

Final Exam

NOTE: The .rmd version of the file is available here: [\(link\)](#)

Instructions

Please prepare responses/solutions for the following questions. On the day of the exam, you will be given a new set of questions. You will use the solutions you've prepared for this exam during the exam.

During the exam, you will also be permitted to access the internet for publicly available content. You will not be allowed to communicate with anyone via the internet or any other means during the exam. This includes, but is not limited to:

- No messaging, emailing, or using social media to contact others.
- No posting questions or seeking answers on forums, chat rooms, chat bots (including large language models like ChatGPT), or any collaborative platforms.
- No sharing or discussing exam content with peers through any online or electronic medium.

You may **NOT** discuss any aspect of the exam or prep questions with anyone other than the instructor or TA. You may **NOT** share code or documents.

Submission instructions

1. Within your course repo, create a folder called `final-exam`
2. Within the folder, create the script file `final-exam.rmd` with your solutions. Create a rendered report in `.pdf` output.
3. Add, commit, and push to your repo on github.com.
4. If you received an email from gradescope, upload your PDF to gradescope.

Questions

All questions, including extra credit are 5 pts.

1. In the examples below, please explain/define what is meant by the word “probability”.

I think there is an 0.8 probability that the defendant committed the crime.

The probability of winning the powerball is less than 0.00001.

In these examples, probability refers to the expression of likelihood that an event will occur. In the first case, we see an expression of belief probability because a personal belief on whether or not the defendant is guilty is being expressed. In the second example, we see a frequency probability which is determined by calculating long run averages that an event occurs.

2. The following table is based on a study of delivery method and postnatal depression ([link](#)). (You do not need to read the publication.) In a cohort of mothers, researchers collected delivery mode and depression

scores (8 weeks postpartum). The data were collected in 1991 and 1992. While some planned vaginal deliveries did result in emergency caesarean section or assisted vaginal delivery, the variable of interest was the **planned** delivery mode.

Suppose that the cell probabilities were provided as a , b , c , and d as in the table below. Complete the rest of the table symbolically.

	Depression score < 13	Depression score ≥ 13	All
Planned vaginal delivery	a	b	$a+b$
row	$a/(a+b)$	$b/(a+b)$	
col	$a/(a+c)$	$b/(b+d)$	
Planned caesarean section delivery	c	d	$c+d$
row	$c/(c+d)$	$d/(c+d)$	
col	$c/(a+c)$	$d/(b+d)$	
All	$a+c$	$b+d$	$a+b+c+d = 1$

3. Now suppose that rather than cell probabilities, conditional probabilities were collected. Define postnatal depression as a depression score ≥ 13 , and let

$$e = P(\text{postnatal depression} | \text{Planned vaginal delivery})$$

$$f = P(\text{postnatal depression} | \text{Planned caesarean section delivery})$$

$$g = \text{incidence of planned caesarean section.}$$

Complete the table symbolically.

	Depression score < 13	Depression score ≥ 13	All
Planned vaginal delivery	$(1-g)-((1-g)*e)$	$(1-g)*e$	$1-g$
row	$(1-g)-((1-g)*e)/(1-g)$	e	
col	$(1-g)-((1-g)e)+g-(gf)/(1-g)-((1-g)*e)$	$((1-g)e)/(1-g)e+(g*f)$	
Planned caesarean section delivery	$g-(g*f)$	$g*f$	g
row	$(g-(g*f))/g$	f	
col	$(g-(g*f))/((1-g)-((1-g)e)+g-(g*f))$	$(gf)/((1-g)e)+(g*f)$	
All	$(1-g)-((1-g)e)+g-(gf)$	$((1-g)e)+(gf)$	

4. (Continuing from the previous problem.) If planned caesarean section is 30% of all deliveries, and the risk of postnatal depression is 0.1 in the planned vaginal delivery group and 0.15 in planned caesarean section delivery groups, what is

$$P(\text{Planned caesarean section delivery} | \text{Depression score} < 13)?$$

We are looking for $(g-(gf))/((1-g)-((1-g)e)+g-(g*f))$, which is a conditional column probability off the table. We are given that $g = .3$, $e = .1$, $f = .15$.

```

g <- .3
e <- .1
f <- .15
prob <- (g-(g*f))/((1-g)-((1-g)*e)+g-(g*f))
prob

```

```
## [1] 0.2881356
```

5. Suppose observational data were collected in which depression rates matched the proportions in question 4. Would the data support the conclusion that caesarean section delivery leads to higher rates of depression? If not, why not? (Hint: Recall chapter 8 of the text “Understanding Uncertainty”.)

6. The Monte Hall problem is a classic game show. Contestants on the show were shown three doors. Behind one randomly selected door was a sportscar; behind the other doors were goats.

At the start of the game, contestants would select a door, say door A. Then, the host would open either door B or C to reveal a goat. At that point in the game, the host would ask the contestant if she would like to change her door selection. Once a contestant decided to stay or change, the host would open the chosen door to reveal the game prize, either a goat or a car.

In this problem, consider a **modified** version of the Monte Hall problem in which the number of doors is **variable**. Rather than 3 doors, consider a game with 4 or 5 or 50 doors. In the modified version of the game, a contestant would select an initial door, say door A. Then, the host would open **one** of the remaining doors to reveal a goat. At that point in the game, the host would ask the contestant if she would like to change her door selection. Once a contestant decided to stay or change, the host would open the chosen door to reveal the game prize, either a goat or a car.

Consider two strategies:

1. Always stay with the first door selected.
2. Always switch to the unopened door.

The function `game` below plays a single game of Monte Hall. The function returns a vector of length two, the first element is the prize under strategy 1 and the second element is the prize under strategy 2. The function has a single input parameter, `N`, which is the number of doors in the game.

Use the `game` function to estimate the probability that strategy 1 results in a goat and strategy 2 results in a car. Let `N=5`.

```

suppressPackageStartupMessages(require(magrittr))
suppressPackageStartupMessages(require(dplyr))

game <- function(N){
  if(N<3) stop("Must have at least 3 doors")
  prize <- sample(c(rep("goat",N-1),"car"), N)
  guess <- sample(1:N,1)
  game <- data.frame(door = 1:N, prize = prize, stringsAsFactors = FALSE) %>%
    mutate(first_guess = case_when(
      door == guess ~ 1
      , TRUE ~ 0
    )) %>%
    mutate(potential_reveal = case_when(
      first_guess == 1 ~ 0
      , prize == "car" ~ 0
      , TRUE ~ 1
    ))
}

```

```

)) %>%
mutate(reveal = 1*(rank(potential_reveal, ties.method = "random") == 3)) %>%
mutate(potential_switch = case_when(
  first_guess == 1 ~ 0
  , reveal == 1 ~ 0
  , TRUE ~ 1
)) %>%
mutate(switch = 1*(rank(potential_switch, ties.method = "random") == 3))
c(game$prize[game$first_guess == 1], game$prize[game$switch == 1])
}
# we want the proportion of games in which goat is the first word and also the proportions for which car
# create a table with proportions
R <- replicate(10000, game(5) |> paste(collapse = ""))
table(R) |> proportions()

```

```

## R
##  cargoat  goatcar goatgoat
##   0.1979   0.2660   0.5361

```

```

# Probability of strategy 1 resulting in a goat: goatcar + goatgoat = 0.2761 + 0.5288 = 0.8049
# Probability of strategy 2 resulting in a car:  goatcar = 0.2761

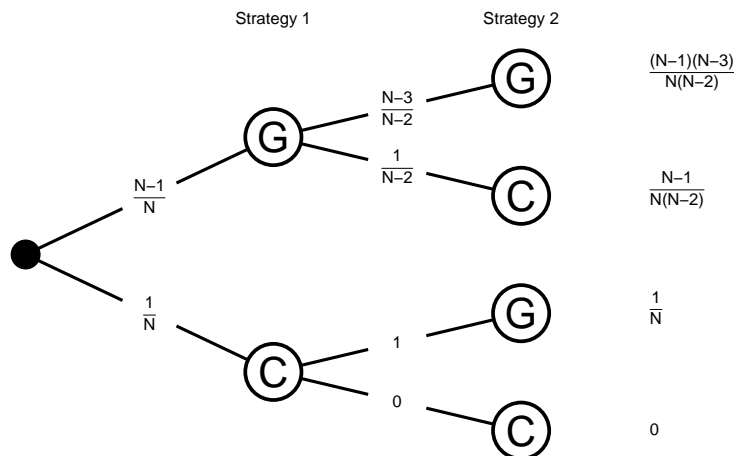
```

7. Consider the following tree, a possible analytic solution proposed by your classmate for the Monte Hall game with N doors. Your classmate argues that at the start of the game, there is only $\frac{1}{N}$ chance of getting the car in the initial guess. Consequently, there is a $\frac{N-1}{N}$ of selecting a goat in the initial guess. The initial guess is the outcome of strategy 1.

If strategy 1 results in a goat, then the outcome of strategy 2 is either a goat or car. As the host as revealed a door with a goat behind it, there are $N-2$ doors to choose from, 1 of which hides a car and $N-3$ of which hide goats. (Or so your classmate argues.)

Likewise, your classmate argues that if strategy 1 results in a car, then the outcome of strategy 2 must be a goat.

Multiplying the probabilities along the pathway, your classmate argues, generates the probability of the path itself.



The joint distribution of the outcomes of strategy 1 and strategy 2 can be represented, then, with the following contingency table.

		Statrategy 2	
Statrategy 1	Car	Car 0	Goat $\frac{1}{N}$
	Goat	$\frac{N-1}{N(N-2)}$	$\frac{(N-1)(N-3)}{N(N-2)}$

Using simulation, check the solution of your classmate for $N=5$. Show the contingency table which results from simulation next to the proposed analytic solution proposed by your classmate. How well does the simulation solution match the proposed solution?

it matches perfectly with the car/car estimation of 0, $1/5 = .2$ for car/goat which sort of matches .11, for goat car $4/5 \cdot 3 = 4/15 = .266$ which sort of matches .36 , and for goat goat = .53 which matches very well with the table, which also gave .53

we need to generate a bunch of replications of the game and create a contingency table based on of res
we need to store data from the replications in a data frame

```
library(gmodels)
a1 <- replicate(100, game(5))
gmodels::CrossTable(a1[1, ], a1[2, ])
```

```
##
##
##      Cell Contents
## |-----|
## |              N |
## | Chi-square contribution |
## |      N / Row Total |
## |      N / Col Total |
## |      N / Table Total |
## |-----|
##
```

```
##
## Total Observations in Table: 100
##
##
##      | a1[2, ]
##      |
##      | a1[1, ] | car | goat | Row Total |
## -----|-----|-----|-----|
##      | car | 0 | 18 | 18 |
##      |      | 6.120 | 3.153 |      |
##      |      | 0.000 | 1.000 | 0.180 |
##      |      | 0.000 | 0.273 |      |
##      |      | 0.000 | 0.180 |      |
## -----|-----|-----|-----|
##      | goat | 34 | 48 | 82 |
##      |      | 1.343 | 0.692 |      |
##      |      | 0.415 | 0.585 | 0.820 |
##      |      | 1.000 | 0.727 |      |
##      |      | 0.340 | 0.480 |      |
## -----|-----|-----|-----|
## Column Total | 34 | 66 | 100 |
##      |      | 0.340 | 0.660 |      |
## -----|-----|-----|-----|
##
##
```

8. Calculate the relative and absolute simulation error of your simulated probability in question 6, supposing that the your classmate's solution in question 7 is correct.

abosulte error = |simultaion - actual|

for car/car and goat/goat error is zero abs error car/goat = .11 - .2 = .09 abs error goat/car = .36 - .266 = 0.094

relative error = abosulte error/ actual rel error car/goat = .09/.2 = .45 rel error goat/car = 0.094/ .266 = 0.3533835

9. Consider a test for a rare genetic condition. Let T+ denote a test result that indicates the condition is present, while T- denotes absence. Let D+ and D- denote the true status of the disease.

Using the following information,

- $P(T+|D+) = .85$, and - sensitivty
- $P(T-|D-) = .95$, and - specifcity
- $P(D+) = 0.001$ - prevelence

calculate the **negative** predictive value of the test, $P(D-|T-)$.

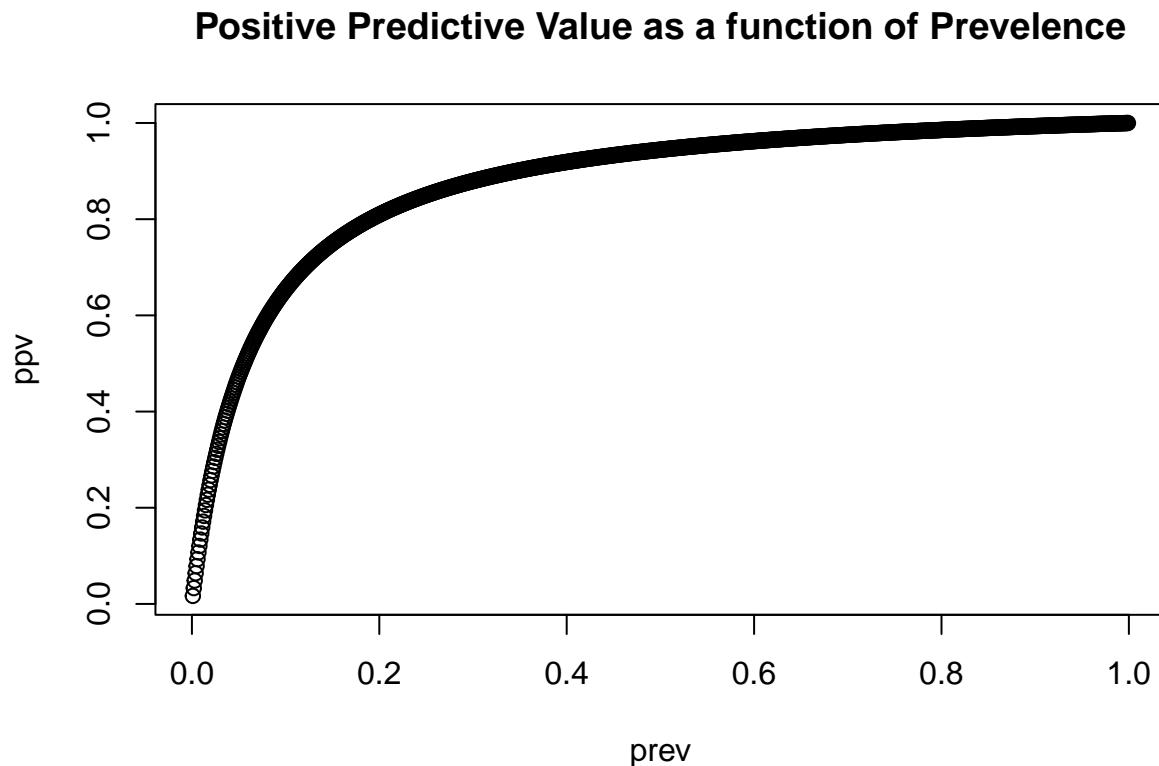
	Disease +		Disease -
Test +			
→ cell	0.00085	0.04995	0.0508
→ row			
→ col	.85	.05	
Test -			
→ cell	0.00015	0.94905	0.9492
→ row		0.999842	

	Disease +	Disease -	
→ col	.15	.95	
	.001	.999	1

the NPV is a conditional row probability where we are given a negative test and determining if that means we don't have the disease. We can use the rules of probability to fill in the table given conditional row probabilities and a marginal probability. The NPV is equal to 0.999842

10. Create a plot that shows how the **positive** predictive value is a function of the prevalence of disease, $P(D+)$. Keep the sensitivity and specificity the same as the previous question.

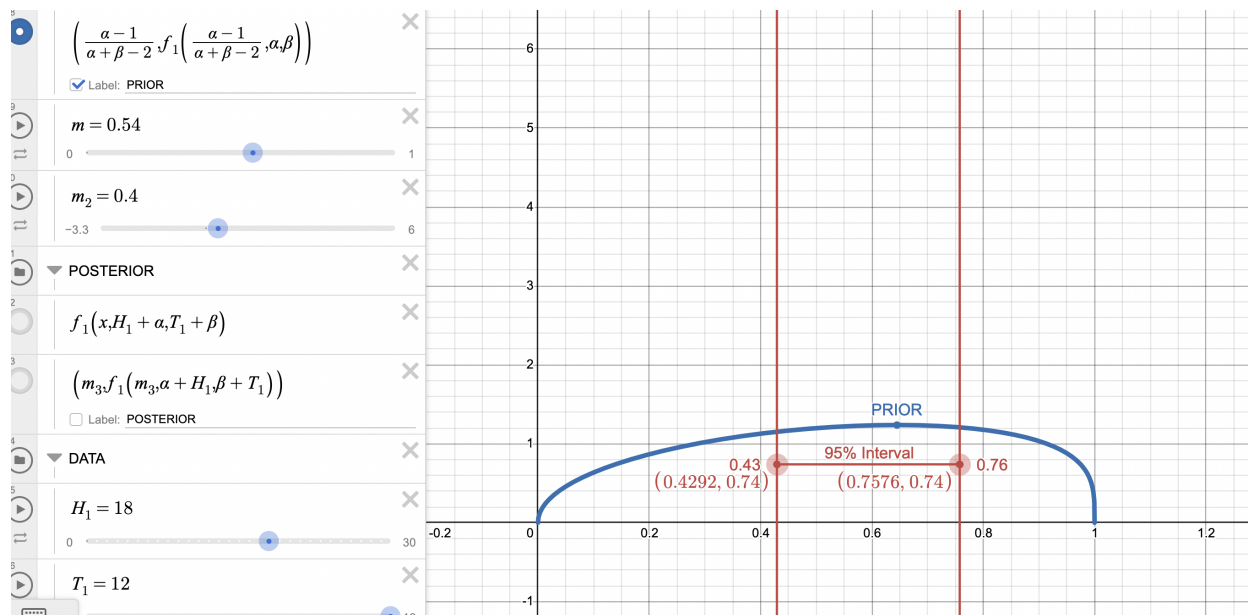
```
# ppv = (D+, T+)/ (T+) because it is a conditional row probability
# (D+, T+) = prev * sensitivity
# (T+) = ((1 - prev) * (1 - specificity)) + (D+, T+)
sen <- .85
spec <- .95
prev <- seq(.001, .999, .001)
cellprob <- prev * sen
tpos <- ((1 - prev) * (1 - spec)) + cellprob
ppv <- cellprob / tpos
plot(prev, ppv, main= "Positive Predictive Value as a function of Prevalence")
```



11. Suppose an upcoming election for UVA student body president is between two candidates. In a survey of 30 students, 18 voiced support for candidate A. Use the Desmos calculator ([link](#)) to fit a probability model

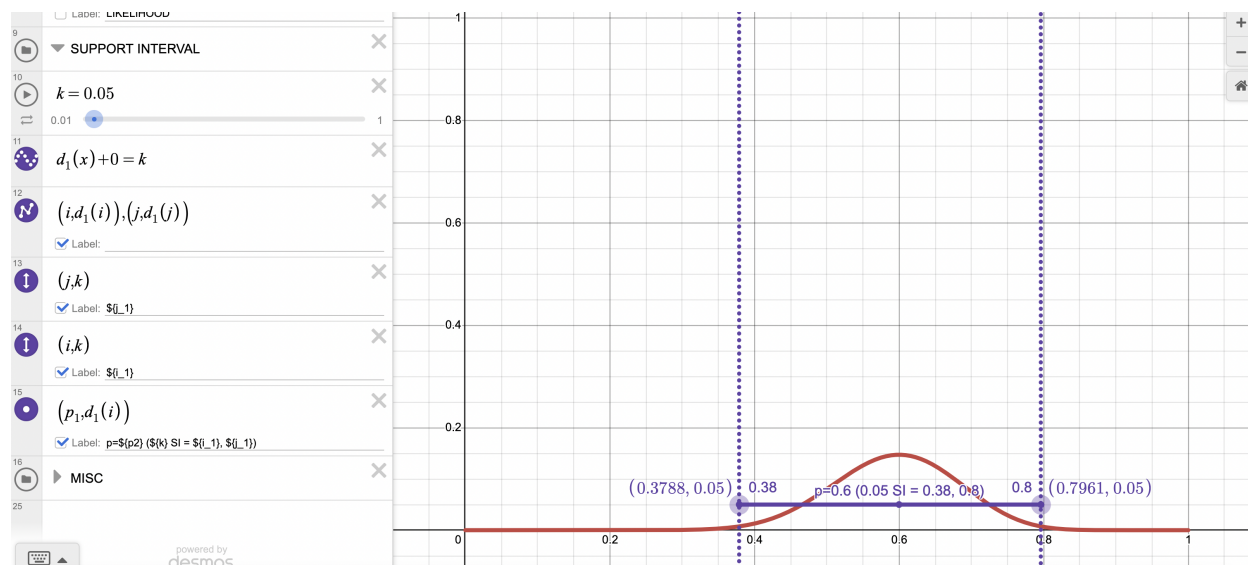
with Bayesian methods for the election, specifically the probability that candidate A is the preferred by the student body. Report the 95% credible interval. (Provide a screen shot of the calculator with your solution.)

First I set my prior knowledge, which I changed slightly from what was given. Then I set my H as 18 and my T as 12 to represent how many people voted for candidates A and B. Then I moved the credible interval bar until it equaled 95%



12. Suppose an upcoming election for UVA student body president is between two candidates. In a survey of 30 students, 18 voiced support for candidate A. Use the Desmos calculator (link) to fit a probability model with Maximum Likelihood for the election, specifically the probability that candidate A is the preferred by the student body. Report the 1/20 support interval. (Provide a screen shot of the calculator with your solution.)

Since Max Likelihood does not require a prior, the first step was to set $n = 30$ and $h = 18$. then I set the support interval to $k = 1/20 = .05$.



13. Suppose diastolic blood pressure (DBP) follows a normal distribution with mean 80 mmHg and SD 15 mmHg. What is the probability that a randomly sampled person's DBP exceeds 104 mmHg?


```
# we want the probability of greater than 104 - complement to that is less than or equal to 104  
# use 1 - pnorm to determine
```

```
1 - pnorm(104, 80, 15)
```

```
## [1] 0.05479929
```

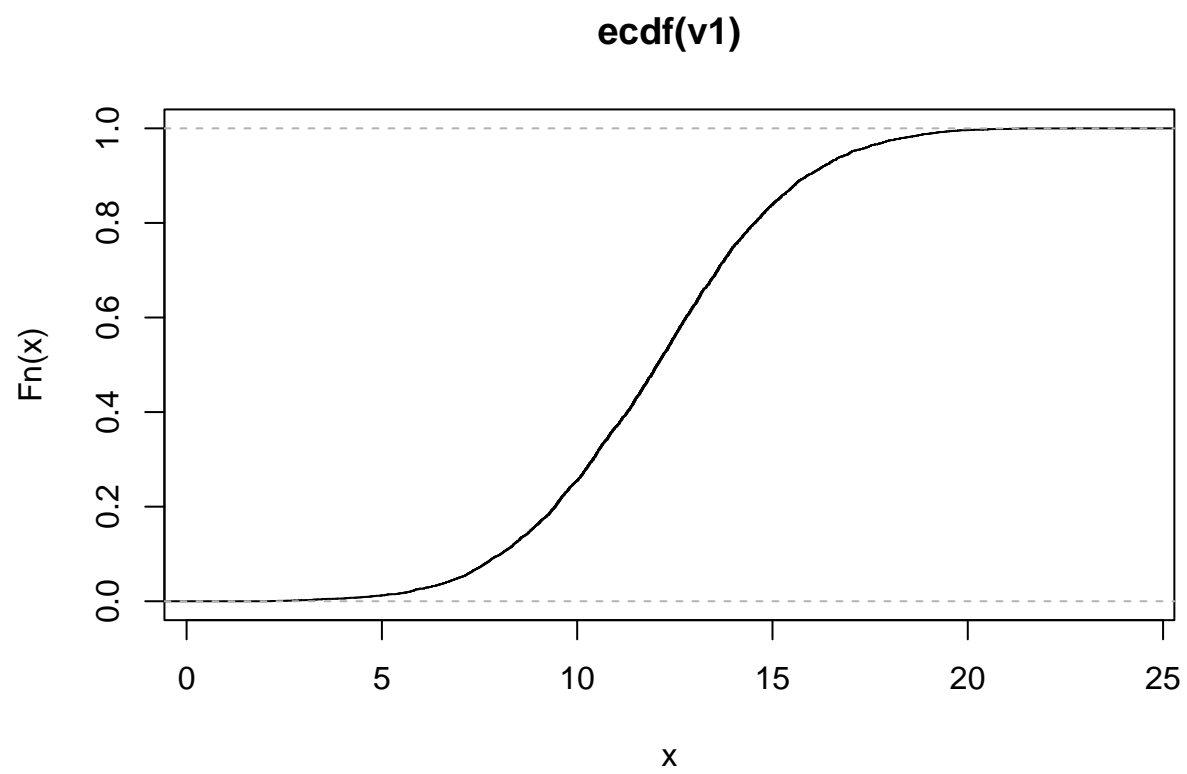
14. Suppose a laptop manufacturer sourced batteries from two different vendors. In testing the batteries, the manufacturer collected the following data on time to battery depletion.

```
d1 <- readRDS(url("https://tgstewart.cloud/battery-data.RDS"))  
head(d1)
```

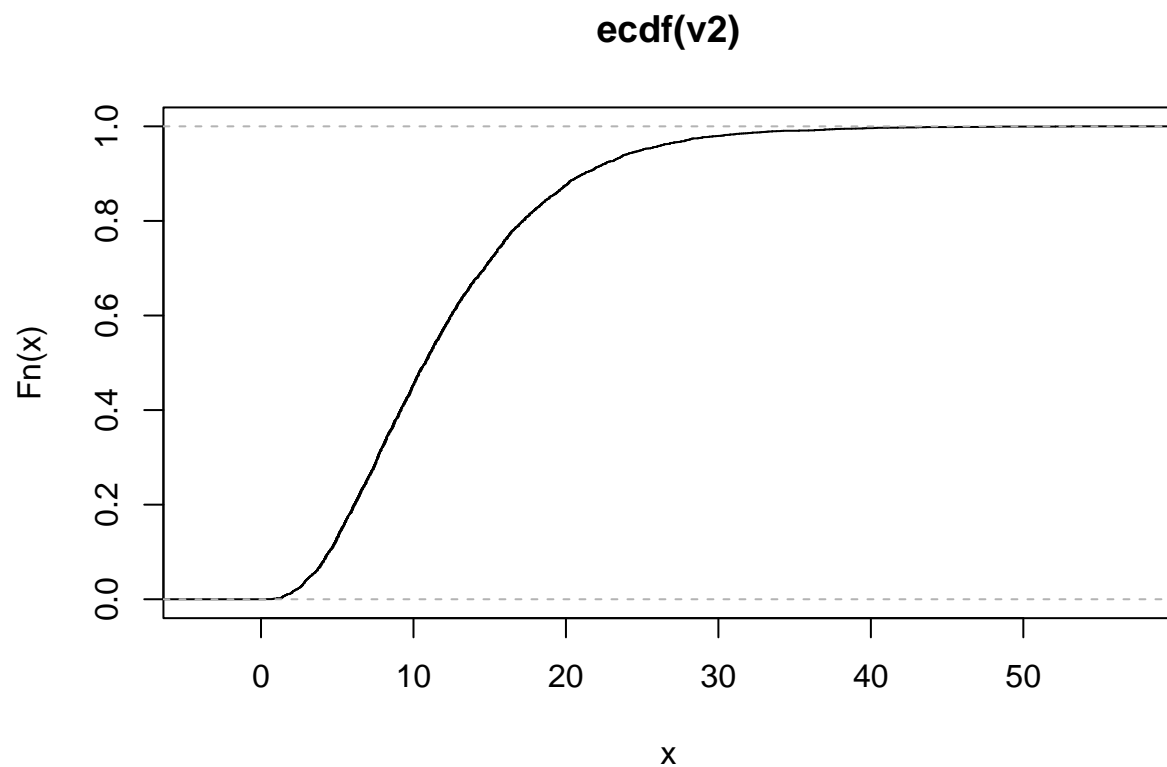
```
##   source    time  
## 1      0 12.957865  
## 2      0 12.972893  
## 3      0 12.025714  
## 4      1 18.317820  
## 5      0 14.463882  
## 6      1  8.110458
```

Using the data, generate a plot of the empirical CDF for time to battery depletion for each vendor. (Generate both eCDFs on the same plot, if possible.)

```
v1 <- d1 %>% filter(source == 0) |> pull(time)  
  
v2 <- d1 %>% filter(source == 1) |> pull(time)  
a2 <- ecdf(v2)  
a1 <- ecdf(v1)  
plot(a1)
```



```
plot(a2)
```

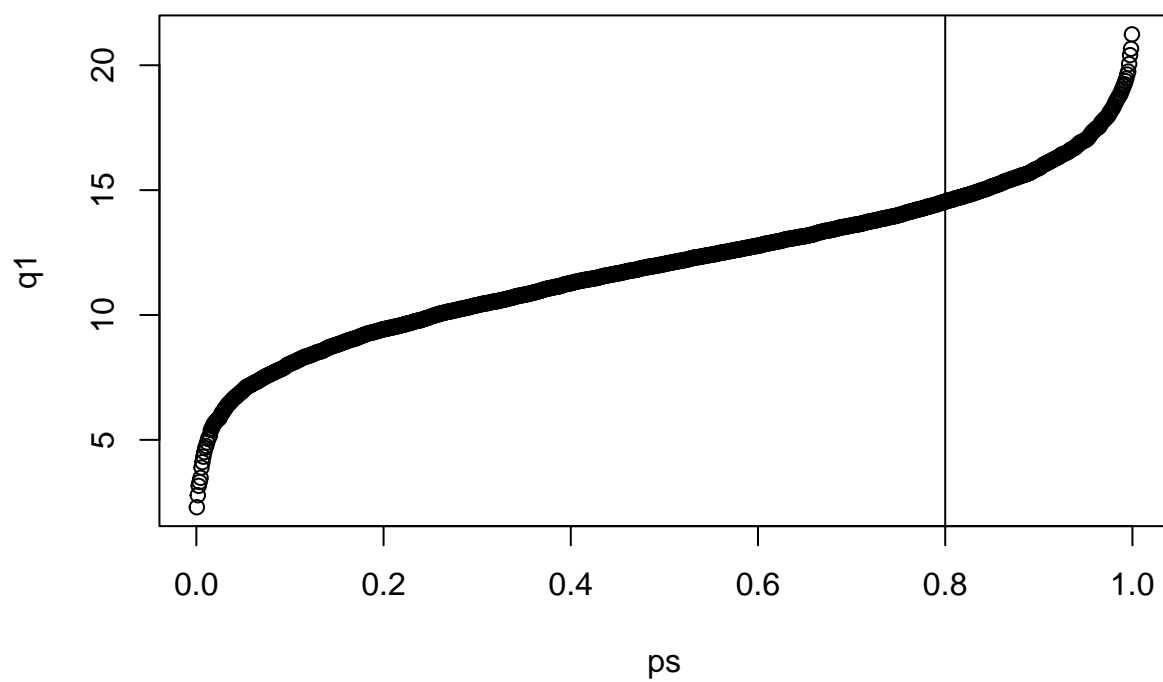


15. Based on the data, what is the 80th percentile for battery life (time to battery depletion) for each vendor?

use the quantile function which is the inverse of the eCDF to get the 80th percentile:

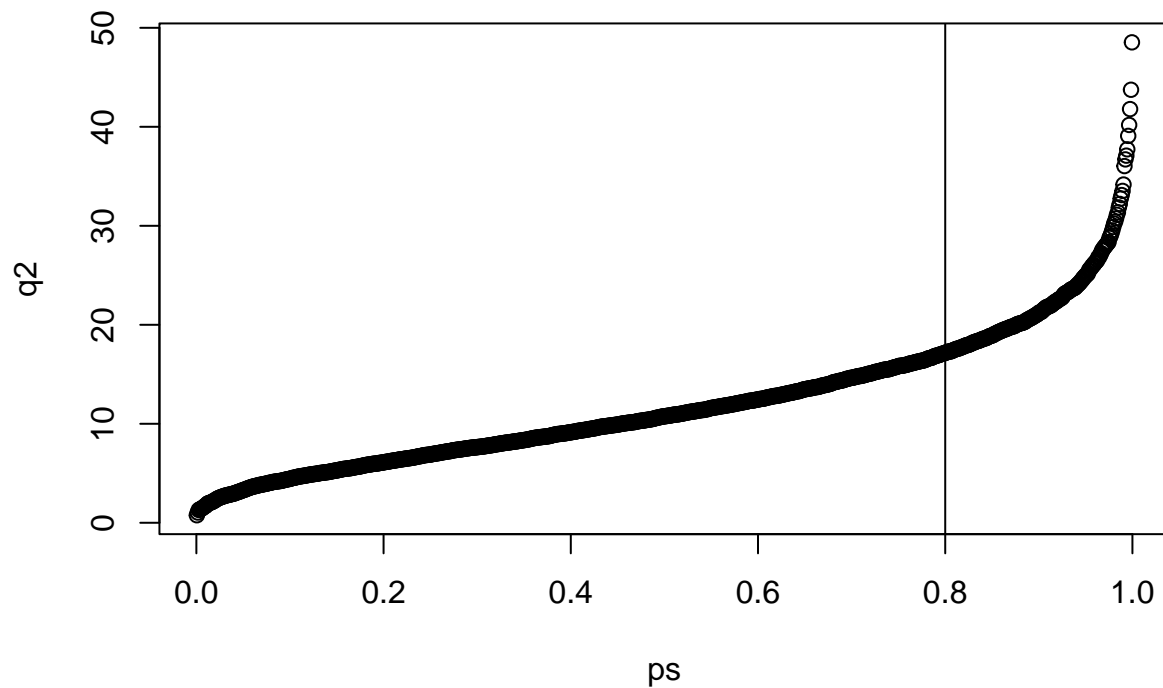
```
ps <- ppoints(1000)
q1 <- quantile(v1, ps)
plot(ps, q1, main = "estimate for 80th percentile vendor 1")
abline(v = .8)
```

estimate for 80th percentile vendor 1



```
q2 <- quantile(v2, ps)
plot(ps, q2, main = "estimate for 80th percentile vendor 2")
abline(v = .8)
```

estimate for 80th percentile vendor 2



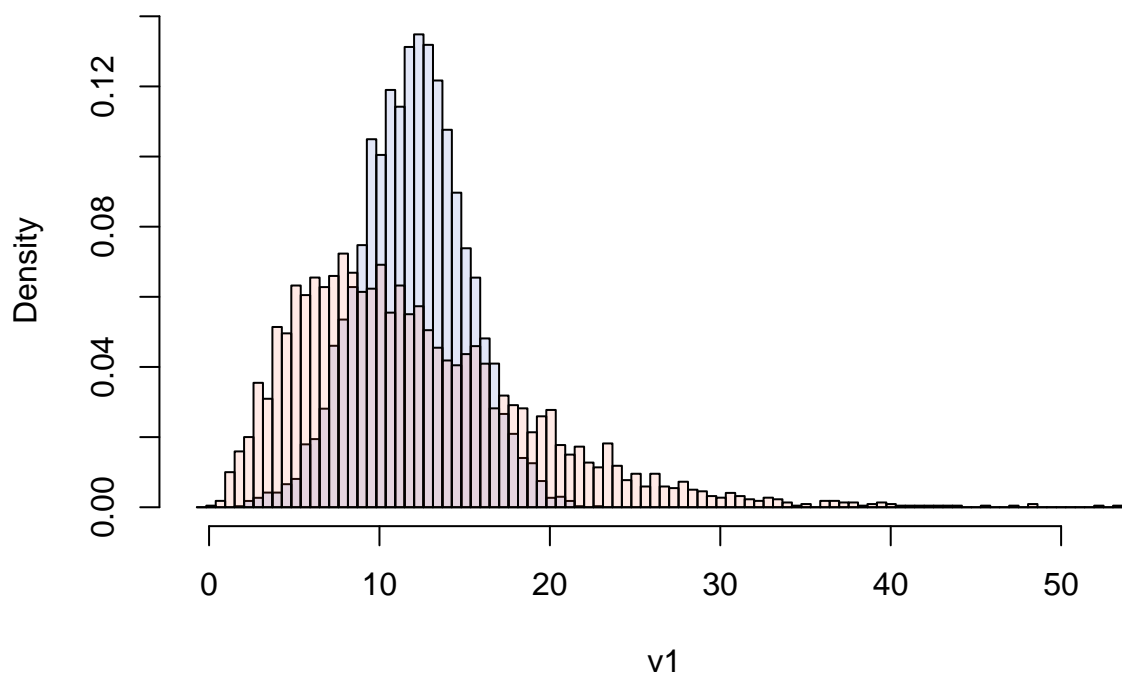
16. Using the data, generate a histogram for time to battery depletion for each vendor. (Generate both histograms on the same plot, if possible.)

```
# Hints for plot
b1 <- d1$time %>% range %>% `+` (c(-1,1))
b2 <- seq(b1[1], b1[2], length=100)

# Source 1
hist(v1, breaks = b2, freq = FALSE, col = "#1338BE20", xlim = b1)

# Source 2
hist(v2, breaks = b2, add=TRUE, col = "#FF573320", freq = FALSE)
```

Histogram of v1



17. The function `rbatlife` was created to mimic the distribution of battery life from the previous problem. It will generate `N` draws from the distribution. Using `rbatlife`, what is the mean battery life for each vendor?

```
rbatlife <- function(N){
  g <- rbinom(N,1,.4)
  o <- rgamma(N,3,scale=4)*g + rnorm(N,12,3)*(1-g)
  data.frame(source = g, time = o)
}
```

use this fucntion with v1 and v2 to along with the mean() fucntion to determine mean battery life for

18. The following code will load the first 500 rows of the NHANES data, a large national survey about nutrition.

```
#install.packages("Hmisc")
library(Hmisc)
```

```
##
## Attaching package: 'Hmisc'
```

```
## The following objects are masked from 'package:dplyr':
##
##   src, summarize
```

```
## The following objects are masked from 'package:base':
##
##     format.pval, units
```

```
suppressPackageStartupMessages(require(dplyr))
Hmisc::getHdata(nhgh)
d1 <- nhgh[1:500,]
head(d1)
```

```
##      seqn      sex      age      re      income tx dx      wt      ht
## 1  51624    male 34.16667 Non-Hispanic White [25000,35000) 0 0  87.4 164.7
## 3  51626    male 16.83333 Non-Hispanic Black [45000,55000) 0 0  72.3 181.3
## 5  51628 female 60.16667 Non-Hispanic Black [10000,15000) 1 1 116.8 166.0
## 6  51629    male 26.08333 Mexican American [25000,35000) 0 0  97.6 173.0
## 7  51630 female 49.66667 Non-Hispanic White [35000,45000) 0 0  86.7 168.4
## 10 51633    male 80.00000 Non-Hispanic White [15000,20000) 0 1  79.1 174.3
##      bmi leg arml armc waist tri sub gh albumin bun SCr
## 1  32.22 41.5 40.0 36.4 100.4 16.4 24.9 5.2      4.8  6 0.94
## 3  22.00 42.0 39.5 26.6  74.7 10.2 10.5 5.7      4.6  9 0.89
## 5  42.39 35.3 39.0 42.2 118.2 29.6 35.6 6.0      3.9 10 1.11
## 6  32.61 41.7 38.7 37.0 103.7 19.0 23.2 5.1      4.2  8 0.80
## 7  30.57 37.5 36.1 33.3 107.8 30.3 28.0 5.3      4.3 13 0.79
## 10 26.04 42.8 40.0 30.2  91.1  8.6 15.2 5.4      4.3 16 0.83
```

```
summary(d1)
```

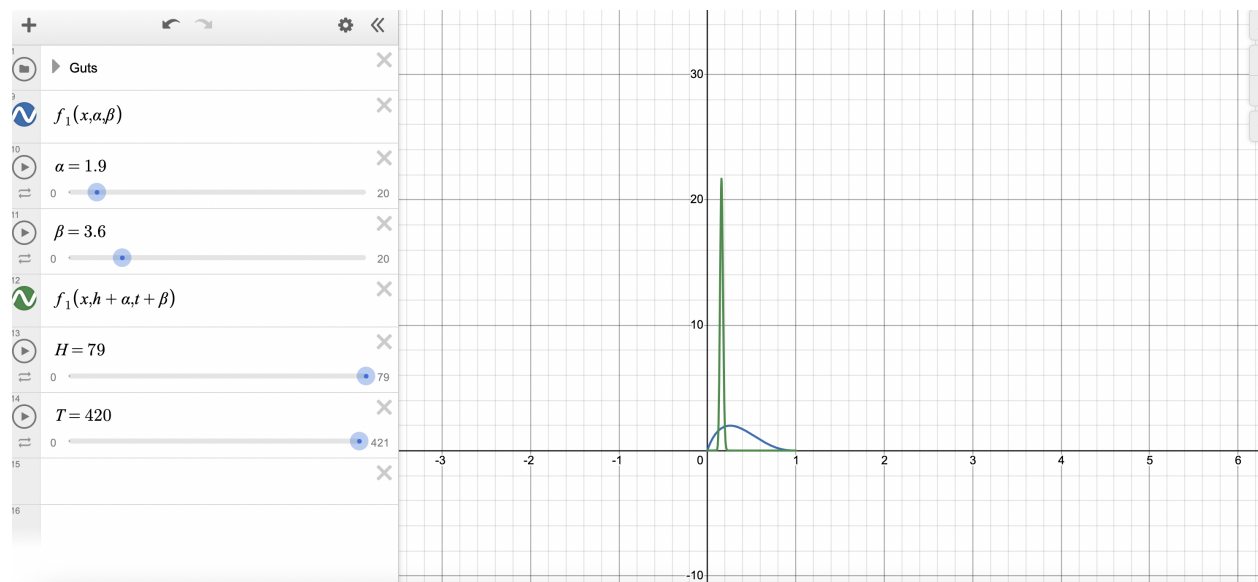
```
##      seqn      sex      age
## Min.   :51624    male :248   Min.   :12.08
## 1st Qu.:51805    female:252   1st Qu.:26.75
## Median :51996                                     Median :43.75
## Mean   :51987                                     Mean  :44.02
## 3rd Qu.:52167                                     3rd Qu.:60.27
## Max.   :52339                                     Max.   :80.00
##
##      re      income      tx
## Mexican American      : 93 [25000,35000) : 62   Min. :0.000
## Other Hispanic        : 50 [35000,45000) : 56   1st Qu.:0.000
## Non-Hispanic White    :235 >= 100000      : 54   Median :0.000
## Non-Hispanic Black    : 87 [75000,100000): 43   Mean   :0.102
## Other Race Including Multi-Racial: 35 [20000,25000) : 40   3rd Qu.:0.000
##      (Other)           :224   Max.    :1.000
##      NA's              : 21
##
##      dx      wt      ht      bmi
## Min.   :0.000   Min.   : 35.10   Min.   :140.3   Min.   :14.59
## 1st Qu.:0.000   1st Qu.: 64.90   1st Qu.:159.6   1st Qu.:23.56
## Median :0.000   Median : 76.55   Median :166.4   Median :27.48
## Mean   :0.158   Mean   : 80.16   Mean   :167.1   Mean   :28.58
## 3rd Qu.:0.000   3rd Qu.: 92.30   3rd Qu.:174.8   3rd Qu.:32.13
## Max.   :1.000   Max.   :230.70   Max.   :192.6   Max.   :81.25
##
##      leg      arml      armc      waist
## Min.   :25.00   Min.   :27.00   Min.   :20.50   Min.   : 60.20
## 1st Qu.:36.00   1st Qu.:34.80   1st Qu.:29.18   1st Qu.: 83.95
```

```
## Median :38.50 Median :36.90 Median :32.20 Median : 95.50
## Mean :38.44 Mean :36.96 Mean :32.58 Mean : 96.68
## 3rd Qu.:41.50 3rd Qu.:39.00 3rd Qu.:35.50 3rd Qu.:107.15
## Max. :50.30 Max. :44.00 Max. :48.60 Max. :147.70
## NA's :13 NA's :11 NA's :12 NA's :13
## tri sub gh albumin
## Min. : 3.10 Min. : 5.10 Min. : 4.000 Min. :2.500
## 1st Qu.:12.80 1st Qu.:14.00 1st Qu.: 5.200 1st Qu.:4.100
## Median :19.00 Median :19.80 Median : 5.500 Median :4.300
## Mean :19.14 Mean :20.24 Mean : 5.761 Mean :4.274
## 3rd Qu.:25.20 3rd Qu.:26.80 3rd Qu.: 5.800 3rd Qu.:4.500
## Max. :39.40 Max. :40.00 Max. :15.500 Max. :5.200
## NA's :31 NA's :67 NA's :5
## bun SCr
## Min. : 2.00 Min. :0.3900
## 1st Qu.: 9.00 1st Qu.:0.7100
## Median :12.00 Median :0.8300
## Mean :12.91 Mean :0.9034
## 3rd Qu.:16.00 3rd Qu.:0.9700
## Max. :56.00 Max. :9.1300
## NA's :5 NA's :5
```

```
table(d1$dx)
```

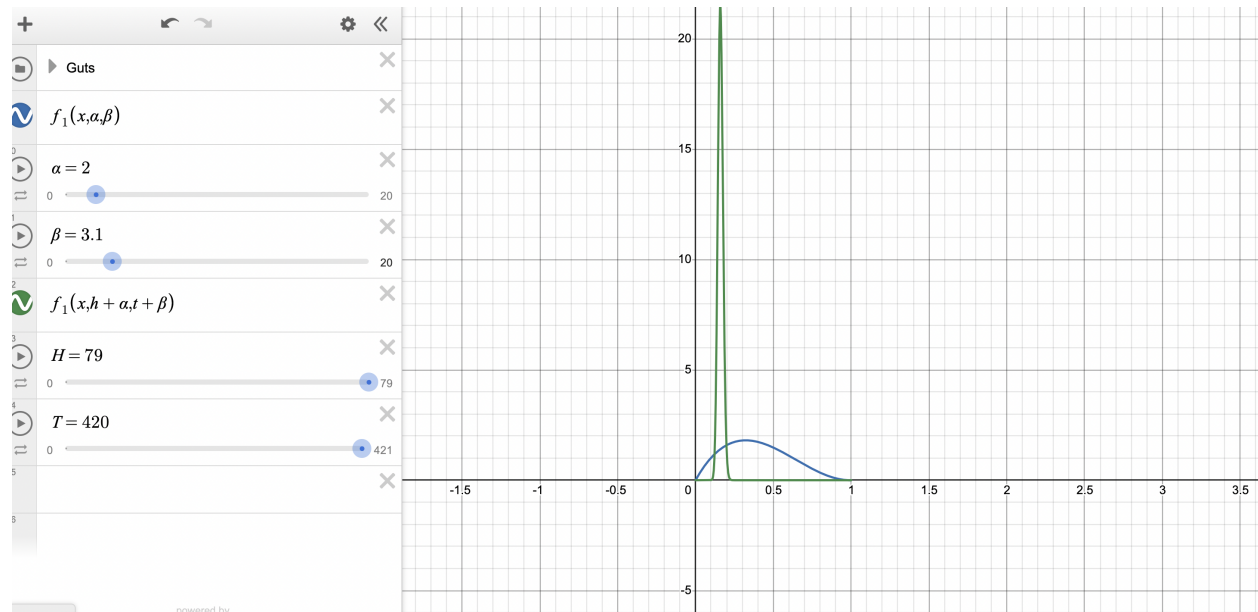
```
##
## 0 1
## 421 79
```

Estimate the prevalence of diabetes (dx) for all respondents using Bayesian updating with a binomial likelihood and beta prior. Use the following Desmos calculator ([link](#)). Change α and β to control the prior. Use H (heads) and T (tails) to plug in the data. Take a screenshot of the posterior distribution and the prior.



19. Reestimate the prevalence of diabetes (dx) with a more informative prior. Take a screenshot of the resulting posterior distribution with the new prior. Explain why the new prior is more informative.

I updated the prior based on what I would assume to be true in society about diabetes, and I would assume not many people have it but many would in the study but is why i increased α by a little bit. The new prior is more informative because it aligns more closely with my prior beliefs and more accurate prior beliefs and therefore some uncertainty is eliminated



20. Suppose the posterior distribution of the mean birthweight of infants whose mothers did not smoke was a normal distribution with mean = 3100 and standard deviation = $\sqrt{10}$. The symmetric density credible interval is calculated by identifying the 0.025 and 0.975 quantiles from the posterior. Calculate the interval.

```
pnorm(0.025,3100,sqrt(10))
```

```
## [1] 0
```

```
qnorm(0.975, 3100, sqrt(10)) - qnorm(.025, 3100, sqrt(10))
```

```
## [1] 12.3959
```

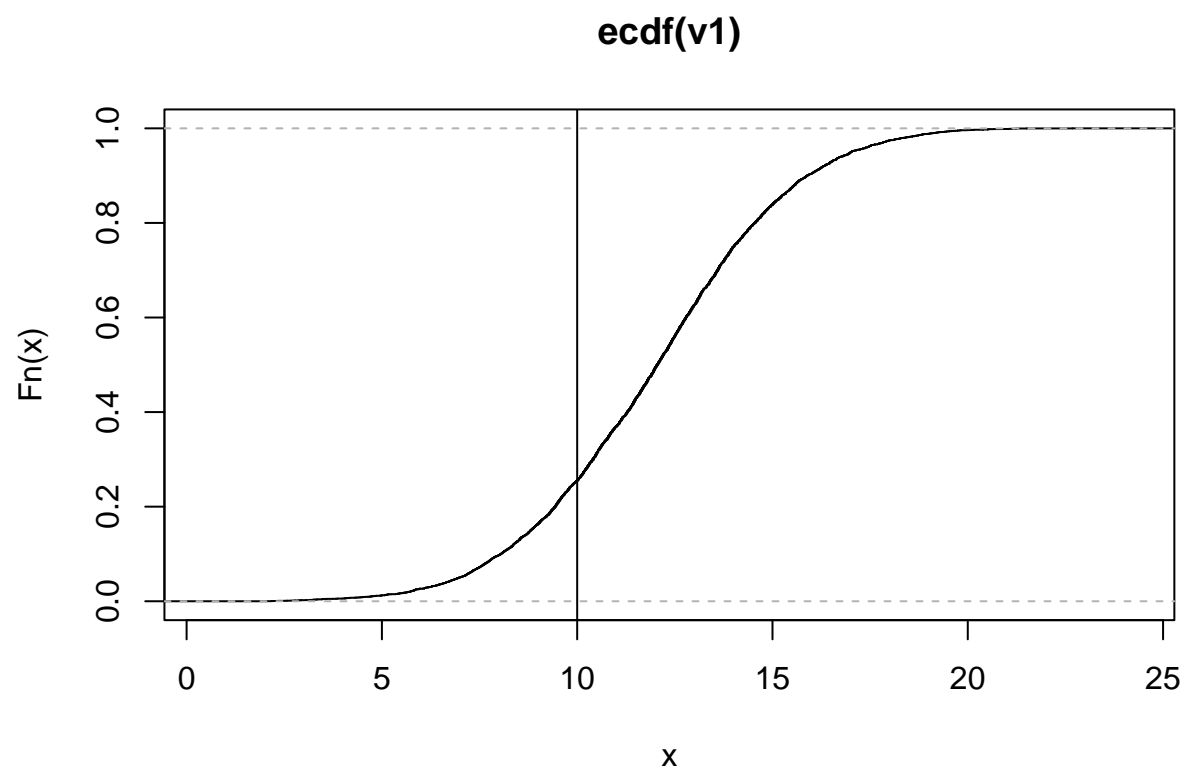
EC1. A creative writing essay was submitted without the author's name. After informing the class of the unnamed essay, two students claimed to be the author. Student A is known to use exclamation points in 10% of sentences. Student B is known to use exclamation points in 5% of sentences. A review of the unnamed essay revealed that 5 of 60 sentences used an exclamation point. With this information, calculate

$$P(\text{Student A authored the essay} \mid 5 \text{ of } 60 \text{ sentences used an exclamation point}).$$

EC2. Continuing the previous question, create a plot with number of exclamation points on the x-axis and the probability that student A authored the essay on the y-axis.

EC3. (Continuing problems 14 - 17) If a battery lasted for 10 or fewer hours, what is the probability it was from source 1?

```
# the probability within source 1 itself would be just over .2 and around .4 for source 2.
plot(a1)
abline(v=10)
```



```
plot(a2)  
abline(v=10)
```

