

Homework 2: Housing Price

DingHuan, 3170102085

1. Loading and cleaning

a.

```
ca_pa <- read.csv("data/calif_penn_2011.csv", header = T)
```

b. The dataframe has 11275 rows and 34 columns.

```
dim(ca_pa)
```

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

c. The results below are hidden given that it is too long to show them all. `apply(ca_pa, c(1,2), is.na)` returns a matrix having the same dimension as `ca_pa`, whose elements are Boolean numbers indicating whether the data in `ca_pa` is NA. `colSums()` sums the columns and returns a named vector indicating how many NA elements are there in each column of `ca_pa`.

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

d.

```
ca_pa <- na.omit(ca_pa)
```

e. There are 670 rows containing NA elements, which is now removed.

```
11275 - nrow(ca_pa)
```

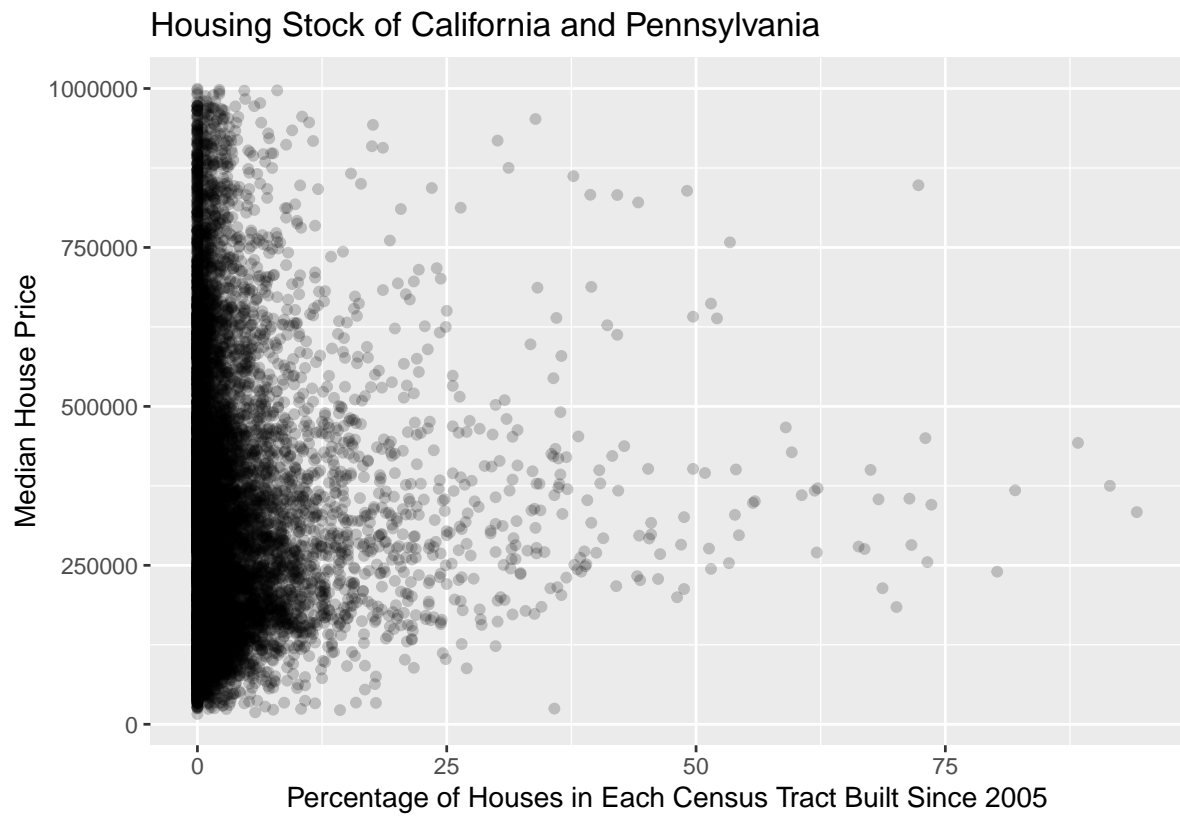
```
## [1] 670
```

f. The answers in (c) and (e) are compatible. Although we know the number of NA elements in each column, we still have no idea whether they are in the same row or not.

2. This Very New House

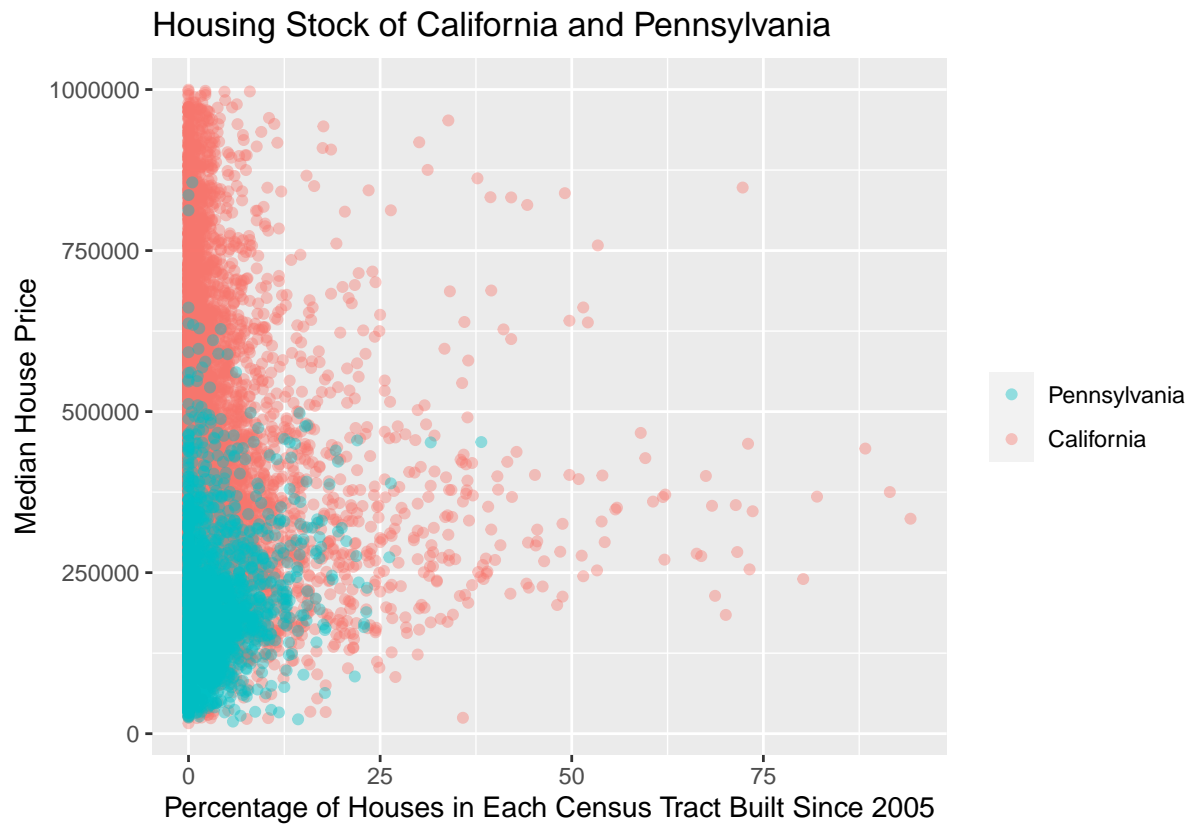
a.

```
library(tidyverse)
ca_pa %>% ggplot(aes(x = Built_2005_or_later, y = Median_house_value)) +
  geom_point(alpha = 0.2) +
  labs(x = "Percentage of Houses in Each Census Tract Built Since 2005",
       y = "Median House Price",
       title = "Housing Stock of California and Pennsylvania")
```



b.

```
ca_pa %>% ggplot() +
  geom_point(aes(x = Built_2005_or_later, y = Median_house_value,
                 color = (STATEFP==42)),
             alpha = 0.4) +
  labs(x = "Percentage of Houses in Each Census Tract Built Since 2005",
       y = "Median House Price",
       title = "Housing Stock of California and Pennsylvania") +
  scale_colour_hue(element_blank(), breaks=c(TRUE, FALSE),
                  labels=c("Pennsylvania", "California"))
```



3. *Nobody Home*

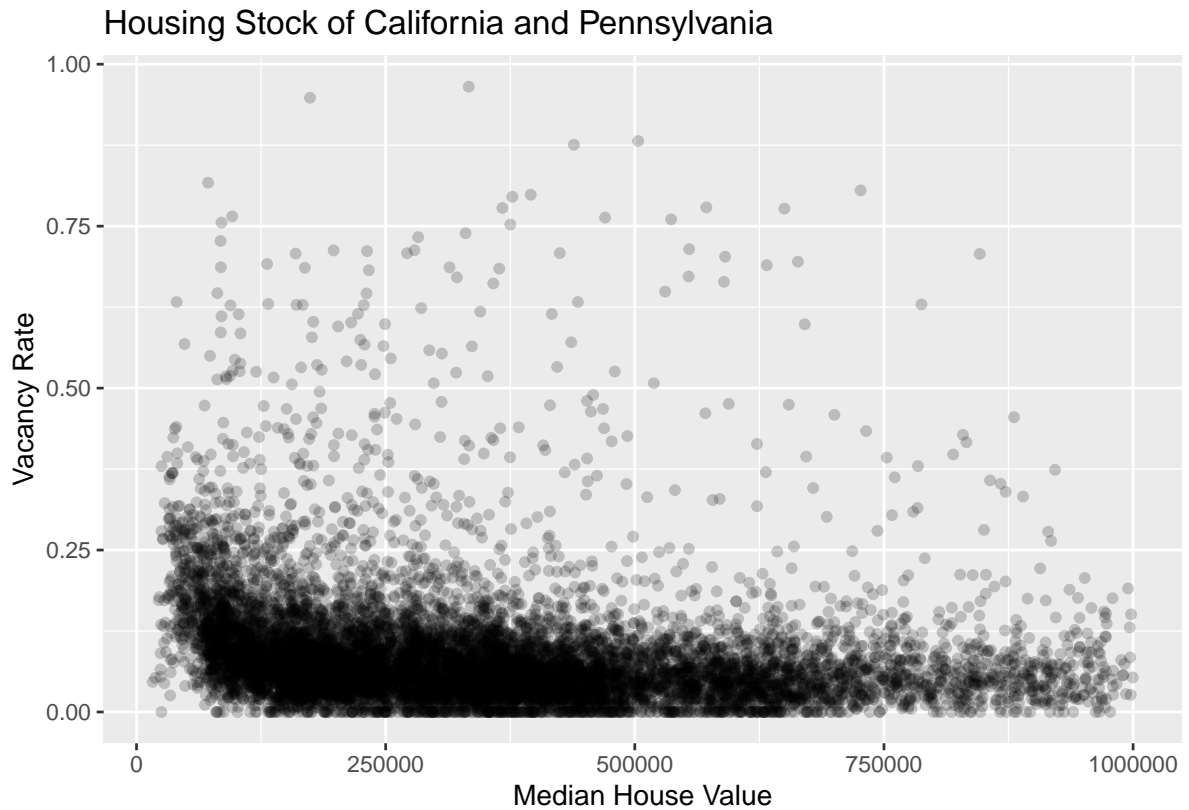
- a. The minimum, maximum, mean, and median vacancy rates are 0.00000, 0.96531, 0.08889 and 0.06767 respectively.

```
ca_pa <- ca_pa %>% mutate(Vacancy_rate = Vacant_units / Total_units)
summary(ca_pa$Vacancy_rate)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.00000 0.03846 0.06767 0.08889 0.10921 0.96531
```

- b.

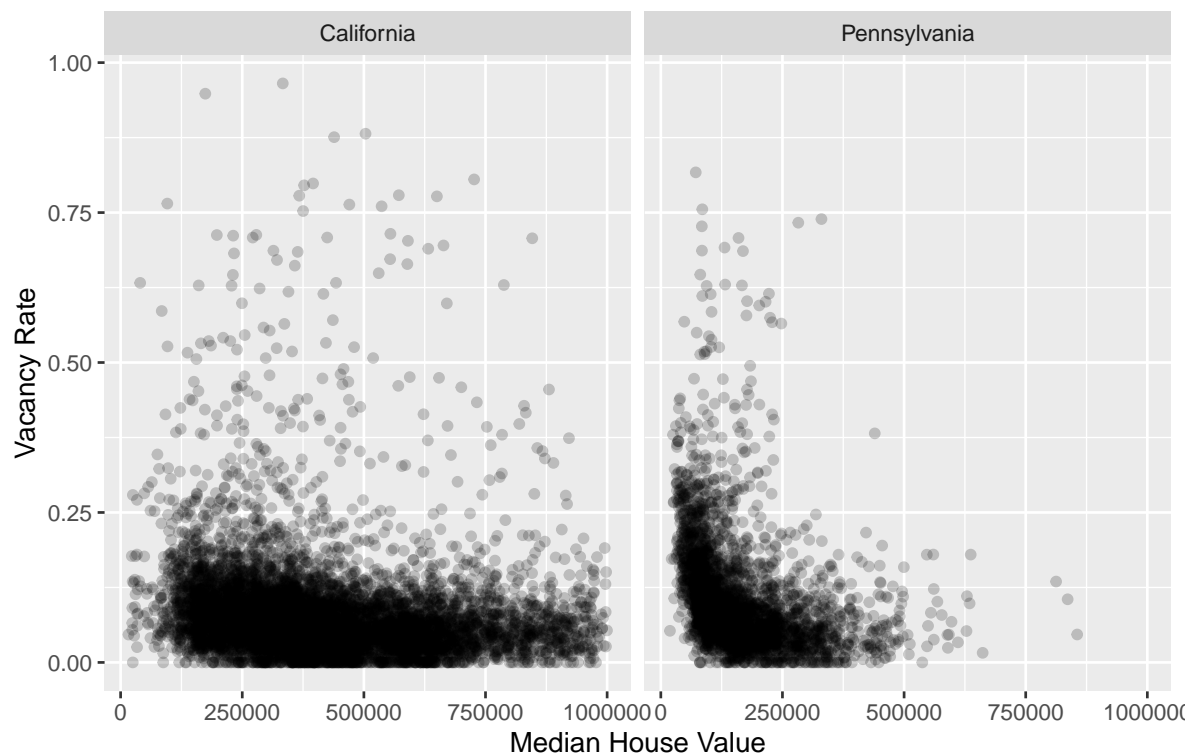
```
ca_pa %>% ggplot(aes(x = Median_house_value, y = Vacancy_rate)) +
  geom_point(alpha = 0.2) +
  labs(x = "Median House Value",
       y = "Vacancy Rate",
       title = "Housing Stock of California and Pennsylvania")
```



- c. It is clear that there are fewer Census tracts with high median house value in Pennsylvania, among which there are more Census tracts have higher vacancy rate. Although the number of Census tracts with high median house value in Pennsylvania is not high, their vacancy rate is much lower than those in California. The distribution of vacancy rate among Census tracts in California seems not to change with median house value.

```
ca_pa %>% ggplot(aes(x = Median_house_value, y = Vacancy_rate)) +
  geom_point(alpha = 0.2) +
  labs(x = "Median House Value",
       y = "Vacancy Rate",
       title = "Housing Stock of California and Pennsylvania") +
  facet_wrap(~ STATEFP,
             labeller = as_labeller(c('6' = "California", '42' = "Pennsylvania")))
```

Housing Stock of California and Pennsylvania



4. a. The first iteration records the row numbers of the county marked as 1 in California to the variable `acca`. The second iteration records the median house value of the Census tracts in `acca` to the variable `accamhv`, and finally calculate thier median value (the median value of the median values of the Census tracts recorded in `accamhv`).

```
acca <- c()
for (tract in 1:nrow(ca_pa)) {
  if (ca_pa$STATEFP[tract] == 6) {
    if (ca_pa$COUNTYFP[tract] == 1) {
      acca <- c(acca, tract)
    }
  }
}
accamhv <- c()
for (tract in acca) {
  accamhv <- c(accamhv, ca_pa[tract,10])
}
median(accamhv)
```

```
## [1] 474050
```

b.

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

```
## [1] 474050
```

c.

```
ca_pa_long <- ca_pa %>%
  gather(key = built_year,
```

```

      value = built_numbers,
      dplyr::starts_with('Built'))
ca_pa_long %>%
  filter((STATEFP == 6 & COUNTYFP %in% c(1,85)) | (STATEFP == 42 & COUNTYFP == 3)) %>%
  group_by(COUNTYFP) %>%
  summarise(APHB_2005 = sum((built_year == "Built_2005_or_later") *
                           built_numbers)/sum(built_numbers)) %>%
  ungroup()

```

```

## # A tibble: 3 x 2
##   COUNTYFP APHB_2005
##   <int>     <dbl>
## 1       1     0.0282
## 2       3     0.0147
## 3      85     0.0320

```

d.

```

p <- ca_pa_long %>% group_by(X) %>%
  summarise(APHB_2005 = sum((built_year == "Built_2005_or_later") *
                           built_numbers)/sum(built_numbers)) %>%
  ungroup()
ca_pa <- dplyr::left_join(ca_pa, p, by = "X")
rm(p)
cor(ca_pa$Median_house_value, ca_pa$APHB_2005)

```

```
## [1] -0.01893763
```

```

ca_pa %>% mutate(STATENAME = ifelse(STATEFP == 6, "California", "Pennsylvania")) %>%
  group_by(STATENAME) %>%
  summarise(Correlation_coefficient = cor(Median_house_value, APHB_2005))

```

```

## # A tibble: 2 x 2
##   STATENAME   Correlation_coefficient
##   <chr>             <dbl>
## 1 California        -0.115
## 2 Pennsylvania       0.268

```

```

ca_pa %>%
  filter((STATEFP == 6 & COUNTYFP %in% c(1,85)) | (STATEFP == 42 & COUNTYFP == 3)) %>%
  mutate(COUNTYNAME = ifelse(COUNTYFP == 1, "Alameda County",
                             ifelse(COUNTYFP == 85, "Santa Clara",
                                    "Allegheny County"))) %>%
  group_by(COUNTYNAME) %>%
  summarise(Correlation_coefficient = cor(Median_house_value, APHB_2005))

```

```

## # A tibble: 3 x 2
##   COUNTYNAME   Correlation_coefficient
##   <chr>             <dbl>
## 1 Alameda County    0.0130
## 2 Allegheny County  0.194
## 3 Santa Clara      -0.173

```

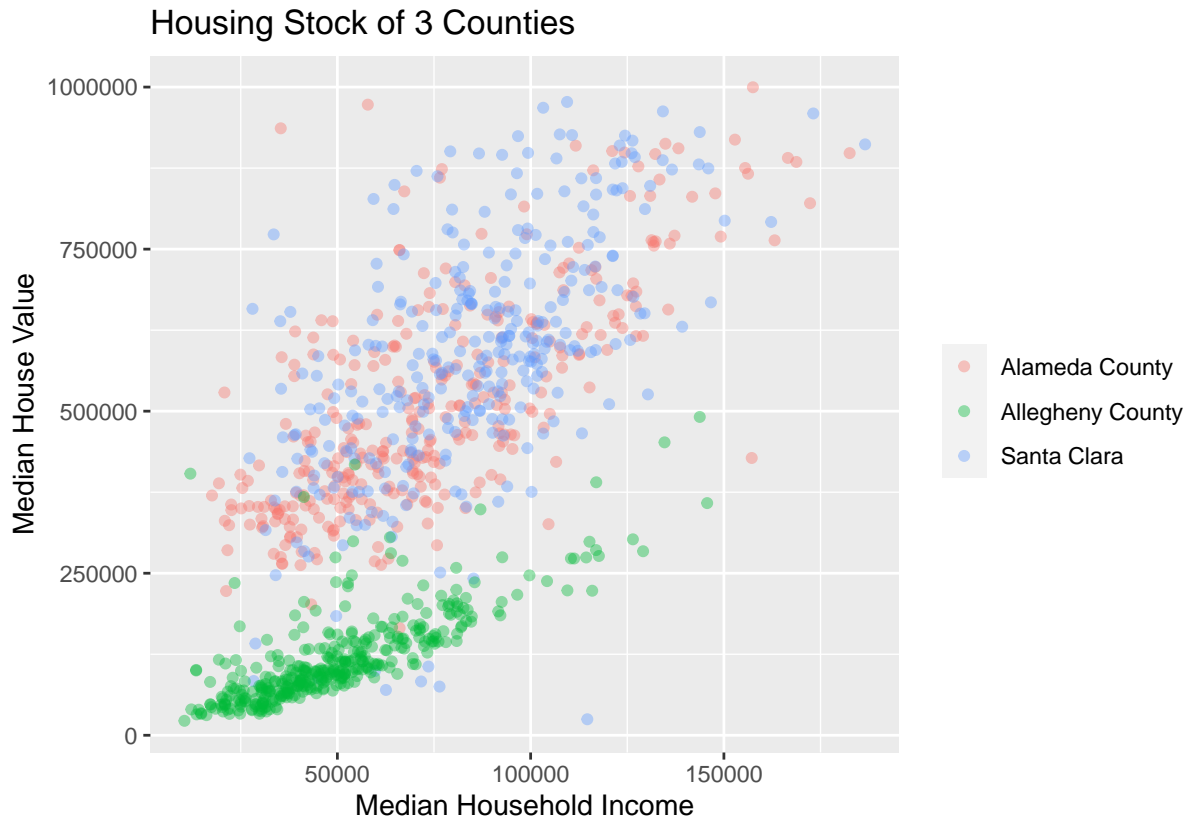
e.

```

ca_pa %>%
  filter((COUNTYFP %in% c(1,85) & STATEFP == 6) | (COUNTYFP==3&STATEFP == 42)) %>%

```

```
mutate(COUNTYNAME = ifelse(COUNTYFP == 1, "Alameda County",
                           ifelse(COUNTYFP == 85, "Santa Clara",
                                   "Allegheny County"))) %>%
ggplot(aes(x = Median_household_income, y=Median_house_value, color = COUNTYNAME)) +
geom_point(alpha = 0.4) +
labs(x = "Median Household Income",
     y = "Median House Value",
     title = "Housing Stock of 3 Counties") +
theme(legend.title=element_blank())
```



5. (MB.CH1.11) The first line create an variable **gender** with 2 levels “female” and “male”, whose first 91 elements are “female” and the remaining elements are “male”. The second line **table(gender)** shows the factor levels. The third line exchanges the order of levels. The function searches **gender** first, finding the same level as elements in **levels=c("male", "female")**, and then changes the numeric order to the new levels. But when it doesn't find the same level, just as line 5, the function removes the old levels and create a new level named **Male** and matches nothing, thus the result of **table(gender)** is 0 for level “Male”. When NA is included in table, like what line8 does, we can see an NA level with 92 elements, which are exactly those whose level “male” are removed in line5.

```
gender <- factor(c(rep("female", 91), rep("male", 92)))
table(gender)
```

```
## 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
```

6. (MB.CH1.12)

a.

```
cutoff <- function(x, value){
  prop = sum(x > value) / length(x)
  return(prop)
}
cutoff(1:100, 10)
```

```
## [1] 0.9

cutoff(1:100, 35)
```

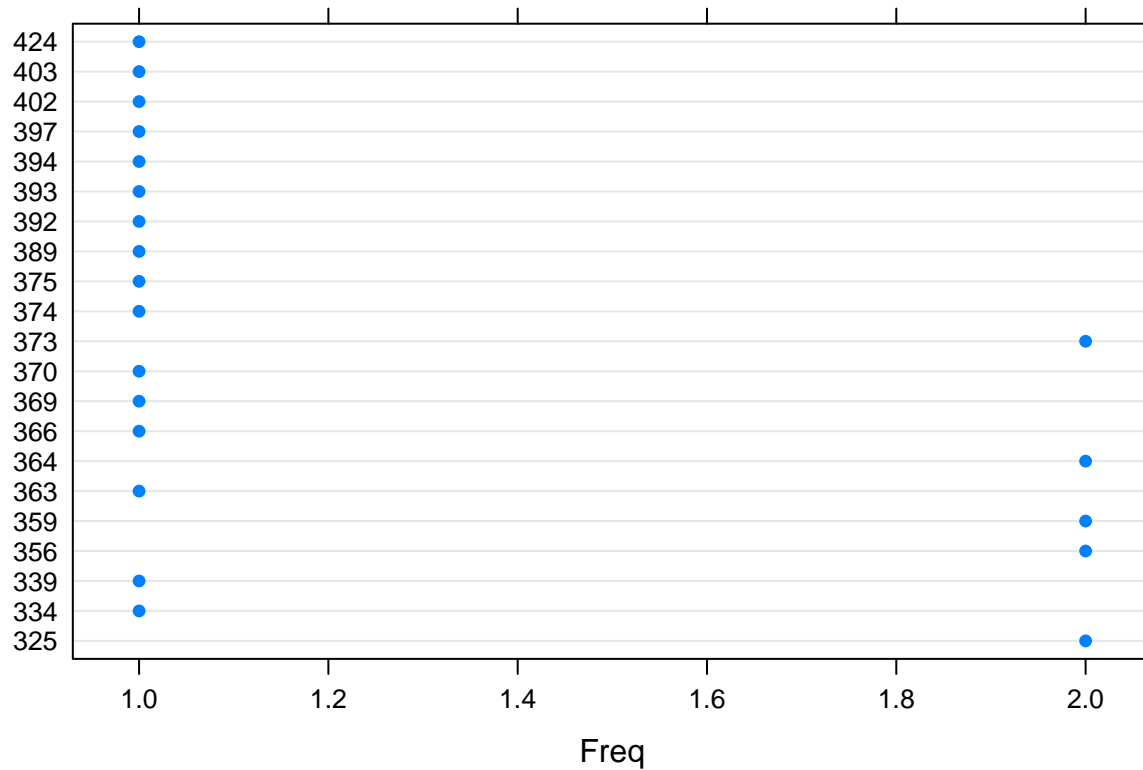
```
## [1] 0.65

cutoff(1:100, 35.5)
```

```
## [1] 0.65
```

b.

```
library(Devore7)
dotplot(ex01.36)
```

```
cutoff(ex01.36$C1, 420)
```

```
## [1] 0.03846154
```

7. (MB.CH1.18)

```
library(MASS)
data(Rabbit)
Treatment <- unstack(Rabbit, Treatment ~ Animal)
Dose <- unstack(Rabbit, Dose ~ Animal)
BPchange <- unstack(Rabbit, BPchange ~ Animal)
Rabbit <- data.frame(Treatment = Treatment[,1], Dose = Dose[,1])
Rabbit <- cbind(Rabbit, BPchange)
Rabbit
```

```
##   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
```