

HÁSKÓLINN Í REYKJAVÍK
REYKJAVIK UNIVERSITY

COMPUTER SCIENCE DEPARTMENT

STATS 1
TÖLFRÆÐI 1

Workbook

Lilja Ýr Guðmundsdóttir

Email: liljag18@ru.is

Contents

1 Week 1	4
1.1 Concepts	4
1.2 R Code	4
1.3 Example Problem Set	4
1.4 Problem Set	10
2 Week 2	11
2.1 Concepts	11
2.2 Example Problem Set	11
2.3 Problem Set	17
2.4 Problem Set solutions	18
3 Week 3	20
3.1 Concepts	20
3.2 Example Problem Set	21
3.3 Problem Set	32
3.4 Problem set Solution	34
4 Week 2	35
4.1 Concepts	35
4.2 Example Problem Set	35
4.3 Problem Set	42
4.4 Problem Set Solutions	43
5 Week 2	44
5.1 Concepts	44
5.2 R code	44
5.3 Example Problem Set	45
5.4 Problem Set	51
5.5 Problem Set Solution	52
6 Week 2	53
6.1 Concepts	53
6.2 R code	53
6.3 Example Problem Set	54
6.4 Problem Set	66
6.5 Problem Set Solutions	68
7 Week 7	70
7.1 Concepts	70
7.2 R code	70
7.3 Example Problem Set	71

7.4	Problem Set	79
7.5	Problem Set Solutions	80
8	Week 8	81
8.1	Concepts	81
8.2	Example Problem Set	82
8.3	Problem Set	88
8.4	Problem Set Solution	90
9	Week 9	91
9.1	Concepts	91
9.2	R code	91
9.3	Example Problem Set	92
9.4	Problem Set	104
9.5	Problem Set Solution	107
10	Week 10	109
10.1	Concepts	109
10.2	R code	109
10.3	Example Problem Set	110
10.4	Problem Set	119
10.5	Problem Set Solution	122
11	Week 11	123
11.1	Concepts	123
11.2	R code	123
11.3	Example Problem Set	124
12	Week 11	136
12.1	Concepts	136
12.2	R code	136
12.3	Example Problem Set	137

Introduction

This is my workbook for Statistics class containing any notes and work I've done over the semester. There are 118 problems and examples solved here (even though the workbook turn in guideline said there were only 117, might have made a counting error). There are also some that aren't correct but I learned from them.

1 Week 1

1.1 Concepts

We calculate the mean, median and standard deviation, easiest to just use R for this.

Example: The Statistics course has 220 students: 88 girls and 132 boys.

The problem solving classes are Tuesday and Wednesday.

On Tuesday 44 students are expected: 11 girls and 33 boys

On Wednesday 176 are expected: 77 girls and 99 boys.

A student name is randomly picked from the list of all registered students.

What is the probability that the randomly selected student is a girl? -> 0.4

What is the probability that the randomly selected student is expected Wednesday? -> 0.8

What is the probability that the student is a girl who is expected on Wednesday? -> 0.35

What is the probability that the student is a girl, or a student expected on Tuesday? -> 0.55

1.2 R Code

```
getwd()
etwd( "/Users/liljayrgudmundsdottir/Documents/HR_Files/5onn/Tolfraedi/datasets/Ch1/table1-2.csv")
data <- read.csv( file = 'table1-2.csv' )
x = data$PM
median(x)
mean(x)
sd(x)
```

1.3 Example Problem Set

Datasets and software: Download the datasets from Canvas the online version of the book and use Ch1\table1-2.csv to calculate the median, mean value, and standard deviation. Make a histogram and a boxplot. Use Excel or R. Use 4, 6, and 8 bins and compare.

```
> getwd()
[1] "/Users/liljayrgudmundsdottir"
> setwd("/Users/liljayrgudmundsdottir/Documents/HR_Files/5onn/Tolfraedi")
> getwd()
[1] "/Users/liljayrgudmundsdottir/Documents/HR_Files/5onn/Tolfraedi"
> setwd("/Users/liljayrgudmundsdottir/Documents/HR_Files/5onn/Tolfraedi/datasets/Ch1")
> getwd()
[1] "/Users/liljayrgudmundsdottir/Documents/HR_Files/5onn/Tolfraedi/datasets/Ch1"
> data <- read.csv(file =table1-2.csv)
Error: unexpected symbol in "data <- read.csv(file =table1-2.csv"
> data <- read.csv(file ='table1-2.csv')
```

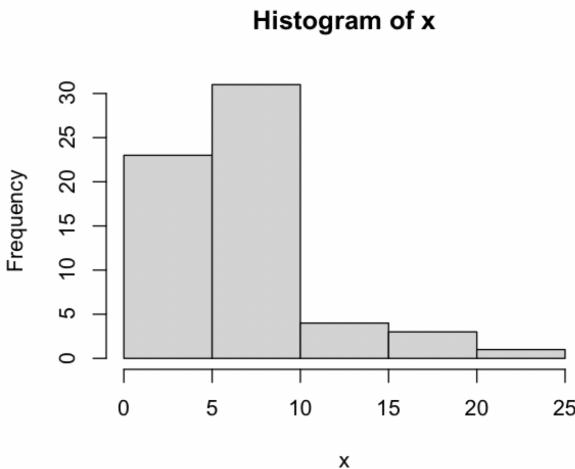
```

> x = data$PM
> x
[1]  7.59  6.28  6.07  5.23  5.54  3.46  2.44  3.01 13.63 13.02 23.38  9.24  3.22  2.06
[15] 4.04 17.11 12.26 19.91  8.50  7.81  7.18  6.95 18.64  7.10  6.04  5.66  8.86  4.40
[29] 3.57  4.35  3.84  2.37  3.81  5.32  5.84  2.89  4.68  1.85  9.14  8.67  9.52  2.68
[43] 10.14  9.20  7.31  2.09  6.32  6.53  6.32  2.01  5.91  5.60  5.61  1.50  6.46  5.29
[57]  5.64  2.07  1.11  3.32  1.83  7.56
> median(x)
[1] 5.75
> mean(x)
[1] 6.596452
> sd(x)
[1] 4.518998

```

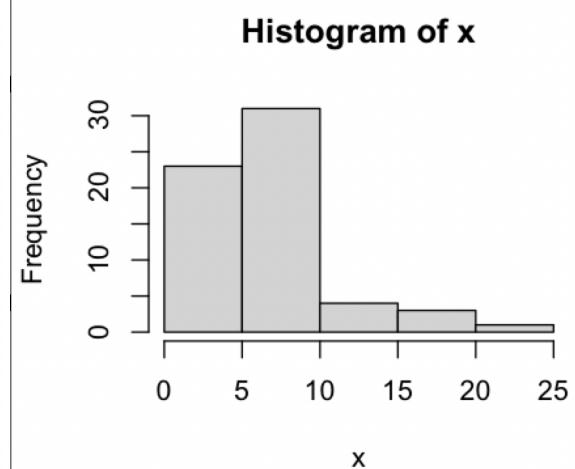
> hist(x, breaks=4)

Quartz 2 [*]



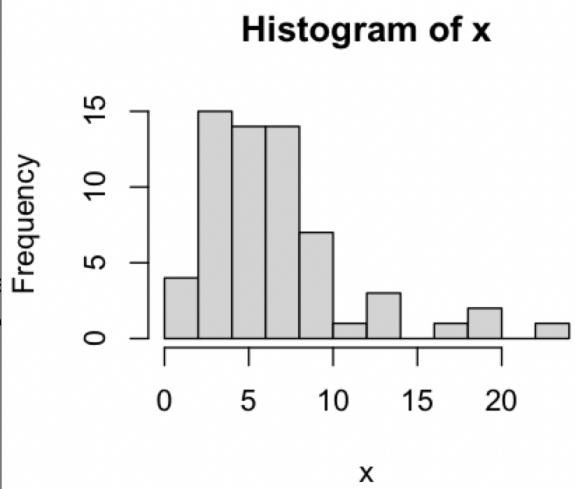
> hist(x, breaks=6)

Quartz 2 [*]



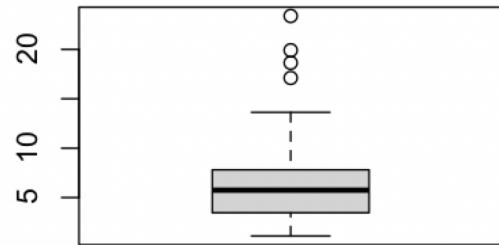
> hist(x, breaks=8)

Quartz 2 [*]



> boxplot(data\$PM)

Quartz 2 [*]



Data summary:

5. Find a sample size for which the median will always equal one of the values in the sample.

When $n = 1$ then the median will always be equal to one value in the dataset. This is true for all odd number amounts of data points.

6. For a list of positive numbers, is it possible for the standard deviation to be greater than the mean? If so, give an example. If not, explain why not.

Yes, depending on the data set, i.e (1,2,3,16) gives a mean of 5.5 and a standard deviation of 7.04. Also works for spread values.

7. Is it possible for the standard deviation of a list of numbers to equal 0? If so, give an example. If not, explain why not.

If dataset = {1} then $sd = \sqrt{\frac{1}{1}(1+1)^2} = 0$	
If dataset = {1, 1} then $\bar{x} = \frac{1}{2}(1+1) = 1$, $sd = \sqrt{\frac{1}{2}((1-1)^2 + (1-1)^2)} = 0$	
If dataset = {2, 2, 2} then $\bar{x} = \frac{1}{3}(2+2+2) = 2$ $sd = \sqrt{\frac{1}{3}((2-2)^2 + (2-2)^2 + (2-2)^2)} = 0$	
Therefore sd of list of values can be 0 if all the values are the same	

Problem 1.2.10 page 24 Edition 5: A sample of 100 cars driving on a freeway during a morning commute was drawn, and the number of occupants in each car was recorded. The results were as follows:

Occupants	1	2	3	4	5
Number of Cars	70	15	10	3	2

- Find the sample mean number of occupants.
- Find the sample standard deviation of the number of occupants.
- Find the sample median number of occupants.
- Compute the first and third quartiles of the number of occupants.
- What proportion of cars had more than the mean number of occupants?
- For what proportion of cars was the number of occupants more than one standard deviation greater than the mean?
- For what proportion of cars was the number of occupants within one standard deviation of the mean?

10. a) Sample mean: $\bar{x} = \frac{152}{100} = 1.52$

$x: 1 \ 2 \ 3 \ 4 \ 5$

$f_i: 70 \ 15 \ 10 \ 3 \ 2$ $x^2 f_i: 70 \ 60 \ 90 \ 48 \ 50$

$$x^2 f_i: 70 \ 30 \ 30 \ 12 \ 10 \Rightarrow \# \text{ of occupants total} = \sum x \cdot f_i = 152$$

$$\begin{aligned} b) \text{sd} &= \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} = \sqrt{\left(\frac{1}{n-1} \left(\sum (x^2 f_i) - \frac{(\sum x \cdot f_i)^2}{n} \right) \right)^{-1}} \\ &= \left(\frac{1}{n-1} \left(318 - \frac{152^2}{n} \right) \right)^{-1/2} = (0.878)^{-1/2} = 0.94 \end{aligned}$$

c) median: 1 occupant

d) 1st quart: $0.25n = 25 \rightarrow 1$ occupant

3rd quart: $0.75n = 75 \rightarrow 2$ occupants

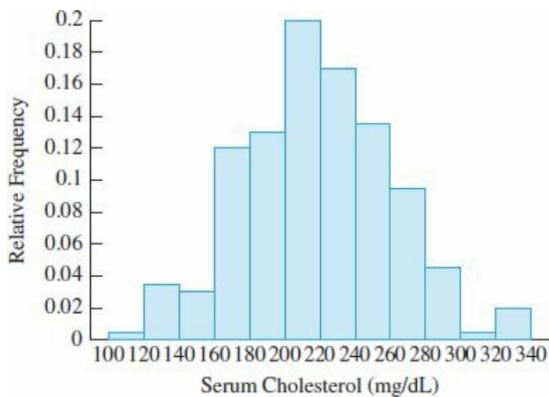
$$e) P(x > 1.25) = \frac{15}{100} + \frac{10}{100} + \frac{3}{100} + \frac{2}{100} = 0.15 + 0.10 + 0.03 + 0.02 = 0.3 \quad \underline{30\%}$$

$$f) P(x > 1.52 + 0.94) = P(x > 2.46) = \frac{10}{100} + \frac{3}{100} + \frac{2}{100} = \underline{15\%}$$

$$g) P(0.58 < x < 2.46) = \frac{70}{100} + \frac{15}{100} = 0.70 + 0.15 = \underline{85\%}$$

7. The figure below is a histogram showing the distribution of serum cholesterol level for a sample of men. Use the histogram to answer the following questions:

- Is the percentage of men with cholesterol levels above 240 mg/dL closest to 30%, 50%, or 70%?
- In which interval are there more men: 240–260 mg/dL or 280–340 mg/dL?



1.3.7
 a) closer to the 30% of total
 b) 240–260 mg/dL

Standard deviation:

Prove Eq.1.3 using Eq.1.2 at page 15. What do you obtain if $n \gg 1$?

Couldn't quite get this one, didn't make sense.

Sample space:

Example 2.1

An electrical engineer has on hand two boxes of resistors, with four resistors in each box. The resistors in the first box are labeled 10Ω (ohms), but in fact their resistances are 9, 10, 11, and 12Ω . The resistors in the second box are labeled 20Ω , but in fact their resistances are 18, 19, 20, and 21Ω . The engineer chooses one resistor from each box and determines the resistance of each.

Let A be the event that the first resistor has a resistance greater than 10, let B be the event that the second resistor has a resistance less than 19, and let C be the event that the sum of the resistances is equal to 28. Find a sample space for this experiment, and specify the subsets corresponding to the events A , B , and C .

$$\begin{aligned} \text{Example 2.1} \\ S &= \{(9, 18), (9, 19), (9, 20), (9, 21), (10, 18), (10, 19), (10, 20), (10, 21), \\ &\quad (11, 18), (11, 19), (11, 20), (11, 21), (12, 18), (12, 19), (12, 20), (12, 21)\} \\ A &: (\Omega_1 > 10) : \{(11, 18), (11, 19), (11, 20), (11, 21), (12, 18), (12, 19), (12, 20), (12, 21)\} \\ B &: (\Omega_2 < 19) : \{(9, 18), (10, 18), (11, 18), (12, 18)\} \\ C &: (\Omega_{\text{sum}} = 28) : \{(10, 18), (9, 19)\} \end{aligned}$$

3. A section of an exam contains four True-False questions. A completed exam paper is selected at random, and the four answers are recorded.
 - a. List all 16 outcomes in the sample space.
 - b. Assuming the outcomes to be equally likely, find the probability that all the answers are the same.
 - c. Assuming the outcomes to be equally likely, find the probability that exactly one of the four answers is “True.”
 - d. Assuming the outcomes to be equally likely, find the probability that at most one of the four answers is “True.”

Problem 2.1.3.

a) $S = \{(T, T, T, T), (T, T, T, F), (T, T, F, T), (T, F, T, T), (F, T, T, T), (F, T, T, F), (F, T, F, T), (F, F, T, T), (F, F, T, F), (F, F, F, T), (F, T, F, F), (T, F, F, F), (T, T, F, F), (T, F, T, F), (T, F, F, T), (F, F, F, F)\}$

b) $\frac{2}{16} = \frac{1}{8} = 0.125 \quad 12.5\% \text{ chance}$

c) $\frac{4}{16} = \frac{1}{4} = 0.25 \quad 25\% \text{ chance}$

d) $\frac{4}{16} + \frac{1}{16} = \frac{5}{16} = 0.3125 \quad 31.25\% \text{ chance}$

1.4 Problem Set

Lilja Ýr Guðmundsdóttir
23 August, 2020

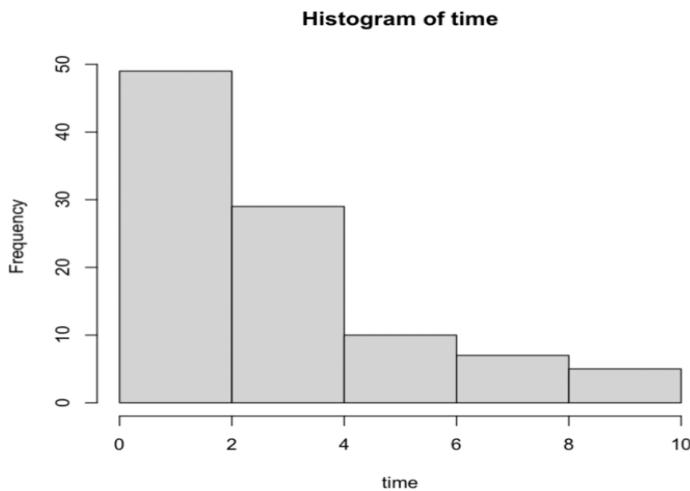
Verkefni 1

I used R to find the correct values (see code on the next page for how).

Results:

	Result
Mean	2.809726
Standard Deviation	2.35351

Histogram with 5 bins:



Lilja Ýr Guðmundsdóttir
23 August, 2020

Code (after navigating to right directory):

```
> data= read.table('dataC.csv', header=T)
> attach(data)
> mean(time)
[1] 2.809726
> sd(time)
[1] 2.35351
> hist(time, breaks=5)
```

2 Week 2

2.1 Concepts

- $P(\bar{A}) = 1 - P(A)$
- $P(A)P(B) = P(A \cup B)$
- $P(A \cap B) = P(A) + P(B) - P(A \cup B)$
- $P(A \cap B) = P(A) + P(B)$ if A and B are independent
- $P(A)P(B) = P(A \cap B)$ if A and B are independent
- Baye's rule: $P(A|B) = \frac{P(A \cap B)}{P(B)}$
- $P(A|B)$ means A given B
- Suppose $P(A|B)=P(A)$ ($\neq 0$ or 1) then In this case $P(B|A) = P(B)$
- If $P(A|B)=P(A)$ then the events A and B are independent
- A disease test can be 99% accurate, but still a large proportion of people tested positive are actually disease-free. This is possible when the disease is rare.
- A parent and an offspring (daughter or son) have 50% DNA in common
- Suppose both parents have genotype (a,b) at a specific location on the genome. The probability for their child to have (b,b) is 0.25

2.2 Example Problem Set

Combined events

2.1.13 Let S be the event that a randomly selected college student has taken a statistics course, and let C be the event that the same student has taken a chemistry course. Suppose $P(S) = 0.4$, $P(C) = 0.3$, and $P(S \cap C) = 0.2$.

- a. Find the probability that a student has taken statistics, chemistry, or both.
- b. Find the probability that a student has taken neither statistics nor chemistry.
- c. Find the probability that a student has taken statistics but not chemistry.

2.1.13 $P(S) = 0.4$, $P(C) = 0.3$, $P(S \cap C) = 0.2$

a) $P(S) + P(C) - P(S \cap C) = 0.5 = P(S \cup C)$

b) $P(S^c \cap C^c) = 1 - P(S \cup C) = 0.5$

c) Find $P(S \cap C^c)$. We have $P(S) = P(S \cap C) + P(S \cap C^c)$
which gives: $P(S \cap C^c) = P(S) + P(C) - P(S \cup C) = 0.4 + 0.3 - 0.5 = 0.2$
Since $P(S) = 0.4$ and $P(S \cap C) = 0.2$ then $P(S \cap C^c) = 0.2$.

2.1.14 Six hundred paving stones were examined for cracks, and 15 were found to be cracked. The same 600 stones were then examined for discoloration, and 27 were found to be discolored. A total of 562 stones were neither cracked nor discolored. One of the

600 stones is selected at random.

- Find the probability that it is cracked, discolored, or both.
- Find the probability that it is both cracked and discolored.
- Find the probability that it is cracked but not discolored.

2.1.14 600 stone total $P(C) = \frac{15}{600}$ $P(D) = \frac{27}{600}$ $P(C \cap D^c) = \frac{562}{600}$

~~a) $P(C \cup D) = P(C) + P(D) - P(C \cap D)$~~

~~$= 0.025 + 0.045 - 0.018 = 0.052$~~

~~b) $P(C \cap D^c) = P(C) - P(C \cap D)$~~

~~$= 0.025 - 0.018 = 0.007$~~

a) $P(C \cup D) = 1 - P(C^c \cap D^c) = 1 - 0.9367 = 0.0633$

b) $P(C \cap D) = P(C) + P(D) - P(C \cup D) = 0.025 + 0.045 - 0.0633 = 0.0067$

c) $P(C \cap D^c) \rightarrow P(C) = P(C \cap D) + P(C \cap D^c) \Rightarrow P(C \cap D^c) = P(C \cap D) - P(C) = \dots$

2.1.16 A system contains two components, A and B. The system will function so long as either A or B functions. The probability that A functions is 0.95, the probability that B functions is 0.90, and the probability that both function is 0.88. What is the probability that the system functions?

2.1.16 $P(A) = 0.95$ $P(B) = 0.90$ $P(A \cup B) = 0.88$

$P(\text{sys functions}) = P(A \cap B) = P(A) + P(B) - P(A)P(B) = P(A) + P(B) - P(A \cup B)$

$= 0.95 + 0.90 - 0.88 = 0.97$

2.1.17 A system contains two components, A and B. The system will function only if both components function. The probability that A functions is 0.98, the probability that B functions is 0.95, and the probability that either A or B functions is 0.99. What is the probability that the system functions?

2.1.17 $P(A) = 0.98$ $P(B) = 0.95$ $P(A \cup B) = 0.99$

$P(\text{SF}) = P(A \cap B) = P(A) + P(B) - P(A \cup B) = 0.98 + 0.95 - 0.99 = 0.94$

Combined events and conditional probabilities

Example 2.4 A target on a test firing range consists of a bull's-eye with two concentric rings around it. A projectile is fired at the target. The probability that it hits the bull's-eye is 0.10, the probability that it hits the inner ring is 0.25, and the probability

that it hits the outer ring is 0.45. What is the probability that the projectile hits the target? What is the probability that it misses the target?

Example 2.4: $P(B) = 0.10$	$P(I) = 0.25$	$P(O) = 0.45$
$P(B \cap I \cap O) = P(B) + P(I) + P(O) = 0.60 = P(H)$		
Misses: $P(M) = 1 - P(H) = 0.20$		

Example 2.6 An extrusion die is used to produce aluminum rods. Specifications are given for the length and the diameter of the rods. For each rod, the length is classified as too short, too long, or OK, and the diameter is classified as too thin, OK, or too thick. In a population of 1000 rods, the number of rods in each class is as follows:

Length	Diameter		
	Too thin	OK	Too thick
Too Short	10	3	5
OK	38	900	4
Too long	2	25	13

A rod is sampled at random from this population. What is the probability that it is too short?

Example 2.6: sample size = 1000
$P(TS) = (10+3+5)/1000 = 0.018.$

Example 2.7 Refer to Example 2.6. If a rod is sampled at random, what is the probability that it is either too short or too thick?

Example 2.7
$P(TS \cup TT) = P(\text{too short}) + P(\text{too thick}) - P(\text{too short and too thick})$
$= \frac{18}{1000} + \frac{22}{1000} - \frac{5}{1000}$

2.2.11 One drawer in a dresser contains 8 blue socks and 6 white socks. A second drawer contains 4 blue socks and 2 white socks. One sock is chosen from each drawer.

What is the probability that they match?

$$2.2.11 \quad ① B=8 \quad W=6$$

$$② B=4 \quad W=2$$

$$P(\text{matching}) = P(BB) + P(WW)$$

$$= \frac{8}{14} \cdot \frac{4}{6} + \frac{6}{14} \cdot \frac{2}{6} = \frac{44}{84} = 0.5238$$

2.2.12 A drawer contains 6 red socks, 4 green socks, and 2 black socks. Two socks are chosen at random. What is the probability that they match?

$$2.2.12 \quad R=6 \quad B=8 \quad W=2 \quad T=16$$

$$P(\text{matching}) = \frac{6}{16} \cdot \frac{5}{15} + \frac{6}{16} \cdot \frac{7}{15} + \frac{2}{16} \cdot \frac{1}{15} = 0.26 \approx \frac{1}{3}$$

2.3.8 A drag racer has two parachutes, a main and a backup, that are designed to bring the vehicle to a stop after the end of a run. Suppose that the main chute deploys with probability 0.99, and that if the main fails to deploy, the backup deploys with probability 0.98.

- What is the probability that one of the two parachutes deploys?
- What is the probability that the backup parachute deploys?

$$2.3.8 \quad P(M) = 0.99 \quad P(B) = 0.98$$

$$\begin{aligned} a) P(\text{either works}) &= P(M \cup B) = P(M) + P(B) - P(M)P(B) \\ &= 0.99 + 0.98 - 0.99 \cdot 0.98 = 0.998 \end{aligned}$$

$$b) P(B) = (1 - P(A))P(B) = (1 - 0.99)0.98 = 0.0098$$

2.3.9 At a certain car dealership, 20% of customers who bought a new vehicle bought an SUV, and 3% of them bought a black SUV. Given that a customer bought an SUV, what is the probability that it was black?

$$2.3.9 P(S) = 20\% \quad P(B) = 39.3\%$$

$$P(\text{black given it's an SUV}) = P(S \cap B) / P(S) = P(B) / P(S) = \underline{0.15}$$

2.3.15 A population of 600 semiconductor wafers contains wafers from three lots. The wafers are categorized by lot and by whether they conform to a thickness specification. The following table presents the number of wafers in each category. A wafer is chosen at random from the population.

Lot	Conforming	Nonconforming
A	88	12
B	165	35
C	260	40

- a. If the wafer is from Lot A, what is the probability that it is conforming?
- b. If the wafer is conforming, what is the probability that it is from Lot A?
- c. If the wafer is conforming, what is the probability that it is not from Lot C?
- d. If the wafer is not from Lot C, what is the probability that it is conforming?

$$2.3.15 T = 600$$

$$a) P(\text{Conforming given from A}) = \frac{88}{88+12} = 0.88$$

$$b) P(\text{Conforming and from A}) = \frac{88}{88+165+260} = 0.172$$

$$c) P(\text{Conforming and not from C}) = \frac{88+165}{88+165+260} = 0.493$$

$$d) P(\text{Conforming given not from C}) = \frac{88+165}{88+12+165+35} = 0.843$$

2.3.24 A lot of 1000 components contains 300 that are defective. Two components are drawn at random and tested. Let A be the event that the first component drawn is defective, and let B be the event that the second component drawn is defective.

- a. Find $P(A)$.
- b. Find $P(B|A)$.
- c. Find $P(A \cap B)$.
- d. Find $P(A^c \cap B)$.
- e. Find $P(B)$.
- f. Find $P(A|B)$.

g. Are A and B independent? Is it reasonable to treat A and B as though they were independent? Explain.

2.3.24	$T = 1000$	$D = 300$	$A \rightarrow \text{first def.}$	$B \rightarrow \text{second def.}$
a)	$P(A) = \frac{300}{1000} = 0.3 = \left(\frac{300}{1000} \cdot \frac{299}{999}\right) + \left(\frac{300}{1000} \cdot \frac{700}{999}\right)$			
b)	$P(B A) = \frac{P(B \cap A)}{P(A)} = \left(\frac{300}{1000} \cdot \frac{299}{999}\right) \div 0.3 = 0.2993$			
c)	$P(A \cap B) = P(B A) \cdot P(A) = 0.2993 \cdot 0.3 = 0.0898$			
d)	$P(A^c \cap B) = \frac{700}{1000} \cdot \frac{300}{999} = 0.2102$			
e)	$P(B) = \left(\frac{300}{1000} \cdot \frac{299}{999}\right) + \left(\frac{700}{1000} \cdot \frac{300}{999}\right) = 0.3$			
f)	$P(A B) = \frac{P(A \cap B)}{P(B)} = \frac{0.0898}{0.3} = 0.2993$			
g)	Yes, independent b/c $P(B A) = P(A B)$			

2.3.27 Each day, a weather forecaster predicts whether or not it will rain. For 80% of rainy days, she correctly predicts that it will rain. For 90% of non-rainy days, she correctly predicts that it will not rain. Suppose that 10% of days are rainy and 90% are non-rainy.

- a. What proportion of the forecasts are correct?
- b. Another forecaster always predicts that there will be no rain. What proportion of these forecasts are correct?

correct given it rains			
2.3.27	$P(C R) = 0.80$	$P(C N) = 0.90$	$P(R) = 0.10$
			$P(N) = 0.90$ no rain
a)	$P(C) = P(C R)P(R) + P(C N)P(N) = 0.89$		
b)	It will always be right when there is no rain so it will be right 90% of the time.		

2.3 Problem Set

Litla Yr Guðmundsdóttir

1. 50 light bulbs {^{5 bulbs}_{bad}} $P(B) = 0.1$ ∴ $P(G) = (1 - P(B)) = 0.9$
- a) $P(2 \text{ good bulbs}) = P(G_1) \cdot P(G_2) = \frac{45}{50} \cdot \frac{44}{49} = \frac{1980}{2450} \approx 0.809 \text{ or } \underline{\sim 81\%}$
- b) $P(1^{\text{st}} \text{ bad}, 2^{\text{nd}} \text{ good}) = P(B_1) \cdot P(G_2) = \frac{5}{50} \times \frac{45}{49} = \frac{45}{490} \approx 0.092 \text{ or } \underline{\sim 9.2\%}$
- c) $P(1 \text{ bulb}, 2^{\text{nd}} \text{ good}) = P(B \cap G_2) \cdot P(G_2) = (P(B) \cdot P(G_2)) + (P(G) \cdot P(G_2))$
 $= \left(\frac{5}{50} \cdot \frac{45}{49}\right) + \left(\frac{45}{50} \cdot \frac{44}{49}\right) = \frac{1980}{2450} + \frac{45}{100} \approx 0.901 \text{ or } \underline{\sim 90\% \text{ chance}}$

Explanations for problem 1

a) You first pick one light bulb out of all 50 so the probability of that one being good is $P(G_1) = \frac{\# \text{ good bulbs}}{\text{total } \# \text{ of bulbs}}$. Then when the second bulb is being picked then there is one less good bulb to pick and one less bulb in total assuming we don't replace. Therefore (∴) The probability of picking another good bulb is $P(G_2) = \frac{\# \text{ good bulbs remaining}}{\text{total } \# \text{ of remaining bulbs}} = \frac{\# \text{ good bulbs total} - 1}{\text{total } \# \text{ of bulbs} - 1}$.

Then just multiply the two values together to get total probability

b) Same principle as in (a) except now we have to factor in the # of bad bulbs. First we pick one light bulb and the probability of that one being bad is $P(B_1) = \frac{\# \text{ of bad bulbs}}{\text{total } \# \text{ of bulbs}}$.

Then we pick another bulb from a pile of bulbs one fewer than before in total but, assuming that the first bulb that was picked was bad, the amount of good bulbs hasn't changed. ∴ the probability of a good bulb is $P(G_2) = \frac{\# \text{ of good bulbs}}{\text{total } \# \text{ of remaining bulbs}} = \frac{\# \text{ of good bulbs}}{\text{total } \# \text{ of bulbs} - 1}$.

Then just multiply the two together like before.

c) This splits into 2 parts:

i) Probability of drawing 2 good bulbs: this was accomplished in (a).

ii) Probability of drawing first 1 bad bulb and then 1 good bulb: this was accomplished in (b).

To get the final probability of drawing a good bulb irrespective of whether the first was good or bad simply add parts (i) and (ii) together.

BRUNNEN, TH

Lilja Ýr Guðmundsdóttir

2 There are 2 versions for problem 1: 1A and 1B

and 2 versions for problem 2: 2A and 2B

a) The probability of the combination 1B and 2B is simple enough to solve with brute force. But for the purpose of this exercise we will solve it with the proper mathematical methods. ∴ the probability of a student getting problem 1B is $P(1B) = \frac{\# \text{ of } 1B}{\# \text{ of problems} 1}$ = $\frac{1}{2}$ and the probability of a student getting problem 2B is $P(2B) = \frac{\# \text{ of } 2B}{\# \text{ of problems} 2} = \frac{1}{2}$.

Then we simply multiply the 2 together and get $\frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4}$. This method works because the students have an equal probability of getting one of the 2 problems in each problem set.

b) If we assume the probability of getting 1A and 1B and getting 2A and 2B always remain the same no matter how many assignments the teacher has given then the probability for each student to get the combination of 1B and 2B stays the same (25% as per (a)).

∴ the probability of 2 students getting that assignment is at 25%. The probability of them getting 1A-2A is also 25% and so on for 1A-2B and 1B-2A.
∴ the probability is always 25%.

□

2.4 Problem Set solutions

1. A box contains 50 light bulbs, some of them bad and some of them good. The probability to pick a bad one is p . Two bulbs are drawn at random. Calculate:
 - The probability that both are good. We have N bulbs in total. Of them $B = pN$ are

bad and $G = (1 - p)N$ are good, and $G + B = N$. $P(\text{both 1 and 2 good}) = P(\text{1 good} \cap \text{2 good}) = P(\text{2 good} \cap \text{1 good}) = P(\text{2 good} | \text{1 good})P(\text{1 good}) = \frac{G-1}{N-1} \frac{G}{N}$.

b) The probability that if the first is bad the second is good.

If we picked one bad, there are still G good inside, but $N-1$ in total. $P(\text{2 good} | \text{1 bad}) = \frac{G}{N-1}$.

c) The probability that the second is good irrespective of whether the first was good or bad.

Use the law of total probability: $P(\text{2 good}) = P(\text{2 good} \cap \text{1 good}) + P(\text{2 good} \cap \text{1 bad}) = P(\text{2 good} | \text{1 good})P(\text{1 good}) + P(\text{2 good} | \text{1 bad})P(\text{1 bad}) = \frac{G-1}{N-1} \frac{G}{N} + \frac{G}{N-1} \frac{B}{N} = \frac{G}{N}$. So we see that if we do not condition on the quality of bulb 1, then the probability that the 2-nd is good is the same as the probability that the 1-st is good.

2. An assignment for the statistics class consists of two problems. Problem 1 has two versions, 1A and 1B, and problem 2 has also two versions, 2A and 2B. The two versions of each problem are randomized. Meaning that each student receives at random 1A or 1B (not both) with equal probabilities, and 2A or 2B (not both) with equal probabilities. Calculate:

a) The probability that a student obtains a specific combination.

All possible combinations for a single student: 1A-2A, 1A-2B, 2A-1B, 2A-2B. The probability of any combination is $1/4$.

b) The probability that two students obtain the same assignment:

All combinations for student 1: 1A-2A, 1A-2B, 2A-1B, 2A-2B

All combinations for student 2: 1A-2A, 1A-2B, 2A-1B, 2A-2B

We observe 4 possibilities for the same assignment out of 16 possibilities in total, so the result is $4/16 = 1/4$.

3 Week 3

3.1 Concepts

Example: The table shows the values of a discrete random variable n and their probabilities p(n).

n	0	1	2	3
p(n)	0.4	0.5	0.2	

What is the mean value of n? -> Mean value of n is 1.

Bayes rule:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|D^c)P(D^c)}$$

- The cumulative distribution function F(a) is the probability that the random variable is less or equal to a
- The variance of a random variable x is not $[E(x)]^2 - E(x^2)$
- Considering a continuous random variable x with values between 0 and 1 the probability that $x=0.5$ is 0
- Consider a continuous random variable x. The area enclosed between the curve representing the probability density function f(x) and the x axis depends on f(x)
- Suppose R and S are two independent random variables, both with the same variance v. The variance of R-S is 2v
- Suppose Z is a random variable with standard deviation s. The standard deviation is $\sqrt{2} * x$
- We want to measure the average mass of coffee beans of a certain type. We measure the mass of 100 beans, one at a time, and take the average. We also calculate the sample standard deviation s. The standard deviation of the average mass is $s/\sqrt{10}$
- In order to reduce the std. of the average mass 10 times we should increase the number of measured coffee beans $\sqrt{10}$ times
- Two independent random variables have values between 0 and 1. Their covariance is 0

3.2 Example Problem Set

Examples of statistics in health sciences

Example 2.26 The proportion of people in a given community who have a certain disease is 0.005. A test is available to diagnose the disease. If a person has the disease, the probability that the test will produce a positive signal is 0.99. If a person does not have the disease, the probability that the test will produce a positive signal is 0.01. If a person tests positive, what is the probability that the person actually has the disease?

D = person has the disease

$$P(D) = 0.005, P(+|D) = 0.99, P(+|D^c) = 0.01 \text{ (false positive)}$$

Using Bayes rule:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|D^c)P(D^c)} = \frac{0.99 * 0.005}{0.99 * 0.005 + 0.001 * 0.995} = 0.332$$

Probability that person has the disease if they tested positive is 33.2%

2.26

P of people in com. who have a ^{certain} disease is 0.005

P of test being positive if person has disease is 0.99

P — if person doesn't have dis. is 0.01

D = person has disease

$P(D) = 0.005$ $P(+|D) = 0.99$ $P(+|D^c) = 0.01$, false pos.

Using Bayes rule:

$$P(D|+) = \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|D^c)P(D^c)}$$

$$= \frac{0.99 \cdot 0.005}{0.99 \cdot 0.005 + 0.001 \cdot 0.995} = 0.332$$

P that person has disease if her tested positive

Problem 2.3.30 A drag racer has two parachutes, a main and a backup, that are designed to bring the vehicle to a stop after the end of a run. Suppose that the main chute deploys with probability 0.99, and that if the main fails to deploy, the backup deploys with probability 0.98.

- a. What is the probability that one of the two parachutes deploys?

$$P(D|+) = \frac{0.99 * 0.05}{0.99 * 0.005 + 0.01 * 0.95} + \frac{0.0495}{0.059} = 0.83$$

- b. What is the probability that the backup parachute deploys?

$P(-|D^c) = 0.99, P(-|D) = 0.01, P(D^c) = 0.95$ which gives:

$$P(D^c|-) = \frac{P(-|D^c)P(D)}{P(-|D^c)P(D^c) + P(-|D)P(D)} = \frac{0.99 * 0.95}{0.99 * 0.95 + 0.01 * 0.05} = 0.99$$

	2.3.30
	$P(O)$ changed to 0.05
a)	$P(D +) = \frac{0.99 \cdot 0.05}{0.99 \cdot 0.05 + 0.01 \cdot 0.95} = \frac{0.0495}{0.0495 + 0.0095} = \frac{0.0495}{0.059} = 0.83$
b)	$P(D^c -) \quad P(- D^c) = 0.99 \quad P(- D) = 0.01 \quad P(D^c) = 0.95$
	$P(D^c -) = \frac{P(- D^c)P(D)}{P(- D^c)P(D^c) + P(- D)P(D)} = \frac{0.99 \cdot 0.95}{0.99 \cdot 0.95 + 0.01 \cdot 0.05}$
	$= \frac{0.9405}{0.9405 + 0.005} = \frac{0.9405}{0.941} = 99\%$

2.3.31 Sickle-cell anemia is an inherited disease in which red blood cells are misshapen and sticky. Sickle cells tend to form clumps in blood vessels, inhibiting the flow of blood. Humans have two genes for sickle-cell anemia, either of which may be S for normal cells or s for sickle cells. A person with two copies of the s gene will have sickle-cell anemia. A person with one s gene and one S gene will not have the disease, but will be a carrier, which means that the s gene may be transmitted to the person's offspring. If two carriers have a child, the probability is 0.25 that the child will have the disease and 0.5 that the child will be a carrier. Outcomes among children are independent.

a. A mother and father who are both carriers have two children. What is the probability that neither child has the disease?

$$P(\neg s) \cdot P(\neg s) = \frac{3}{4} \cdot \frac{3}{4} = 9/16$$

b. What is the probability that both children are carriers?

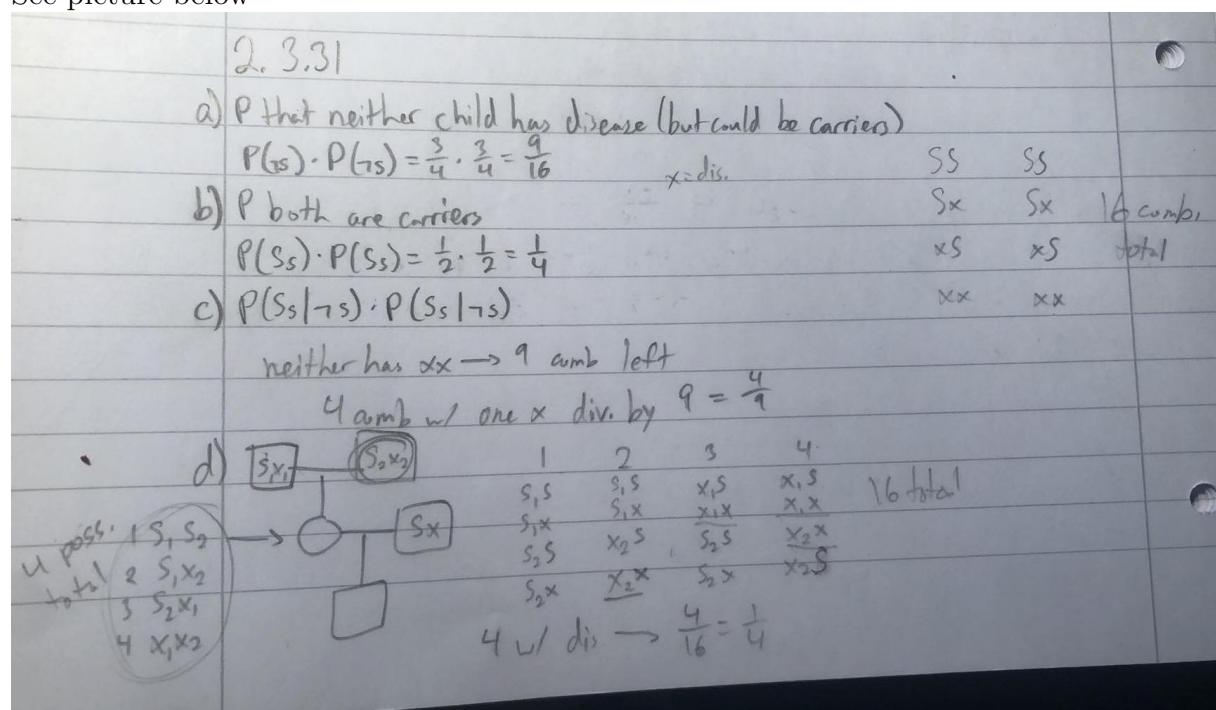
$$P(Ss) \cdot P(Ss) = 1/2 \cdot 1/2 = 1/4$$

c. If neither child has the disease, what is the probability that both are carriers?

$P(Ss|\neg s) \cdot P(Ss|\neg s)$ = neither has ss \Rightarrow 9 combos left, 4 combos with one s divided by 9 gives 4/9

d. A woman who is the child of two carriers has a child by a man who is a carrier. What is the probability that this child has the disease?

See picture below



Reliability analysis

2.3.34 A system consists of four components connected as shown in the diagram in the image below. Assume A, B, C, and D function independently. If the probabilities that A, B, C, and D fail are 0.10, 0.05, 0.10, and 0.20, respectively, what is the probability that the system functions?

$$P(\text{virki}) = P(l_1 \cup l_2) = P((A \cap B) \cup (C \cup D)) = P(A \cap B) + P(C \cup D) - P((A \cap B)(C \cup D))$$

$$P(A \cap B) = P(A)P(B) = 0.855$$

$$P(C \cup D) = P(C) + P(D) - P(C \cap D) = P(C) + P(D) - P(C)P(D) = 0.98$$

$$P(\text{virki}) = 0.855 + 0.98 - 0.855 \cdot 0.98 = 0.9971$$

Rel. analysis
2.3.34

work independently

$P(A^c) = 0.10$ $P(B^c) = 0.05$
 $P(C^c) = 0.10$ $P(D^c) = 0.20$

$$P(\text{virki}) = P(l_1 \cup l_2) = P((A \cap B) \cup (C \cup D)) = P(A \cap B) + P(C \cup D) - P((A \cap B)(C \cup D))$$

$$P(A \cap B) = P(A)P(B) = 0.855$$

$$P(C \cup D) = P(C) + P(D) - P(C \cap D) = P(C) + P(D) - P(C)P(D) = 0.98$$

$$P(\text{virki}) = 0.855 + 0.98 - P(A \cap B)P(C \cup D) = 0.85 + 0.98 - 0.85 \cdot 0.98 = \underline{\underline{0.9971}}$$

Discrete random variables

2.4.2 Computer chips often contain surface imperfections. For a certain type of computer chip, the probability mass function of the number of defects X is presented in the following table.

x	0	1	2	3	4
p(x)	0.4	0.3	0.15	0.10	0.05

a. Find $P(X \leq 2)$.

$$P(X \leq 2) = P(0) + P(1) + P(2) = 0.4 + 0.3 + 0.15 = 0.85$$

b. Find $P(X > 1)$.

$$P(X > 1) = P(2) + P(3) + P(4) = 0.15 + 0.10 + 0.05 = 0.3$$

c. Find μ_x .

$$\mu_x = \sum x_i P(x_i) = 0 \cdot 0.4 + 1 \cdot 0.3 + 2 \cdot 0.15 + 3 \cdot 0.10 + 4 \cdot 0.05 = 1.1$$

d. Find σ_x^2 .

$$\sigma_x^2 = \sum (x_i - \mu)^2 P(x_i) = \sum x_i^2 P(x_i) - \mu^2 = \mu_{x^2} - \mu_x^2 = 1.24 \text{ (full equation in image below).}$$

2.4.2	x	0	1	2	3	4
	p(x)	0.4	0.3	0.15	0.10	0.05
a)	$P(X \leq 2) = P(0) + P(1) + P(2) = 0.4 + 0.3 + 0.15 = 0.85$					
b)	$P(X > 1) = P(2) + P(3) + P(4) = 0.15 + 0.10 + 0.05 = 0.3$					
c)	$\mu_x = \sum x_i P(x_i) = 0 \cdot 0.4 + 1 \cdot 0.3 + 2 \cdot 0.15 + 3 \cdot 0.10 + 4 \cdot 0.05 = 1.1$					
d)	$\sigma_x^2 = \sum (x_i - \mu)^2 P(x_i) = \sum x_i^2 P(x_i) - \mu^2 = \mu_{x^2} - \mu_x^2$					
	$= 0^2 \cdot 0.4 + 1^2 \cdot 0.15 + 2^2 \cdot 0.15 + 3^2 \cdot 0.10 + 4^2 \cdot 0.05 - 1.1^2 = 2.45 - 1.21 = 1.24$					

2.4.7 A computer sends a packet of information along a channel and waits for a return signal acknowledging that the packet has been received. If no acknowledgment is received within a certain time, the packet is re-sent. Let X represent the number of times the packet is sent. Assume that the probability mass function of X is given by (see image below) where c is a constant.

SEE IMAGE BELOW FOR ANSWERS.

- Find the value of the constant c so that $p(x)$ is a probability mass function.
- Find $P(X = 2)$.
- Find the mean number of times the packet is sent.
- Find the variance of the number of times the packet is sent.
- Find the standard deviation of the number of times the packet is sent.

$2.4.7$ $p(x) = \begin{cases} cx, & x=1,2,3,4,5 \\ 0 & \text{otherwise} \end{cases}$
a) $x=0, 1, 2, 3, 4, 5 \Rightarrow 0, c, 2c, 3c, 4c, 5c = 1$ $\rightarrow 0 + c + 2c + 3c + 4c + 5c = 1 \Rightarrow 15c = 1 \Rightarrow c = \frac{1}{15}$
b) $P(2) = \frac{2}{15}$
c) $\mu_x = \sum x_i P(x_i) = \sum x_i \times \frac{x_i}{15} = \sum \frac{x_i^2}{15} = \frac{1+2^2+3^2+4^2+5^2}{15} = \frac{11}{3}$
d) $\text{var} = \sum (x_i - \mu)^2 P(x_i) = \sigma^2 = \sum x_i^2 P(x_i) - \mu^2 = 1 \cdot \frac{1}{15} + \frac{2^2}{15} + \frac{3^2}{15} + \frac{4^2}{15} + \frac{5^2}{15} - \left(\frac{11}{3}\right)^2$ $= \frac{1}{15} + \frac{8}{15} + \frac{27}{15} + \frac{64}{15} + \frac{125}{15} - \frac{121}{9} = \frac{225}{15} - \frac{121}{9} = \frac{45}{3} - \frac{121}{9} = \frac{135}{9} - \frac{121}{9} = \frac{14}{9}$
e) $\sigma = \sqrt{\sigma^2} = \sqrt{\frac{14}{9}} \approx 1.247 \quad \text{sd}$

2.4.10 Candidates for a job are interviewed one by one until a qualified candidate is found. Thirty percent of the candidates are qualified.

a. What is the probability that the first candidate is qualified?

The $P(C_1)$ being qualified is $0.3 \rightarrow 30\%$

b. What is the probability that the first candidate is unqualified and the second candidate is qualified?

$$P(C_1Q^c) \cdot P(C_2Q) = 0.7 \cdot 0.3 = 0.21$$

c. Let X represent the number of candidates interviewed up to and including the first qualified candidate. Find $P(X = 4)$.

$$P(X=4) = 0.7 \cdot 0.7 \cdot 0.7 \cdot 0.3 = 0.1029$$

d. Find the probability mass function of X .

$$P(x) = (1-p)^{x-1} p \dots$$

2.4.10

$$P(Q^c) = 0.3 \quad 30\% \text{ of cand. are qual.}$$

a) The $P(C_1)$ being qual. is $0.3 \rightarrow 30\%$

$$b) P(C_1Q^c) \cdot P(C_2Q) = 0.7 \cdot 0.3 = 0.21$$

$$c) P(X=4) = 0.7 \cdot 0.7 \cdot 0.7 \cdot 0.3 = 0.1029$$

$$d) P(x) = (1-p)^{x-1} p = (0.7)^{x-1} \cdot 0.3, x=1,2,3,\dots$$

2.4.8 After manufacture, computer disks are tested for errors. Let X be the number of errors detected on a randomly chosen disk. The following table presents values of the cumulative distribution function $F(x)$ of X .

x	0	1	2	3	4
$F(x)$	0.41	0.72	0.83	0.95	1.00

a. What is the probability that two or fewer errors are detected?

$$P(x \leq 2) = 0.83$$

b. What is the probability that more than three errors are detected?

$$P(x > 3) = 1 - P(X \leq 3) = 1 - 0.95 = 0.05$$

c. What is the probability that exactly one error is detected?

$$P(x=1) = 0.72 - 0.41 = 0.31$$

d. What is the probability that no errors are detected?

$$P(x=0) = 0.41$$

e. What is the most probable number of errors to be detected?

Most probable = 0

2.4.8	x	0	1	2	3	4
	$F(x)$	0.41	0.72	0.83	0.95	1.00
a)	$P(x \leq 2) = 0.83$					
b)	$P(x > 3) = 1 - P(x \leq 3) = 1 - 0.95 = 0.05$					
c)	$P(x=1) = 0.72 - 0.41 = 0.31$					
d)	$P(x=0) = 0.41$					
e)	most probable = 0					

Continuous random variables

2.4.15 The lifetime of a transistor in a certain application is random with probability density function in the image below.

SEE IMAGE FOR CALCULATIONS.

- Find the mean lifetime.
- Find the standard deviation of the lifetimes.
- Find the cumulative distribution function of the lifetime.
- Find the probability that the lifetime will be less than 12 months.

$f(t) = \begin{cases} 0.1e^{-0.1t} & t > 0 \\ 0 & t \leq 0 \end{cases}$
$\text{a) mean } \mu = \int_{-\infty}^{\infty} x f(x) dx = \int_0^{\infty} t f(t) dt = 0.1 \int_0^{\infty} t e^{-0.1t} dt = 0.1 \frac{1}{0.1^2} = 10$
$\text{b) } \sigma^2 = \int_{-\infty}^{\infty} x^2 f(x) dx - \mu^2 = 0.1 \int_0^{\infty} t^2 f(t) dt - \mu^2 = 0.1 \frac{2}{0.1^3} - 10^2 = 200 - 100 = 100, \sigma = \sqrt{100} = 10$
$\text{c) } F(L) = P(0 < t < L) = \int_0^L f(t) dt = 0.1 \int_0^L e^{-0.1t} dt = 1 - e^{-0.1L}$
$\text{d) } F(12) = 1 - e^{-0.1 \cdot 12} = 0.6988$

2.4.14 Elongation (in percent) of steel plates treated with aluminum are random with probability density function

a. What proportion of steel plates have elongations greater than 25%?

SEE IMAGE FOR CALCULATIONS.

b. Find the mean elongation.

c. Find the variance of the elongations.

d. Find the standard deviation of the elongations.

e. Find the cumulative distribution function of the elongations.

f. A particular plate elongates 28%. What proportion of plates elongate more than this?

$$2.4.14 \quad f(x) = \begin{cases} \frac{x}{250}, & 20 < x < 30 \\ 0 & \text{otherwise} \end{cases}$$

$$a) P(x > 25) = \int_{25}^{30} \frac{x}{250} dx = \left[\frac{x^2}{500} \right]_{25}^{30} = \frac{900}{500} - \frac{625}{500} = 0.55$$

$$b) \mu_x = \int_{20}^{30} x \cdot \frac{x}{250} dx = \left[\frac{x^3}{750} \right]_{20}^{30} = \frac{2700}{750} - \frac{800}{750} = 25.33$$

$$c) \sigma^2 = \int_{20}^{30} x^2 \cdot \frac{x}{250} dx - \mu^2 = \left[\frac{x^4}{1000} \right]_{20}^{30} - \mu^2 = 8.22$$

$$d) \sigma = \sqrt{\sigma^2} = 2.8674$$

$$e) F(x) = \int_{-\infty}^x f(t) dt. \quad \text{If } x < 20: F(x) = \int_{-\infty}^x 0 dt = 0$$

$$\text{If } 20 \leq x \leq 30: F(x) = \int_{20}^x \frac{t}{250} dt = \frac{x^2}{500} - \frac{400}{500}$$

$$\text{If } x \geq 30: F(x) = \int_{-\infty}^{20} 0 dt + \int_{20}^{30} \frac{t}{250} dt + \int_{30}^x 0 dt = 1$$

$$f) P(x > 28) = 1 - F(28) = 1 - 0.768 = 0.232.$$

2.4.24 Particles are a major component of air pollution in many areas. It is of interest to study the sizes of contaminating particles. Let X represent the diameter, in micrometers, of a randomly chosen particle. Assume that in a certain area, the probability density function of X is inversely proportional to the volume of the particle; that is, assume that (in function in image below) c is a constant.

- Find the value of c so that $f(x)$ is a probability density function.
- Find the mean particle diameter.
- Find the cumulative distribution function of the particle diameter.
- Find the median particle diameter.
- The term PM_{10} refers to particles 10 μm or less in diameter. What proportion of the contaminating particles are PM_{10} ?
- The term PM_{25} refers to particles 2.5 μm or less in diameter. What proportion of the contaminating particles are PM_{25} ?
- What proportion of the PM_{10} particles are PM_{25} ?

2.4.24 $X = \text{diameter of random particle}$

$$f(x) = \begin{cases} \frac{c}{x^3}, & x \geq 1 \\ 0, & x < 1 \end{cases}$$

a) val of c so $f(x)$ is a probability density function

$$\int_1^\infty \frac{c}{x^3} dx = 1 \Rightarrow \left[-\frac{c}{2x^2} \right]_1^\infty = 1 \Rightarrow 0 - \frac{-c}{2(1)^2} = 1 \Rightarrow c = 2$$

b) mean particle diameter

$$\mu_x = \int_{-\infty}^\infty x f(x) dx = \int_1^\infty \frac{2x}{x^3} dx = \int_1^\infty \frac{2}{x^2} dx = \left[-\frac{2}{x} \right]_1^\infty = 2$$

→ c) cumulative distribution function of particle diameter

$$x < 1: F(x) = \int_{-\infty}^x f(t) dt = \int_{-\infty}^x 0 dt = 0$$

$$x \geq 1: F(x) = \int_{-\infty}^x f(t) dt = \int_{-\infty}^1 0 dt + \int_1^x \frac{2}{t^3} dt = 0 + \left[-\frac{2}{t^2} \right]_1^x = 1 \dots$$

d) median particle diameter

$$\int_{-\infty}^{x_m} f(x) dx = 0.5 \Rightarrow \int_{-\infty}^{x_m} \frac{2}{x^3} dx = 0.5 \Rightarrow \left| -\frac{2}{x^2} \right|_{-\infty}^{x_m} = 0.5 \Rightarrow \left(-\frac{2}{x_m} - 0 \right) = 0.5$$

$$\Rightarrow \text{according to online integral calc} \Rightarrow x_m = 1.4$$

$$e) P(X=10) = 1 - \frac{1}{10^2} \Rightarrow P = 1 - \frac{1}{100} = 0.99$$

$$f) P(X=2.5) = 1 - \frac{1}{2.5^2} = 0.84$$

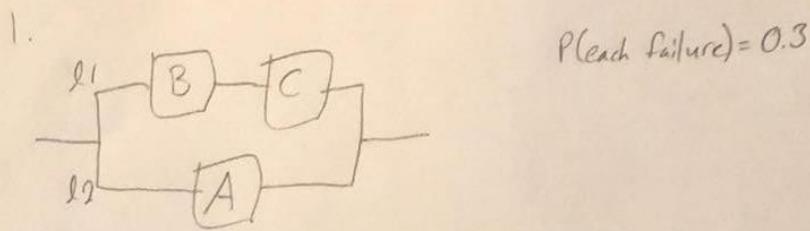
→ g)

3.3 Problem Set

1. The system shown in the figure consists of three elements A, B, C, which work independently. Each of them has a probability of failing $p=0.3$. Calculate:
 - a) The probability that all elements work.
 - b) The probability that the system works.
2. Calculate the standard deviation of the random variable n represented in the table.

n	0	1	2	3
$p(n)$	0.4	0.3	0.2	0.1

3. Consider a random variable $x \in (0, \infty)$ with the probability density function $f(x) = e^{-x}$.
 - a) Find the cumulative distribution function
 - b) Find the probability that $x > 1$.



a) $P(B) = P(C) = P(A) = 1 - P(\text{Fail}) = \underline{\underline{0.7}}$

b) $P(\text{vict}) = P(l_1 \cup l_2) = P((C \cap B) \cup (A))$
 $= P(A \cap B) + P(A) - P((A \cap B) \cap (A))$

$P(C \cap B) = P(C)P(B) = 0.7 \cdot 0.7 = 0.49$

$P(A) = 0.7$

$P(\text{vict}) = 0.49 + 0.7 - P(C \cap B)P(A) = 0.49 + 0.7 - 0.49 \cdot 0.7$
 $= \underline{\underline{0.847}}. \quad \square$

2.

n	0	1	2	3
$P(n)$	0.4	0.3	0.2	0.1

probability density for $\hat{X}(x)$

i) $\mu_x = \sum x_i P(x_i) = 0 \cdot 0.4 + 1 \cdot 0.3 + 2 \cdot 0.2 + 3 \cdot 0.1 = 0.3 + 0.4 + 0.3 = 1$

ii) $\sigma_x^2 = \sum x_i^2 P(x_i) - \mu_x^2 = 0^2 \cdot 0.4 + 1^2 \cdot 0.3 + 2^2 \cdot 0.2 + 3^2 \cdot 0.1 - 1^2$
 $= 0 + 0.3 + 0.8 + 0.9 - 1 = 2 - 1 = 1$

$\therefore \text{the } sd = \sigma_x = \sqrt{1} = \underline{\underline{1}}. \quad \square$

3. $x \in (0; \infty)$; probability density function $f(x) = e^{-x}$.

a) $F(x) = \int_0^x f(t) dt = \int_0^x e^{-t} dt = \left[-e^{-t} \right]_0^x = \underline{\underline{1 - e^{-x}}}$

b) $P(x > 1) = \int_1^\infty x f(t) dt = \int_1^\infty t e^{-t} dt = \left[-t e^{-t} \right]_1^\infty + \int_1^\infty e^{-t} dt = 1 - \int_1^\infty e^{-t} dt = 1 - \left[-e^{-t} \right]_1^\infty = 1 - 1 + 2e^{-1} = \underline{\underline{2e^{-1}}}. \quad \square$

3.4 Problem set Solution

1. The system shown in the figure consists of three elements A, B, C, which work independently. Each of them has a probability of failing $p=0.3$. Calculate:

a) The probability that all elements work.

$P(\text{one element works})=1-p$. $P(\text{all elements work})= P(A \cap B \cap C) = P(A)P(B)P(C)=(1-p)^3=0.343$ Where $P(A)$ means $P(A \text{ works})$

b) The probability that the system works.

$P(A \cup (B \cap C)) = P(A) + P(B \cap C) - P(A \cap B \cap C) = (1-p) + (1-p)^2 - (1-p)^3 = 0.847$. The three versions differ only by a permutation of elements, but the numerical results are identical.

2. Calculate the standard deviation of the random variable n represented in the table.

n	0	1	2	3
Version A $p(n)$	0.4	0.3	0.2	0.1
Version B $p(n)$	0.1	0.2	0.3	0.4
Version C $p(n)$	0.2	0.3	0.1	0.4

$$\sigma^2 = \sum_i n^2 p(n) - \left[\sum_i np(n) \right]^2 = 2 - 1 \text{ (A)} \text{ or } 5 - 2^2 \text{ (B)} \text{ or } 4.3 - 1.7^2 \text{ (C)} = 1 \text{ (A, B)} \text{ or } 1.41 \text{ (C)}$$

$$\sigma = 1 \text{ (A, B)} \text{ or } 1.187 \text{ (C)}$$

3. Consider a random variable $x \in (0, \infty)$ with the probability density function $f(x)=e^{-x}$.

a) Find the cumulative distribution function

$$F(x) = \int_0^x f(x') dx' = 1 - e^{-x}$$

b) Find the probability that $x > 1$.

$$P(x > 1) = 1 - P(x < 1) = 1 - F(1) = e^{-1} \approx 0.368$$

4 Week 2

4.1 Concepts

- You measure a certain physical quantity several times and you calculate the mean value. The uncertainty of the mean value is smaller than the uncertainty of each measured value
- The bias of measurements can be reduced by reducing the systematic errors.
- Consider a cube with side equal to L. If the uncertainty of L is 0.3% the uncertainty of the volume is 0.9%
- Consider a measured variable x and the function $U(x) = c\sqrt{x}$ where c is a constant number. The relative uncertainty of U is two times smaller than the relative uncertainty of x.
- The uncertainty of a function U of two measured variables x and y can be calculated with the formula based on partial derivatives if x and y are independent.

4.2 Example Problem Set

Linear functions of random variables

1.2.8 In a certain company, every worker received a \$50-per-week raise. How does this affect the mean salary? The standard deviation of the salaries?

1.2.8 \$50 per week raise affects mean and sd how?	
i) mean:	$\bar{x} = \frac{\sum x_i}{\# \text{empl}} \rightarrow \bar{x}' = \frac{\sum (x_i + 50)}{\# \text{empl}} = \frac{\sum (x_i) + 50}{\# \text{empl}} \Rightarrow \text{increase by } 50$
ii) sd:	$\sigma = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} \Rightarrow \sigma' = \sqrt{\frac{1}{n-1} \sum ((x_i + 50) - (\bar{x} + 50))^2} = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x} + 50 - 50)^2} = \sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2} = \sigma$ stays the same

1.2.9 In another company, every worker received a 5% raise. How does this affect the mean salary? The standard deviation of the salaries?

1.2.9 5% increase affects mean and sd how?	
i) mean:	$\bar{x}' = \frac{\sum 1.05 x_i}{\# \text{empl}} = \frac{0.05 \sum x_i}{\# \text{empl}} + \bar{x} = 0.05 \bar{x} + \bar{x} \therefore \text{mean increases by } 5\%$
ii) sd:	$\sigma' = \sqrt{\frac{1}{n-1} \sum (0.05 x_i - 0.05 \bar{x})^2} = \sqrt{\frac{1}{n-1} \cdot (0.05)^2 \sum (x_i - \bar{x})^2} = 0.05 \sigma \therefore \text{sd increases by } 5\%$

2.5.10 A gas station earns \$2.60 in revenue for each gallon of regular gas it sells, \$2.75 for each gallon of midgrade gas, and \$2.90 for each gallon of premium gas. Let X_1 , X_2 , and X_3 denote the numbers of gallons of regular, midgrade, and premium gasoline sold in a day. Assume that X_1 , X_2 , and X_3 have means $\mu_1 = 1500$, $\mu_2 = 500$, and $\mu_3 = 300$, and standard deviations $\sigma_1 = 180$, $\sigma_2 = 90$, and $\sigma_3 = 40$, respectively.

$2.5.10$ $\$2.60$ for reg, $\$2.75$ for mid, $\$2.90$ for premium $X_1, X_2, X_3 \rightarrow \# \text{ of gallons sold of each}$ $M_1 = 1500, M_2 = 500, M_3 = 300 \rightarrow \text{mean of each}; \sigma_1 = 180, \sigma_2 = 90, \sigma_3 = 40 \rightarrow \text{sd of each}$ a) mean daily revenue daily revenue is $2.60X_1 + 2.75X_2 + 2.90X_3$ $M_{\text{tot}} = C_1M_1 + C_2M_2 + C_3M_3 = 2.60 \cdot 1500 + 2.75 \cdot 500 + 2.90 \cdot 300 = \6145
b) X_1, X_2, X_3 are independent, what is sd of daily revenue? $\sigma = \sqrt{180^2 + 90^2 + 40^2} = \sqrt{32400 + 8100 + 1600} = \sqrt{42100} \approx 205.18$

3.2.6 A cylindrical hole is bored through a steel block, and a cylindrical piston is machined to fit into the hole. The diameter of the hole is 20.00 ± 0.01 cm, and the diameter of the piston is 19.90 ± 0.02 cm. The clearance is one-half the difference between the diameters. Estimate the clearance and find the uncertainty in the estimate.

$3.2.6$ $d_{\text{hole}} = 20.00 \pm 0.01$ $d_{\text{piston}} = 19.90 \pm 0.02$ $c_l = \frac{1}{2}(d_{\text{hole}} - d_{\text{piston}}) = \frac{1}{2}(20.00 - 19.90) = 0.05 = \text{estimate of clearance}$ uncertainty: $\sigma = \sqrt{\left(\frac{1}{2}0.01\right)^2 + \left(\frac{1}{2}0.02\right)^2} = \sqrt{0.000025 + 0.0001} \approx 0.0112 \approx 0.01$

Uncertainty of sample mean

3.2.3 The length of a rod is to be measured by a process whose uncertainty is 3 mm. Several independent measurements will be taken, and the average of these measurements will be used to estimate the length of the rod. How many measurements must be made so that the uncertainty in the average will be 1 mm?

$\text{Uncertainty of sample mean}$ $3.2.3$ uncertainty = 3 mm, # of samples to make uncertainty, $\delta = 1$ mm $\delta_x = \frac{\sigma}{\sqrt{n}} \rightarrow \sqrt{n} = \frac{\sigma}{\delta_x} \Rightarrow n = \left(\frac{\sigma}{\delta_x}\right)^2 = \left(\frac{3}{1}\right)^2 = 9 \text{ samples}$

Bias and uncertainty

3.2.14 The volume of a rock is measured by placing the rock in a graduated cylinder partially filled with water and measuring the increase in volume. Eight independent measurements are made. The average of the measurements is 87.0 mL, and the standard deviation is 2.0 mL.

- Estimate the volume of the rock, and find the uncertainty in the estimate.
- Eight additional measurements are made, for a total of 16. What is the uncertainty, approximately, in the average of the 16 measurements?
- Approximately how many measurements would be needed to reduce the uncertainty to 0.4 mL?

*→ null still adj taken
→ shows 2 grams more than it should*

Bias and uncertainty

3.2.14 uncertainty 3g and bias 2g on scale

on average

- a) single measurement: uncertainty = 3g , bias = 2g
- b) 4 indep. meas. : $\sigma_x = \frac{\sigma}{\sqrt{n}} = \frac{3}{\sqrt{4}} = 1.5g$, bias = 2g (unchanged)
- c) 400 indep. meas. : $\sigma_x = \frac{3}{\sqrt{400}} = 0.15g$, bias = 2g (still unchanged)
- d) more meas. ~~does~~ uncertainty decreases → makes sense b/c more data
- e) more meas. bias stays the same no matter how many measurements

3.2.15 A certain scale has an uncertainty of 3 g and a bias of 2 g.

- a. A single measurement is made on this scale. What are the bias and uncertainty in this measurement?
- b. Four independent measurements are made on this scale. What are the bias and uncertainty in the average of these measurements?
- c. Four hundred independent measurements are made on this scale. What are the bias and uncertainty in the average of these measurements?
- d. As more measurements are made, does the uncertainty get smaller, get larger, or stay the same?
- e. As more measurements are made, does the bias get smaller, get larger, or stay the same?

3.2.15 V of rock measured by putting it in grad. cyl. of water
and measuring increase in volume. 8 indep. meas. taken

$$\bar{x} = 87.0 \text{ mL} \quad \text{sd} = 2.0 \text{ mL}$$

a) estimate V of rock, $n=8$, $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{2.0}{\sqrt{8}} \approx 0.7 \rightarrow 87.0 \pm 0.7 \text{ mL}$

b) 8 more meas., uncertainty: $\sigma'_{\bar{x}} = \frac{\sigma}{\sqrt{n_2}} = \frac{2.0}{\sqrt{16}} = 0.5$

c) \sim # meas. to reduce uncert. to 0.4 mL: $0.4 = \frac{2.0}{\sqrt{n}} \Rightarrow n = \left(\frac{2.0}{0.4}\right)^2 = 25$

Propagation of errors

3.3.2 Given that X and Y are related by the given equation, and that $X = 3.0 \pm 0.1$, estimate Y and its uncertainty.

- a. $XY = 1$
- b. $Y/x = 2$
- c. $\sqrt{XY} = 3$
- d. $Y\sqrt{X} = 4$

$$3.3.2 \quad X = 3.0 \pm 0.1$$

$$a) XY = 1 \Rightarrow Y = \frac{1}{X} = 0.3 \pm 0.1$$

$$\frac{dy}{dx} = \frac{-1}{x^2} \quad \text{using } \frac{d}{dx} \frac{1}{x} = -\frac{1}{x^2} \quad \sigma = \sqrt{\left(\frac{-1}{3^2}\right)^2 (0.1)^2} = 0.01$$

$$b) \frac{Y}{X} = 2 \Rightarrow Y = 2X = 6 \pm 0.2$$

$$\frac{dy}{dx} = 2 \quad \sigma = \sqrt{2^2 (0.1)^2} = 0.2$$

$$c) \sqrt{XY} = 3 \Rightarrow Y = 9X = 27$$

$$\frac{dy}{dx} = 9 \cdot 1 \quad \sigma = \sqrt{9^2 \cdot 0.1^2} = 0.9$$

$$d) Y\sqrt{X} = 4 \Rightarrow Y = \frac{4}{\sqrt{X}} = 2.31 \pm 0.0385$$

$$\frac{dy}{dx} = \frac{-4}{2x^{3/2}} = -0.385 \quad \sigma = \sqrt{(-0.385)^2 (0.1)^2} = 0.0385$$

3.3.3 The volume of a cone is given by $V = \pi r_r h / 3$, where r is the radius of the base and h is the height. Assume the height is 6 cm, measured with negligible uncertainty, and the radius is $r = 5.00 \pm 0.02$ cm. Estimate the volume of the cone, and find the uncertainty in the estimate.

Propagation of errors

$$3.3.2 \quad V_{\text{cone}} = \frac{\pi r^2 h}{3}, \quad h = 6 \text{ cm}, \quad r = 5.00 \pm 0.02 \text{ cm}$$

$$V_{\text{cone}} = \frac{\pi 5.00^2 \cdot 6}{3} = 157.08 = 157 \text{ cm}^3 \quad \frac{dV}{dr} = \frac{2\pi r h}{3} = \frac{10 \cdot 6 \cdot \pi}{3} = 20\pi$$

$$\text{uncertainty} = \sigma = \left| \frac{dV}{dr} \sigma_r \right| = |20\pi \cdot 0.02| = 0.4\pi = 1.25$$

$$V = 157.1 \pm 1.3 \text{ cm}^3$$

3.4.1 Find the uncertainty in U, assuming that $X = 10.0 \pm 0.5$, $Y = 5.0 \pm 0.1$, and

- $U = XY^2$
- $U = X^2 + Y^2$
- $U = (X + Y^2)/2$

	<p>3.4.1 $X = 10.0 \pm 0.5$ $Y = 5.0 \pm 0.1$</p>
a)	$U = XY^2$ $\frac{\partial U}{\partial X} = Y^2 = 25.0$ $\frac{\partial U}{\partial Y} = 2XY = 2 \cdot 10 \cdot 5 = 100.0$ $\sigma_U = \sqrt{\left(\frac{\partial U}{\partial X}\right)^2 \sigma_x^2 + \left(\frac{\partial U}{\partial Y}\right)^2 \sigma_y^2} = \sqrt{25^2 \cdot 0.5^2 + 100^2 \cdot 0.1^2} = \sqrt{156.25 + 100} = 16$ $U = (10.0 \pm 16) = 250 \pm 16$
b)	$U = X^2 + Y^2$ $\frac{\partial U}{\partial X} = 2X = 20.0$ $\frac{\partial U}{\partial Y} = 2Y = 10.0$ $\sigma = \sqrt{\left(\frac{\partial U}{\partial X}\right)^2 \sigma_x^2 + \left(\frac{\partial U}{\partial Y}\right)^2 \sigma_y^2} = \sqrt{20^2 \cdot 0.5^2 + 10^2 \cdot 0.1^2} = \sqrt{100 + 1} \approx 10.1$ $U = (10^2 + 5^2) \pm 10.1 = 125.0 \pm 10.1$
c)	$U = \frac{x+y^2}{2}$ $\frac{\partial U}{\partial X} = \frac{1}{2}$ $\frac{\partial U}{\partial Y} = Y = 5.0$ $\sigma = \sqrt{\left(\frac{\partial U}{\partial X}\right)^2 \sigma_x^2 + \left(\frac{\partial U}{\partial Y}\right)^2 \sigma_y^2} = \sqrt{\left(\frac{1}{2}\right)^2 \cdot 0.5^2 + 5^2 \cdot 0.1^2} = \sqrt{0.0625 + 0.25} = 0.56$ $U = \left(\frac{10.0 + 5.0^2}{2}\right) \pm 0.56 = 17.5 \pm 0.6$

3.4.3 From a fixed point on the ground, the distance to a certain tree is measured to be $s = 55.2 \pm 0.1$ m and the angle from the point to the top of the tree is measured to be $\theta = 0.50 \pm 0.02$ radians. The height of the tree is given by $h = s \tan \theta$.

	<p>3.4.3 $s = 55.2 \pm 0.1$ m</p>
	$\theta = 0.50 \pm 0.02$ radians
a)	$\frac{\partial h}{\partial s} = \tan \theta = 0.546$ $\frac{\partial h}{\partial \theta} = \frac{s}{\cos^2(\theta)} = 55.2$ $\sigma = \sqrt{0.546^2 \cdot 0.1^2 + 55.2^2 \cdot 0.02^2} = \sqrt{0.00298 + 1.219} = 1.105$ $h = 55.2 \cdot \tan(0.5) \pm 1.105 = 30.16 \pm 1.11$
b)	<p>greater reduction of uncertainty: θ uncert. reduced to 0.01 radians b/c the number it is multiplied by in σ calc. is bigger (55.2)</p>

Relative uncertainty

Prove formula 3.14 at page 192 using the formula 3.12 at page 186 with 3 variables X, Y, Z. Use the relative uncertainty in problems 3.3.5 page 184 and 3.4.2 page 193.

$$\text{Formula 3.14: } \frac{\sigma_u}{u} = \sqrt{\left(m_1 \frac{\sigma_{x_1}}{x_1}\right)^2 + \dots + \left(m_n \frac{\sigma_{x_n}}{x_n}\right)^2}$$

$$\text{Formula 3.12: } \sigma_u \approx \sqrt{\left(\frac{\partial u}{\partial x_1}\right)^2 \sigma_{x_1}^2 + \dots + \left(\frac{\partial u}{\partial x_n}\right)^2 \sigma_{x_n}^2}$$

don't know where to go with this, can't take ∂u out of square root because it is a differ...

4.3 Problem Set

Verkefni 4

1. $x \in (0,1)$ (RV) $\mu = 0.5$ $\sigma_x = \frac{1}{\sqrt{12}}$ $y \in (a,b) \rightarrow y = a + (b-a)x$
- a) mean of y : we multiply and get $\mu_y = (b-a)\mu_x + a = a + (b-a)\mu_x$
 $\Rightarrow \mu_y = a + 0.5(b-a)$ (given by rules defined by Linear Comb. of RVs).
- b) sd of y : we get by those same rules used in (a) that σ_y is:
 $\sigma_y^2 = (b-a)^2 \sigma_x^2 \Rightarrow \sigma_y = (b-a) \sigma_x \Rightarrow \sigma_y = (b-a) \frac{1}{\sqrt{12}} \Rightarrow \sigma_y = \frac{b-a}{\sqrt{12}}$. \square
2. Same RV as in (1) but a sample of it of size n . How large does n have to be for the uncertainty to be ± 0.01 ?
We obtain uncertainty of sample size by using the equation $\sigma_x' = \frac{\sigma_x}{\sqrt{n}}$.
In order to get $\sigma_x' = \pm 0.01$ then: $n = \left(\frac{\sigma_x}{\sigma_x'}\right)^2 = \left(\frac{\frac{1}{\sqrt{12}}}{0.01}\right)^2 = \frac{50\sqrt{3}}{3} \approx 28.86$
- Therefore the sample size n has to be at least 29 (to be sure the uncertainty is around ± 0.01). \square
3. Cyl. silver wire w/ $D = 4.0 \pm 0.1$ mm
resistance: $R = \rho \frac{L}{A}$, $\rho = 1.59 \times 10^{-8} \Omega \text{m}$, $L = 200 \pm 1$ → length of wire
 A is cross sectional area.
- a) relative uncertainty of A
To calculate A we do $A = \frac{D}{2}\pi = \frac{4.0}{2}\pi = 2\pi$. Now let's find the uncertainty.
We find that $\sigma_A = \left| \frac{\partial A}{\partial D} \right| \sigma_D = \left| \frac{1}{2}\pi \right| \sigma_D = \frac{1}{2}\pi \cdot 0.1 = 0.157 \approx 0.2$.
- b) relative uncertainty of resistance R .
To calculate R we use the given equation $R = \rho \frac{L}{A}$.
We find that $\sigma_R = \sqrt{\left(\left| \frac{\partial R}{\partial A} \right| \sigma_A \right)^2 + \left(\left| \frac{\partial R}{\partial L} \right| \sigma_L \right)^2} = \sqrt{\left(-\frac{\rho L}{A^2} \cdot \sigma_A \right)^2 + \left(\frac{\rho}{A} \sigma_L \right)^2}$
 $\Rightarrow \sigma_R = \sqrt{\left(1.59 \times 10^{-8} \cdot \frac{200}{(2\pi)^2} \cdot \frac{1}{20} \right)^2 + \left(\frac{1.59 \times 10^{-8}}{2\pi} \cdot 1 \right)^2}$
 $\Rightarrow \sigma_R = \sqrt{1.6 \times 10^{-16} + 6.4 \times 10^{-18}} = 1.289 \times 10^{-8} \approx 1.3 \times 10^{-8}$. \square

4.4 Problem Set Solutions

1. A random variable $x \in (0,1)$ with uniform distribution has mean value $\mu=0.5$ and standard deviation $\sigma=1/\sqrt{12}$. If we need a random variable y with uniform distribution in another interval, $y \in (a,b)$, we can obtain it with the linear transformation $y=a+(b-a)x$. Use the properties of the linear combinations of random variables to calculate

a) the mean of y : $E(y)=a+(b-a)E(x)=a+(b-a)0.5=(b+a)/2$

b) the standard deviation of y : $\text{Var}(y)=(b-a)^2 \text{Var}(x)=(b-a)^2/12 \Rightarrow \sigma_y=(b-a)/\sqrt{12}$

NOTE: Do not use integrals.

2. Consider a sample of the previous random variable x of size n . How large should be n such that the sample mean has an uncertainty ± 0.01 ?

The uncertainty of the sample mean is $\sigma_{\bar{x}} = \frac{\sigma}{n} = 0.01 \Rightarrow n = (\frac{\sigma}{0.01})^2 = \frac{1}{12*10^{-4}} = 833$

3. Consider a cylindrical Silver wire with diameter $D=4.0 \pm 0.1 \text{ mm}$. The electrical resistance of the wire is $R=\Omega$ where $\rho=1.59 \times 10^8 \Omega \text{ m}$ (ohm meter) is the resistivity of the material, $L=200 \pm 1 \text{ mm}$ is the length of the wire, and A is the cross sectional area.

a) relative uncertainty of A : $A = \frac{\pi D^2}{4} \Rightarrow \frac{\sigma_A}{A} = 2 \frac{\sigma_D}{D} = 2 \frac{0.1}{4} = 0.05$

b) relative uncertainty of R : $R = \rho \frac{L}{A} \Rightarrow \frac{\sigma_R}{R} = \sqrt{\left(\frac{\sigma_L}{L}\right)^2 + \left(\frac{\sigma_A}{A}\right)^2} = \sqrt{\left(\frac{1}{200}\right)^2 + (0.05)^2} = 0.05025$

Acceptable result for any given L (200, 250, or 300 nm): $\frac{\sigma_R}{R} \approx 0.05$

COMMENTS: The reason to ask for relative uncertainties, $\frac{\sigma_A}{A}$ and $\frac{\sigma_R}{R}$, is that in this example they are much easier to calculate than the absolute uncertainties σ_A and σ_R . We see that only a couple of ratios are needed, which can be obtained in a couple of minutes. However, many students still tried to do the full math, calculated all derivatives from scratch, or did other unnecessary calculations, and wasted their time. Other students did not show the relative uncertainties, like they did not know what they are. This concept has been very thoroughly covered in videos 4 and 5 of the Lecture Package 1 and in the last group of problems of week 4, but also in the data analyses of the physics labs already done by most of the students. That is why the course material must be studied before going over the assignment, and if so the 60 min time is more than enough for solutions and for the upload.

(Problems similar to 1 and 2 were also covered over the week 4.)

5 Week 2

5.1 Concepts

- 56 teams of 5 can be formed out of 8 individuals
- A random variable x has a binomial distribution with $n=8$ and $p=0.8$. The probability that $x=5$ is 0.147.
- The probability that x is less or equal to 5 if $x \sim Bin(8, 0.8)$ is 0.203
- The mean value of $x \sim Bin(n, p)$ is np
- The variance of $x \sim Bin(n, p)$ is $np(1-p)$
- The temperature in an office space is expected to be $T = 22$ Celsius degrees with an (absolute) uncertainty $S = 1$ Celsius degree, such that the temperature is reported as $T \pm S$. The relative uncertainty is S/T
- In the Poisson distribution $P(k, \lambda)$ the meaning of the parameter lambda (λ) is the mean value of k
- If the mean value of a Poisson random variable k is 15 then the probability that $k=15$ is 0.102
- Hvað veldur stórum stökkum á rafmagnsfjármálamörkuðum? Misræmi í framboði og eftirspurn, og að það að erfitt er að geyma rafmagn sem er búið að framleiða.
- Hvernig þarf að breyta Poisson ferlinu til þess að það geti orðið gott líkan fyrir þyrringamyndun á rafmagnsmörkuðum? Gera tíðnina slembna, þannig að þegar stórt stökk kemur þá eykst tíðnin tímabundið, en er þó miðsækin til lengri tíma litið.

Estimate proportion and uncertainty:

$$\sigma_{p_A} = \sqrt{\frac{p_A(1-p_A)}{n}} \text{ where } p_A = \frac{x}{n} \text{ and } n \text{ is the total number of things in the sample.}$$

5.2 R code

```
dbinom(2, 4, 0.6)          # P(X=2), Bin(4, 0.6)
1-pbinom(2, 8, 0.2)        # P(X>2), Bin(8, 0.2)
pbinom(2, 5, 0.4)          # P(X<=2), Bin(5, 0.4)
sum(dbinom(3:5, 6, 0.7))   # P(3<=X<=5), Bin(6, 0.7)
num = runif(10)             # generate 10 random numbers
```

5.3 Example Problem Set

Computer generated random numbers

Week 5										
LP2 - random var generator										
1) $x \in (0, 1)$	$\mu = 0.5$	$\sigma = \frac{1}{\sqrt{12}}$	$f(x) = \begin{cases} \frac{1}{b-a} = 1, & 0 < x < 1 \\ 0, & \text{otherwise} \end{cases}$							
???										
2) num = runif(10) / 1000	1000	10000								
$\bar{x} = 0.51$	$\sigma = 0.30$	$\bar{x} = 0.47$	$\sigma = 0.28$	$\bar{x} = 0.50$	$\sigma = 0.29$	$\bar{x} = 0.49$	$\sigma = 0.29$			
100000	1000000									
$\bar{x} = 0.50$	$\sigma = 0.29$	$\bar{x} = 0.50$	$\sigma = 0.288$							
				$\bar{x} \xrightarrow{n \rightarrow \infty} \mu$	$s \xrightarrow{n \rightarrow \infty} \sigma$					
3)?										$n = \text{num of rand numbers}$
4) 10: $\bar{x} = 0.51$	0.56	0.34	0.47	0.53	0.65	0.39	0.59	0.45	0.44	
$\sigma = 0.30$	0.33	0.30	0.25	0.28	0.30	0.24	0.27	0.25	0.28	
100: $\bar{x} = 0.47$	0.53	0.48	0.52	0.54	0.51	0.45	0.46	0.45	0.50	
$\sigma = 0.30$	0.28	0.28	0.30	0.31	0.27	0.30	0.28	0.26	0.29	
1000: $\bar{x} = 0.52$	0.49	0.51	0.49	0.52	0.51	0.51	0.49	0.51	0.50	
$\sigma = 0.29$	0.29	0.29	0.29	0.28	0.28	0.28	0.29	0.29	0.29	
10000: $\bar{x} = 0.50$	0.50	0.50	0.49	0.49	0.49	0.50	0.49	0.50	0.49	
$\sigma = 0.28$	0.29	0.28	0.28	0.28	0.28	0.29	0.29	0.28	0.29	
100000: $\bar{x} = 0.49$	0.50	0.49	0.49	0.49	0.49	0.5	0.50	0.49	0.49	
$\sigma = 0.28$	0.28	0.28	0.28	0.28	0.28	0.28	0.29	0.28	0.28	
1000000: $\bar{x} = 0.50$	0.49	0.50	0.50	0.49	0.49	0.50	0.49	0.50	0.49	
$\sigma = 0.28$	0.28	0.28	0.28	0.28	0.28	0.28	0.28	0.28	0.28	
5)?										

Counting combinations

- 2.2.4** A group of 18 people have gotten together to play baseball. They will divide themselves into two teams of 9 players each, with one team wearing green uniforms and the other wearing yellow uniforms. In how many ways can this be done?

2.2.4

18 people into two teams of 9 for baseball game

$$18! / 9!9! = 48620$$

- 2.2.6** A college math department consisting of 10 faculty members must choose a department head, an assistant department head, and a faculty senate representative. In how many ways can this be done?

2.2.6

10 faculty choose 1 head, 1 assis., 1 rep

$$10 \cdot 9 \cdot 8 = 720$$

- 2.2.10** A company has hired 15 new employees, and must assign 6 to the day shift, 5 to the graveyard shift, and 4 to the night shift. In how many ways can the assignment be made?

2.2.10

15 emp, 6 day shift, 5 grave shift, 4 night shift

$$15! / 6!5!4! = 630,630$$

Bernoulli (or binary) variables

4.1.6 Two dice are rolled. Let $X = 1$ if the dice come up doubles and let $X = 0$ otherwise. Let $Y = 1$ if the sum is 6, and let $Y = 0$ otherwise. Let $Z = 1$ if the dice come up both doubles and with a sum of 6 (that is, double 3), and let $Z = 0$ otherwise.

- Let p_X denote the success probability for X . Find p_X .
- Let p_Y denote the success probability for Y . Find p_Y .
- Let p_Z denote the success probability for Z . Find p_Z .
- Are X and Y independent?
- Does $p_Z = p_X p_Y$?
- Does $Z = XY$? Explain.

$2 \text{ dice rolled } X=1 \text{ if double, } X=0 \text{ otherwise}$		$\rightarrow Z=0 \text{ otherwise}$
$Y=1 \text{ if sum}=6, Y=0 \text{ otherwise}$		$Z=1 \text{ both doubles + sum}=6 \text{ (double 3)}$
a)	$p_X = \text{success prob for } X: \frac{6}{36} = \frac{1}{6}$	
b)	$p_Y = \text{success prob for } Y: \frac{5}{36}$	
c)	$p_Z = \text{success prob for } Z = \frac{1}{36}$	
d)	$\text{are } X \text{ and } Y \text{ indep? } \rightarrow \text{no}$	
e)	$p_Z = p_X p_Y? \text{ no } p_X p_Y = \frac{5}{216}$	
f)	$Z = XY? \rightarrow \text{yes but } p(Z) \neq p(X)p(Y)$	

Binomial distribution:

Example 4.7 Find the probability mass function of the random variable X if $X \sim \text{Bin}(10, 0.4)$. Find $P(X = 5)$.

Ex. 4.7	$\text{prob mass funct. of RV } X \text{ if } X \sim \text{Bin}(10, 0.4). \text{ Find } P(X=5)$	
	$p(x) = \begin{cases} \frac{10!}{x!(10-x)!} (0.4)^x (0.6)^{10-x}, & x \in \{1, 10\} \text{ whole number} \\ 0 & \text{else} \end{cases} \quad \text{prob mass funct.}$	
	$P(X=5) = p(5) = \frac{10!}{5!(10-5)!} (0.4)^5 (0.6)^{10-5} = 0.2007$	

Example 4.8 A fair die is rolled eight times. Find the probability that no more than 2 sixes come up.

Ex 4.8

die rolled 8 times, find prob of no more than 2 6's coming up
 each die is Bernoulli trial w/ success prob of $\frac{1}{6}$. X denotes # of
 sixes in 8 rolls. Find $P(X \leq 2) = P(X=0) + P(X=1) + P(X=2)$
 $pbinom(2, 8, \frac{1}{6}) = 0.865$

4.2.3 Find the following probabilities:

- $P(X = 2)$ when $X \sim \text{Bin}(4, 0.6)$
- $P(X > 2)$ when $X \sim \text{Bin}(8, 0.2)$
- $P(X \leq 2)$ when $X \sim \text{Bin}(5, 0.4)$
- $P(3 \leq X \leq 5)$ when $X \sim \text{Bin}(6, 0.7)$

a) $P(X = 2)$ $X \sim \text{Bin}(4, 0.6)$	$dbinom(2, 4, 0.6) = 0.3456$
b) $P(X > 2)$ $X \sim \text{Bin}(8, 0.2)$	$1 - pbinom(2, 8, 0.2) = 0.20308$
c) $P(X \leq 2)$ $X \sim \text{Bin}(5, 0.4)$	$pbinom(2, 5, 0.4) = 0.68256$
d) $P(3 \leq X \leq 5)$ $X \sim \text{Bin}(6, 0.7)$	$\sum(dbinom(3:5, 6, 0.7)) = 0.8118$

4.2.9 Several million lottery tickets are sold, and 60% of the tickets are held by women. Five winning tickets will be drawn at random.

- What is the probability that three or fewer of the winners will be women?
- What is the probability that three of the winners will be of one gender and two of the winners will be of the other gender?

4.2.9 sev. mill. lotto tickets sold, 60% held by women

5 winning tickets drawn at random

a) $P(3 \text{ or fewer women win})$ X is # of women among 5 winners
 $X \sim \text{Bin}(5, 0.6) = 0.6528$

b) 3 winners of one gender, 2 of another

$$P(X=2) + P(X=3) = dbinom(2, 5, 0.6) + dbinom(3, 5, 0.6) = 0.5760$$

4.2.11 A quality engineer samples 100 steel rods made on mill A and 150 rods made on mill B. Of the rods from mill A, 88 meet specifications, and of the rods from mill B, 135 meet specifications.

- Estimate the proportion of rods from mill A that meet specifications, and find the uncertainty in the estimate.
- Estimate the proportion of rods from mill B that meet specifications, and find the uncertainty in the estimate.
- Estimate the difference between the proportions, and find the uncertainty in the estimate.

4.2.11

100 parts ordered from A, 12 defective; 200 parts ord. from B, 10 defective

a) prop. of def. parts in A + uncertainty in estimate

$$\hat{p}_A = \frac{x}{100} = \frac{12}{100} = 0.12 \quad \sigma_{\hat{p}_A} = \sqrt{\frac{p_A(1-p_A)}{100}} = \sqrt{\frac{0.12(1-0.12)}{100}} = 0.032$$

b) same but for B

$$\hat{p}_B = \frac{x}{200} = \frac{10}{200} = 0.05 \quad \sigma_{\hat{p}_B} = \sqrt{\frac{0.05(1-0.05)}{200}} = 0.015$$

c) diff. in prop. + uncertainty in estimate

$$\hat{p}_A - \hat{p}_B = 0.07 ; \text{ uncertainty } \sqrt{\sigma_{\hat{p}_A}^2 + \sigma_{\hat{p}_B}^2} = 0.036$$

4.2.24 One design for a system requires the installation of two identical components. The system will work if at least one of the components works. An alternative design requires four of these components, and the system will work if at least two of the four components work. If the probability that a component works is 0.9, and if the components function independently, which design has the greater probability of functioning?

4.2.24

Sys will work if at least 1 of 2 comp work OR
at least 2 of 4 comp work

$P(\text{comp works}) = 0.9$; comp work independently

Which has a greater prob. of functioning?

Design w/ 4 comp

5.4 Problem Set

Verkefni 5

Litja Yr Gud

lifbag 18

1) $n=100$ rand num $(0, 1)$

a) `for(i in 1:5){`

`num = runif(100)`

`print(mean(num))`

}

Exp:	1	2	3	4	5
\bar{x} :	0.4959	0.4663	0.4670	0.4968	0.5205

`result <- numeric(5)`

`for(i in 1:5){`

`num = runif(100)`

`result[i] = mean(num)`

}

`sd(result)`

Exp:	1	2	3	4	5
\bar{x} :	0.5322	0.5183	0.4995	0.4886	0.4983

sd of exp $\bar{x} = 0.0175 = 0.018$

but the sd for each set of random numbers
is ~ 0.2919 (do same for loop as for mean)

b) The exact μ and σ for uniform distribution are $\mu=0.50$ and $\sigma=0.28$

In part (a) a larger sample size should be used to get closer to
the theoretical value (gets pretty accurate at 10000).

2) Prob of defective battery: $p=0.04$

a) Prob of no defective battery in sample of 20:

If $X = \#$ of def. batteries then $P(X=0) = dbinom(0, 20, 0.04) = 0.44$

b) Prob of at most 2 def. batteries in sample of 20:

$P(X \leq 2) = pbinom(2, 20, 0.04) = 0.956$

c) mean value and sd of defective batteries in sample of 100

mean or μ is acquired with $\mu = np = 100 \cdot 0.04 = 4$

sd or σ is acquired with $\sigma = \sqrt{npq} = \sqrt{\mu \cdot q} = \sqrt{4 \cdot q} = 2\sqrt{q}$ because it

depends on the number of defective batteries q

5.5 Problem Set Solution

1.a) Do this computer experiment: Generate with your favourite software a sample of $n=100$ random numbers with uniform distribution in the interval $(0,1)$ and obtain their sample mean \bar{X} . Call this Experiment 1 and put the value of \bar{X} in a small table, as shown below. Do five more such experiments and complete the table. Use at least three significant digits for each number to be able to compare them properly. Finally, calculate the standard deviation of these six values for \bar{X} seen as a set of six random numbers. Mention briefly the computer functions or commands used.

Experiment	1	2	3	4	5	6	st.dev of \bar{x}
\bar{x}	0.5010	0.5056	0.5291	0.4746	0.4387	0.5188	0.03306

The table shows an example of results which can be obtained with R or Excel using the models posted in Canvas.

b) What is the theoretically expected standard deviation (or uncertainty) of the sample mean of $n=100$ random numbers $x \in (0,1)$ with uniform distribution? (Remember the exact μ and σ for the uniform distribution.) What should be done at point a) to obtain the st.dev. of \bar{x} as close as possible to the theoretically expected value

The theoretical standard deviation: $\frac{\sigma}{\sqrt{n}} = \frac{1}{\sqrt{12}} \frac{1}{\sqrt{100}} = 0.0289$. See also Verkefni 4 problem 2 and Miniproject (Phase 1) point 5.

It is assumed that the student practiced such exercises before the assignment and understood the idea of randomness and the effects of the sample size. It is also assumed that the student distinguishes two types of samples in this problem: 6 samples of $n=100$ numbers with uniform distribution, and one sample of 6 numbers \bar{X} representing their sample means. In order to obtain the st.dev. of the \bar{x} 's as close as possible to 0.0289 one has to increase the sample of \bar{x} 's, or the number of "experiments", from 6 to a lot more. To how many more? The answer will be given in Week 7.

2. Consider a brand of electric batteries which are defective with probability $p = 0.04$. Calculate:

- a) The probability that there is no defective battery in a sample of 20. $dbinom(0, 20, p) = (1-p)^{20} = 0.442$ (or BINOM.DIST in Excel with FALSE option)
 - b) The probability that at most two batteries are defective in a sample of 20. $pbinom(2, 20, p) = 0.956$ (or BINOM.DIST in Excel with TRUE option)
 - c) The mean value and the standard deviation of the number of defective batteries in a sample of 100. $mean = np = 4$ $std = \sqrt{np(1-p)} = 1.96$
-

6 Week 2

6.1 Concepts

- The area under the normal probability distribution function with parameters mu and sigma is equal to 1
- When sigma increases the normal probability distribution function becomes wider
- Consider $x \sim N(0, 1)$. The probability that $x < -1$ is 16%
- Consider $x \sim N(0, 1)$. The probability that $-1 < x < 1$ is 68%
- Consider a sample of n random numbers x_1, x_2, \dots, x_n with a distribution D. If n is large enough their sample mean follows a normal distribution irrespective of the distribution D
- A Poisson distribution with a large parameter lambda can be approximated with a normal distribution having the standard deviation sqrt (lambda)
- The mean value of a random variable having a lognormal distribution with parameters $\mu = 0$ and sigma=1 is 1.65
- The distribution of the body mass index is lognormal
- Consider x a random variable having an exponential distribution with parameter lambda=3. The probability that x is less or equal to the mean value is 0.63
- Hvaða líkindadreifingar eru oftast notaðar í jarðvarmamati með Monte-Carlo hermun? Jöfn dreifing og þríhyrningslag dreifing

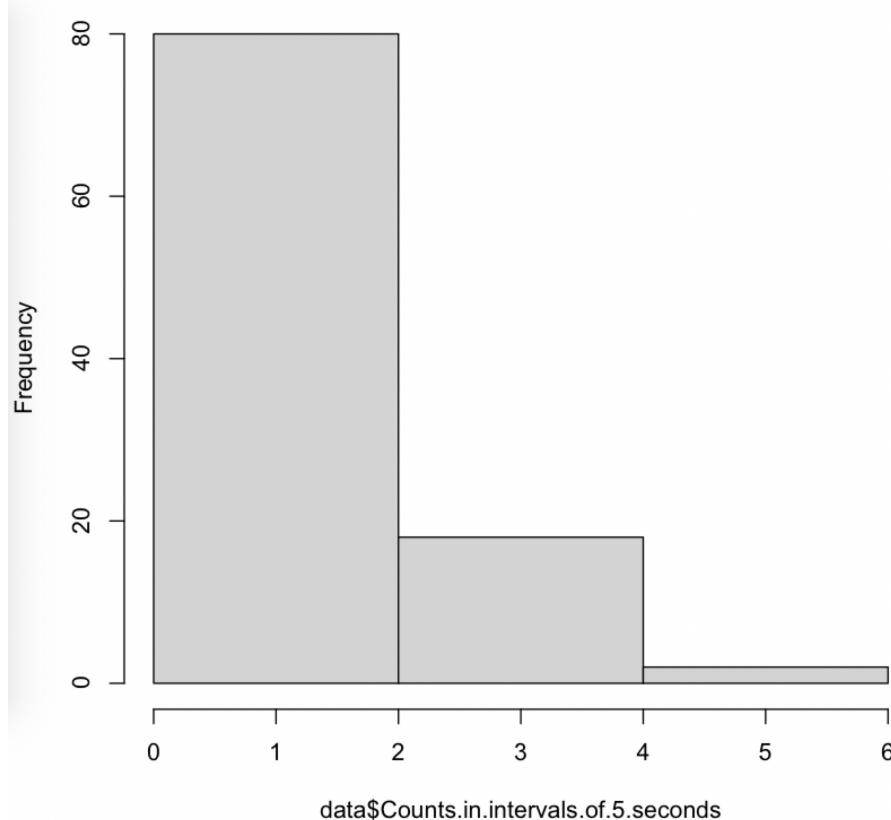
6.2 R code

```
hist(data$Counts.in.intervals.of.5.seconds, 4)w = table(data$Counts.in.i  
for (i in w) print(ppois(i, 1.47))  
dpois(1,6) + dpois(0,6)  #P(X<=1)  
dbinom(1, 0.2)  #same as dpois(1, 0.2)  
1 - pnorm(-0.85,0,1)    #area under normal curve right of z--0.85  
1-ppois(5,5)      #poisson distribution  
n=10  
z=vector()  
z=0  
k=100  
for (i in 1:k){z[i]=mean(runif(n,0,1))}  
curve(dnorm(x, mean=mean(z), sd=sd(z)), col="red", add=T)  
  
qnorm(0.25,mean,sd)      #first quartile
```

6.3 Example Problem Set

Poisson distribution: Poisson distribution of radiation: Put the data on background radiation uploaded in Canvas (Modules/Other_Materials/Background_radiation.csv) on a histogram, calculate the frequencies predicted by the Poisson distribution with λ equal to the observed mean value, create another histogram. Compare the experimental and the expected histograms. See video

Histogram of data\$Counts.in.intervals.of.5.seconds



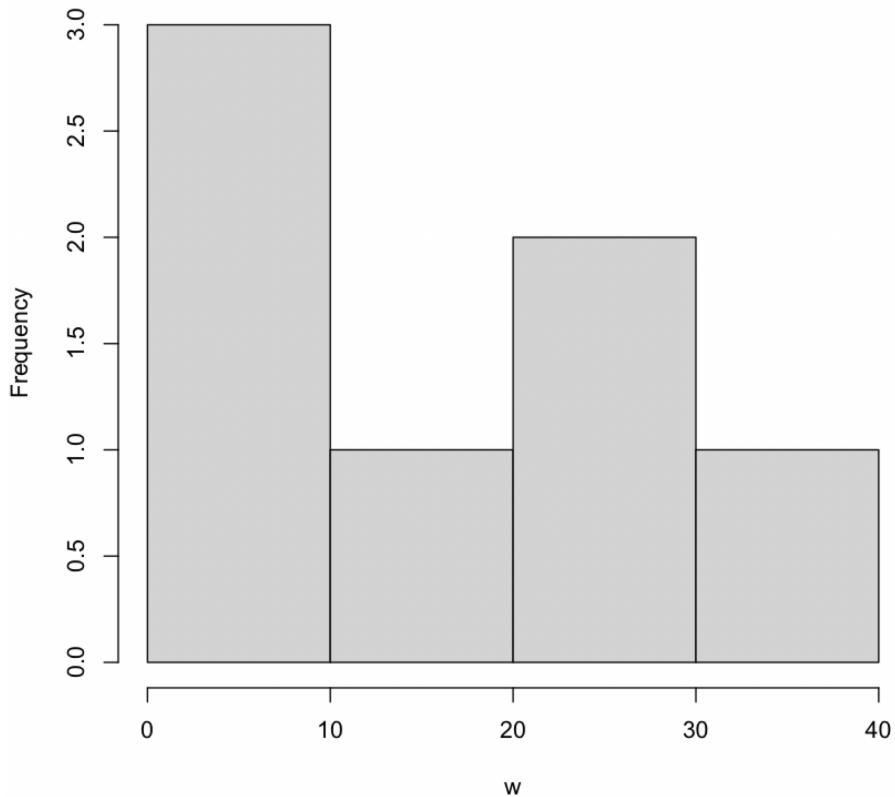
```
hist(data$Counts.in.intervals.of.5.seconds, 4)
w = table(data$Counts.in.intervals.of.5.seconds)
> hist(data$Counts.in.intervals.of.5.seconds, 4)
> w = table(data$Counts.in.intervals.of.5.seconds)
> w
```

0	1	2	3	4	5	6
24	33	23	15	3	1	1

$\lambda = \text{mean} = 1.47$
for (i in w) print(ppois(i, 1.47))

```
> for (i in w){ print(ppois(i, 1.47)) }
[1] 1
[1] 1
[1] 1
[1] 1
[1] 0.9380662
[1] 0.5679159
[1] 0.5679159
```

Histogram of w



4.3.3 In a certain city, the number of potholes on a major street follows a Poisson distribution with a rate of 3 per mile. Let X represent the number of potholes in a two-mile stretch of road. Find:

- a. $P(X = 4)$
 - b. $P(X \leq 1)$
 - c. $P(5 \leq X < 8)$
 - d. μ_X
 - e. σ_X

4.3.3 3 potholes per mile $X = \#$ of potholes in 2 mile stretch
Poisson distribution of road

$$a) P(X=4) = e^{-6} \frac{6^4}{4!} = 0.1339 \quad b/c \quad 6 \text{ per 2 mi} \rightarrow X \sim \text{Poisson}(6)$$

$$P(X=x) = e^{-\lambda} \frac{\lambda^x}{x!} \text{ where } \lambda=6. \quad \text{dpois}(4, 6)$$

$$b) P(X \leq 1) = P(X=0) + P(X=1) = 0.0174 = dpois(0, 6) + dpois(1, 6)$$

$$c) P(5 \leq x < 8) = P(x=5) + P(x=6) + P(x=7) = 0.4589$$

$$d) \mu_x = \lambda = 6$$

$$\text{e)} \sigma_x = \sqrt{\lambda} = 2.45$$

4.3.6 One out of every 5000 individuals in a population carries a certain defective gene. A random sample of 1000 individuals is studied.

- a. What is the probability that exactly one of the sample individuals carries the gene?
 - b. What is the probability that none of the sample individuals carries the gene?
 - c. What is the probability that more than two of the sample individuals carry the gene?
 - d. What is the mean of the number of sample individuals that carry the gene?
 - e. What is the standard deviation of the number of sample individuals that carry the gene?

4.3.6 1 out of 5000 has def. gene. sample = 1000 people

$$a) P(X=1) = dpois(1, np = 1000 \cdot \frac{1}{5000}) = dpois(1, 0.2) = 0.164$$

$$d\text{binom}(1, 1000, 0.0002) = 0.164 \leftarrow \text{the same!}$$

$$b) P(X=0) = dpois(0, 0.2) = 0.819$$

$$c) P(X > 2) = 1 - d_{Pois}(2, 0.2) = 0.00115$$

$$d) \mu_x = \lambda = 0,2$$

$$e) \sigma_x = \sqrt{\lambda} = 0,45$$

4.3.11 A microbiologist wants to estimate the concentration of a certain type of bacterium in a waste-water sample. She puts a 0.5 mL sample of the wastewater on a microscope slide and counts 39 bacteria. Estimate the concentration of bacteria, per mL, in this wastewater, and find the uncertainty in the estimate.

$4.3.11 \quad 0.5 \text{ mL sample} \rightarrow 39 \text{ bacteria}$ estimate concentration per mL + uncertainty of estimate $\hat{\lambda} = \frac{39}{0.5} = 78 \quad \sigma_{\hat{\lambda}} = \sqrt{\frac{\lambda}{t}} = \sqrt{\frac{78}{0.5}} = \sqrt{156} = 12.49 \quad \lambda = 78 \pm 12$

4.3.13 The number of defective components produced by a certain process in one day has a Poisson distribution with mean 20. Each defective component has probability 0.60 of being repairable.

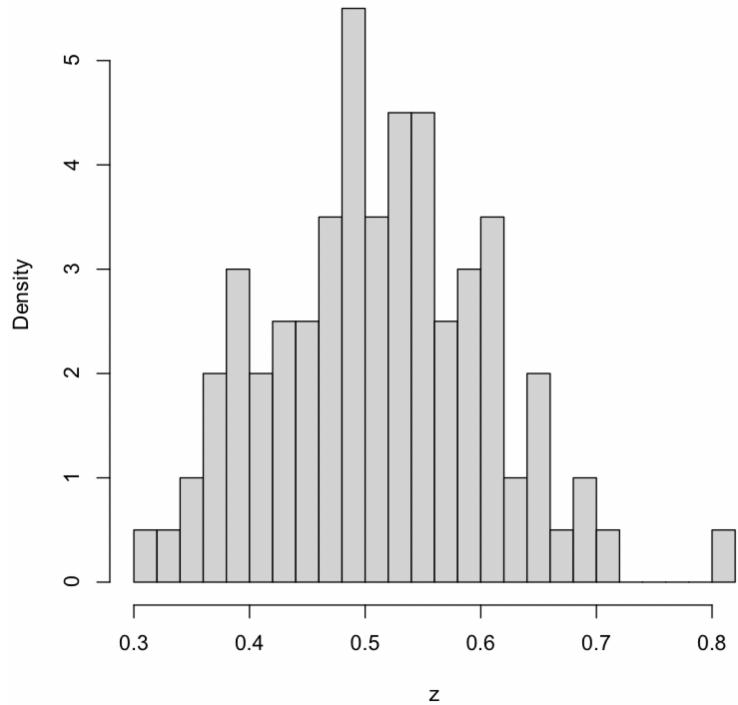
- a. Find the probability that exactly 15 defective components are produced.
- b. Given that exactly 15 defective components are produced, find the probability that exactly 10 of them are repairable.
- c. Let N be the number of defective components produced, and let X be the number of them that are repairable. Given the value of N , what is the distribution of X ?
- d. Find the probability that exactly 15 defective components are produced, with exactly 10 of them being repairable.

$4.3.13 \quad \text{Poisson dist. w/ } \mu_x = 20 = \lambda \quad p = 0.60 \text{ of being repairable}$ a) $P(X=15) = dpois(15, 20) = 0.0516$ b) $P(N=10) = dbinom(10, 15, 0.60) = 0.186$ - of those 15, 10 are ^{Prob.} rep. c) $N = \# \text{ of def. comp.}, \quad X = \# \text{ of rep. comp.}$ $dbinom(X, N, 0.6) \rightarrow \text{Given } N, X \sim Bin(N, 0.6)$ d) exactly 15 def. comp., 10 rep. $\rightarrow P(N=15 \cap X=10)$ $= P(N=15) P(X=10 N=15) = dpois(15, 20) \cdot dbinom(10, 15, 0.6) = 0.0096$

Normal distribution:

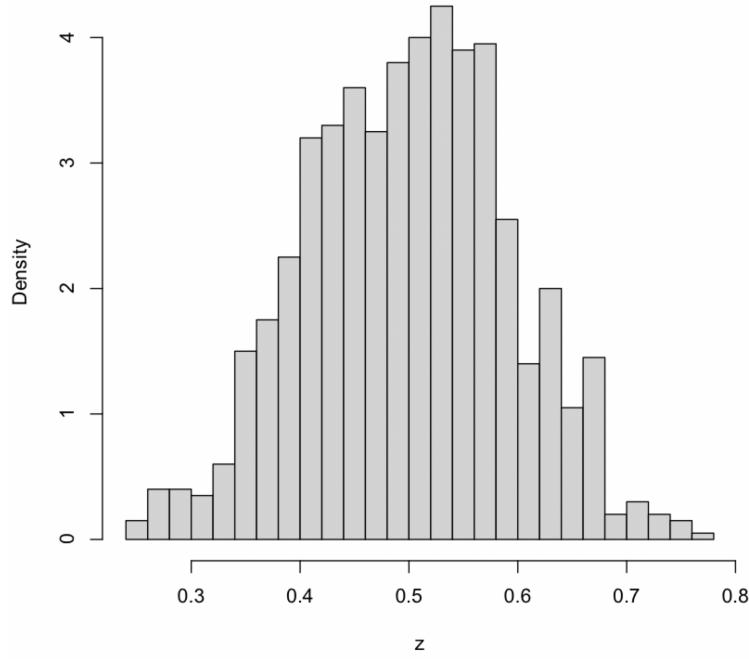
Miniproject Phase 2. Prove the central limit theorem with a numerical experiment.

Histogram of z



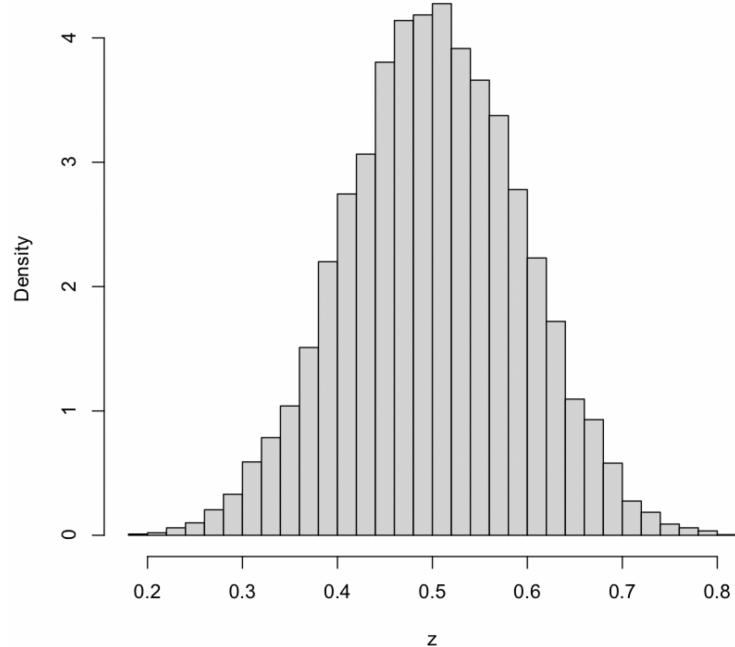
$k=100$

Histogram of z



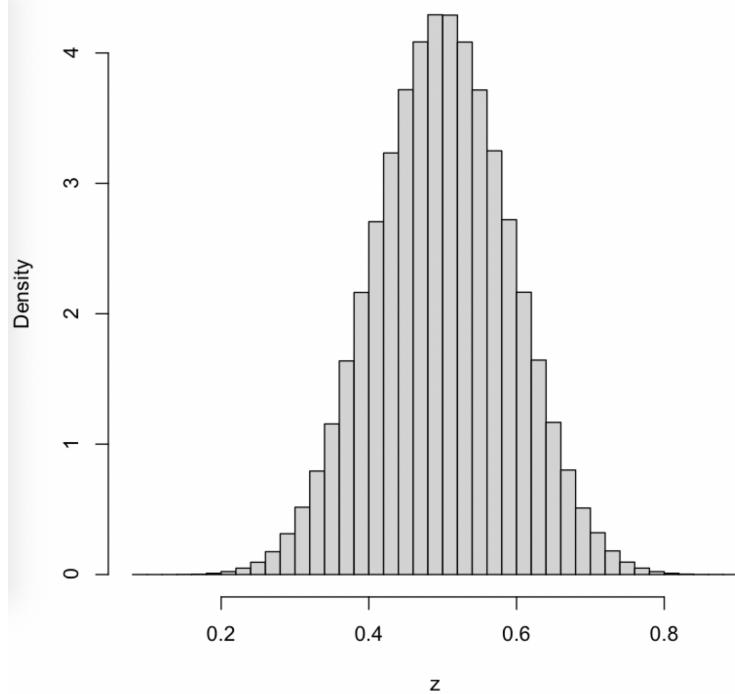
$k=1000$

Histogram of z



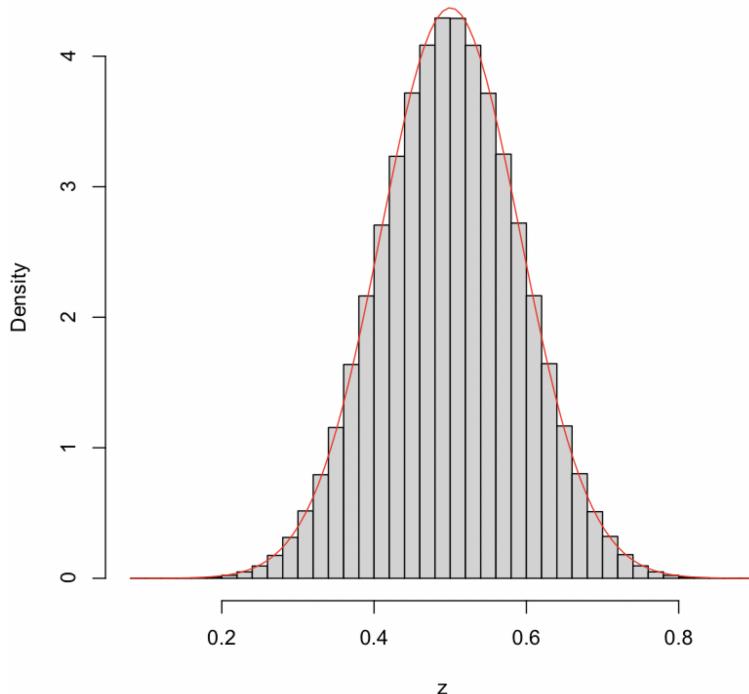
$k=10000$

Histogram of z



$k=1000000$

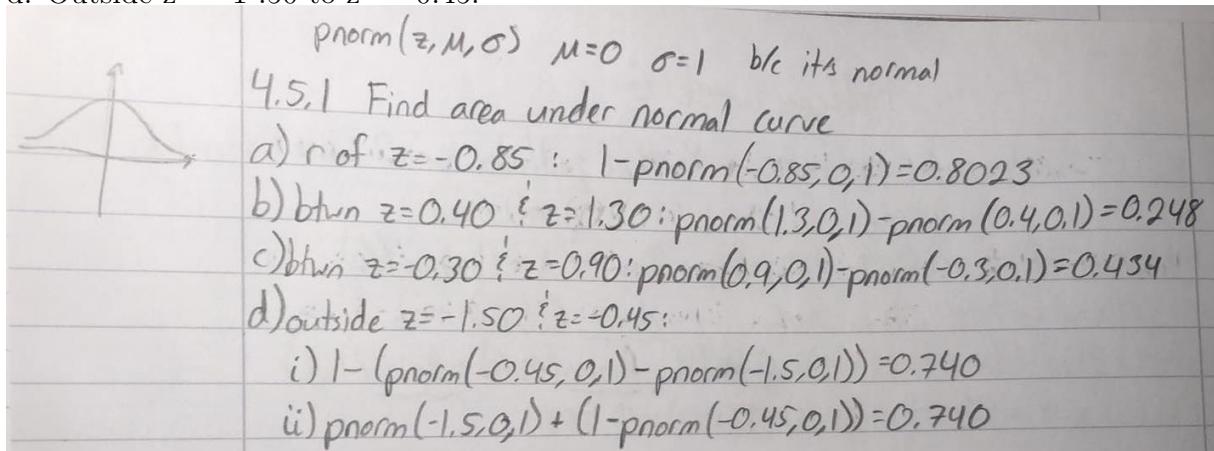
Histogram of z



$k=1000000$ with red curve - $\text{curve(dnorm(x, mean=mean(z), sd=sd(z)), col="red", add=T)}$
 This shows the central limit theorem because with larger sample sizes (k) the more normalized the data looks.

4.5.1 Find the area under the normal curve

- a. To the right of $z = -0.85$.
- b. Between $z = 0.40$ and $z = 1.30$.
- c. Between $z = -0.30$ and $z = 0.90$.
- d. Outside $z = -1.50$ to $z = -0.45$.



4.5.3 Let $Z \sim N(0, 1)$. Find a constant c for which

- a. $P(Z \geq c) = 0.1587$
- b. $P(c \leq Z \leq 0) = 0.4772$
- c. $P(-c \leq Z \leq c) = 0.8664$
- d. $P(0 \leq Z \leq c) = 0.2967$
- e. $P(|Z| \leq c) = 0.1470$

$qnorm(p, \mu, \sigma)$

4.5.3 $Z \sim (0, 1)$, find constant c

a) $P(Z \geq c) = 0.1587 \rightarrow qnorm(1 - 0.1587, 0, 1) = 0.998$

b) $P(c \leq Z \leq 0) = 0.4772 : qnorm((1 - (0.4772 + 0.5)), 0, 1) = -1.999$ the other half

c) $P(-c \leq Z \leq c) = 0.8664 : qnorm((0.5 - \frac{0.8664}{2}), 0, 1) = -1.500$

d) $P(0 \leq Z \leq c) = 0.2967 : qnorm((0.5 + 0.2967), 0, 1) = 0.82$

e) $P(|Z| \leq c) = 0.1470 : qnorm(0.1470/2, 0, 1) = -1.450 \quad c = 1.45$

4.5.7 Scores on a standardized test are approximately normally distributed with a mean of 480 and a standard deviation of 90.

- a. What proportion of the scores are above 700?
- b. What is the 25th percentile of the scores?
- c. If someone's score is 600, what percentile is she on?
- d. What proportion of the scores are between 420 and 520?

$$4.5.7 \quad \mu = 480 \quad \sigma = 90$$

a) what proportion are above 700?

$$z = (700 - 480) / 90 = 2.44 \quad \text{and } pnorm(2.44, 0, 1) = 0.0073$$

b) 25th % of scores?

$$z \approx -0.7 \quad (\text{table A2}) \approx -0.67 \quad (\text{same table})$$

$$25^{\text{th}} \% \approx 480 - 0.67(90) = 419.7$$

c) score = 600, what percentile

$$z = (600 - 480) / 90 = 1.33 \quad \text{and } pnorm(1.33, 0, 1) = 0.9082 = 91^{\text{st}} \%$$

d) proportion of scores between 420 and 520

$$z_1 = (420 - 480) / 90 = -0.67, \quad z_2 = (520 - 480) / 90 = 0.44$$

$$pnorm(0.44, 0, 1) - pnorm(-0.67, 0, 1) = 0.4186$$

4.5.9 The lifetime of a lightbulb in a certain application is normally distributed with mean $\mu = 1400$ hours and standard deviation $\sigma = 200$ hours.

- What is the probability that a lightbulb will last more than 1800 hours?
- Find the 10th percentile of the lifetimes.
- A particular lightbulb lasts 1645 hours. What percentile is its lifetime on?
- What is the probability that the lifetime of a light-bulb is between 1350 and 1550 hours?
- Eight lightbulbs are chosen at random. What is the probability that exactly two of them have lifetimes between 1350 and 1550 hours?

$$4.5.9 \quad \mu = 1400 \text{ hrs} \quad \sigma = 200 \text{ hrs}$$

$$a) P(X > 1800) \quad z = (1800 - 1400)/200 = 2.00 \quad \text{of } z: 1 - \text{pnorm}(2.0, 1) = 0.0228$$

$$b) 10^{\text{th}} \text{ perc. } \sim -1.28 \approx z \quad 1400 - 1.28(200) = 1144$$

$$c) 1645 \text{ hrs} = ? \% \quad z = (1645 - 1400)/200 = 1.23 \quad \text{of } z: \text{pnorm}(1.23, 0, 1) = 0.8907$$

$$d) 8 \text{ random bulbs } P(\text{that 2 have lifetime btwn 1350 and 1550 hrs})$$

$$z_1 = (1350 - 1400)/200 = -0.25 \quad z_2 = (1550 - 1400)/200 = 0.75$$

$$\text{pnorm}(0.75, 0, 1) - \text{pnorm}(-0.25, 0, 1) = 0.3721$$

4.5.13 A cylindrical hole is drilled in a block, and a cylindrical piston is placed in the hole. The clearance is equal to one-half the difference between the diameters of the hole and the piston. The diameter of the hole is normally distributed with mean 15 cm and standard deviation 0.025 cm, and the diameter of the piston is normally distributed with mean 14.88 cm and standard deviation 0.015 cm. The diameters of hole and piston are independent.

- Find the mean clearance.
- Find the standard deviation of the clearance.
- What is the probability that the clearance is less than 0.05 cm?
- Find the 25th percentile of the clearance.
- Specifications call for the clearance to be between 0.05 and 0.09 cm. What is the probability that the clearance meets the specification?
- It is possible to adjust the mean hole diameter. To what value should it be adjusted so as to maximize the probability that the clearance will be between 0.05 and 0.09 cm?

$$4.5.13 \quad \mu_h = 15 \quad \sigma_{d_h} = 0.025 \quad \mu_p = 14.88 \quad \sigma_{d_p} = 0.015$$

a) mean clearance

$$\mu_c = \mu_{0.5h - 0.5p} = 0.5\mu_h - 0.5\mu_p = 0.06$$

b) sd of clearance

$$\sigma_c = \sqrt{0.5^2 \sigma_h^2 + (-0.5)^2 \sigma_p^2} = 0.01458$$

c) $P(\text{clearance} < 0.05)$

$$z\text{-score} = (0.05 - 0.06) / 0.01458 = -0.09 \quad \text{lof } z = \text{pnorm}(0.69, 0, 1) = 0.245 = P(c < 0.05)$$

$$d) 25^{\text{th}} \% z \approx 0.06 - 0.67(0.01458) = 0.0502 \text{ cm}$$

$$e) P(0.05 \leq c \leq 0.09) \quad z_1 = (0.05 - 0.06) / 0.01458 = -0.69 \quad z_2 = (0.09 - 0.06) / 0.01458 = 2.06$$

$$\text{pnorm}(2.06, 0, 1) - \text{pnorm}(-0.67, 0, 1) = 0.729$$

f) maximize clearance prob $P(0.05 < c < 0.09)$

maximized at $\mu_c = 0.07$ (middle of range), adjust μ_h (hole size)

$$\mu_c = 0.5\mu_h - 0.5\mu_p = 0.07 = 0.5\mu_h - 0.5(14.88) \rightarrow \mu_h = 15.02$$

$$z_1 = (0.05 - 0.07) / 0.01458 = -1.37 \quad z_2 = (0.09 - 0.07) / 0.01458 = 1.37$$

$$\text{pnorm}(1.37, 0, 1) - \text{pnorm}(-1.37, 0, 1) = 0.8294$$

4.5.21 Let $X \sim N(\mu, \sigma^2)$, and let $Z = (X - \mu)/\sigma$. Use Equation (4.25) to show that $Z \sim N(0, 1)$.

4.5.21

$X \sim N(\mu, \sigma^2)$, $Z = \frac{X - \mu}{\sigma}$, Show that $Z \sim N(0, 1)$ eq. 4.25

eq 4.25: $aX + b \sim N(a\mu + b, a^2\sigma^2)$

let $a = \frac{1}{\sigma}$ and $b = -\frac{\mu}{\sigma}$, $Z = aX + b$. Now $a\mu + b = \frac{1}{\sigma}\mu - \frac{\mu}{\sigma} = 0$ and
 $a^2\sigma^2 = \left(\frac{1}{\sigma}\right)^2\sigma^2 = 1$ $\therefore Z \sim N(0, 1)$ (eq. 4.25)

6.4 Problem Set

Lilja Ýr Guðmundsdóttir - liljag18

Verkefni 6

Here are the solutions, the histograms for question 1 are on the next page.

Verkefni 6

Lilja Ýr Guð
liljag18

1. 100 data samples, $n=5$ random num: w/ distr. bwn. (0,1)
calculate mean \bar{x} and hist w/ 100 \bar{x}

a) $n=5$
 $z=\text{vector}()$
 $z=0$
 $k=100000$
 $\text{for}(i \text{ in } 1:k) \{$
 $z[i] = \text{mean}(\text{runif}(n, 0, 1))$
 $\}$

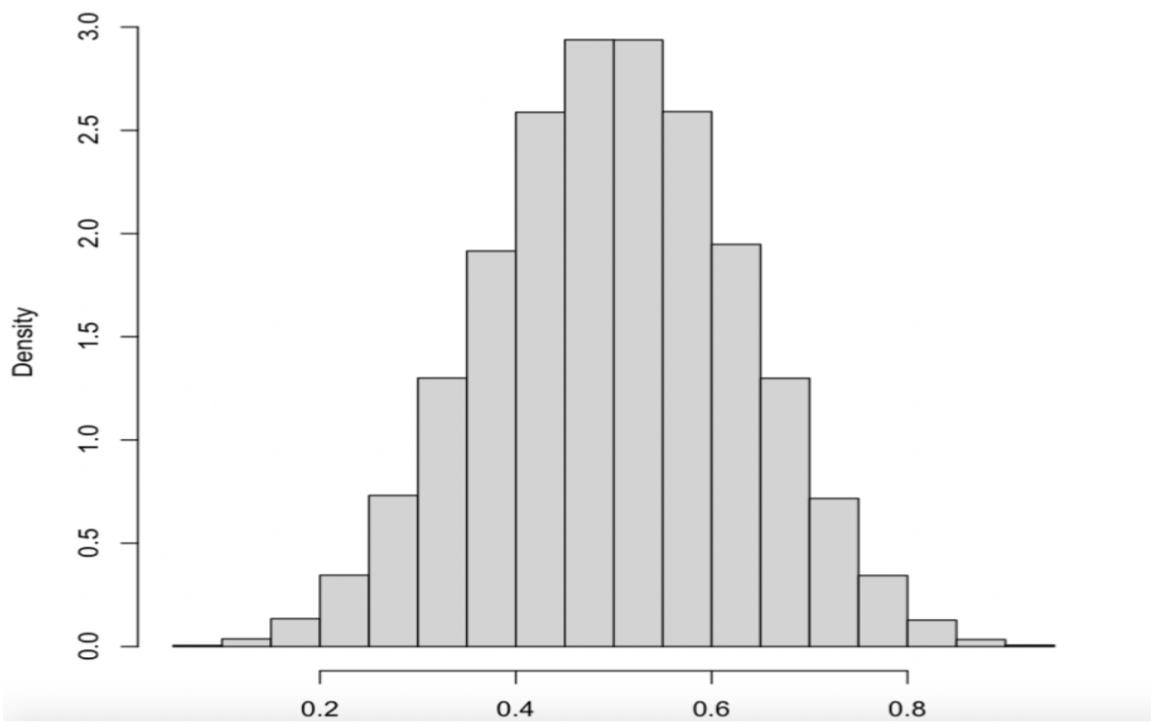
b) $\text{hist}(z, \text{prob}=T, \text{breaks}=20)$

b) The histogram for $n=5$ is wider while for sample $n=10$ it gets more narrow. This is because the sample size is bigger leading to less weight in any outliers since there is a higher probability of the samples chosen being similar to one another. □

2 Poisson distr. of 5 cars per minute. Find probability of 5 cars arriving in one minute.
We need to find $P(X=5)$. We find lambda by $\lambda=5$ b/c there are 5 cars per minute and we are calculating for one minute. There are two ways to solve this:
i) By hand: $P(X=5) = e^{-5} \frac{5^5}{5!} \approx 0.175$
ii) With R: $P(X=5) = \text{dpois}(5, 5) = 0.175$. □

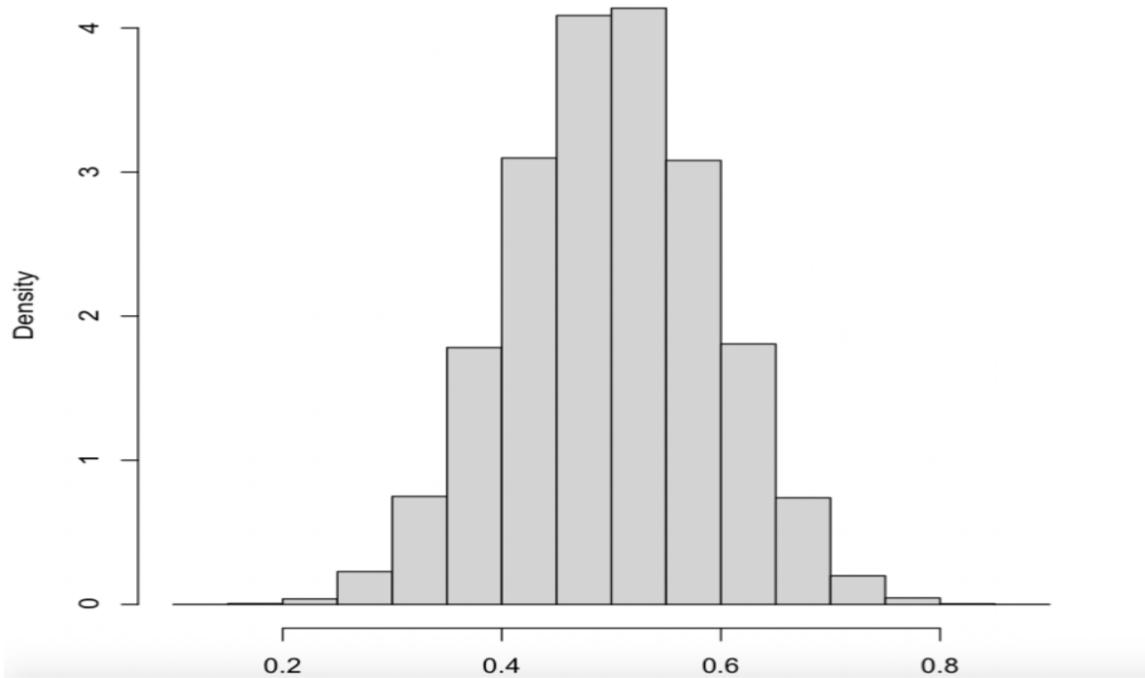
3. Lifetime of battery w/ normal distr. $\mu = 300$ hrs, $\sigma = 50$ hrs
a) $P(\mu < x < \mu + \sigma) = P(300 < x < 350)$ We get (assuming $\mu=300$ and $\sigma=50$):
 $z_{\mu} = \frac{300-300}{50} = 0$ and $z_{\mu+\sigma} = \frac{350-300}{50} = 1$, so the prob is found using R with command $\text{pnorm}(1, 0, 1) - \text{pnorm}(0, 0, 1) = 0.341 = P(\mu < x < \mu + \sigma)$
b) first quartile or first 25th percentile
According to table A2 in the book the $z \approx -0.67$. Then we get
 $25^{\text{th}} \% = 300 - 0.67(50) = 266.5$. □

Histogram of z



Here is the histogram for n=5.

Histogram of z



Here is the histogram for n=10.

6.5 Problem Set Solutions

1. a) Do the following computer experiment: -Use your favourite computer tools to generate 100 data samples, each of $n=5$ random numbers with uniform distribution between $(0,1)$. For each sample calculate the sample mean \bar{x} . Make a histogram with these 100 calculated \bar{x} 's.

-Repeat this experiment using $n=10$ and make a new histogram with the new 100 sample means.

-Include both histograms in the uploaded solution using at least 5 bins and mention briefly the computer functions or commands used.

a) Example of solution: Here I will denote x_{bar} with $z_n=5$;

```
k=100; z1=0
for ( i in 1:k) {z1 [ i]=mean( runif(n ,0 ,1)) }
hist (z1 , xlim=c(0 ,1))
```

```
n=10; k=100; z2=0
for ( i in 1:k) {z2 [ i]=mean( runif(n ,0 ,1)) }
hist (z2 , xlim=c(0 ,1))
```

b) Which histogram is wider and which one is narrower? Explain why.

The width of the histogram depends on the standard deviation of the 100 numbers.

The larger the std. the wider the histogram, which can be seen as an approximation of a normal distribution.

For $n=5$ I got $sd(z_1)=0.118$ and the "true" value is $1/\sqrt{12}/\sqrt{5}=0.129$

For $n=10$ I got $sd(z_2)=0.0893$ and the "true" value is $1/\sqrt{12}/(10)=0.0913$.

So the histogram for z_1 is wider than the histogram for z_2 , because $sd(z_1) > sd(z_2)$.

Comment 1: Here 2 points are given only if the explanation includes some numbers describing the std. in each case, or the argument that the true std decreases like $1/\sqrt{12}/\sqrt{n}$.

Comment2: We know that for larger k , like 1000000, $sd(z)$ calculated with R will be very close to the "true" or "exact" std., both for $n=5$ and $n=10$. Try this exercise if you did not do it yet, to consolidate your understanding.

2. The number of cars arriving at an intersection follows a Poisson distribution with a mean rate of 5 per minute. Find the probability that more than 5 cars arrive in one minute.

$1 - \text{ppois}(5, 5) = 0.384$

3. The lifetime of a battery is normally distributed with mean $\mu=300$ hours and standard deviation $\sigma=50$ hours. Calculate:

a) The probability that a battery lasts between μ and $\mu + \sigma$ hours.

$m=300; s=50$

$\text{pnorm}(m+s, m, s) - \text{pnorm}(m, m, s) = 0.3413$ b)

b) The first quartile of the lifetimes.

$\text{qnorm}(0.25, m, s) = 266.28$

7 Week 7

7.1 Concepts

- An estimator is biased if its mean value is different from the exact value
- An estimated value of the probability parameter p of a binomial distribution is the observed number of successes divided by the number of trials n . The mean squared error of this estimate is $p(1-p)/n$
- The likelihood function is the probability of the observed data as function of a parameter.
- The maximum likelihood estimator of the parameter λ of the Poisson distribution is equal to the sample mean of the observed data.
- Consider x a random variable with normal distribution. The maximum likelihood estimator of the mean value is equal to the sample mean of the observed data.
- A confidence interval is reported as plus or minus two standard deviations. Considering that the random variable has a normal distribution, this interval corresponds to confidence level of approximately 95%
- Decrease the estimated standard deviation to narrow down a 95% confidence interval in the case of a normal distribution
- A random variable with a chi square distribution is always positive
- A random variable has a Chi square distribution with 3 degrees of freedom. Its mean value is 3
- Consider two independent random variables z_1 and z_2 with standard normal distribution. The random variable equal to z_1 squared plus z_2 squared has Chi square distribution with 2 degrees of freedom
-

$$\text{Mean: } E(Y) = e^{\mu + \frac{\sigma^2}{2}}$$

$$\text{SD: variance} = V(Y) = \frac{1}{\lambda^2} = e^{2\mu+2\sigma^2} - e^{2\mu+\sigma^2} = sd^2$$

$$\text{Median: } m = e^\mu$$

7.2 R code

```
1-pexp(1,lambda)      # Returns the cumulative probability distribution
                        #for value x
qexp(p, r)            #Returns the inverse cumulative probability
                        #distribution for probability p
plnorm(mean,6,6)      #used to compute the log normal value of the
                        #cumulative probability density function
```

7.3 Example Problem Set

Lognormal distribution:

4.6.1 The lifetime (in days) of a certain electronic component that operates in a high-temperature environment is lognormally distributed with $\mu = 1.2$ and $\sigma = 0.4$.

- Find the mean lifetime.
- Find the probability that a component lasts between three and six days.
- Find the median lifetime.
- Find the 90th percentile of the lifetimes.

●	<p>4.6.1 Lognormally distr. w/ $\mu = 1.2$, $\sigma = 0.4$</p> <p>a) mean lifetime</p> <p>$Y = \text{lifetime of randomly chosen component}$</p> <p>use eq. 4.31: $E(Y) = e^{\mu + \frac{\sigma^2}{2}} = \text{mean} = e^{1.2 + \frac{0.4^2}{2}} = 3.6$</p> <p>b) prob. comp. lasts between 3 and 6 days</p> $P(3 < Y < 6) = P(\ln(3) < \ln(Y) < \ln(6)) \Rightarrow \ln(Y) \sim N(1.2, 0.4^2)$ <p>$z\text{-score: } z_1 = \frac{\ln(3) - 1.2}{0.4} = -0.25, z_2 = \frac{\ln(6) - 1.2}{0.4} = 1.48$</p> <p>area between the two is $0.9306 - 0.4013 = 0.5293$</p> <p>$\therefore P(3 < Y < 6) = 0.5293$</p> <p>c) median lifetime</p> <p>median $= e^\mu = e^{1.2} = 3.32$ ($P(Y \leq m) = P(\ln(Y) \leq \ln(m)) = 0.5 = P(\ln(Y) < 1.2)$)</p> <p>d) 90th percentile of the lifetimes</p> <p>$P(Y \leq p_{90}) = 0.90 = P(\ln(Y) \leq \ln(p_{90}))$, z score for 90th% = 1.28</p> <p>$z\text{-score find decimal perc. in table A2 (bulk data part - loc of } \sim 0.9)$</p> <p>$\ln(p_{90}) = 1.2 + 1.28 \cdot 0.4 = 1.712 \therefore p_{90} = e^{1.712} = 5.54$</p> <p>$1.28 = (\ln(p_{90}) - 1.2) / 0.4$</p>
------------------------------------	--

4.6.3 The body mass index (BMI) of a person is defined to be the person's body mass divided by the square of the person's height. The article "Influences of Parameter Uncertainties within the ICRP 66 Respiratory Tract Model: Particle Deposition" (W. Bolch, E. Farfan, et al., Health Physics, 2001;378–394) states that body mass index (in kg/m^2) in men aged 25–34 is lognormally distributed with parameters $\mu = 3.215$ and $\sigma = 0.157$.

- Find the mean BMI for men aged 25–34.
- Find the standard deviation of BMI for men aged 25–34.
- Find the median BMI for men aged 25–34.
- What proportion of men aged 25–34 have a BMI less than 22?
- Find the 75th percentile of BMI for men aged 25–34.

4.6.3 BMI = $\frac{\text{body mass}}{\text{height}^2} \text{ kg/m}^2$ logn. distr. $\mu = 3.215, \sigma = 0.157$

a) mean

$$E(Y) = e^{\mu + \frac{\sigma^2}{2}} = e^{3.215 + \frac{0.157^2}{2}} = 25.2$$

b) sd

$$\text{variance} = V(Y) = e^{2\mu + 2\sigma^2} - e^{2\mu + \sigma^2} = sd^2 \therefore$$

$$sd = \sqrt{e^{2 \cdot 3.215 + 2 \cdot 0.157^2} - e^{2 \cdot 3.215 + 0.157^2}} = \sqrt{15.86} = 3.9828$$

c) median

$$m = e^\mu = e^{3.215} = 24.903$$

d) $P(Y < 22) = P(\ln(Y) < \ln(22)) = P(\ln(Y) < 3.0910)$

z score of 3.0910 is $\frac{3.0910 - 3.215}{0.157} = -0.79$, loc of -0.79
area to the left is 0.2148 (read table A2 for left column)
 $P(Y < 22) = 0.2148$

e) 75th percentile

$$\text{z-score } 75^{\text{th}} \% = 0.68 \rightarrow 0.68 = \frac{\ln(p_{75}) - 3.215}{0.157} \rightarrow \ln(p_{75}) = 3.32 \rightarrow p_{75} = 27.71$$

4.6.4 The article "Stochastic Estimates of Exposure and Cancer Risk from Carbon Tetrachloride Released to the Air from the Rocky Flats Plant" (A. Rood, P. McGavran, et al., Risk Analysis, 2001:675-695) models the increase in the risk of cancer due to exposure to carbon tetrachloride as lognormal with $\mu = -15.65$ and $\sigma = 0.79$.

- Find the mean risk.
- Find the median risk.
- Find the standard deviation of the risk.
- Find the 5th percentile.
- Find the 95th percentile.

4.6.4 logn. dist. $\mu = -15.65$ and $\sigma = 0.79$

- mean
 $E(Y) = e^{\mu + \frac{\sigma^2}{2}} = e^{-15.65 + \frac{0.79^2}{2}} = 2.18 \times 10^{-7}$
- median
 $m = e^\mu = e^{-15.65} = 1.597 \times 10^{-7}$
- sd
 $sd = \sqrt{e^{2\mu + 2\sigma^2} - e^{2\mu + \sigma^2}} = 2.031 \times 10^{-7}$
- 5th %
 $Z\text{-score: } -1.65 \rightarrow -1.65 = \frac{\ln(p_{5\%}) - (-15.65)}{0.79} \Rightarrow \ln(p_{5\%}) = -1.65 \cdot 0.79 - 15.65$
 $\Rightarrow \ln(p_{5\%}) = -16.95 \Rightarrow p_{5\%} = 4.337 \times 10^{-8}$
- 95th %
 $Z\text{-score: } 1.65 \rightarrow 1.65 = \frac{\ln(p_{95\%}) - (-15.65)}{0.79} \Rightarrow \ln(p_{95\%}) = 1.65 \cdot 0.79 - 15.65$
 $\Rightarrow \ln(p_{95\%}) = -14.3465 \Rightarrow p_{95\%} = 5.88 \times 10^{-7}$

Exponential distribution:

Example 4.59 The lifetime of a particular integrated circuit has an exponential distribution with mean 2 years. Find the probability that the circuit lasts longer than three years.

4.59 exp. distr. w/ mean 2 yrs $P(Y > 3)$ $\mu_T = 2 \rightarrow \lambda = \frac{1}{\mu_T} = 0.5$

$$P(Y > 3) = 1 - P(Y \leq 3) = 1 - (1 - e^{-0.5(3)}) = e^{-1.5} = 0.223$$

Example 4.60 Refer to Example 4.59. Assume the circuit is now four years old and is still functioning. Find the probability that it functions for more than three additional years. Compare this probability with the probability that a new circuit functions for more than three years, which was calculated in Example 4.59.

4.60 4 yrs old $P(Y>3+4)$ vs $P(Y>3)$

$$P(Y>7|Y>4) = \frac{P(Y>7 \text{ and } Y>4)}{P(Y>4)} = \frac{P(Y>7)}{P(Y>4)} = \frac{e^{-0.5(7)}}{e^{-0.5(4)}} = e^{-1.5} = 0.223$$

P the circuit lasts an additional 3 yrs (after 4 yrs) is the same P as it lasting 3 years

4.7.2 The time between requests to a web server is exponentially distributed with mean 0.5 seconds.

- What is the value of the parameter λ ?
- What is the median time between requests?
- What is the standard deviation?
- What is the 80th percentile?

$$4.7.2 \mu_s = 0.5$$

$$a) \lambda = \frac{1}{0.5} = 2 s^{-1}$$

$$b) m = \frac{\ln(2)}{\lambda} = \frac{\ln(2)}{2} = 0.347 \quad qexp(0.5, \lambda) = qexp(\mu, \lambda)$$

$$c) sd^2 = \sigma^2 = \frac{1}{\lambda^2} = \mu_s^2 \rightarrow \sigma = \mu_s = \frac{1}{\lambda} = 0.5$$

$$d) 80^{\text{th}} \% P(T \leq x_{80}) = 1 - e^{-\lambda x_{80}} = 0.80 \therefore e^{-2x_{80}} = 0.20 \quad qexp(0.8, \lambda) \\ \Rightarrow \ln(e^{-2x_{80}}) = \ln(0.20) \Rightarrow -2x_{80} = -1.609 \Rightarrow x_{80} = 0.8047$$

$$1 - pexp(1, \lambda) \quad e) > 1s \text{ b/w req.} \Rightarrow P(t > 1) = 1 \Rightarrow 1 - P(t \leq 1) = 1 - (1 - e^{-\lambda t}) = e^{-2} = 0.135$$

$$f) P(t > 2 | t > 1) = \frac{P(t > 3 \text{ and } t > 2)}{P(t > 3)} = \frac{P(t > 3)}{P(t > 2)} = \frac{e^{-6}}{e^{-4}} = e^{-2} = 0.135$$

$$(1 - pexp(3, \lambda)) / (1 - pexp(2, \lambda))$$

4.7.6 Refer to Exercise 2.

- Find the probability that there will be exactly 5 requests in a 2-second time interval.
- Find the probability that there will be more than 1 request in a 1.5-second time interval.
- Find the probability that there will be no requests in a 1-second time interval.
- Find the probability that the time between requests is greater than 1 second.
- Find the probability that the time between requests is between 1 and 2 seconds.
- If there have been no requests for the past two seconds, what is the probability that more than one additional second will elapse before the next request?

4.7.6

a) P of 5 req in a 2 sec interval $\mu_T = 0.5 \text{ s}$, $\lambda = \frac{\text{(req/s)} \cdot \lambda}{\text{req!}} = \frac{2}{2} = 1$

$\frac{2}{2} \text{ s between each request} \rightarrow 0.405$

$P(X=5) = \frac{(1)^5}{5!} e^{-1} = \frac{1}{120} e^{-1} = 0.0067$

$P(5) = \frac{(2 \cdot 2)^5}{5!} e^{-2} = 0.1563$

b) $P(X > 1) = \frac{(1.5 \cdot 2)^2}{2!} e^{-1.5} = 0.1493$

c) $P(\text{no req in one s interv.}) = e^{-\lambda t} = e^{-2 \cdot 1} = 0.1353$

d) $1 - e^{-\lambda t} = 1 - e^{-2} = 0.8647$

e) $e^{-2} - e^{-2 \cdot 2} = 0.117$

f) $P(t > 2+1 | t > 2) = \frac{P(t > 3)}{P(t > 2)} = \frac{e^{-6}}{e^{-4}} = e^{-2} = (1 - p_{\text{exp}}(3, \lambda)) / (1 - p_{\text{exp}}(2, \lambda))$

4.7.8 Someone claims that the waiting time, in minutes, between hits at a certain website has the exponential distribution with parameter $\lambda = 1$.

- Let X be the waiting time until the next hit. If the claim is true, what is $P(X \geq 5)$?
- Based on the answer to part (a), if the claim is true, is five minutes an unusually long time to wait?
- If you waited five minutes until the next hit occurred, would you still believe the claim? Explain.

4.7.8 exp. distribution w/ parameter $\lambda = 1$

a) $X = \text{waiting time until next hit. } P(X \geq 5) = 1 - P(X < 5) = 1 - (1 - e^{-5}) = 0.00674$

b) Yes it is not unusual

c) No, this time is too long, the claim looks incorrect

4.7.9 A certain type of component can be purchased new or used. Fifty percent of all new components last more than five years, but only 30% of used components last more than five years. Is it possible that the lifetimes of new components are exponentially distributed? Explain.

4.7.9 new 50% last > 5 yrs, used 30% last > 5 yrs. Is this exp. distr?
 No, if lifetimes = exp. distr. then the proportion of used comp. lasting
 > 5 yrs would be the same as new comp. lasting > 5 yrs b/c of lack of mem. prop.

4.7.11 The number of traffic accidents at a certain intersection is thought to be well modeled by a Poisson process with a mean of 3 accidents per year.

- Find the mean waiting time between accidents.
- Find the standard deviation of the waiting times between accidents.
- Find the probability that more than one year elapses between accidents.
- Find the probability that less than one month elapses between accidents.
- If no accidents have occurred within the last six months, what is the probability that an accident will occur within the next year?

4.7.11 mean = 3 accidents/year

a) $\mu_T = \frac{1}{3}$ of a year

b) $\sigma_T = \frac{1}{\sqrt{3}}$ of a year

c) $P(T > 1) = 1 - P(T \leq 1) = 1 - (1 - e^{-3(1)}) = 0.0498$

d) $P(T < \frac{1}{12}) = 1 - e^{-3(\frac{1}{12})} = 0.2212$

e) Prob: $P(T \leq 1.5 | T > 0.5) = 1 - P(T > 1.5 | T > 0.5)$, by lack of mem. prop..

$P(T > 1.5 | T > 0.5) = P(T > 1)$ so $P(T \leq 1.5 | T > 0.5) = 1 - P(T > 1) = 1 - e^{-3(1)} = 0.9502$.

Estimates

Show the equivalence of Eq.(4.53) and Eq.(4.54) from the textbook.

Formula 4.53 $MSE_{\hat{\theta}} = (\mu_{\hat{\theta}} - \theta)^2 + \sigma_{\hat{\theta}}^2$

Formula 4.54 $MSE_{\hat{\theta}} = \mu_{(\hat{\theta} - \theta)^2}$

There is some relationship between mean (μ) and standard deviation (σ) but the θ makes this all too difficult to figure out.

4.9.6 Let X_1, \dots, X_n be a random sample from a population with the Poisson(λ) distribution. Find the MLE of λ .

4.9.6 X_1, \dots, X_n = random sample from pop w/ Poisson(λ) distribution.

Find MLE of λ .

MLE identical to sample mean

4.9.10 Let X_1, \dots, X_n be a random sample from a $N(\mu, \sigma^2)$ population. Find the MLEs of μ and of σ . (Hint: The likelihood function is a function of two parameters, μ and σ . Compute partial derivatives with respect to μ and σ and set them equal to 0 to find the values $\hat{\mu}$ and $\hat{\sigma}$ that maximize the likelihood function.)

4.9.10

MLE of μ identical to the sample mean, $\text{MLE } \sigma^2 = \frac{1}{n} \sum_i (x_i - \bar{x})^2$

4.9.3 Let X_1 and X_2 be independent, each with unknown mean μ and known variance $\sigma^2 = 1$.

a. Let $\hat{\mu}_1 = \frac{X_1 + X_2}{2}$. Find the bias, variance, and mean squared error of $\hat{\mu}_1$.

b. Let $\hat{\mu}_2 = \frac{X_1 + 2X_2}{3}$. Find the bias, variance, and mean squared error of $\hat{\mu}_2$.

c. Let $\hat{\mu}_3 = \frac{X_1 + X_2}{4}$. Find the bias, variance, and mean squared error of $\hat{\mu}_3$.

d. For what values of μ does $\hat{\mu}_3$ have smaller mean squared error than $\hat{\mu}_1$?

d. For what values of μ does $\hat{\mu}_3$ have smaller mean squared error than $\hat{\mu}_2$?

4.9.3 $\sigma^2 = 1$

$$a) \hat{\mu}_1 = \frac{X_1 + X_2}{2}, E(\hat{\mu}_1) = \frac{\mu X_1 + \mu X_2}{2} = \frac{\mu + \mu}{2} = \mu$$

$$\text{bias: } E(\hat{\mu}_1) - \mu = 0 \quad \text{variance: } V(\hat{\mu}_1) = \frac{\sigma^2 + \sigma^2}{4} = \frac{\sigma^2}{2} = \frac{1}{2}$$

mean² error: sum of variance and square of bias: $MSE(\hat{\mu}_1) = \frac{1}{2} + 0^2 = \frac{1}{2}$

$$b) \hat{\mu}_2 = \frac{X_1 + 2X_2}{3}, E(\hat{\mu}_2) = \frac{\mu X_1 + 2\mu X_2}{3} = \frac{\mu + 2\mu}{3} = \mu$$

$$\text{bias: } E(\hat{\mu}_2) - \mu = 0 \quad \text{variance: } V(\hat{\mu}_2) = \frac{\sigma^2 + 2\sigma^2}{9} = \frac{3\sigma^2}{9} = \frac{1}{3}$$

$$MSE(\hat{\mu}_2) = \frac{1}{3} + 0^2 = \frac{1}{3}$$

$$c) \hat{\mu}_3 = \frac{X_1 + X_2}{4}, E(\hat{\mu}_3) = \frac{\mu X_1 + \mu X_2}{4} = \frac{\mu}{2} \quad \text{bias: } E(\hat{\mu}_3) - \mu = -\frac{\mu}{2}$$

$$\text{variance: } \frac{\sigma^2 + \sigma^2}{16} = \frac{\sigma^2}{8} \quad MSE(\hat{\mu}_3) = \frac{\sigma^2}{8} + \left(-\frac{\mu}{2}\right)^2 = \frac{2\mu^2 + 1}{8}$$

d) $\hat{\mu}_3$ has smaller MSE than $\hat{\mu}_1$, whenever $\frac{2\mu^2+1}{8} < \frac{1}{2}$. Solving for μ gives us $-1.2247 < \mu < 1.2247$

e) $MSE(\hat{\mu}_3) < MSE(\hat{\mu}_2)$ when $\frac{2\mu^2+1}{8} < \frac{5}{9}$, solving gives,
 $-1.3123 < \mu < 1.3123$

Example 4.68 Let X_1, \dots, X_n be a random sample from a $N(\mu, \sigma^2)$ population. The sample variance is $s^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / (n - 1)$. It can be shown that S^2 has mean $\mu_{s^2} = \sigma^2$ and variance $\sigma_{s^2}^2 = 2\sigma^2 / (n - 1)$. Consider the estimator $\hat{\sigma}^2 = \sum_{i=1}^n (X_i - \bar{X})^2 / n$ in which the sum of the squared deviations is divided by n rather than $n - 1$. Compute the bias, variance, and mean squared error of both s^2 and $\hat{\sigma}^2$. Show that $\hat{\sigma}^2$ has smaller mean squared error than s^2 .

Example 4.68

$\mu_{s^2} = \sigma^2$, s^2 is unbiased for σ^2 so mean squared error = variance!

$$MSE_{s^2} = 2\sigma^4 / (n-1), \quad \hat{\sigma}^2 = \frac{n-1}{n} s^2.$$

It follows that $MSE_{\hat{\sigma}^2} = \frac{n-1}{n} MSE_{s^2} = \frac{n-1}{n} \sigma^2 \quad \therefore \text{Bias of } \hat{\sigma}^2 = \frac{n-1}{n} \sigma^2 - \sigma^2 = \frac{-\sigma^2}{n}$

$$\text{Variance given by: } \sigma_{\hat{\sigma}^2}^2 = \frac{(n-1)^2}{n^2} \sigma_{s^2}^2 = \frac{2(n-1)}{n^2} \sigma^4.$$

$$\therefore MSE_{\hat{\sigma}^2} = \left(\frac{-\sigma^2}{n} \right)^2 + \frac{2(n-1)}{n^2} \sigma^4 = \frac{2n-1}{n^2} \sigma^4$$

To show $\hat{\sigma}^2$ has smaller MSE than s^2 we subtract:

$$MSE_{s^2} - MSE_{\hat{\sigma}^2} = \frac{2\sigma^4}{n-1} - \frac{(2n-1)\sigma^4}{n^2} = \frac{3n-1}{n^2(n-1)} > 0 \quad (\text{since } n > 1).$$

7.4 Problem Set

Verkefni 7

Lilja Ýr Guð
11/18/18

1. The mechanical tensile strength of a certain type of steel screws has a lognormal distribution w/ params $\mu=6 \text{ MPa}$ and $\sigma=0.6 \text{ MPa}$
- Calculate the probability that a randomly selected screw has tensile strength larger than the mean.
We have to find $P(X > \mu) = P(X > 6) = P(\ln(X) > \ln(6)) = P(\ln(X) > 1.792)$
We could solve this by hand but why do that when R is so much easier. By running `pnorm(6, 6, 0.6)` I get that $x = 1.1602 \times 10^{-12}$ (very little chance). □
 - Calculate the probability that a randomly selected screw has a tensile strength larger than the median tensile strength.
We first have to find the median: $m = e^{\mu} = e^6 = 403.43$
This is easy to calculate by doing $z = \frac{\ln(403.43) - 6}{0.6} = 0$. This z-score gives us an area to the right that is 0.5 and therefore $P(X > \text{median}) = 50\%$ which makes sense. □
 - Is the mean value of a lognormal distribution always larger than the median?
Explain.

We know that the mean value is μ and the median is e^{μ} . Therefore $\mu > e^{\mu} \Rightarrow \ln(\mu) > \mu \Rightarrow 0 > \mu$. This shows that mean can never be bigger than the median because then μ would have to be less than 0 but $\ln(\text{negative number})$ is not possible. So no. □

2. The number of cars observed by a speed detector mounted on a street is modelled w/ a Poisson distribution w/ mean $\mu=5$ car per min.

- Calculate prob the waiting time between 2 cars is smaller than mean waiting time.
 $P(t < \sigma) = 1 - e^{-\lambda t} = 1 - e^{-\frac{5}{6}} = 1 - e^{-0.83} = 1 - 0.36 = 0.63$. (where $\lambda = \frac{1}{\sigma}$) □
- Find the variance of the waiting time between two cars.

According to the text book then the $\sigma = \mu = 5$ therefore the variance is $\sigma^2 = 25$. □

7.5 Problem Set Solutions

1. The mechanical tensile strength of a certain type of steel screws has a lognormal distribution with parameters $\mu=6$ MPa (mega Pascal) and $\sigma=0.6$ MPa.

a) Calculate the probability that a randomly selected screw has a tensile strength larger than the mean tensile strength. $\text{mean} = e^{\mu+\sigma^2/2} = 483$

$$1 - \text{plnorm}(\text{mean}, 6, 6) = 0.382$$

b) Calculate the probability that a randomly selected screw has a tensile strength larger than the median tensile strength.

By definition the median is the value of an RV x such that $P(x < \text{median}) = P(x > \text{median}) = 0.5$. Alternatively, obtain the median 403.4 and check.

c) Is the mean value of a lognormal distribution always larger than the median? Explain. If x has a lognormal distribution, $\ln(x)$ has a normal distribution whose median is μ , and so the median of x is e^μ . We see that $e^{\mu+\sigma^2/2} > e^\mu$ because $e^{\sigma^2/2} > 1$. So yes, the mean of the lognormal distribution is always larger than the median. This inequality is also expected due to the asymmetry of the distribution, which has a long tail to large values of x , or it is skewed to the right. The 2 points are given either for the proof with exponential functions, or for some explanations of the tail effect indicating that the student understood it.

2. The number of cars observed by a speed detector mounted on a street is modelled with a Poisson distribution with a mean of 5 cars per minute.

a) Calculate the probability that the waiting time between two cars is smaller than the mean waiting time.

The waiting time has an exponential distribution with parameter $\lambda=5$ whose mean value is $T=1/\lambda$. Using the cumulative distribution $F(t)=1-e^{-\lambda t}$ we obtain $P(t < T) = F(T) = 1 - e^{-1} = 0.632$ or

$$\text{pexp}(1/5, 5) = 0.632$$

b) Find the variance of the waiting time between two cars.

$$x \sim \text{Exp}(\lambda) \Rightarrow \text{Var}(x) = \frac{1}{\lambda^2} = 0.04$$

8 Week 8

8.1 Concepts

- If the p-value is less than 0.05 the null hypothesis is rejected
- In a test of a population mean the one sided p-value is one half of the two sided p-value
- The z-score is calculated assuming H_0 true
- The p-value is the probability that the disagreement between the null and alternative hypotheses is greater than observed
- Let's think about a case when the alternative hypothesis is true. If we could increase the sample size the p-value should decrease
- The mean value of a random variable with t distribution with n degrees of freedom is 0
- In a t test with 5 degrees of freedom you obtain $t=2$. Calculate the corresponding 2-sided p-value. Use point (not comma) as decimal symbol. (For example 0.2 not 0,2) -> 0.1019
- If in the previous t test the t number increases, the p-value decreases
- You perform a 1-sided t test with n degrees of freedom and obtain $t=0$. The p-value is 0.5
- A group of eight individuals with high cholesterol levels were given a new drug that was designed to lower cholesterol levels. Cholesterol levels, in mg/dL, were measured before and after treatment for each individual. The best method to test the effect of the new drug using this data is a t test with paired data

8.2 Example Problem Set

Miniproject phase 3: Computer generated random numbers with normal and chi square distributions

1. See first the recent video on this subject.
2. Use the function rnorm(n,0,1) in R or NORM.INV(RAND(),0,1) in Excel to generate a vector of n=100,1000, ... random numbers with standard normal distribution N(0,1). Observe their histogram and how it evolves with increasing n.
3. Generate now several vectors of n random numbers with N(0,1), z1, z2, z3, ... and make histograms of $x_1 = z_1^2$, $x_2 = z_1^2 + z_2^2$, $x_3 = z_1^2 + z_2^2 + z_3^2$, What distribution you obtain? Compare with random numbers generated with the functions rchisq(n,k) or CHISQ.INV(RAND(),k).
4. Consolidate your understanding of the Chi square distribution by practicing such exercises. Make good graphical representations, optimise the number of bins, add on the graphs probability density functions for comparison. Calculate sample mean values using random numbers and compare with the true values. Observe median values. Imagine yourselves other experimental tests.

Program

2. tried rnorm(n, 0, 1) w/ n = 100, 1000, 10000, 100000, 1000000

As n got bigger the histogram gained a more bellshaped curve

(Another method is $x=0$; $n=1000$; $x=runit(n, 0, 1)$)

$z=0$; $z=qnorm(x, 0, 1)$

- Chi² distrib. with 1 DOF (degree of freedom)


```
n=1000; z1=rnorm(n,0,1); hist(z1, breaks=50); hist(z1**2, breaks=50)
hist(z1**2, breaks=50, prob=T, xlim=c(0,5));
curve(dchisq(x,1), col="red", add=T)
```
- Chi² distr. w/ 2+ DOF


```
n=1000; z1=rnorm(n,0,1); z2=rnorm(n,0,1); z3=rnorm(n,0,1); z4=...
z1=0; z2=z1**2 + z2**2 & z3**2 + z4**2
hist(z2, breaks=50)
curve(dchisq(x,2), col="red", add=T)
```

Curve looks good when n = 1000000
- Random #s w/ rchisq fun


```
n=1000; z2=0; DOF=3; z2=rchisq(n,DOF); hist(z2, prob=T)
```

Confidence intervals for population mean:

5.1.13 The sugar content in a one-cup serving of a certain breakfast cereal was measured for a sample of 140 servings. The average was 11.9 g and the standard deviation was 1.1 g.

- Find a 95% confidence interval for the mean sugar content.
- Find a 99% confidence interval for the mean sugar content.
- What is the confidence level of the interval (11.81, 11.99)?
- How large a sample is needed so that a 95% confidence interval specifies the mean to within ± 0.1 ?
- How large a sample is needed so that a 99% confidence interval specifies the mean to within ± 0.1 ?

population mean:

5.1.13 140 1-cup servings, $\bar{x} = 11.9$ g, $sd = 1.1$ g

a) 95% confidence interval for mean sugar content

population mean = μ , confidence interval = $\bar{x} \pm z_{\frac{\alpha}{2}} \sigma_{\bar{x}}$

We want 95% $\therefore 1 - \alpha = 0.95 \Rightarrow \alpha = 0.05$ and $z_{0.05} = 1.96$

$$\sigma_{\bar{x}} \approx \frac{1.1}{\sqrt{140}} = 0.09297$$

$$\therefore \text{conf. interv.} = 11.9 \pm (1.96)(0.09297) = 11.9 \pm 0.1822$$

b) 99% conf. interval

pop. mean = μ , conf. interv. = $\bar{x} \pm z_{\frac{\alpha}{2}} \sigma_{\bar{x}}$

$1 - \alpha = 0.99 \rightarrow \alpha = 0.01 \therefore z_{0.005} = -2.57$

$$\sigma_{\bar{x}} = \frac{1.1}{\sqrt{140}} = 0.09297 \therefore \text{conf. interv.} = 11.9 \pm (-2.57)(0.09297)$$

$$= 11.9 \pm 0.2389$$

c) conf. level of interv. (11.81, 11.99)

area to the right

$$11.99 = 11.9 + z_{\frac{\alpha}{2}} \left(\frac{1.1}{\sqrt{140}} \right) \rightarrow z_{\frac{\alpha}{2}} = 0.968, \quad \alpha = 2 \cdot 0.8340 \sqrt{2 \cdot 0.1660}$$

$$\alpha = 0.332 \therefore \text{conf. level} = 66.8\%$$

d) size of sample so that 95% int. specifies mean within ± 0.1

$$z = 1.96 \text{ (from (a))} \quad 0.1 = 1.96 \left(\frac{1.1}{\sqrt{n}} \right) \text{ solve for } n:$$

$$\frac{5}{98} = \frac{1.1}{\sqrt{n}} \Rightarrow \sqrt{n} = 1.1 \cdot \frac{98}{5} \Rightarrow n = 464.8 \approx 465$$

e) size of sample so that 99% int. specifies mean within ± 0.1

$$z = 2.57 \text{ (from (b))} \quad 0.1 = 2.57 \left(\frac{1.1}{\sqrt{n}} \right) \text{ solve for } n:$$

$$\frac{10}{257} = \frac{1.1}{\sqrt{n}} \Rightarrow \sqrt{n} = \frac{2827}{100} \Rightarrow n = 799.19 \approx 800$$

5.1.23 Based on a large sample of capacitors of a certain type, a 95% confidence interval for the mean capacitance, in μF , was computed to be (0.213, 0.241). Find a 90% confidence interval for the mean capacitance of this type of capacitor.

5.1.23

95% conf. inter. (in μF) is (0.213, 0.241). Find 90% conf. inter.

Sample mean is midpoint of range $\therefore \bar{x} = \frac{0.213 + 0.241}{2} = 0.227$

Then we get $0.241 = 0.227 + 1.96\left(\frac{s}{\sqrt{n}}\right)$ (from (a))

This means that $\frac{s}{\sqrt{n}} = \frac{1}{140}$. Using this we get for 90%:

$$0.227 \pm 1.645 \frac{s}{\sqrt{n}} = 0.227 \pm 0.01175$$

Confidence intervals for the variance:

Example 5.28 A simple random sample of 15 pistons is selected from a large population whose diameters are known to be normally distributed. The sample standard deviation of the piston diameters is $s = 2.0 \text{ mm}$. Find a 95% confidence for the population variance σ^2 .

Ex. 5.28

15 random pistons selected w/ normally distr. diameters; sample sd = 2.0 mm

Find 95% conf. for pop. variance σ^2

Quantity $((n-1)s^2)/\sigma^2$ has chi-square distr. w/ $n-1=14$ DOF

Use Table A.7 p. 1123

upper and lower 0.025 pts are $\chi^2_{14, 0.025} = 5.629$ and $\chi^2_{14, 0.975} = 26.119$

These values contain 95% of area under curve between them :

$$5.629 < \frac{(n-1)s^2}{\sigma^2} < 26.119, \quad s^2 = 4 \text{ and } n = 15 \therefore$$

$$\frac{5.629}{(n-1)s^2} < \sigma^2 < \frac{26.119}{(n-1)s^2} \Rightarrow 2.144 < \sigma^2 < 9.948.$$

5.8.6 Following are weights, in pounds, of 12 two-month-old baby girls. Assume that the population is normally distributed.

12.23	12.32	11.87	12.34	11.48	12.66
8.51	14.13	12.95	10.30	9.34	8.63

- a. Find the sample standard deviation s .
- b. Construct a 95% confidence interval for population standard deviation σ .

5.8.6

weight of 1yr old babies normally distributed

a) sample sd of s

using $r \rightarrow r \leftarrow c(\text{all inputs}) ; sd(r) = 1.798172$

b) 95% conf. interv. for pop. sd σ

$$s = 1.798 \quad n = 12 \quad \alpha = 1 - 0.95 = 0.05 \quad X^2_{12, 0.025} = 23.337 \quad X^2_{12, 0.975} = 4.404$$

interval is: $\left(\frac{(12-1)(1.798)^2}{23.337}, \frac{(12-1)(1.798)^2}{4.404} \right) = (1.524, 8.076)$

Z-test for population mean:

6.1.6 A certain type of stainless steel powder is supposed to have a mean particle diameter of $\mu = 15 \mu\text{m}$. A random sample of 87 particles had a mean diameter of $15.2 \mu\text{m}$, with a standard deviation of $1.8 \mu\text{m}$. A test is made of $H_0: \mu = 15$ versus $H_1: \mu \neq 15$

- a. Find the P-value.
- b. Do you believe it is plausible that the mean diameter is $15 \mu\text{m}$, or are you convinced that it differs from $15 \mu\text{m}$? Explain your reasoning.

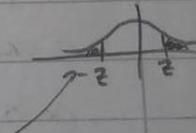
6.1.6 $\mu = 15 \mu\text{m}$, random sample of n particles, $\bar{x} = 15.2 \mu\text{m}$, $s = 1.8 \mu\text{m}$

Test is made of $H_0: \mu = 15$ vs. $H_1: \mu \neq 15$

a) Find P-value

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{15.2 - 15}{1.8/\sqrt{87}} = 0.1928803 \quad \therefore z = 1.0364$$

$$p_val = 2 * pnorm(-|z|, 0, 1) = 0.30 \quad (\text{calc. area twice})$$



b) H_0 more plausible b/c p-value $>> 0.05$, we can't say definitely that it is better, need more data

6.1.7 When it is operating properly, a chemical plant has a mean daily production of at least 740 tons. The output is measured on a simple random sample of 60 days. The sample had a mean of 715 tons/day and a standard deviation of 24 tons/day. Let μ represent the mean daily output of the plant. An engineer tests $H_0: \mu \geq 740$ versus $H_1: \mu < 740$.

- Find the P-value.
- Do you believe it is plausible that the plant is operating properly or are you convinced that the plant is not operating properly? Explain your reasoning.

6.1.7

$\bar{x}_{\text{daily}} = 740$ tons, sample 60 days, sample $\bar{x} = 715$ tons/day, $sd = 24$ tons/day

Test: $H_0: \mu \geq 740$ vs. $H_1: \mu < 740$

a) p-value

$$z = \frac{\bar{x} - \mu}{\sigma} = \frac{715 - 740}{24} = -8.07. \text{ Since alternate hyp. is of form } \mu < \mu_0, \text{ p-value is area to the left of } z \therefore P \approx 0$$

$$2 \Phi_{\text{norm}}(z, 0, 1) = 7.1 \times 10^{-16} \approx 0$$

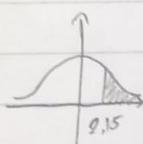
b) Operating properly?

If $\bar{x} = 740$ tons then prob of observing a sample mean as small as 715 tons is nearly 0. \therefore the mean daily output is less than 740 tons

Z-test for population proportions:

6.3.5 In a survey of 500 residents in a certain town, 274 said they were opposed to constructing a new shopping mall. Can you conclude that more than half of the residents in this town are opposed to constructing a new shopping mall?

6.3.5



500 residents, 274 opposed, can you conclude $\frac{1}{2}$ pop are against it?

$$n = 500, k = 274, p = 0.5 \quad \hat{p} = \frac{k}{n} = 0.548$$

$$H_0: p \leq 0.5 \quad H_1: p > 0.5 \quad z = \frac{\hat{p} - p}{\sigma_p}, \sigma_p = \sqrt{\frac{p(1-p)}{n}} = 0.022, z = 2.15$$

$$p\text{-val} = 1 - \Phi_{\text{norm}}(z, 0, 1) = 0.015 < 0.05 \therefore \text{statistically significant}$$

$\therefore > 50\%$ against it (H_1).

6.3.8 A grinding machine will be qualified for a particular task if it can be shown to produce less than 8% defective parts. In a random sample of 300 parts, 12 were defective. On the basis of these data, can the machine be qualified?

6.3.8

Sample 300, 12 defective, need to be <8% defi. is machine qual.

$$n=300, k=12, p=0.08 \quad \hat{p} = \frac{k}{n} = 0.04$$

$$H_0: p \leq 0.08 \quad H_1: p > 0.08 \quad z = \frac{\hat{p} - p}{\sigma} = -2.55$$

p-val = pnorm(z, 0, 1) = 0.0053. This observed fraction of defect parts (0.04) is much less than acceptance lvl so machine is good

8.3 Problem Set

Verkefni 8

liljag18

Verkefni 8

liljag18
Lilja Ýr Guð.

1. A sample of 50 Aluminum plates has $\bar{x}_{\text{thick}} = 5.20 \text{ mm}$ with $sd = 0.50 \text{ mm}$

- a) Find the 95% confidence interval for population mean thickness.

We have to find the confidence interval $\bar{x} \pm z_{\alpha/2} \cdot \sigma_x$ for population mean μ .

First off we do $1-\alpha = 0.95 \Rightarrow \alpha = 0.05$ which we use to find z

using a table in the book. This gives us $z_{0.025} = 1.96$. Using all of this we get

$$\sigma_x = \frac{0.50}{\sqrt{50}} \text{ giving us a confidence interval of } 5.20 \pm 1.96 \left(\frac{0.50}{\sqrt{50}} \right) = 5.20 \pm 0.14. \quad \square$$

- b) What sample size is needed to obtain the 95% confidence interval $(5.10, 5.30)$?

We have the same z as in (a) since we are using the same confidence interval level $\therefore z = 1.96$. The sample mean is easily found

with $\bar{x} = \frac{5.10 + 5.30}{2} = 5.20$. Now we have the equation

$$z = \frac{(\bar{x}-\mu)\sqrt{n}}{\sigma} \Rightarrow n = \left(\frac{\bar{x}-\mu}{z \cdot \sigma} \right)^2. \text{ We plug and chug (with } \bar{x}-\mu = 5.30 - 5.10 = 0.2) \text{ and get } n = \left(\frac{1.96 \cdot 0.5}{0.2} \right)^2 = 25. \quad \square$$

- c) What is the 99% confidence interval in the second case?

Much like we did in (a) we get $1-\alpha = 0.99 \Rightarrow \alpha = 0.01$. This gives us $z_{0.005} = -2.57$. Then we do the same as in (a) and get

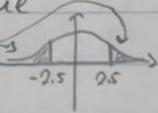
$$\bar{x} \pm z_{0.005} \left(\frac{sd}{\sqrt{n}} \right) = 5.20 \pm 2.57 \left(\frac{0.2}{\sqrt{25}} \right) = 5.20 \pm 0.10. \quad \square$$

Verkefni 8

liljag18

2. An instrument designed to measure thin layers of paint is tested. To determine the accuracy of the instrument, it was used to take 100 measurements of a layer with a known thickness of 100 micron. The mean of the measurements was 115 micron w/ a $sd = 60$ micron.

a) Perform a z test to compare the mean of the measurements w/ the real thickness and calculate the 2 sided p-value

$$\text{We get that } z = \frac{(\bar{x} - \mu)}{sd} = \frac{(115 - 100)}{60} = \frac{5}{60} = 2.5$$


$$\text{Then we use r and get } p\text{-value} = 2 \text{pnorm}(-z, 0, 1)$$

$$= 0.0124 \quad \square$$

Two tests are $H_0: \mu = 100$ and $H_1: \mu \neq 100$.

b) Do you believe the instrument is properly calibrated? Explain.

No because the probability of getting these results is 0.012 which is really low (and $p\text{-value} < 0.05$ so we reject H_0 in favor of H_1 , which essentially means the instrument is not properly calibrated.) \square

8.4 Problem Set Solution

1. A sample of 50 Aluminum plates has a mean thickness 5.20 mm with a standard deviation 0.50 mm.

a) Find the 95% confidence interval for the population mean thickness.

$$x=5.20 \quad \sigma_x = 0.50/\sqrt{50} = 0.0707 \quad 95\%CI = [x - 1.96\sigma_x, x + 1.96\sigma_x] = [5.06, 5.34]$$

b) What sample size is necessary to obtain the 95% confidence interval [5.10, 5.30]?

With the new sample size n' and the new std. σ'_x

$$x - 5.10 = 1.96\sigma'_x = \frac{0.5}{\sqrt{n}} \Rightarrow n' = \left(\frac{1.96 \cdot 0.5}{0.1}\right)^2 = 96.$$

c) What is the 99% confidence interval in the second case?

$$\sigma'_x = \frac{0.5}{\sqrt{96}} = 0.051 \Rightarrow 99\% CI = [x - 2.58\sigma'_x, x + 2.58\sigma'_x] = [5.07, 5.33].$$

2. An instrument designed to measure thin layers of paint is tested. To determine the accuracy of the instrument, it was used to take 100 measurements of a layer with a known thickness of 100 micron. The mean of the measurements was 115 micron with a standard deviation of 60 micron.

a) Perform a z test to compare the mean of the measurements with the real thickness and calculate the 2-sided p-value.

$$z\text{-score: } z = \frac{115-100}{60/\sqrt{100}} = 2.5 \quad p\text{-value} = 2 * \text{pnorm}(-z, 0, 1) = 0.0124.$$

b) Do you believe the instrument is properly calibrated? Explain.

The mean measurement is significantly different than the expected value. The instruments needs recalibration.

9 Week 9

9.1 Concepts

- The p-value returned by a chi square test is always one sided.
- In a chi-square test with 6 DOF the calculated test statistic chi squared is equal to 16. The p_value is between 0.013 and 0.014
- The electroencephalography data discussed in the presentation of today was analysed using t test with paired data
- The data table for a chi square test of independence has 3 rows and 3 columns. It has 4 degrees of freedom
- In a chi-square test the expected values in each category are calculated using the null hypothesis
- In a study on gender balance the proportions of girls and boys are tested in 100 universities with a chi square test. A result can be considered significant if the p-value is lower than 0.0005
- The probability distribution used in a Fisher test is the hypergeometric distribution
- A false positive result corresponds to a p-value < 0.05
- The statistical power of a test increases by increasing the sample size
- What is the odds ratio attributed to the genetic variant associated with myocardial infarction discovered at Decode Genetics in 2007? Between 1.2 and 1.3

9.2 R code

```
1-pnorm(z,0,1)      #calculate p_value
setwd("/Users/liljayrgudmundsdottir/Documents/HR_Files/5onn/Tolfraedi/
datasets/Ch6")
data <- read.csv('ex6-7-2.csv')
disk_tablets = data$Disk
oval_tablets = data$Oval
t.test(disk_tablets, oval_tablets)      #problem 6.7.2
data_A = data$Brand.A
data_B = data$Brand.B
ttest(data_A, data_B)
t.test(data_A, data_B)      #problem 6.7.14
data_60 = data$X60.degrees
data_61 = data$X61.degrees
t.test(data_60, data_61)      #problem 6.7.15
t.test(FW, SW, var.equal=TRUE)      #problem 6.7.6
t.test(B, G, alt="two.sided", paired=TRUE)      #problem 6.8.1
x=c(467,191,42); y=c(445,171,34); z=c(254,129,93)
chisq.test(rbind(x, y, z))      #6.2.10 chi square
```

9.3 Example Problem Set

Z-test for difference of population means:

6.5.4 The article "Wired: Energy Drinks, Jock Identity, Masculine Norms, and Risk Taking" (K. Miller, Journal of American College Health, 2008:481–489) reports that in a sample of 413 male college students, the average number of energy drinks consumed per month was 2.49 with a standard deviation of 4.87, and in a sample of 382 female college students, the average was 1.22 with a standard deviation of 3.24. Can you conclude that the mean number of energy drinks is greater for male students than for female students?

6.5.4

sample = 413 \bar{x} drinks/month = 2.49, $sd = 4.87$ $H_0: \mu_x - \mu_y = 0$

sample = 382 \bar{y} drinks/month = 1.22, $sd = 3.24$ $H_1: \mu_x - \mu_y \neq 0$

$\bar{x} - \bar{y} = 2.49 - 1.22 = 1.27$, drawn from a normal population

with mean $\mu_x - \mu_y$ and variance $\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}$

H_0 gives $\mu_x - \mu_y = 0$. Now $sd' = \sqrt{4.87^2/413 + 3.24^2/382} = 0.2914$

Null distr. of $\bar{x} - \bar{y}$ is $\bar{x} - \bar{y} \sim N(0, 0.2914^2)$

The z-score is $z = \frac{(\bar{x} - \bar{y}) - 0}{0.2914} = \frac{1.27 - 0}{0.2914} = 4.36$

$1 - pnorm(z, 0, 1) = 6.5 \times 10^{-6}$ \therefore men drink larger # of drinks

6.5.9 Are low-fat diets or low-carb diets more effective for weight loss? This question was addressed in the article “Comparison of the Atkins, Zone, Ornish, and LEARN Diets for Change in Weight and Related Risk Factors Among Overweight Premenopausal Women: The A TO Z Weight Loss Study: A Randomized Trial” (C. Gardner, A. Kiazand, et al., Journal of the American Medical Association, 2007:969–977). A sample of 77 subjects went on a low-carbohydrate diet for six months. At the end of that time the sample mean weight loss was 4.7 kg with a sample standard deviation of 7.2 kg. A second sample of 79 subjects went on a low-fat diet. Their sample mean weight loss was 2.6 kg with a standard deviation of 5.9 kg.

- Can you conclude that the mean weight loss is greater for those on the low carbohydrate diet?
- Can you conclude that the mean weight loss on the low-carbohydrate diet is more than 1 kg greater than that of the low-fat diet?

6.5.9

$$n_x = 77 \quad \bar{X} = 4.7 \quad s_x = 7.2; \quad n_y = 79 \quad \bar{Y} = 2.6 \quad s_y = 5.9$$

$$a) H_0: \mu_x - \mu_y \leq 0 \quad vs. \quad H_1: \mu_x - \mu_y > 0$$

$$\text{standard deviation: } \sqrt{7.2^2/77 + 5.9^2/79} = 1.055$$

$$z\text{-score: } z = (4.7 - 2.6 - 0)/1.055 = 1.9898 > 0 \rightarrow \text{alt hyp}$$

alt hyp: $\mu_x - \mu_y > \Delta$ making p-val the area to the right of z

$$P = 1 - \text{pnorm}(z, 0, 1) = 0.0233 < 0.05$$

Mean weight loss is greater for those on low-carb diet

$$b) H_0: \mu_x - \mu_y \leq 1 \quad vs. \quad H_1: \mu_x - \mu_y > 1$$

sd: same as in (a), values haven't changed = 1.055

$$z\text{-score: } z = (4.7 - 2.6 - 1)/1.055 = 1.0426 > 1 \rightarrow \text{alt hyp}$$

$$P = 1 - \text{pnorm}(z, 0, 1) = 0.1486 > 0.05$$

Cannot conclude weight loss had 1kg diff. ($\bar{X} = \bar{Y} + 1 \not\in \Delta$)

t-test:

6.7.2 In a study of the relationship of the shape of a tablet to its dissolution time, 6 disk-shaped ibuprofen tablets and 8 oval-shaped ibuprofen tablets were dissolved in water. The dissolve times, in seconds, were as follows:

Disk: 269.0 249.3 255.2 252.7 247.0 261.6

Oval: 268.8 260.0 273.5 253.9 278.5 289.4 261.6 280.2

Can you conclude that the mean dissolve times differ between the two shapes?

```
> t.test(disk_tablets, oval_tablets)
```

Welch Two Sample t-test

```
data: disk_tablets and oval_tablets
t = -2.7757, df = 11.961, p-value = 0.01683
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-26.667185 -3.207815
sample estimates:
mean of x mean of y
255.8000 270.7375
```

Yes we can conclude that the mean dissolve times differ because p_value is 0.01683 << 0.05. Unequal Variance.

6.7.14 The data file Datasets.zip/datasets/Ch5/ex5-6-11.csv shows measurements of the sodium content (in grams) in samples of two brands of chocolate bar. Can you conclude that the mean sodium content is higher for brand B than for brand A?

```
> t.test(data_A, data_B)
```

Welch Two Sample t-test

```
data: data_A and data_B
t = -3.1601, df = 19.552, p-value = 0.005024
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-6.466653 -1.319501
sample estimates:
mean of x mean of y
33.00000 36.89308
```

Yes we can conclude that because p_value is << 0.05.

6.7.15 The data file Datasets.zip/datasets/Ch5/ex5-6-12.csv shows the permeability coefficient (in $10^{-6} \text{cm}^3(\text{STP})/\text{cm} \cdot \text{s} \cdot \text{MPa}$) of CO₂ measured polyethylene at both 60°C and 61°C. Can you conclude that the mean permeability coefficient at 60°C differs from that at 61°C?

```
> t.test(data_60, data_61)
```

Welch Two Sample t-test

```
data: data_60 and data_61
t = -1.0236, df = 14.983, p-value = 0.3223
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-7.133668 2.505097
sample estimates:
mean of x mean of y
60.40000 62.71429
```

No there is not a difference because the p_value is >> 0.05.

6.7.6 Two weights, each labeled as weighing 100 g, are each weighed several times on the same scale. The results, in units of μg above 100 g, are as follows:

First weight:	53	88	89	62	39	66
Second weight:	23	39	28	2	49	

Since the same scale was used for both weights, and since both weights are similar, it is reasonable to assume that the variance of the weighing does not depend on the object being weighed. Can you conclude that the weights differ?

```
> t.test(FW, SW, var.equal=TRUE)
```

Two Sample t-test

```
data: FW and SW
t = 3.3306, df = 9, p-value = 0.008791
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
12.17926 63.75407
sample estimates:
mean of x mean of y
66.16667 28.20000
```

There they differ because p_value << 0.05. Equal Variance.

6.8.1 The article “Improved Bioequivalence Assessment of Topical Dermatological Drug Products Using Dermatopharmacokinetics” (B. N’Dri-Stempfer, W. Navidi, R. Guy, and A. Bunge, Pharmaceutical Research, 2009:316–328) described a study comparing the amounts of econazole nitrate absorbed into human skin for several formulations of antifungal ointment. Both a brand name and generic drug were applied to the arms of 14 subjects, and the amounts absorbed, in $\mu\text{g}/\text{cm}^2$, were measured. Following are the results. Can you conclude that the mean amount absorbed differs between the brand name and the generic drug?

Brand Name	Generic	Difference
2.23	1.42	0.81
1.68	1.95	-0.27
1.96	2.58	-0.62
2.81	2.25	0.56
1.14	1.21	-0.07
3.20	3.01	0.19
2.33	2.76	-0.43
4.06	3.65	0.41
2.92	2.89	0.03
2.92	2.85	0.07
2.83	2.44	0.39
3.45	3.11	0.34
2.72	2.64	0.08
3.74	2.82	0.92

```
#both types of ointment tested on the same person -> paired samples
#For T tests the alternative hypothesis can be:
#One sided, H1: mu > mu0, in R this is denoted as: alternative = "greater"
#One sided, H1: mu < mu0, in R this is denoted as: alternative = "less"
#Two sided, H1: mu not equal to mu0, in R this is denoted as:
#alternative = "two.sided"
#In this problem „can you conclude that the means differ?”
# -> mu not equal to mu0 -> two sided
#Now we will take some examples:
```

```
> t.test(B, G, alt="two.sided", paired=TRUE)
```

Paired t-test

```
data: B and G
t = 1.4593, df = 13, p-value = 0.1682
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.0826935 0.4269792
sample estimates:
mean of the differences
0.1721429
```

They don't differ because pvalue >> 0.05. Paired Samples.

6.8.5 Two formulations of a certain coating, designed to inhibit corrosion, are being tested. For each of eight pipes, half the pipe is coated with formulation A and the other half is coated with formulation B. Each pipe is exposed to a salt environment for 500 hours. Afterward, the corrosion loss (in μm) is measured for each formulation on each pipe.

Pipe	A	B
1	197	204
2	161	182
3	144	140
4	162	178
5	185	183
6	154	163
7	136	156
8	130	143

Can you conclude that the mean amount of corrosion differs between the two formulations?

```
> t.test(A, B, alt="two.sided", paired=TRUE)
```

Paired t-test

```
data: A and B
t = -3.0151, df = 7, p-value = 0.01952
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-17.842571 -2.157429
sample estimates:
mean of the differences
-10
```

They differ because p_value is less than 0.05.

6.8.10 A group of eight individuals with high cholesterol levels were given a new drug that was designed to lower cholesterol levels. Cholesterol levels, in mg/dL, were measured before and after treatment for each individual, with the following results:

Subject	Before	After
1	283	215
2	299	206
3	274	187
4	284	212
5	248	178
6	275	212
7	293	192
8	277	196

- a. Can you conclude that the mean cholesterol level after treatment is less than the mean before treatment? b. Can you conclude that the reduction in mean cholesterol level after treatment is greater than 75 mg/dL?

```
> t.test(B, A)

Welch Two Sample t-test

data: B and A
t = 10.984, df = 13.789, p-value = 3.357e-08
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 63.85366 94.89634
sample estimates:
mean of x mean of y
 279.125    199.750
```

They differ because p_value is less than 0.05.

Chi-square test:

6.10.2 At an assembly plant for light trucks, routine monitoring of the quality of welds yields the following data:

		Number of Welds		
		High Quality	Moderate Quality	Low Quality
Day Shift		467	191	42
Evening Shift		445	171	34
Night Shift		254	129	17

Can you conclude that the quality varies among shifts?

- State the appropriate null hypothesis.
- Compute the expected values under the null hypothesis.
- Compute the value of the chi-square statistic.
- Find the P_value. What do you conclude?

6.10.2

a) Null hypothesis: The probability distributions do not depend on shift

b) Use $E_{ij} = \frac{O_{i\cdot} O_{\cdot j}}{O_{\cdot\cdot}}$

$E_{11} = \frac{700 \cdot 1166}{1750} = 466.4$	$E_{12} = \frac{700 \cdot 491}{1750} = 196.4$	$E_{13} = \frac{700 \cdot 93}{1750} = 37.2$	$E_{21} = \frac{650 \cdot 1166}{1750} = 433.1$	$E_{22} = \frac{650 \cdot 491}{1750} = 182.4$	$E_{23} = \frac{650 \cdot 93}{1750} = 34.5$	$E_{31} = \frac{400 \cdot 1166}{1750} = 266.5$	$E_{32} = \frac{400 \cdot 491}{1750} = 112.2$	$E_{33} = \frac{400 \cdot 93}{1750} = 21.3$	467	191	42	700	
445	171	34	1166	491	93	1750	400	1750	445	171	34	650	
254	129	17	254	129	17	400	254	129	17	1166	491	93	1750

c) $\chi^2 = \sum_{i=1}^3 \sum_{j=1}^3 \frac{(O_{ij} - E_{ij})^2}{E_{ij}} = 85.163$ (Using chisq.test(rbind(x,y,z)))

$x = c(467, 191, 42)$; $y = c(445, 171, 34)$; $z = c(254, 129, 17)$

d) $P\text{-value} \approx 2.2 \times 10^{-16}$

6.10.3 The article “Inconsistent Health Perceptions for US Women and Men with Diabetes” (M. McCollum, L. Hansen, et al., Journal of Women’s Health, 2007:1421–1428) presents results of a survey of adults with diabetes. Each respondent was categorized by gender and income level. The numbers in each category (calculated from percentages given in the article) are presented in the following table.

Poor	Near Poor	Low Income	Middle Income	High Income	
Men	156	77	253	513	604
Women	348	152	433	592	511

Can you conclude that the proportions in the various income categories differ between men and women?

6.10.3

T

156	77	253	513	604	1603
348	152	433	592	511	2036
T	504	229	686	1105	3639

$$E: \begin{array}{ccccc} 222.01 & 100.88 & 302.19 & 486.76 & 491.16 \\ 281.99 & 128.12 & 383.81 & 618.24 & 623.84 \end{array}$$

$$(2-1)(5-1) = 4 \text{ degrees of freedom}$$

Computer gives us $\chi^2 = 108.35$ and $p = 2.2 \times 10^{-16} << 0.05$

Proportions in var. income categories differ between M and W

6.10.4 The article “Analysis of Time Headways on Urban Roads: Case Study from Riyadh” (A. Al-Ghamdi, Journal of Transportation Engineering, 2001:289–294) presents a model for the time elapsed between the arrival of consecutive vehicles on urban roads. Following are 137 arrival times (in seconds) along with the values expected from a theoretical model.

Time	Observed	Expected
0-2	18	23
2-4	28	18
4-6	14	16
6-8	7	13
8-10	11	11
10-12	11	9
12-18	10	20
18-22	8	8
>22	30	19

Can you conclude that the theoretical model does not explain the observed values well?

```
> chisq.test(rbind(O,E))
```

Pearson's Chi-squared test

```
data: rbind(O, E)
X-squared = 10.72, df = 8, p-value = 0.2181
```

p_val is 0.2181 >> 0.05 so no.

6.10.12 The article “Determination of Carboxyhemoglobin Levels and Health Effects on Officers Working at the Istanbul Bosphorus Bridge” (G. Kocasoy and H. Yalin, Journal of Environmental Science and Health, 2004:1129-1139) presents assessments of health outcomes of people working in an environment with high levels of carbon monoxide (CO). Following are the numbers of workers reporting various symptoms, categorized by work shift. The numbers were read from a graph.

Shift	Morning	Evening	Night
Influenza	16	13	18
Headache	24	33	6
Weakness	11	16	5
Shortness of Breath	7	9	9

Can you conclude that the proportions of workers with the various symptoms differ among the shifts?

```
> data <- read.csv('ex6-10-12.csv')
> data
      X Morning Evening Night
1 Influenza     16      13    18
2 Headache      24      33     6
3 Weakness      11      16     5
4 Shortness of Breath    7      9     9
> I = c(16,13,18)
> H = c(24,33,6)
> W = c(11,16,5)
> S = c(7,9,9)
> chisq.test(rbind(W,I,S,H))

Pearson's Chi-squared test

data: rbind(W, I, S, H)
X-squared = 17.572, df = 6, p-value = 0.007397
```

There are differences between the shifts.

9.4 Problem Set

Problem 1:

Verkefni 9

Lilja Ýr Guð.
litjog 18

Problem 1

a) Formulate the null and alternative hypothesis and decide what type of t-test is appropriate.

$H_0: \mu_d = \mu_o - \mu_n = 0$ $H_1: \mu_d = \mu_o - \mu_n \neq 0$, 2 sided test + paired t-test
where μ_o = Old column and μ_n = New column. \square

b) Perform the test and conclude

2-sided test is performed since value can be greater than or less than 0. t-test will also be a paired t-test because the columns have same # of elements and data is the difference between the pairs.

R command: `t.test(data$Old, data$New, alt="two-sided", paired=TRUE)`
(see picture for result)

Given the found results hypothesis H_0 is accepted and the two samples cannot be considered different. (because $p\text{-val}=0.5 > 0.05$) \square

c) Suppose the value of t (the test statistic) would be the double of what you obtained. What would be the p-value?

When t-value is doubled from (b) we get $t = -1.26492$, $df = 15$ (Dof don't change). Then the p-val can be calculated in R:
`p(t(-1.26492, 15)) = 0.1126`. \square

```

> data = read.csv("ex6-S-18.csv", header = TRUE)
> data
   Old  New
1 16.3 15.9
2 15.9 16.2
3 15.8 16.0
4 16.2 15.8
5 16.1 16.1
6 16.0 16.1
7 15.7 15.8
8 15.8 16.0
9 15.9 16.2
10 16.1 15.9
11 16.3 15.7
12 16.1 16.2
13 15.8 15.8
14 15.7 15.8
15 15.8 16.2
16 15.7 16.3
> t.test(data$Old,data$New,alt="two.sided",paired=TRUE)

```

Paired t-test

```

data: data$Old and data$New
t = -0.63246, df = 15, p-value = 0.5366
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.2185059  0.1185059
sample estimates:
mean of the differences
-0.05

```

```

> t2 = 2*(-0.63246)
> pt(t2, 15)
[1] 0.1125996

```

```

data = read.csv("ex6-5-18.csv", header = TRUE)
t.test(data$Old, data$New, alt="two.sided", paired = TRUE)
t2 = 2*(-0.63246) # t = -0.63246 from the test
pt(t2, 15) #15 degrees of freedom

```

Problem 2:

Problem 2

The table shows the preferences for running shoes (from 3 unknown brands) of 630 individuals grouped by age: under 35 and over 35. Can you say that the preferred brand depends on age?

Age group	Adidas	Reebok	Nike
Under 35	174	132	90
Over 35	84	72	82

- a) Formulate the null and alternative hypothesis, perform a chi squared test, and conclude.

H_0 : same probability for people to buy from each brand regardless (or independent) of age.

We put it all in google sheets (excel for mac essentially)

First image is sum of row and column values (O_{11}, O_{12}) with the bottom right cell containing sum of all values (O).

The next image uses formula $E_{ij} = \frac{O_{1j} \cdot O_{i1}}{O}$ $\frac{(O_{ij} - E_{ij})^2}{E_{ij}}$
Then we have a difference table to calculate E_{ij} , in the image below that marked difference table.

Then we sum it all up (all the values in the difference table) to get the value: 8.826 (chisq val)

Degrees of freedom are: $(\text{rows}-1) \times (\text{cols}-1) = 1 \cdot 2 = 2$

Using R (using chisq val) we get

$$p\text{-val} = pchisq(8.826, df=2, \text{lower.tail}=FALSE) \\ = 0.0121187 < 0.05$$

This shows that the null hyp. is rejected and age does matter when buying shoes from these brands. □

- b) What is the probability that the value of χ^2 is greater than the obtained value assuming random data?

No idea because I don't know what obtained value is being referred to here. □

Sum of rows and columns				Sum I:
	174	132	90	
	84	72	78	234
Sum J:	258	204	168	630

	J = 1	J = 2	J = 3	
I = 1	162.1714286	128.2285714		105.6
I = 2	95.82857143	75.77142857		62.4

Difference table				
0.8627604953	0.1109243697	2.304545455		
1.460056223	0.1877181642	3.9		
		SUM:	8.826004707	

9.5 Problem Set Solution

1. The data file ex6-S-18.csv (from the data collection of the textbook or Canvas/Modules/Other materials/Datasets.zip) contains the diameters of two samples of bearings. One sample is produced with an old technology and the other sample with a new technology. The project manager is asking you for a t-test to find out if the diameters of the two samples can be considered different or not.

a) Formulate the null and the alternative hypotheses and decide what type of t-test is appropriate.

H_0 : the diameters of the two samples are equal H_1 : the diameters of the two samples are different

t-test with unequal variances.

b) Perform the test and conclude.

p-value = 0.485 For an (incorrect) paired test p-value = 0.536

p-value > 0.05 $\Rightarrow H_0$ is accepted, the diameters cannot be considered different

c) Suppose the value of t (the test statistic) would be the double of what you obtained. What would be the p-value?

For a correct t-test $t = -0.70711; 2 * pt(2 * t, 29.7) \approx 2 * pt(2 * t, 30) \approx 0.168$

In case of an (incorrect) paired test $t = -0.63246$; $2 * \text{pt}(2 * t, 15) \approx 0.225$

2. The table shows the preferences for running shoes, from three known brands, of 630 individuals grouped by age: under 35 and over 35. Can you say that the preferred brand depends on age?

Age group	Adidas	Reebok	Nike
Under 35	174	132	90
Over 35	84	72	78

a) Formulate the null and alternative hypotheses, perform a chi square test, and conclude.

H_0 : the preferred brand does not depend on the age group H_1 : the preferred brand depends on the age group.

`x=c(174,132,90); y=c(84,72,78)`

`chisq.test(rbind(x,y))`

`p-value = 0.01212`

$p\text{-value} < 0.05 \Rightarrow$ The preferred brand and the age are not independent variables. The preferred brand (seems to) depend on the age group.

b) What is the probability that the value of X^2 (the test statistic) is greater than the obtained value, assuming random data?

The obtained value of X^2 is 8.826. According to the definition, the p-value is $P(X^2 > 8.826) = 0.01212$. So no additional calculation is needed. But, optionally, or alternatively, it is a good idea to check: $P(X^2 > 8.826) = 1 - \text{pchisq}(8.826, 2) = 0.01212$

10 Week 10

10.1 Concepts

- The covariance of two random variables is always positive, negative, or zero
- The correlation coefficient of two random variables is always between -1 and 1
- The correlation coefficient of two independent variables is expected to be zero
- The 95% confidence interval of the correlation coefficient of two random variables is [-0.15,0.23]. Based on this information we can conclude that the two variables are independent
- The correlation coefficient of two random variables estimated from a data sample is $r=-0.6$. These two variables can be considered correlated only if the p-value is below 0.05
- The residuals are the differences between the response data y_i and the response predicted by the linear fit
- In a linear regression of the response y to the independent variable x the slope coefficient is 0.5 and the intercept coefficient is 1.5. The predicted response for $x=0$ is 1.5
- In a regression analysis of the response y to the independent variable x the slope coefficient was 0.5. If now x is considered the response and y the independent variable, the slope will depend on the correlation coefficient
- If the slope coefficient of the response y to the independent variable x is negative then the correlation coefficient of x and y is negative
- The regression coefficients are obtained by minimizing the sum of all squared residuals

10.2 R code

```
mat = rbind(c(17,9), c(16,22))
chisq.test(mat)
fisher.test(mat)
cor.test(x,y)      # testing correlation of x and y
cov(x,y)/(sd(x)*sd(y))      # gives us r used in formula 7.4
cov(dom, fre)
```

10.3 Example Problem Set

Fisher test of a 2x2 contingency table. Experimental approach.

A prestigious university organizes an admission exam for the Master's program. The admitted students are 17 girls and 9 boys. The rejected students are 16 girls and 22 boys. Do you see a gender bias in this process? Make a 2x2 contingency table and perform a Fisher exact test. Change the numbers and repeat the test several times. Observe the relation between data and p-value. Increase the difference between numbers, change the sample size, etc. What if you have all numbers equal, or two pairs of equal numbers? Understand by doing. Try also a chi square test and compare.

Fisher test of a 2x2 contingency table

admitted: 17 girls 9 boys

rejected: 16 girls 22 boys

17	9	26
16	22	38
33	31	64

mat = rbind(c(17, 9), c(16, 22))

chisq.test(mat)

fisher.test(mat)

test other values

Correlation coefficients and correlation tests:

Example 7.3 In a study of reaction times, the time to respond to a visual stimulus (x) and the time to respond to an auditory stimulus (y) were recorded for each of 10 subjects. Times were measured in ms. The results are presented in the following table.

x	161	203	235	176	201	188	228	211	191	178
y	159	206	241	163	197	193	209	189	169	201

Find a 95% confidence interval for the correlation between the two reaction times.

Ex. 7.3

i tolfr. db folder

`data = read.csv("exm7-3.csv")`

`x = data$x; y = data$y`

`cor.test(x,y)`

`cov(x,y)/(sd(x)*sd(y))`

`cor.test(x,y, alt = "g")`

`cor.tst(x,y, alt = "l")`

`cor.tst(x,y)`

`x = runif(100,0,1); y = runif(100,0,1)`

`cor.test(x,y)`

`cov(x,y)/(sd(x)*sd(y))`

see video

```
> x=data$x  
> y=data$y  
> cor.test(x,y)
```

Pearson's product-moment correlation

```
data: x and y  
t = 3.991, df = 8, p-value = 0.004  
alternative hypothesis: true correlation is not equal to 0  
95 percent confidence interval:  
 0.3830003 0.9549492  
sample estimates:  
 cor  
 0.8158796
```

```
> cov(x,y)/(sd(x)*sd(y))  
[1] 0.8158796  
> cor.test(x,y,alt='g')
```

Pearson's product-moment correlation

```
data: x and y  
t = 3.991, df = 8, p-value = 0.002  
alternative hypothesis: true correlation is greater than 0  
95 percent confidence interval:  
 0.4797593 1.0000000  
sample estimates:  
 cor  
 0.8158796
```

```
> cor.test(x,y,alt='l')
```

Pearson's product-moment correlation

```
data: x and y  
t = 3.991, df = 8, p-value = 0.998  
alternative hypothesis: true correlation is less than 0  
95 percent confidence interval:  
 -1.0000000 0.9431764  
sample estimates:  
 cor  
 0.8158796
```

7.1.1 Compute the correlation coefficient for the following data set

x	1	2	3	4	5	6	7
y	2	1	4	3	7	5	6

7.1.1

Corr. coef. for : $x \mid 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7$
 $y \mid 2 \ 1 \ 4 \ 3 \ 7 \ 5 \ 6$

$$\text{corr test}(x, y) = 0.82142 \quad p\text{-val} = 0.02345$$

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} = 0.82142$$

7.1.2 For each of the following data sets, explain why the correlation coefficient is the same as for the data set in Exercise 1.

7.1.2

a)	x	1	2	3	4	5	6	7	}
	y	5	4	7	6	10	8	9	
b)	x	11	21	31	41	51	61	71	}
	y	5	4	7	6	10	8	9	
c)	x	53	43	73	63	103	83	93	}
	y	4	6	8	10	12	14	16	

linear transformation of
x or y or both has
been performed, that's
why corr. cof. is same as 7.1.1

7.1.6 In a study of ground motion caused by earthquakes, the peak velocity (in m/s) and peak acceleration (in m/s^2) were recorded for five earthquakes. The results are presented in the following table (seen in the image).

7.1.6					
Velocity 1.54 1.60 0.95 1.30 2.92					
Acceleration 7.64 8.04 8.04 6.37 3.25					
$\text{read.table("ex7-1-6.csv", header=T, sep=",") = results}$					
Velocity = results\$Velocity; Acceleration = results\$Acceleration					
plot(V, A)					
cor.test(V, A) % gives -0.802, p-val=0.102					
See video, look at it better					

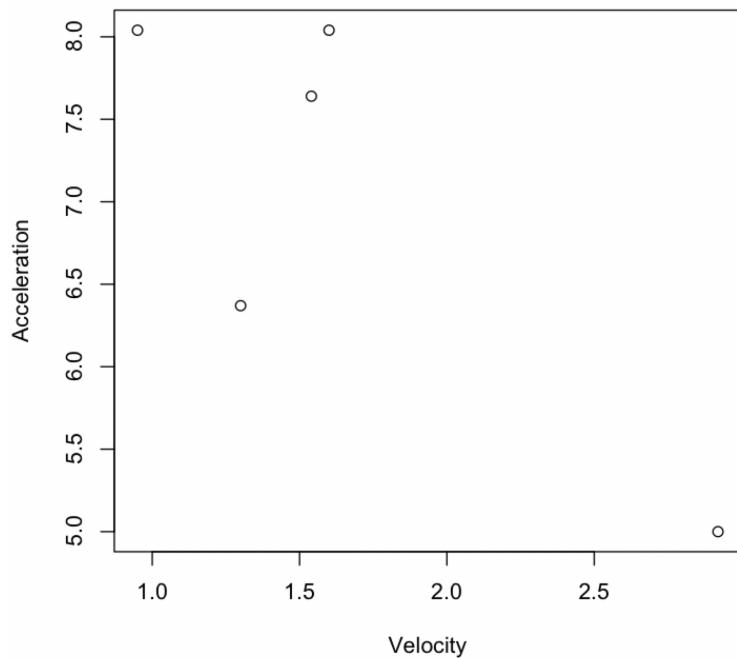
a. Compute the correlation coefficient between peak velocity and peak acceleration.

> `cor.test(V, A)`

Pearson's product-moment correlation

```
data: V and A
t = -2.3317, df = 3, p-value = 0.102
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
-0.9864057 0.2725283
sample estimates:
cor
-0.8027545
```

b. Construct a scatterplot for these data.



c. Is the correlation coefficient an appropriate summary for these data? Explain why or why not.

I'd say it is because there is a weak correlation between the Velocity and Acceleration as can be seen in the scatterplot.

d. Someone suggests converting the units from meters to centimeters and from seconds to minutes. What effect would this have on the correlation?

This would have no affect because it is just transforming the data a little bit and shouldn't affect the correlation.

7.1.9 Tire pressure (in kPa) was measured for the right and left front tires on a sample of 10 automobiles. Can you conclude that the correlation coefficient is positive? Calculate the p-value corresponding to this hypothesis. Use the data set ex7-1-9.csv.

```
cor.test(r,l,alt="great")
```

```
> setwd("/Users/liljayrgudmundsdottir/Documents/HR_Files/5onn/Tolfraedi/datasets/Ch7")
> data <- read.csv('ex7-1-9.csv')
> data
   Right Left
1    184 185
2    206 203
3    193 200
4    227 213
5    193 196
6    218 221
7    213 216
8    194 198
9    178 180
10   207 210
> r = data$Right
> l = data$Left
> cor.test(r,l,alt="great")

Pearson's product-moment correlation

data: r and l
t = 7.1965, df = 8, p-value = 4.638e-05
alternative hypothesis: true correlation is greater than 0
95 percent confidence interval:
0.7786318 1.0000000
sample estimates:
cor
0.9306981
```

General properties of covariance and correlation coefficients:

- 1 - Show that $\text{Cov}(x,y) = E(xy) - E(x)E(y)$
- 2 - Show that if x and y are independent random variables $\text{Cov}(x,y) = 0$
- 3 - Show that the correlation coefficient $\rho \in [-1,1]$. This is Problem 2.6.29

$\rightarrow 1 - \text{show that } \text{Cov}(x,y) = E(xy) - E(x)E(y)$ $(\text{Cov}(x,y)) = E((X-\mu_x)(Y-\mu_y)) = E(XY - \mu_x Y - \mu_y X + \mu_x \mu_y)$ $= E(XY) - E(\mu_x Y) - E(\mu_y X) + \mu_x \mu_y = E(XY) - \mu_x E(Y) - \mu_y E(X) + \mu_x \mu_y$ $E(X+Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (x+y) f_{xy}(x,y) dx dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f_{xy}(x,y) dx dy + \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f_{xy}(x,y) dx dy$ $= E(X) + E(Y)$ $E(XY) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_{xy}(x,y) dx dy = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f_x(x) f_y(y) dx dy$ $= (\int_{-\infty}^{\infty} x f_x(x) dx) (\int_{-\infty}^{\infty} y f_y(y) dy) = E(X)E(Y)$ $E(XY) - \mu_x E(Y) - \mu_y E(X) + \mu_x \mu_y = E(XY) - \mu_x \mu_y - \mu_x \mu_y + \mu_x \mu_y$ $= E(XY) - \mu_x \mu_y = E(XY) - E(X)E(Y).$	$2 - \text{Show that if } x \text{ and } y \text{ are independent random variables } \text{Cov}(x,y) = 0$ Assume we have $E(XY) = E(X)E(Y)$ and $\text{Cov}(X,Y) = E(XY) - E(X)E(Y)$ from above. Then $\text{Cov}(X,Y) = E(XY) - E(X)E(Y) = E(XY) - E(XY) = 0$
$\Rightarrow 3 - \text{Show that the correlation coefficient } \rho \in [-1,1] \text{ (problem 2.6.29)}$	

Joint distributions

2.6.2 In a certain community, levels of air pollution may exceed federal standards for ozone or for particulate matter on some days. For a particular summer season, let X be the number of days on which the ozone standard is exceeded and let Y be the number of days on which the particulate matter standard is exceeded. Assume that the joint probability mass function of X and Y is given in the following table:

x	y			$p_X(x)$
	0	1	2	
0	0.10	0.11	0.05	0.26
1	0.17	0.23	0.08	0.48
2	0.06	0.14	0.06	0.26
$p_Y(y)$	0.33	0.48	0.19	

a. Find the marginal probability mass function $p_X(x)$.

See values in the table above for $p_X(x)$. Now when we add them all together we get $\sum_i p_X(x) = 0.26 + 0.48 + 0.26 = 1$ which fulfills the property of the marginal probability mass function where the sum of all possible x and y have to add up to 1.

b. Find the marginal probability mass function $p_Y(y)$.

See values in the table above for $p_Y(y)$. Now when we add them all together we get

$\sum_j p_Y(y) = 0.33 + 0.48 + 0.19 = 1$ which fulfills the property of the marginal probability mass function where the sum of all possible x and y have to add up to 1.

c. Find μ_X .

We have $\mu_X = E(x) = \sum_i x_i p(x_i) = 0 * 0.26 + 1 * 0.48 + 2 * 0.26 = 1$.

d. Find μ_Y .

We have $\mu_Y = E(y) = \sum_j y_j p(y_j) = 0 * 0.33 + 1 * 0.48 + 2 * 0.19 = 0.86$.

e. Find σ_X .

We have $E(x^2) = \sum_i x_i^2 p(x_i) = 0^2 * 0.26 + 1^2 * 0.48 + 2^2 * 0.26 = 1.52$. This gives us $\sigma^2 = E(x^2) - (E(x))^2 = 1.52 - 1 = 0.52$

f. Find σ_Y .

We have $E(y^2) = \sum_j y_j^2 p(y_j) = 0^2 * 0.33 + 1^2 * 0.48 + 2^2 * 0.19 = 1.24$. This gives us $\sigma^2 = E(y^2) - (E(y))^2 = 1.24 - 0.86 = 0.38$

g. Find $Cov(X, Y)$.

First we have to find μ_{XY} because $Cov(X, Y) = \mu_{XY} - \mu_X \mu_Y$.

$$\mu_{XY} = 0 * 0 * 0.10 + 1 * 0 * 0.17 + 2 * 0 * 0.06 + 0 * 1 * 0.11 + 1 * 1 * 0.23 + 2 * 1 * 0.14 + 0 * 2 * 0.05 + 1 * 2 * 0.08 + 2 * 2 * 0.06 = 1 * 1 * 0.23 + 2 * 1 * 0.14 + 1 * 2 * 0.08 + 2 * 2 * 0.06 = 0.91$$

This gives us $Cov(X, Y) = 0.91 - 1 * 0.86 = 0.05$

h. Find $\rho_{X,Y}$.

$$\text{We get } \rho_{X,Y} = \frac{Cov(X,Y)}{\sigma_X \sigma_Y} = \frac{0.05}{0.52 * 0.38} = 0.25$$

i. Are X and Y independent? Explain.

For x and y to be considered independent $Cov(X, Y) = 0$ and $\rho = 0$. As you can see from previous calculations, neither is equal to 0. Correlation coefficient gives info on how much x and y depend on each other.

10.4 Problem Set

Lilja Ýr Guðmundsdóttir - liljag18

Verkefni 10

Verkefni 10

Lilja Ýr Guðst.
liljag18

1. Two binary random variables, x and y , with values 0 or 1, have the joint probability distribution $P(x=0, y=0) = P(x=1, y=1) = 0.2$, $P(x=0, y=1) = P(x=1, y=0) = 0.3$. See also table

	$y=0$	$y=1$
$x=0$	0.2	0.3
$x=1$	0.3	0.2

- a) Calculate covariance of x and y .

We put $X = [0.2, 0.3]$ and $Y = [0.3, 0.2]$ in R
and then do $\text{cov}(X, Y) = \underline{-0.005}$

- b) Calculate the correlation coefficient of x and y .

Again we use R to calculate. Using X and Y from (a) there

are two ways we can go about calculating this coefficient:

i) $\text{cor.test}(X, Y)$ = returns error "not enough finite observations"

since there needs to be 3 or more values and we only have 2.

ii) $\text{cov}(X, Y) / (\text{sd}(X) * \text{sd}(Y)) = \underline{-1} \quad \square$

Verkefni 10

2. A team of engineers performed an experiment in which accelerometers were placed on a bridge for the purpose of measuring the damping ratios of the vibrations at 18 different frequencies. The data file ex7-2-9.csv shows the measurements. The frequencies are measured in Hertz (Hz) meaning the number of vibrations per second.

- a) Calculate the covariance and the correlation coefficient of the two variables (damping ratio & frequency).

We get the data and split into arrays dam and fre as in picture. The covariance is then $\text{cov}(\text{dam}, \text{fre}) = -0.13738$ and the correlation coefficient (found with $\text{cor.test}(\text{dam}, \text{fre})$) is -0.91533. □

- b) Are the two variables correlated? Perform a correlation test and explain.

As per the correlation test done in (a) we got a coefficient of -0.91533 and the p-value is 1.016×10^{-7} and the 95 percent confidence interval is -0.968 to -0.783 (as seen in the picture). Since both values on the interval are negative we look at the p-value. Seeing as the value is less than $0.05 (1.016 \times 10^{-7} < 0.05)$ then there is a correlation between the variables. □

- c) Suppose we change the units of frequencies from Hz to number of vibrations per minute. Will the correlation coefficient change? Explain.

The correlation coefficient will not change because the coefficient is essentially the relationship between the z-scores of the two variables. So even though the numbers in one variable changes (in this case frequency by a factor of $\times 60$) and the standard deviation and mean change, the z-score does not and therefore the coefficient stays the same. □

Verkefni 10

Image for problem 2:

```

> data = read.csv("ex7-2-9.csv", header = TRUE)
> data
   Damping.Ratio Frequency
1       0.3        2.72
2       0.3        2.84
3       0.3        3.77
4       0.4        2.07
5       0.4        2.20
6       0.4        2.34
7       0.4        2.61
8       0.5        1.80
9       0.5        1.93
10      0.5        1.53
11      0.6        0.77
12      0.6        1.26
13      0.6        1.66
14      0.7        0.89
15      0.7        1.00
16      0.7        0.66
17      0.8        1.13
18      0.8        0.37
> dam = data$Damping.Ratio
> dam
[1] 0.3 0.3 0.3 0.4 0.4 0.4 0.4 0.5 0.5 0.5 0.6 0.6 0.6 0.6 0.7 0.7 0.7 0.7 0.8 0.8
> fre = data$Frequency
> fre
[1] 2.72 2.84 3.77 2.07 2.20 2.34 2.61 1.80 1.93 1.53 0.77 1.26 1.66 0.89 1.00 0.66
[17] 1.13 0.37
> cov(dam, fre)
[1] -0.1373758
> cor.test(dam, fre)

Pearson's product-moment correlation

data: dam and fre
t = -9.0917, df = 16, p-value = 1.016e-07
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
-0.9683741 -0.7831137
sample estimates:
cor
-0.9153284

>

```

10.5 Problem Set Solution

1. Two binary random variables, x and y, with values 0 or 1, have the joint probability distribution $P(x=0, y=0)=P(x=1, y=1)=0.2$, $P(x=0, y=1)=P(x=1, y=0)=0.3$. See also the table

	y=0	y=1
x=0	0.2	0.3
x=1	0.3	0.2

- a) Calculate the covariance of x and y. $\text{Cov}(x,y) = E(x^*y) - E(x)*E(y) = 0.2 - (0.5)^2 = -0.05$
 b) Calculate the correlation coefficient of x and y. $\rho(x,y) = \text{Cov}(x,y) / (\sigma_x * \sigma_y) = -0.05 / (0.25) = -0.2$

2. A team of engineers performed an experiment in which accelerometers were placed on a bridge for the purpose of measuring the damping ratios of the vibrations at 18 different frequencies. The data file ex7-2-9.csv (from the data collection of the textbook or Canvas / Modules / Other materials / Datasets.zip) shows the measurements. The frequencies are measured in Hertz (Hz), meaning the number of vibrations per second.
- a) Calculate the covariance and the correlation coefficient of the two variables (damping ratio and frequency).

$$\text{cov}(dr,f) = -0.1373758, \text{cor}(dr,f) = -0.9153284$$

- b) Are the two variables correlated? Perform a correlation test and explain.

`cor.test(dr,f)`

data: dr and f

t = -9.0917, df = 16, p-value = 1.016e-07

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

-0.9683741 -0.7831137

sample estimates:

cor -0.9153284

The p-value = 1.016e-07 < 0.05 indicates that the two variables are correlated.

- c) Suppose we change the units of frequencies from Hz to number of vibrations per minute. Will the correlation coefficient change? Explain.

No, the correlation coefficient does not change by scaling one variable. (Nor by performing other linear transformations of the variables.)

11 Week 11

11.1 Concepts

Under assumptions 1 through 4 (page 549),

- The quantities $\hat{\beta}_0$ and $\hat{\beta}_1$ are normally distributed random variables.
- The means of $\hat{\beta}_0$ and $\hat{\beta}_1$ are the true values β_0 and respectively.
- The standard deviations of $\hat{\beta}_0$ and $\hat{\beta}_1$ are estimated with

$$s_{\hat{\beta}_0} = s \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (7.36)$$

and

$$s_{\hat{\beta}_1} = \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} \quad (7.37)$$

where $s = \sqrt{\frac{(1 - r^2) \sum_{i=1}^n (y_i - \bar{y})^2}{n - 2}}$ is an estimate of the error standard deviation σ . Page 552

- The R squared coefficient represents the proportion of the response variance explained by the regression
- In a regression analysis the null hypothesis (H_0) about the slope coefficient is the true slope is zero
- In a regression analysis the alternative hypothesis (H_1) about the slope coefficient is the true slope is not zero
- If the p-value of the slope coefficient is larger than 0.05 we conclude that the data has no real trend
- The distribution used to calculate the p-value of the regression coefficient is t distribution with the number of DOF equal to the sample size -2

11.2 R code

```
reg=lm(L~W)
summary( reg )
attributes( reg )
reg$coefficients
reg$fitted.values
reg$residuals
confint( reg ) #confidence intervals
abline( reg , col="red" )
sum=0          # remember to reset sum to 0 every time
for( i in T) {sum = sum + (i-means(T))**2}
sum=0
for( i in c(1,2,3,4,5,6,7,8,9)){sum = sum + (T[ i ]-mean(T))*(C[ i ]-mean(C))}
```

```

cor.test(T, M, method = "pearson")
mreg=lm(S ~ M + T)      # multiple regression where S is on the y axis
confint(mreg)

```

11.3 Example Problem Set

Simple linear regression:

Examples 7.6-7.8 using table7-1.csv and converting the data in SI units. (Video) In each case verify the numbers obtained for the regression coefficients with the formulas obtained in the lecture material. (equivalent to Eqs.7.14 and 7.15 in the textbook). Flip the variables length and force, repeat the regression, and verify the result that the product of the former and present slopes = r^2 Show that the correlation coefficient does not depend on the choice American or European units.

Convert to SI units: L=L*2.54; W=W*4.448 Use European units cm, Newtons (SI)

```

> L = data$Measured.Length..y.
> L
[1] 5.06 5.01 5.12 5.13 5.14 5.16 5.25 5.19 5.24 5.46 5.40 5.57 5.47 5.53 5.61 5.59 5.61 5.75
[19] 5.68 5.80
> L
[1] 5.06 5.01 5.12 5.13 5.14 5.16 5.25 5.19 5.24 5.46 5.40 5.57 5.47 5.53 5.61 5.59 5.61 5.75 5.68 5.80
> W = data$Weight..x.
> W
[1] 0.0 0.2 0.4 0.6 0.8 1.0 1.2 1.4 1.6 1.8 2.0 2.2 2.4 2.6 2.8 3.0 3.2 3.4 3.6 3.8
> L=L*2.54; W=W*4.448 #Use European units cm, Newtons (SI)
> L
[1] 12.8524 12.7254 13.0048 13.0302 13.0556 13.1064 13.3350 13.1826 13.3096 13.8684 13.7160 14.1478
[13] 13.8938 14.0462 14.2494 14.1986 14.2494 14.6050 14.4272 14.7320
> W
[1] 0.0000 0.8896 1.7792 2.6688 3.5584 4.4480 5.3376 6.2272 7.1168 8.0064 8.8960 9.7856
[13] 10.6752 11.5648 12.4544 13.3440 14.2336 15.1232 16.0128 16.9024

```

$$\text{Formula 7.14: } \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Formula 7.15: } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

```

> reg=lm(L~W)
> reg

```

Call:
`lm(formula = L ~ W)`

Coefficients:
`(Intercept) W`
`12.6993 0.1168`

Use the formulas to check regression coefficient:

```

> beta1=(mean(W*L)-mean(W)*mean(L))/(mean(W**2)-mean(W)**2)
> beta1
[1] 0.1168492
> beta0=mean(L)-beta1*mean(W)
> beta0
[1] 12.69927

```

Now we flip L and W:

```

> beta1=(mean(W*L)-mean(L)*mean(W))/(mean(L**2)-mean(L)**2)
> beta1
[1] 8.124132
> beta0=mean(W)-beta1*mean(L)
> beta0
[1] -102.7421
> reg=lm(W~L)
> reg

```

Call:
`lm(formula = W ~ L)`

Coefficients:

(Intercept)	L
-102.742	8.124

Some regression stuff including r^2 :

```

> summary(reg)

Call:
lm(formula = L ~ W)

Residuals:
    Min      1Q  Median      3Q     Max 
-0.24432 -0.08652 -0.01359  0.09554  0.30509 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 12.699274  0.062926 201.81 < 2e-16 ***
W           0.116849  0.006365 18.36 4.2e-13 ***  
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.146 on 18 degrees of freedom
Multiple R-squared:  0.9493, Adjusted R-squared:  0.9465 
F-statistic: 337 on 1 and 18 DF,  p-value: 4.204e-13

> attributes(reg)
$names
[1] "coefficients"   "residuals"      "effects"       "rank"        "fitted.values" "assign"
[7] "qr"              "df.residual"   "xlevels"      "call"        "terms"        "model" 

$class
[1] "lm"

> reg$coefficients
(Intercept)          W
12.6992743 0.1168492
> reg$fitted.values
   1   2   3   4   5   6   7   8   9   10  11 
12.69927 12.80322 12.90717 13.01112 13.11507 13.21902 13.32297 13.42692 13.53087 13.63482 13.73876 
   12  13  14  15  16  17  18  19  20 
13.84271 13.94666 14.05061 14.15456 14.25851 14.36246 14.46641 14.57036 14.67431 
> reg$residuals
   1    2    3    4    5    6    7    8 
0.153125714 -0.077823308 0.097627669 0.019078647 -0.059470376 -0.112619398 0.012031579 -0.244317444 
   9   10   11   12   13   14   15   16 
-0.221266466 0.233584511 -0.022764511 0.305086466 -0.052862556 -0.004411579 0.094839398 -0.059909624 
   17   18   19   20 
-0.113058647 0.138592331 -0.143156692 0.057694286 
> confint(reg) #confidence intervals
   2.5 %  97.5 % 
(Intercept) 12.5670722 12.8314764 
W           0.1034767  0.1302216

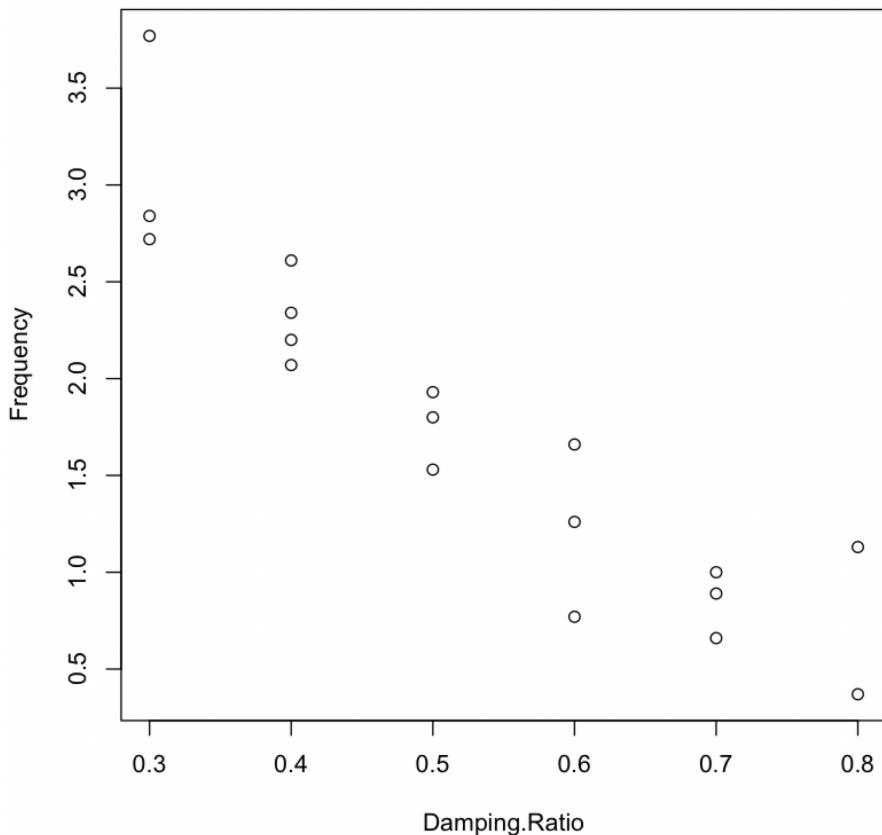
```

Not sure how to prove this though.

7.2.11 Structural engineers use wireless sensor networks to monitor the condition of dams and bridges. The article “Statistical Analysis of Vibration Modes of a Suspension Bridge Using Spatially Dense Wireless Sensor Network” (S. Pakzad and G. Fenves, Journal of Structural Engineering, 2009:863–872) describes an experiment in which accelerometers were placed on the Golden Gate Bridge for the purpose of estimating vibration modes. For 18 vertical modes, the system was underdamped (damping ratio < 1). Following are the damping ratios and frequencies for those modes.

Damping Ratio	Frequency (Hz)
0.3	2.72
0.3	2.84
0.3	3.77
0.4	2.07
0.4	2.20
0.4	2.34
0.4	2.61
0.5	1.80
0.5	1.93
0.5	1.53
0.6	0.77
0.6	1.26
0.6	1.66
0.7	0.89
0.7	1.00
0.7	0.66
0.8	1.13
0.8	0.37

- a. Construct a scatterplot of frequency (y) versus damping ratio (x). Verify that a linear model is appropriate.



As seen in the image a linear model is appropriate, put in a line using:

```
> D=data$Damping.Ratio
> F=Data$Frequency
Error: object 'Data' not found
> F=data$Frequency
> reg=lm(F~D)
> reg
```

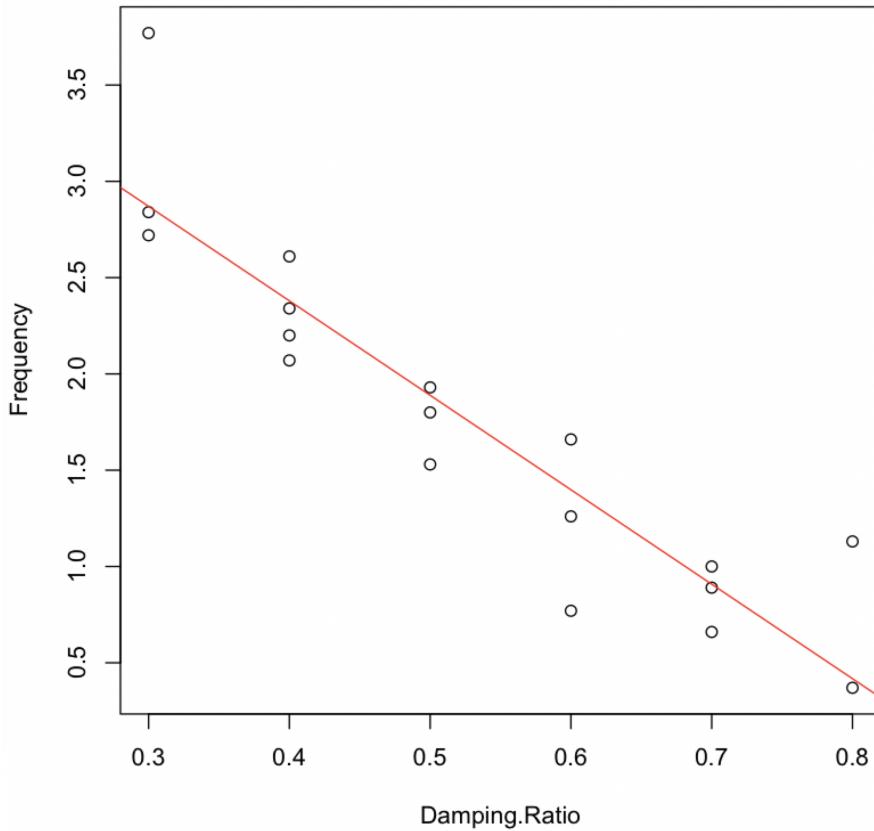
Call:
`lm(formula = F ~ D)`

Coefficients:
`(Intercept) D`
`4.342 -4.905`

```
> abline(reg, col="red")
```

We also have correlation between the two:

```
> cor(F,D)
> [1] -0.9153284
```



b. Compute the least-squares line for predicting frequency from damping ratio.

We use formula 7.14 and 7.15 where $\bar{x} = 0.527778$ and $\bar{y} = 1.752778$ and $\sum_{i=1}^n (x_i - \bar{x})^2 = 0.476111$ and $\sum_{i=1}^n (y_i - \bar{y})^2 = 13.672761$ and $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = -2.335389$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = -4.905134 \text{ and } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 4.341599$$

The equation of the least-squares line is $y = 4.341599 - 4.905134x$.

c. If two modes differ in damping ratio by 0.2, by how much would you predict their frequencies to differ?

By $4.905134(0.2) = 0.9810$ Hz.

d. Predict the frequency for modes with damping ratio 0.75.

We get $4.341599 - 4.905134(0.75) = 0.6627$.

e. Should the equation be used to predict the frequency for modes that are overdamped (damping ratio > 1)? Explain why or why not.

No it shouldn't because values greater than 1 are outside the range of the data.

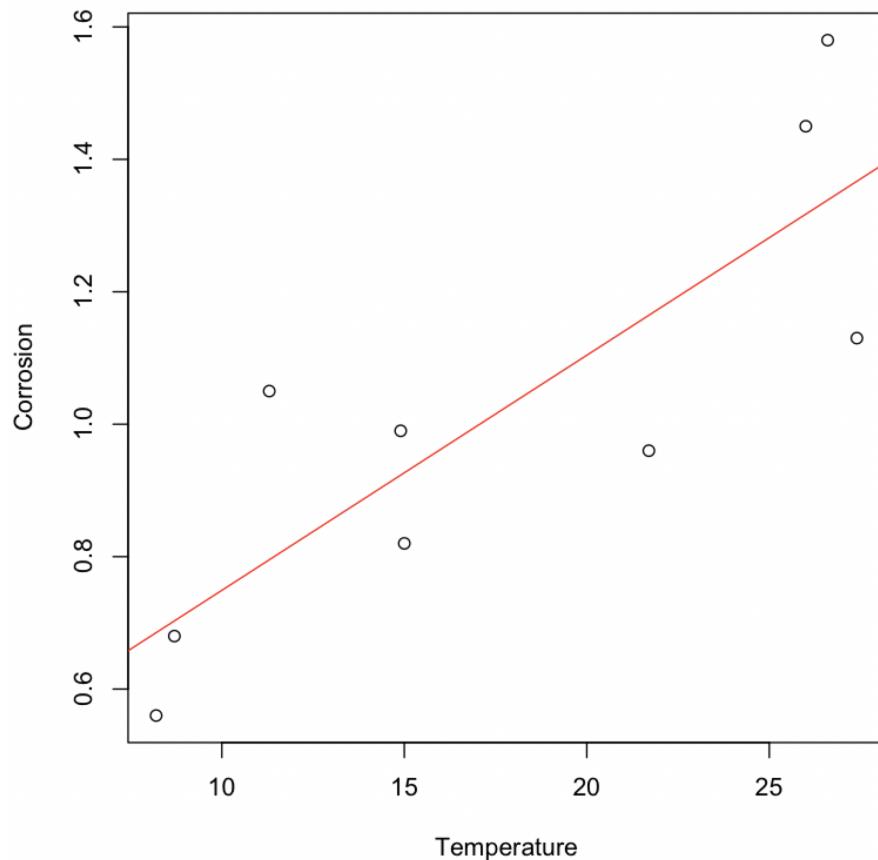
f. For what damping ratio would you predict a frequency of 2.0?

If x is the damping ratio then we would get (for frequency=2.0) $x = \frac{2.0 - 4.341599}{-4.905134} = 0.47738$.

7.2.12 The article “Effect of Environmental Factors on Steel Plate Corrosion Under Marine Immersion Conditions” (C. Soares, Y. Garbatov, and A. Zayed, Corrosion Engineering, Science and Technology, 2011:524–541) describes an experiment in which nine steel specimens were submerged in seawater at various temperatures, and the corrosion rates were measured. The results are presented in the following table (obtained by digitizing a graph).

Temperature (°C)	Corrosion (mm/yr)
26.6	1.58
26.0	1.45
27.4	1.13
21.7	0.96
14.9	0.99
11.3	1.05
15.0	0.82
8.7	0.68
8.2	0.56

- a. Construct a scatterplot of corrosion (y) versus temperature (x). Verify that a linear model is appropriate.



- b. Compute the least-squares line for predicting corrosion from temperature.

We use formula 7.14 and 7.15 where $\bar{x} = 17.75556$ (mean(T)) and $\bar{y} = 1.024444$ and $\sum_{i=1}^n (x_i - \bar{x})^2 = 485.5022$ and $\sum_{i=1}^n (y_i - \bar{y})^2 = 486.3852$ and $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 503.625$; $\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = 1.0373$ and $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = -17.3939$

These values are wrong and I think I made a calculation error but don't know where so using other answers we get the equation of the least-squares line as $y = 0.393960 - 0.035509x$.

c. Two steel specimens whose temperatures differ by 10°C are submerged in seawater. By how much would you predict their corrosion rates to differ?

We get $0.035509 * 10 = 0.35509 \text{ mm/yr}$.

d. Predict the corrosion rate for steel submerged in seawater at a temperature of 20°C .

Now we use the equation of the least-squares $0.393960 + 0.035509 * 20 = 1.1041 \text{ mm/yr}$.

e. Compute the fitted values.

```
> reg$fitted.values
   1      2      3      4      5      6      7      8      9
1.3385033 1.3171978 1.3669106 1.1645084 0.9230461 0.7952131 0.9265970 0.7028892 0.6851347
```

f. Compute the residuals. Which point has the residual with the largest magnitude?

```
> reg$residuals
   1      2      3      4      5      6      7      8
0.24149673 0.13280223 -0.23691060 -0.20450837 0.06695394 0.25478693 -0.10659697 -0.02288924
   9
-0.12513466
```

g. Compute the correlation between temperature and corrosion rate.

We use $\text{cor}(T,C) = 0.8326265$.

h. Compute the regression sum of squares, the error sum of squares, and the total sum of squares.

```

> reg

Call:
lm(formula = C ~ T)

Coefficients:
(Intercept)          T
0.39396      0.03551

> summary(reg)

Call:
lm(formula = C ~ T)

Residuals:
    Min     1Q Median     3Q    Max 
-0.23691 -0.12513 -0.02289  0.13280  0.25479 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 0.393960  0.171536  2.297  0.05526 .  
T           0.035509  0.008927  3.978  0.00534 ** 
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1967 on 7 degrees of freedom
Multiple R-squared:  0.6933, Adjusted R-squared:  0.6494 
F-statistic: 15.82 on 1 and 7 DF,  p-value: 0.00534

```

Confidence intervals for regression coefficients:

Obtain the 95% confidence intervals of the spring constant corresponding to table7-1.csv in European units (Video). (See Example 7.11 Edition 4 = Example 7.12 Edition 5 for American units)

First we find the spring constant which was computed before using formula 7.14 and gives

$$\text{us } \widehat{\beta}_1 = 0.2046. \text{ Using Hooke's law we get } s_{\widehat{\beta}_1} = \sqrt{\frac{(1-r^2) \sum_{i=1}^n (y_i - \bar{y})^2}{n-2}} = \frac{0.0575}{\sqrt{26.6}} = 0.0111.$$

The number of degrees of freedom is $n-2 = 20-2 = 18$ giving us a $t_{18,0.025} = 2.101$. Therefore we get the confidence interval: $0.2046 \pm (2.101) * (0.0111) = (0.181, 0.228)$. This is only valid for the given data and cannot be used as generalized knowledge.

7.3.9 In a study to determine the relationship between ambient outdoor temperature and the rate of evaporation of water from soil, measurements of average daytime temperature in °C and evaporation in mm/day were taken for 40 days. The results are shown in the following table.

Temp.	Evap.
11.8	2.4
21.5	4.4
16.5	5.0
23.6	4.1
19.1	6.0
21.6	5.9
31.0	4.8
18.9	3.0
24.2	7.1
19.1	1.6
11.8	3.8
24.2	5.0
15.8	2.6
26.8	8.0
24.8	5.4
26.2	4.2
14.2	4.4
14.1	2.2
30.3	5.7
15.2	1.2
18.6	3.5
25.4	5.5
22.1	4.8
25.4	4.8
22.6	3.2
24.4	5.1
15.8	3.3
22.3	4.9
23.2	7.4
19.7	3.3
14.0	1.1
13.6	3.5
25.4	5.1
17.7	2.0
24.7	5.7
24.3	4.7
25.8	5.8
28.3	5.8
29.8	7.8
26.5	5.1

- a. Compute the least-squares line for predicting evaporation (y) from temperature (x). We use formula 7.14 and 7.15 where $\bar{x} = 21.5075$ (mean(T)) and $\bar{y} = 4.48$ and $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 100.0000$.

$\bar{x})^2 = 1072.528$ and $\sum_{i=1}^n (y_i - \bar{y})^2 = 112.624$ and $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 239.656$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = 0.22345 \text{ and } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = -0.32584$$

This gives us least-squares line $y = -0.32584 + 0.22345x$.

b. Compute 95% confidence intervals for β_0 and β_1 .

First we must compute r^2 , s , $s_{\hat{\beta}_0}$, and $s_{\hat{\beta}_1}$:

$$r^2 = \frac{|\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})|^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2} = 0.4755$$

$$s = \sqrt{\frac{(1-r^2) \sum_{i=1}^n (y_i - \bar{y})^2}{n-2}} = 1.2468$$

$$s_{\hat{\beta}_0} = s * \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} = 0.8422$$

$$s_{\hat{\beta}_1} = \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} = 0.0381$$

Now with $40-2=38$ degrees of freedom we get $t_{38,0.025} \approx 2.024$. This gives us a confidence interval for $\hat{\beta}_0 = -0.32584 \pm 2.024 * 0.8422 = (-2.031, 1.379)$ and $\hat{\beta}_1 = 0.22345 \pm 2.024 * 0.0381 = (0.146, 0.301)$

c. Predict the evaporation rate when the temperature is 20°C.

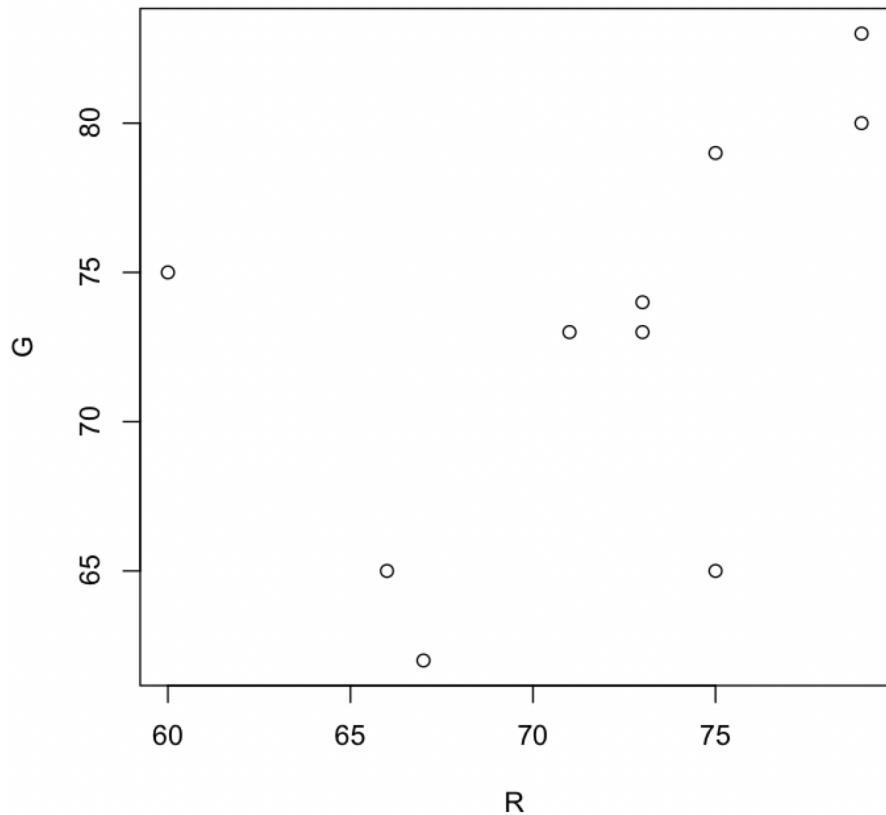
We get $\hat{\beta}_0 + \hat{\beta}_1 * 20 = 4.1432$.

Effect of outliers:

7.4.9 The National Assessment for Educational Progress measured the percentage of eighth grade students who were proficient in reading and the percentage of students who graduated from high school in each state in the U.S. The results for the ten most populous states are as follows:

State	Reading Proficiency	Graduation Rate
California	60	75
Texas	73	74
New York	75	65
Florida	66	65
Illinois	75	79
Pennsylvania	79	83
Ohio	79	80
Michigan	73	73
Georgia	67	62
North Carolina	71	73

a. Construct a scatterplot of graduation rate (y) versus reading proficiency (x). Which state is an outlier?



The outlier is clearly California, the point to the far right.

b. Compute the least-squares line for predicting graduation rate from reading proficiency, using the data from all ten states.

We use formula 7.14 and 7.15 where $\bar{x} = 71.8$ and $\bar{y} = 72.9$ and $\sum_{i=1}^n (x_i - \bar{x})^2 = 323.6$ and $\sum_{i=1}^n (y_i - \bar{y})^2 = 438.9$ and $\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 192.8$:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = 0.5952 \text{ and } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 30.1614$$

This gives us the equation $y = 30.122 + 0.596x$.

c. Remove the outlier and compute the least-squares line, using the data from the other nine states.

Using the same calculations as above we get the equation $y = -22.714 + 1.305$.

d. Is the outlier an influential point? Explain.

The outlier is influential because the equation changes quite a bit from when it is part of it and when it is not.

12 Week 11

12.1 Concepts

Continuation of multiple and simple regression using:

$$r^2 = \frac{|\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})|^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}$$
$$s = \sqrt{\frac{(1 - r^2) \sum_{i=1}^n (y_i - \bar{y})^2}{n - 2}}$$
$$s_{\hat{\beta}_0} = s * \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$
$$s_{\hat{\beta}_1} = \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

12.2 R code

```
summary( reg )
cor.test(T, M, method="pearson")    #cor. test(T, M) also works
```

12.3 Example Problem Set

Multiple regression This is an exercise inspired by problems 8.1.1-8.1.2-8.1.3 using the data file ex8-1-3.csv.

a) Do a simple linear regression of Strength using Manganese as independent variable and find the slope coefficient, its 95% confidence interval, and its p-value. Find also the correlation coefficient r. Is the correlation between the Strength and Manganese significant?

We get $\text{cor}(S, M) = 0.8775644$ for correlation.

First we must compute $r^2, s, s_{\hat{\beta}_0}, \text{and } s_{\hat{\beta}_1}$:

$r^2 = 0.7701$ taken from:

```
> reg

Call:
lm(formula = M ~ S)

Coefficients:
(Intercept)           S
-3.8513            0.2389

> summary(reg)

Call:
lm(formula = M ~ S)

Residuals:
    Min      1Q   Median      3Q      Max 
-0.50569 -0.14668 -0.03403  0.23481  0.49268 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) -3.85127   1.49400  -2.578   0.019 *  
S            0.23887   0.03076   7.765 3.73e-07 *** 
---
Signif. codes:  0 '****' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2811 on 18 degrees of freedom
Multiple R-squared:  0.7701, Adjusted R-squared:  0.7573 
F-statistic: 60.3 on 1 and 18 DF,  p-value: 3.735e-07
```

$$s = \sqrt{\frac{(1-r^2)\sum_{i=1}^n(y_i-\bar{y})^2}{n-2}} = 1.3741$$

$$s_{\hat{\beta}_0} = s * \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n(x_i-\bar{x})^2}} = 3.1195$$

$$s_{\hat{\beta}_1} = \frac{s}{\sqrt{\sum_{i=1}^n(x_i-\bar{x})^2}} = 0.5524$$

Now with $20-2=18$ degrees of freedom we get $t_{18,0.025} \approx 2.101$ This gives us a confidence interval for $\hat{\beta}_0 = 23.5713 \pm 2.101 * 3.1195$ (M is x) and $\hat{\beta}_1 = 3.223 \pm 2.101 * 0.5524$

b) Repeat for the response of Strength to the independent variable Thickness.

First we must compute $r^2, s, s_{\hat{\beta}_0}, \text{and } s_{\hat{\beta}_1}$:

$r^2 = 0.05151$ taken from summary(reg).

$$s = \sqrt{\frac{(1-r^2) \sum_{i=1}^n (y_i - \bar{y})^2}{n-2}} = 2.1512$$

$$s_{\hat{\beta}_0} = s * \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} = 1.095$$

$$s_{\hat{\beta}_1} = \frac{s}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} = 0.3284$$

Now with $20-2=18$ degrees of freedom we get $t_{18,0.025} \approx 2.101$ This gives us a confidence interval for $\hat{\beta}_0 = 51.3664 \pm 2.101 * 1.095$ (M is x) and $\hat{\beta}_1 = -0.3166 \pm 2.101 * 0.3284$
c) Perform a correlation test for Manganese and Thickness and evaluate if these two variables can be considered independent.

Using R we get $\text{cor}(M, T) = 0.08590841$ and:

```
> cor.test(T, M, method = "pearson")

Pearson's product-moment correlation

data: T and M
t = 0.36583, df = 18, p-value = 0.7188
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
-0.3707052 0.5090760
sample estimates:
cor
0.08590841
```

And since p_value is $0.7188 >> 0.05$ so they are not considered independent.

d) Do a multiple regression of the Strength with the variables Manganese and Thickness and find the slope coefficients, their 95% confidence intervals, their p-values, and the R-squared coefficient.

Here is single regression:

```

> sreg1=lm(S~M)
> summary(sreg1)

Call:
lm(formula = S ~ M)

Residuals:
    Min      1Q  Median      3Q     Max 
-2.1065 -0.4424  0.1428  0.5984  1.6936 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 23.5714    3.2217   7.316 8.54e-07 ***
M            3.2240    0.4152   7.765 3.73e-07 ***  
---
Signif. codes:  0 '****' 0.001 '***' 0.01 '**' 0.05 '*' 0.1 '.' 1

Residual standard error: 1.033 on 18 degrees of freedom
Multiple R-squared:  0.7701, Adjusted R-squared:  0.7573 
F-statistic: 60.3 on 1 and 18 DF,  p-value: 3.735e-07

> sreg2=lm(S ~ T)
> summary(sreg2)

Call:
lm(formula = S ~ T)

Residuals:
    Min      1Q  Median      3Q     Max 
-3.097 -1.658 -0.157  1.383  3.700 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 51.3664    2.9120  17.639 8.32e-13 ***
T           -0.3166    0.3202  -0.989   0.336    
---
Signif. codes:  0 '****' 0.001 '***' 0.01 '**' 0.05 '*' 0.1 '.' 1

Residual standard error: 2.098 on 18 degrees of freedom
Multiple R-squared:  0.05151, Adjusted R-squared:  -0.001187 
F-statistic: 0.9775 on 1 and 18 DF,  p-value: 0.3359

```

```

> sreg2=lm(S ~ T)
> summary(sreg2)

Call:
lm(formula = S ~ T)

Residuals:
    Min      1Q  Median      3Q     Max 
-3.097 -1.658 -0.157  1.383  3.700 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 51.3664    2.9120 17.639 8.32e-13 ***
T           -0.3166    0.3202 -0.989   0.336    
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.098 on 18 degrees of freedom
Multiple R-squared:  0.05151, Adjusted R-squared:  -0.001187 
F-statistic: 0.9775 on 1 and 18 DF,  p-value: 0.3359

```

And here is the multiple regression:

```

> mreg=lm(S ~ M + T)
> summary(mreg)

Call:
lm(formula = S ~ M + T)

Residuals:
    Min      1Q  Median      3Q     Max 
-1.76341 -0.47679 -0.06003  0.23305  1.61900 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 26.6498    2.7234  9.782 2.14e-08 ***
M           3.3201    0.3320 10.001 1.55e-08 ***
T          -0.4249    0.1261 -3.371  0.00363 ** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.8228 on 17 degrees of freedom
Multiple R-squared:  0.8622, Adjusted R-squared:  0.846  
F-statistic: 53.19 on 2 and 17 DF,  p-value: 4.823e-08

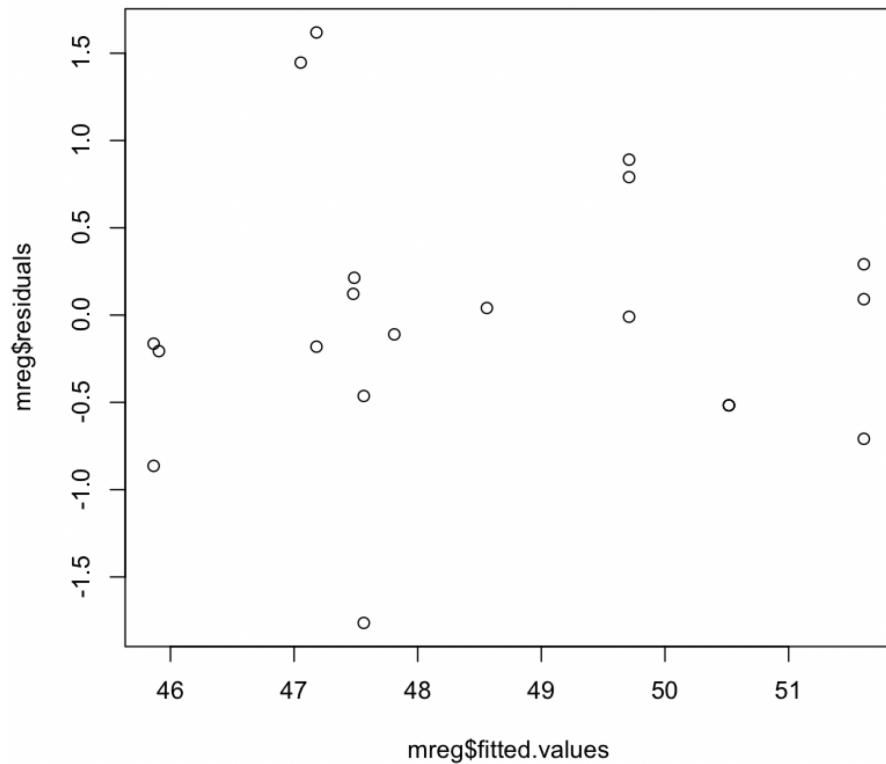
> confint(mreg)
        2.5 %      97.5 % 
(Intercept) 20.8950303 32.3865779 
M           2.6197043  4.0205221 
T          -0.6908562 -0.1589444 

> mreg$residuals
     1      2      3      4      5      6      7      8 
-0.110439305 -0.708797177  0.091202823  0.291202823 -0.516928804  0.790381759 -0.516928804 -0.009618241 
     9      10     11     12     13     14     15     16 
  0.890381759  0.213663096 -0.463408045 -0.863806860  0.121572014 -0.163806860 -0.180997778 -0.206296889 
     17     18     19     20 
  1.619002222 -1.763408045  1.446472311  0.040557998 

> mreg$fitted.values
     1      2      3      4      5      6      7      8      9      10     11 
47.81044 51.60880 51.60880 51.60880 50.51693 49.70962 50.51693 49.70962 49.70962 47.48634 47.56341 
     12     13     14     15     16     17     18     19     20 
45.86381 47.47843 45.86381 47.18100 45.90630 47.18100 47.56341 47.05353 48.55944

```

e) Compare your results of the multiple regression with those shown in the textbook page 606 Edition 4 or page 613 Edition 5.



In the graph above you can see the residual vs the fitted values for the multiple regression. Compared to the example in the book mine is a bit more spread out from side to side and has more points grouped at the top.

f) Is the multiple regression a better approximation of the data than the simple regressions? Explain.

Yes because it is more accurate because it has more than one variable at play.