{4}$.

We introduce two new discrete random variables:

- X' : the number of dollars we win playing one round of the game.
- Y' : the number of dollars we lose playing one round of the game.

We define a third random variable W' to be $X' - Y'$. W' is the number of dollars we win or lose playing one round of the game.

Now we can use the definition of expected value [5] to obtain the expected value of our profit or loss playing the game.

$$E(W') = P(W) E(X') - (1 - P(W)) E(Y') \quad (11)$$

Now we consider how to calculate $E(X')$. We show above that we have a $\frac{1}{4}$ probability of winning the game. We remind the reader that Orloff and Bloom give us that we are playing with fair dice in this problem. We show in the listing of possible outcomes of the game above that there are 6 outcomes where we win. Since the dice are fair, each of the 6 outcomes are equally likely, but we win the same number of dollars for some of the outcomes. So the expected number of dollars we can win playing one round of the game is:

$$\frac{4 + 6 + 6 + 8 + 8 + 8}{6} = \frac{40}{6} = 6\frac{2}{3}. \quad (12)$$

According to the rules of the game that Orloff and Bloom give, the expected number of dollars we can lose playing one round of the game is 1 dollar.

Now we have values for all of the quantities on the right hand side of equation 11. Hence

$$E(W') = \frac{1}{4} \times \$6\frac{2}{3} - \left(1 - \frac{1}{4}\right) \times \$1 \approx \$0.917. \quad (13)$$

Therefore, if we play 60 rounds of the game, we expect to win approximately $60 \times \$0.917 = \55.02 .

References

- [1] John Hancock. *MIT Introduction to Statistics 18.05 Reading 5a Questions*. Available at <https://github.com/jhancock1975/mit-intro-to-stats-18.05/blob/prob-set-3/reading-5a-questions/reading5aQuestions.pdf> (2016/3/16).
- [2] Jeremy Orloff and Jonathan Bloom. *18.05 Problem Set 3, Spring 2014*. Available at https://ocw.mit.edu/courses/mathematics/18-05-introduction-to-probability-and-statistics-spring-2014/assignments/MIT18_05S14_ps3.pdf (Spring 2014).
- [3] Jeremy Orloff and Jonathan Bloom. *Code for Studio 3 (ZIP)*.
- [4] Jeremy Orloff and Jonathan Bloom. *Conditional Probability, Independence and Bayes Theorem Class 3, 18.05, Spring 2014*. Available at https://ocw.mit.edu/courses/mathematics/18-05-introduction-to-probability-and-statistics-spring-2014/readings/MIT18_05S14_Reading3.pdf (Spring 2014).

- [5] Jeremy Orloff and Jonathan Bloom. *Discrete Random Variables Class 4, 18.05, Spring 2014 Jeremy Orloff and Jonathan Bloom*. Available at https://ocw.mit.edu/courses/mathematics/18-05-introduction-to-probability-and-statistics-spring-2014/readings/MIT18_05S14_Reading4a.pdf (Spring 2014).
- [6] Jeremy Orloff and Jonathan Bloom. *Variance of Discrete Random Variables Class 5, 18.05, Spring 2014 Jeremy Orloff and Jonathan Bloom*. Available at https://ocw.mit.edu/courses/mathematics/18-05-introduction-to-probability-and-statistics-spring-2014/readings/MIT18_05S14_Reading5a.pdf (Spring 2014).
- [7] StackExchange.com user Paul Gessler. *TikZ: How to set a node on an exact position on a line?* Available at <http://tex.stackexchange.com/questions/147052/tikz-how-to-set-a-node-on-an-exact-position-on-a-line> (2015/2/08).
- [8] StackExchange user Peter Grill. *Draw a Plot with Point - Answer*.
- [9] StackOverflow.com user Shane. *Counting the number of elements with the values of x in a vector*.
- [10] R Core Team and contributors worldwide. *Standard Deviation*.
- [11] Texample.net user Till Tantau. *Example: Venn diagram*. Available at <http://www.texample.net/tikz/examples/venn-diagram/> (2006/11/08).