

## Homework 2

The data set `calif_penn_2011.csv` contains information about the housing stock of California and Pennsylvania, as of 2011. Information is aggregated into “Census tracts”, geographic regions of a few thousand people which are supposed to be fairly homogeneous economically and socially.

1. *Loading and cleaning*

- a. Load the data into a dataframe called `ca_pa`.

```
setwd("D:/github/Rlab/")
ca_pa<-read.csv("data/calif_penn_2011.csv",header=T)
ca_pa = as.data.frame(ca_pa)
```

- b. How many rows and columns does the dataframe have?

```
dim(ca_pa)#11275 rows and 34 cloumns
```

```
## [1] 11275    34
```

- c. Run this command, and explain, in words, what this does:

```
colSums(apply(ca_pa,c(1,2),is.na))
```

```
##              X              GEO.id2
##              0              0
##      STATEFP      COUNTYFP
##              0              0
##      TRACTCE      POPULATION
##              0              0
##      LATITUDE      LONGITUDE
##              0              0
##      GEO.display.label      Median_house_value
##              0              599
##      Total_units      Vacant_units
##              0              0
##      Median_rooms      Mean_household_size_owners
##              157              215
##      Mean_household_size_renters      Built_2005_or_later
##              152              98
##      Built_2000_to_2004      Built_1990s
##              98              98
##      Built_1980s      Built_1970s
##              98              98
##      Built_1960s      Built_1950s
##              98              98
##      Built_1940s      Built_1939_or_earlier
##              98              98
##      Bedrooms_0      Bedrooms_1
##              98              98
##      Bedrooms_2      Bedrooms_3
##              98              98
##      Bedrooms_4      Bedrooms_5_or_more
```

```
##           98           98
##           Owners       Renters
##           100          100
## Median_household_income Mean_household_income
##           115          126
```

find how many na in each column d. The function `na.omit()` takes a dataframe and returns a new dataframe, omitting any row containing an NA value. Use it to purge the data set of rows with incomplete data.

```
ca_pa = na.omit(ca_pa)
```

e. How many rows did this eliminate?

```
11275-dim(ca_pa)[1]
```

```
## [1] 670
```

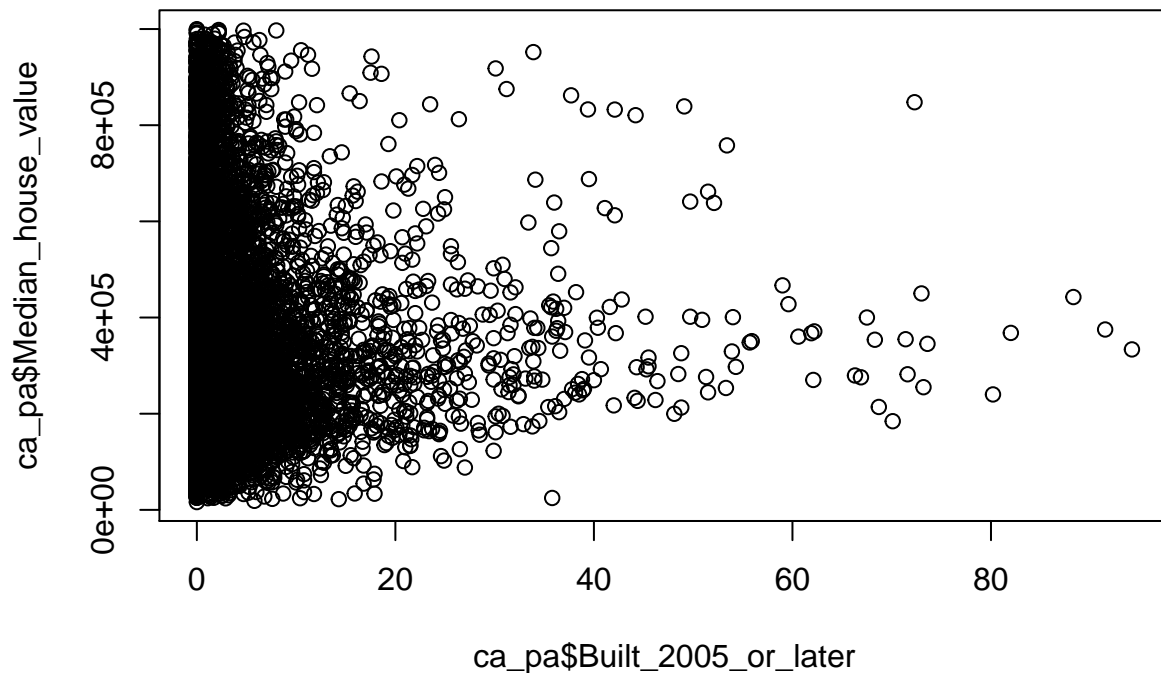
f. Are your answers in (c) and (e) compatible? Explain.

Yes, since some rows have more than 1 NA, thus the number of NAs could be bigger than the number of rows that have NAs.

## 2. *This Very New House*

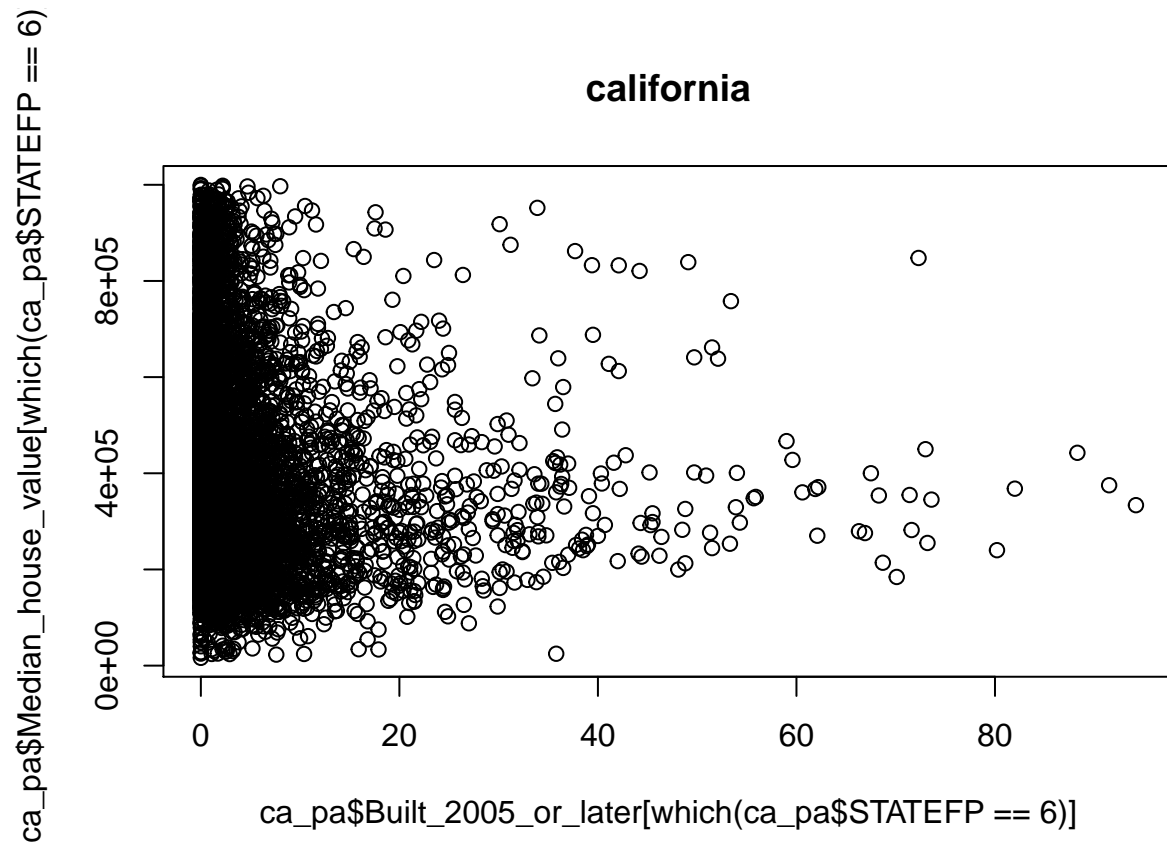
- The variable `Built_2005_or_later` indicates the percentage of houses in each Census tract built since 2005. Plot median house prices against this variable.

```
plot(x=ca_pa$Built_2005_or_later,y=ca_pa$Median_house_value)
```

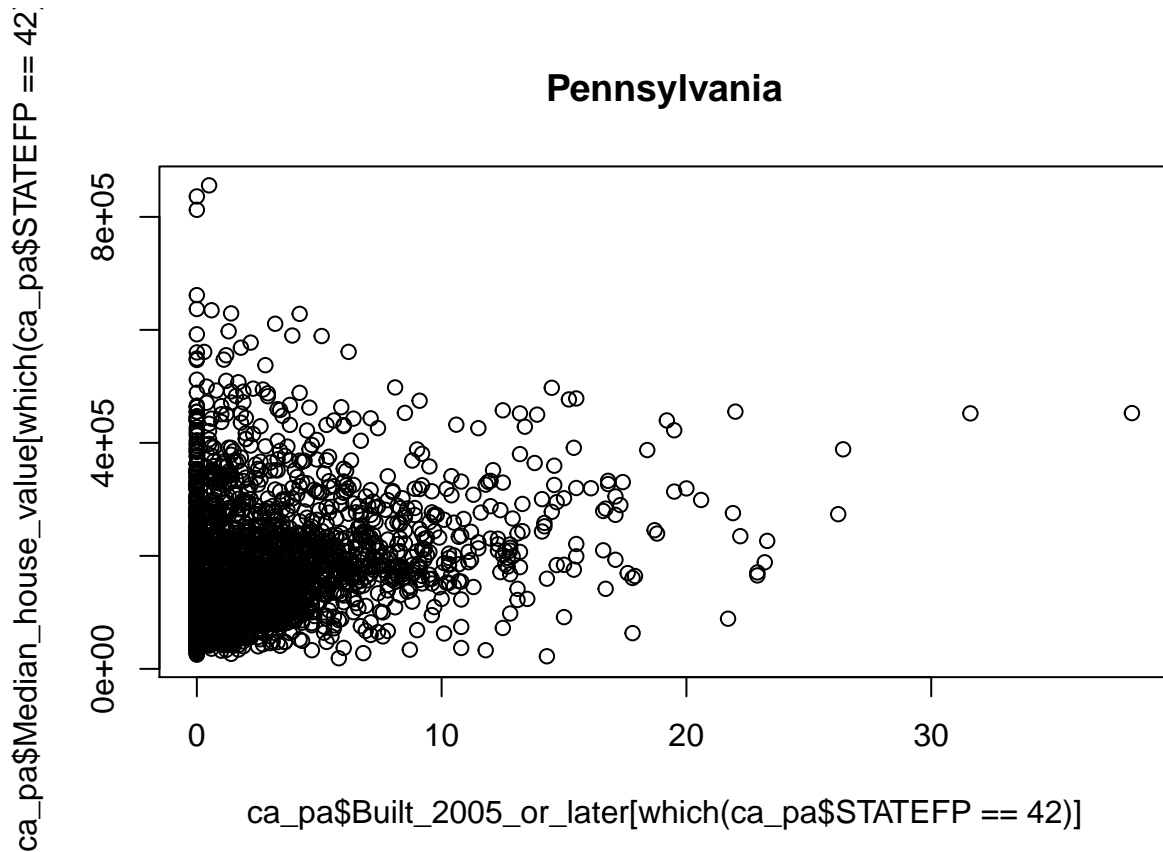


- Make a new plot, or pair of plots, which breaks this out by state. Note that the state is recorded in the `STATEFP` variable, with California being state 6 and Pennsylvania state 42.

```
plot(x=ca_pa$Built_2005_or_later[which(ca_pa$STATEFP==6)],
     y=ca_pa$Median_house_value[which(ca_pa$STATEFP==6)])
title("california")
```



```
plot(x=ca_pa$Built_2005_or_later[which(ca_pa$STATEFP==42)],
     y=ca_pa$Median_house_value[which(ca_pa$STATEFP==42)])
title("Pennsylvania")
```



### 3. Nobody Home

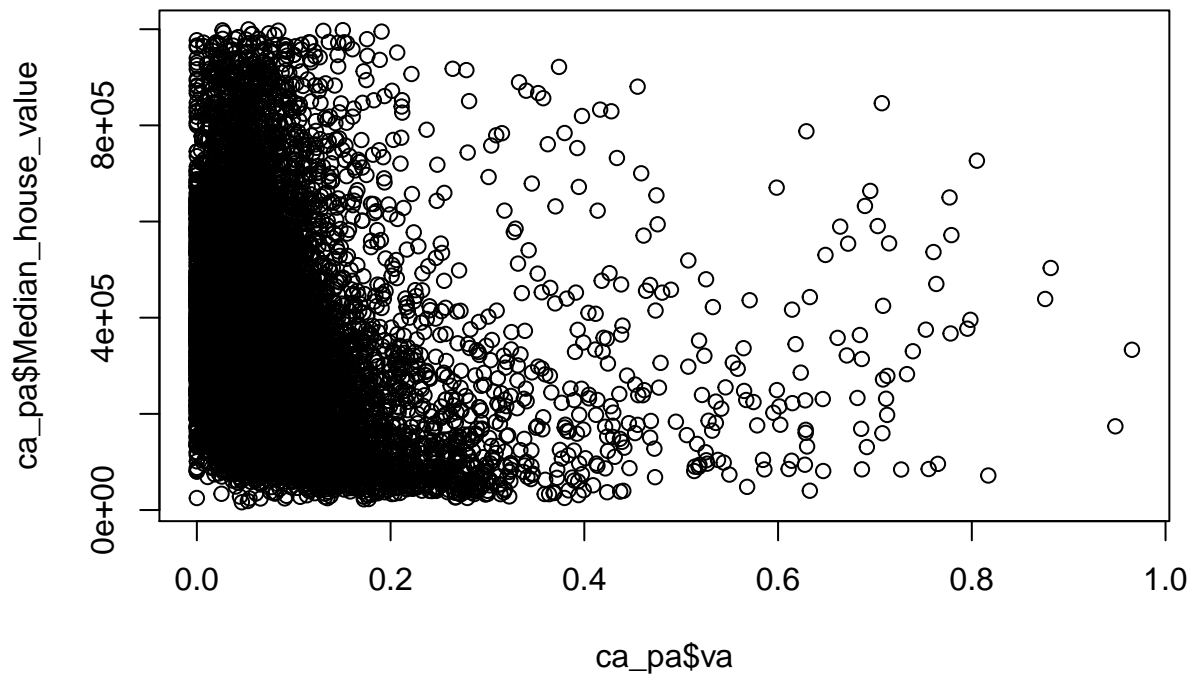
The vacancy rate is the fraction of housing units which are not occupied. The dataframe contains columns giving the total number of housing units for each Census tract, and the number of vacant housing units.

a. Add a new column to the dataframe which contains the vacancy rate. What are the minimum, maximum, mean, and median vacancy rates?

```
ca_pa = ca_pa %>%
  mutate("vacancy rate" = Vacant_units/Total_units)
```

b. Plot the vacancy rate against median house value.

```
plot(x=ca_pa$va,y=ca_pa$Median_house_value)
```

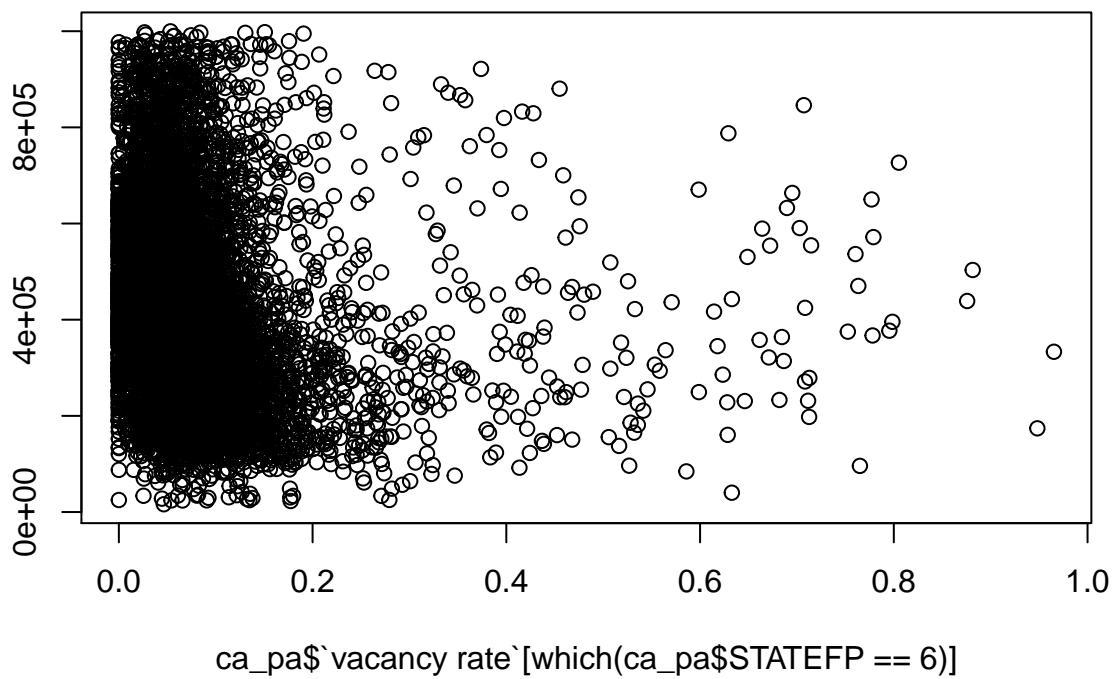


c. Plot vacancy rate against median house value separately for California and for Pennsylvania. Is there a difference?

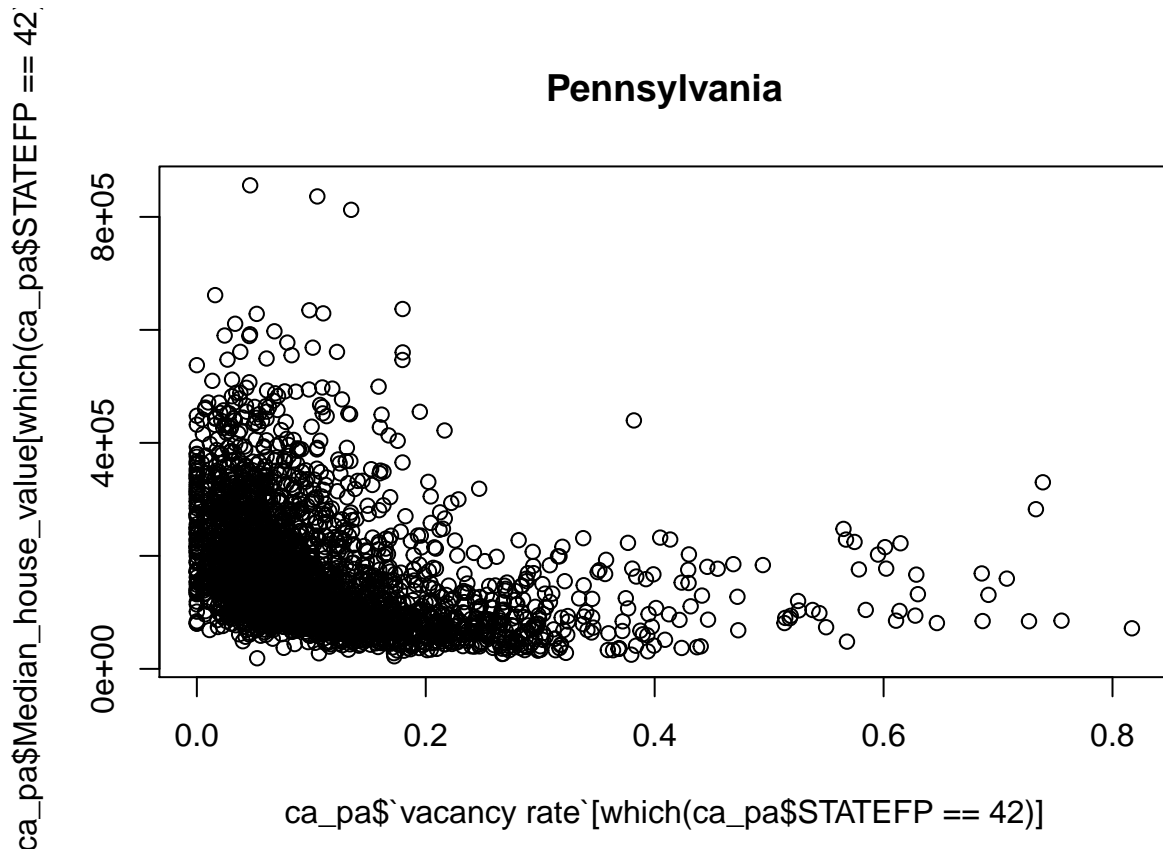
```
plot(x=ca_pa$`vacancy rate`[which(ca_pa$STATEFP==6)],
     y=ca_pa$Median_house_value[which(ca_pa$STATEFP==6)])
title("california")
```

ca\_pa\$Median\_house\_value[which(ca\_pa\$STATEFP == 6)]

## california



```
plot(x=ca_pa$`vacancy rate`[which(ca_pa$STATEFP==42)],  
     y=ca_pa$Median_house_value[which(ca_pa$STATEFP==42)])  
title("Pennsylvania")
```



The houses in Pennsylvania have relatively low price, and there is a more extinct trend that places with high vacancy rate tend to have lower house price.

4. The column COUNTYFP contains a numerical code for counties within each state. We are interested in Alameda County (county 1 in California), Santa Clara (county 85 in California), and Allegheny County (county 3 in Pennsylvania).
  - a. Explain what the block of code at the end of this question is supposed to accomplish, and how it does it. # the code get the median value of median house value in Alameda county. it first get all the observations of Alameda county into vector acca, then find the median value of the median house value in acca.
  - b. Give a single line of R which gives the same final answer as the block of code. Note: there are at least two ways to do this; you just have to find one.

```
median(ca_pa$Median_house_value[which(ca_pa$STATEFP==6&ca_pa$COUNTYFP==1)])
```

- c. For Alameda, Santa Clara and Allegheny Counties, what were the average percentages of housing built since 2005?

```
mean(ca_pa$Built_2005_or_later[which((ca_pa$STATEFP==6&ca_pa$COUNTYFP==1) |
                                       (ca_pa$STATEFP==6&ca_pa$COUNTYFP==85) |
                                       (ca_pa$STATEFP==42&ca_pa$COUNTYFP==3))])
```

- d. The `cor` function calculates the correlation coefficient between two variables. What is the correlation between median house value and the percent of housing built since 2005 in (i) the whole data, (ii) all of California, (iii) all of Pennsylvania, (iv) Alameda County, (v) Santa Clara County and (vi) Allegheny County?





```
## gender
## female    male
##      91      92

gender <- factor(gender, levels=c("male", "female"))
table(gender)
```

```
## gender
##    male female
##      92      91

gender <- factor(gender, levels=c("Male", "female"))
# Note the mistake: "Male" should be "male"
table(gender)
```

```
## gender
##    Male female
##         0      91

table(gender, exclude=NULL)
```

```
## gender
##    Male female  <NA>
##         0      91      92

rm(gender) # Remove gender
```

Explain the output from the successive uses of table().

1 factor() is used to encode a vector as a factor, table() list the values and their numbers in the factor 2 the values in the factor is sorted into increasing order of xtable() will list the value according to their order of appearance. 3 “Male” is not in the vector, thus it’s number is 0 4 when forming the factor using “factor(gender, levels=c(“Male”, “female”))”, “male” is excluded automatically, now set excluded=NULL, table() will show a more factor whose name is and number is the number of values excluded before.

MB.Ch1.12. Write a function that calculates the proportion of values in a vector x that exceed some value cutoff.

(a) Use the sequence of numbers 1, 2, . . . , 100 to check that this function gives the result that is expected.

```
cal<-function(x,k){
  if(typeof(x)=="list") total = length(x)
  else total = length(x)

  y = length(which(x>k))
  return(y/total)
}
x = seq(from=1,to=100,by=1)
cal(x,60)
```

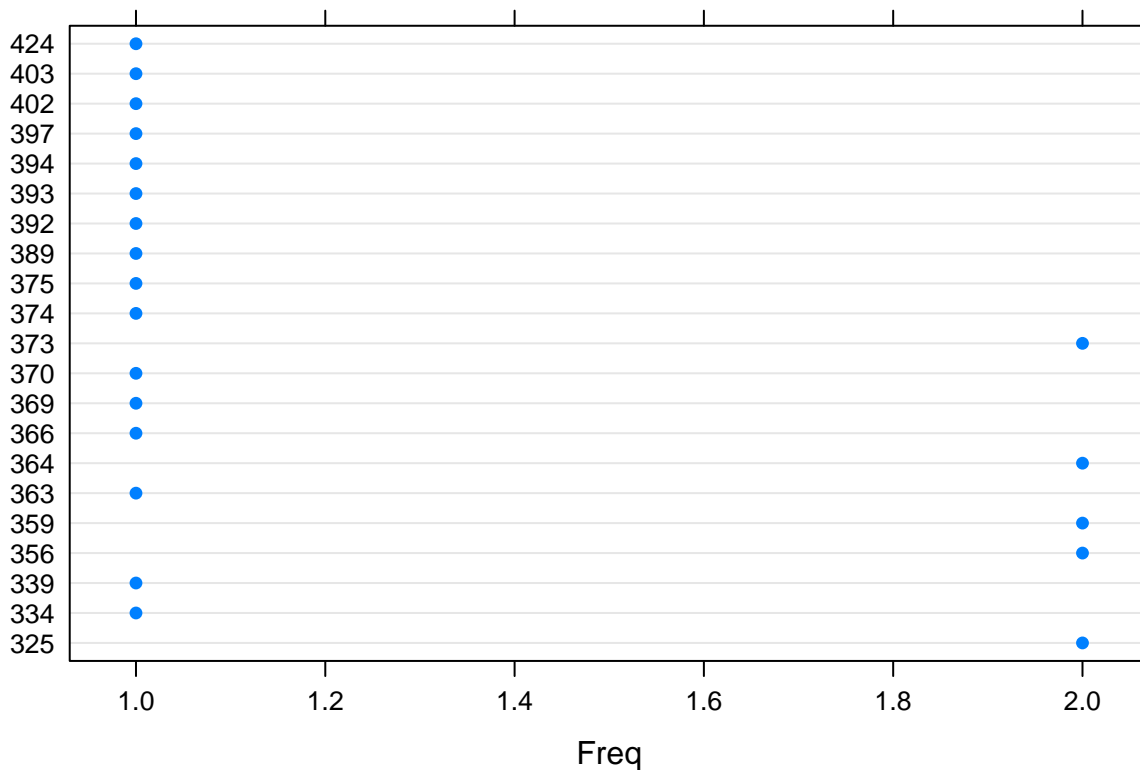
```
## [1] 0.4
```

(b) Obtain the vector ex01.36 from the Devore6 (or Devore7) package. These data give the times required for individuals to escape from an oil platform during a drill. Use dotplot() to show the distribution of times. Calculate the proportion of escape times that exceed 7 minutes.

```
library(Devore7)
```

```
## Loading required package: MASS
##
```

```
## Attaching package: 'MASS'
## The following object is masked from 'package:DAAG':
##
##      hills
## The following object is masked from 'package:dplyr':
##
##      select
dotplot(ex01.36)
```



```
cal(ex01.36,420)
```

```
##      C1
## 0.03846154
```

MB.Ch1.18. The Rabbit data frame in the MASS library contains blood pressure change measurements on five rabbits (labeled as R1, R2, . . . ,R5) under various control and treatment conditions. Read the help file for more information. Use the unstack() function (three times) to convert Rabbit to the following form:

```
Treatment Dose R1 R2 R3 R4 R5
1 Control 6.25 0.50 1.00 0.75 1.25 1.5
2 Control 12.50 4.50 1.25 3.00 1.50 1.5
....
```

```
library(MASS)
Rabbit = Rabbit
```

```

a1 = unstack(Rabbit,BPchange~Animal)
a2 = unstack(Rabbit,Treatment~Animal)
a3 = unstack(Rabbit,Dose~Animal)
data = cbind(a2$R1,a3$R1,a1)
names(data)<-c("Treatment", "Dose", "R1" , "R2" , "R3" , "R4" , "R5" )
data

```

| ##    | Treatment | Dose   | R1    | R2    | R3    | R4    | R5   |
|-------|-----------|--------|-------|-------|-------|-------|------|
| ## 1  | Control   | 6.25   | 0.50  | 1.00  | 0.75  | 1.25  | 1.5  |
| ## 2  | Control   | 12.50  | 4.50  | 1.25  | 3.00  | 1.50  | 1.5  |
| ## 3  | Control   | 25.00  | 10.00 | 4.00  | 3.00  | 6.00  | 5.0  |
| ## 4  | Control   | 50.00  | 26.00 | 12.00 | 14.00 | 19.00 | 16.0 |
| ## 5  | Control   | 100.00 | 37.00 | 27.00 | 22.00 | 33.00 | 20.0 |
| ## 6  | Control   | 200.00 | 32.00 | 29.00 | 24.00 | 33.00 | 18.0 |
| ## 7  | MDL       | 6.25   | 1.25  | 1.40  | 0.75  | 2.60  | 2.4  |
| ## 8  | MDL       | 12.50  | 0.75  | 1.70  | 2.30  | 1.20  | 2.5  |
| ## 9  | MDL       | 25.00  | 4.00  | 1.00  | 3.00  | 2.00  | 1.5  |
| ## 10 | MDL       | 50.00  | 9.00  | 2.00  | 5.00  | 3.00  | 2.0  |
| ## 11 | MDL       | 100.00 | 25.00 | 15.00 | 26.00 | 11.00 | 9.0  |
| ## 12 | MDL       | 200.00 | 37.00 | 28.00 | 25.00 | 22.00 | 19.0 |