

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`.

```
ca_pa<-read.csv("C:/Users/lenovo/Documents/github/Rcourse2020/data/calif_penn_2011.csv")
```

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

```
row1<-nrow(ca_pa)
row1
```

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

```
col1<-ncol(ca_pa)
col1
```

```
## [1] 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
```

这条命令的结果是 `ca_pa` 每一列中 NA 值的个数。

- 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?

```
row2<-nrow(ca_pa)
row1-row2
```

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

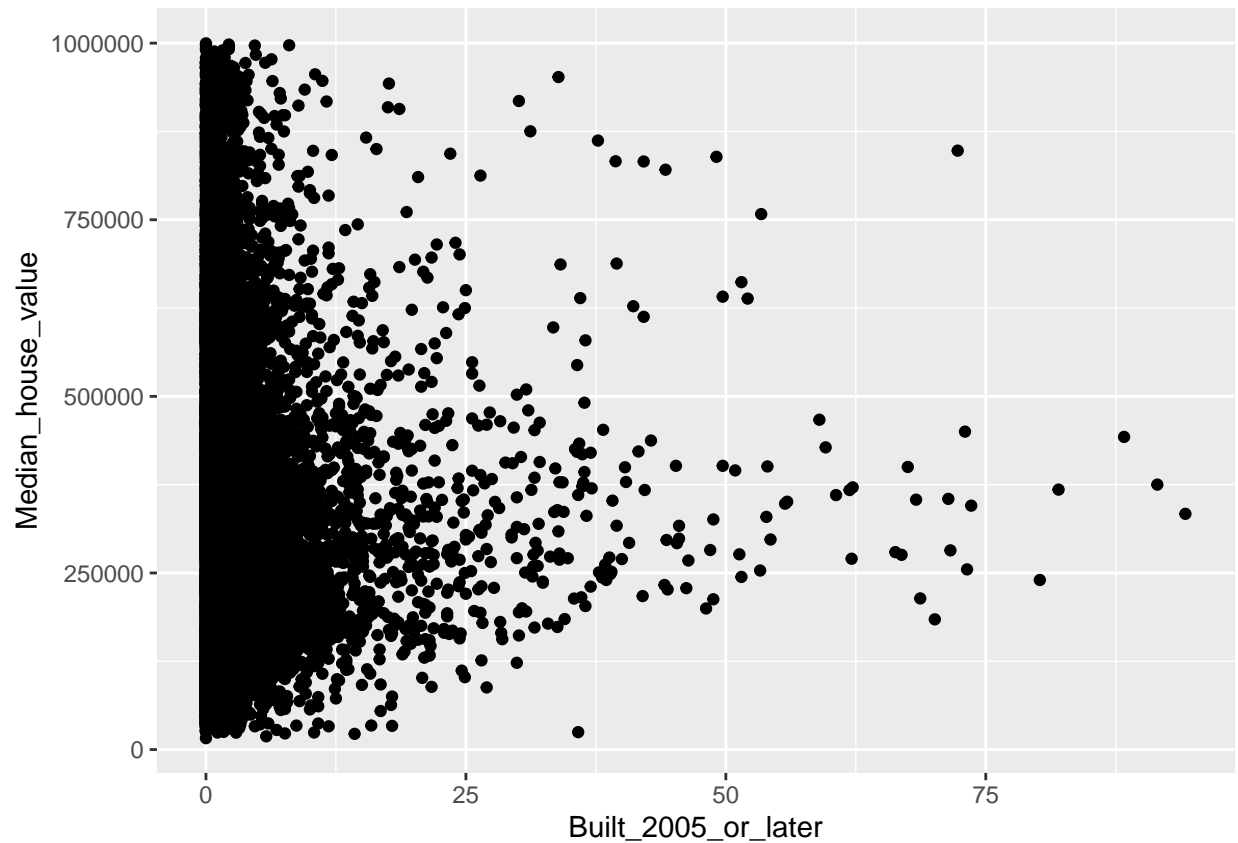
- f. Are your answers in (c) and (e) compatible? Explain.

(c) 中命令的结果是 `ca_pa` 每一列中 NA 值的个数, (e) 中命令的结果是 `ca_pa` 任何列有 NA 值的行的个数。两者结果不冲突。

2. *This Very New House*

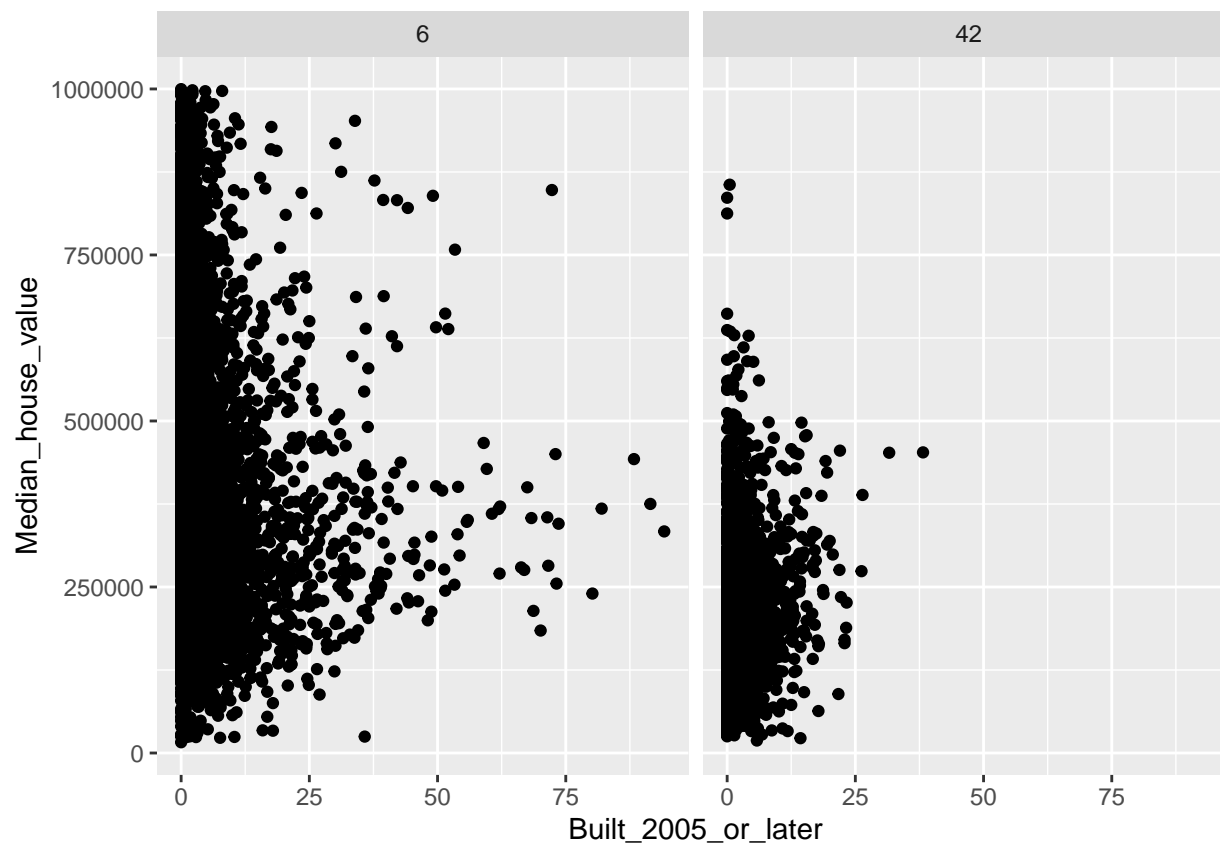
- a. The variable `Built_2005_or_later` indicates the percentage of houses in each Census tract built since 2005. Plot median house prices against this variable.

```
p1<-ggplot(data = ca_pa)+
  geom_point(aes(x=Built_2005_or_later,y=Median_house_value))
p1
```



b. 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.

```
p2<-ggplot(data = ca_pa)+
  geom_point(aes(x=Built_2005_or_later,y=Median_house_value))+
  facet_wrap(~ STATEFP)
p2
```



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.

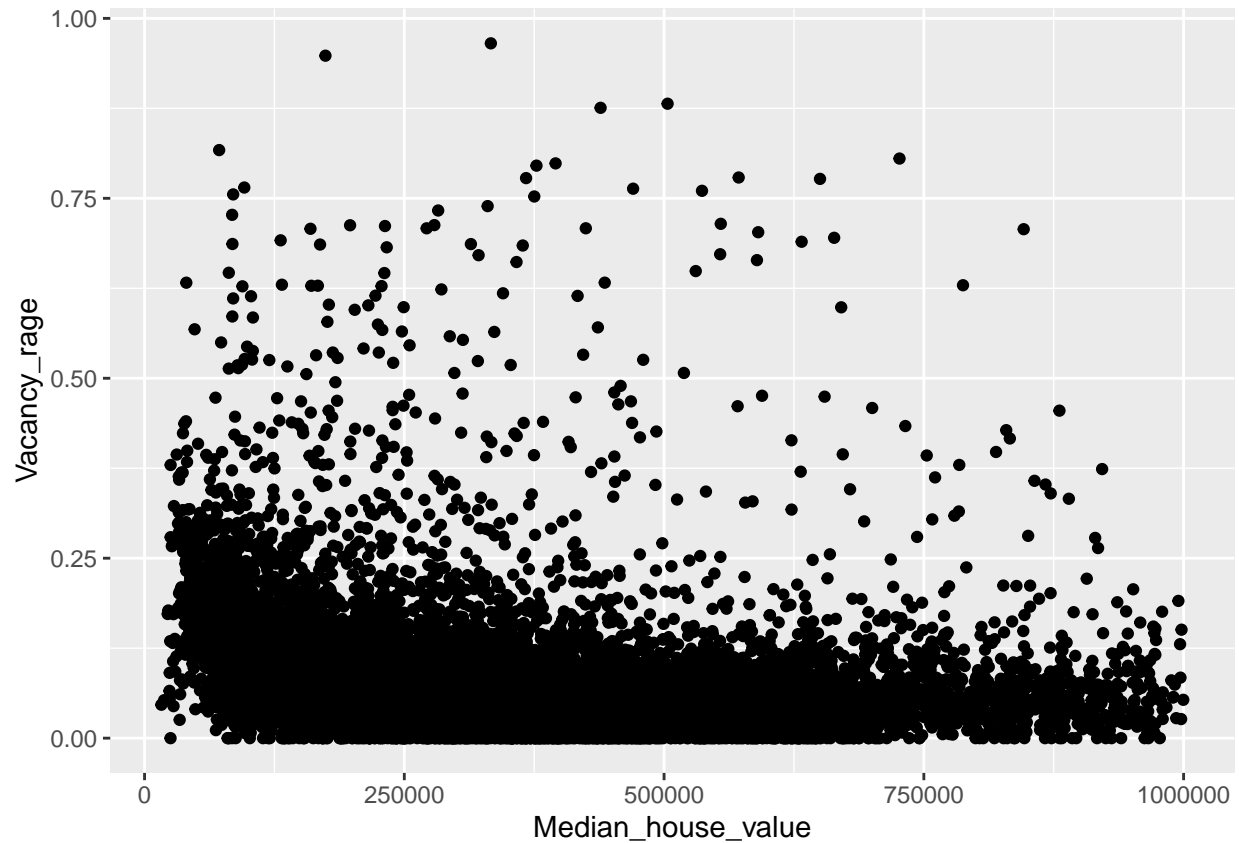
- a. Add a new column to the dataframe which contains the vacancy rate. What are the minimum, maximum, mean, and median vacancy rates?

```
ca_pa<-cbind(ca_pa,"Vacancy_rate"=ca_pa[, "Vacant_units"]/ca_pa[, "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
```

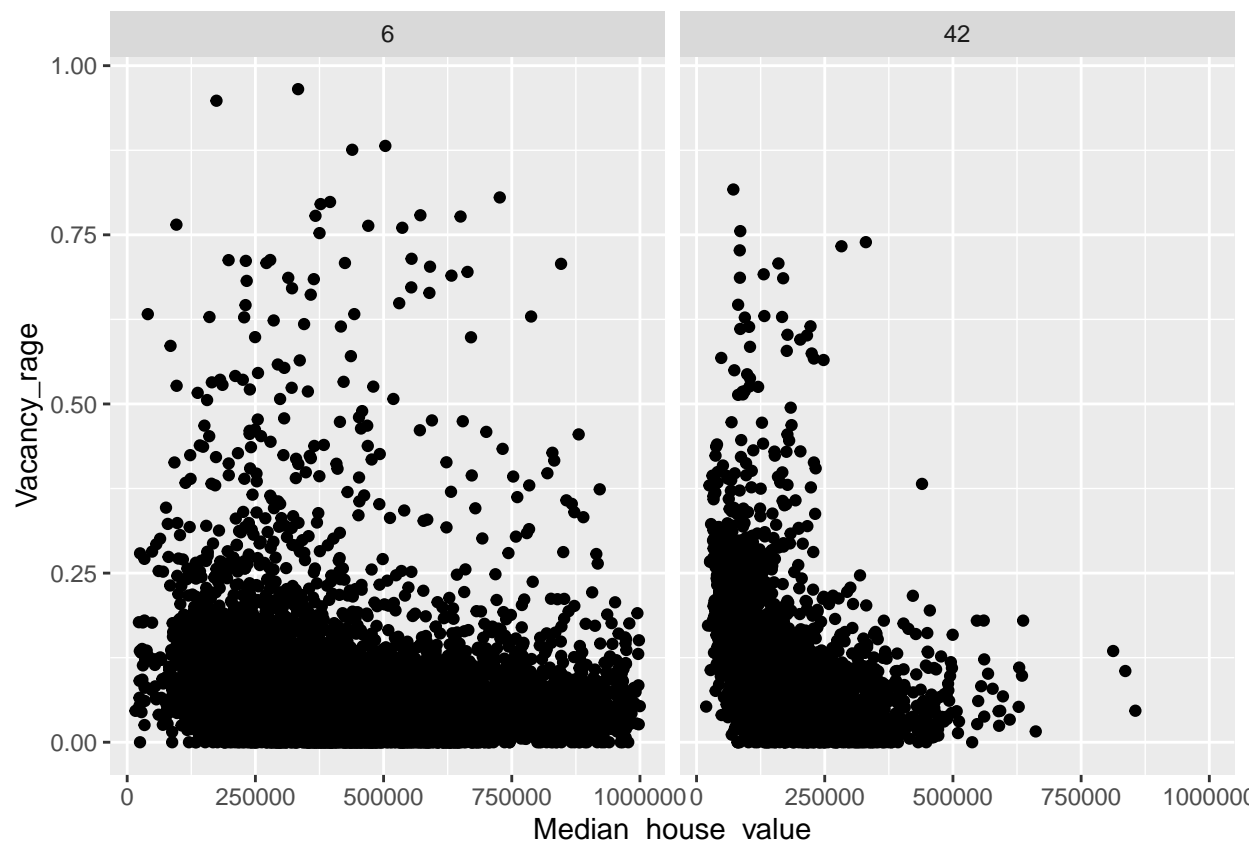
- b. Plot the vacancy rate against median house value.

```
p3<-ggplot(data = ca_pa)+
  geom_point(aes(x=Median_house_value,y=Vacancy_rate))
p3
```



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

```
p3<-ggplot(data = ca_pa)+
  geom_point(aes(x=Median_house_value,y=Vacancy_rate))+
  facet_wrap(~ STATEFP)
p3
```



两个图有区别，California 房屋价值对房屋空置率的影响比 Pennsylvania 的低。不过总的来说，房屋价值越高，空置率越低。

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.

```

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)

```

这段代码想要求出 Alameda County 的房屋价值的中位数。方法：首先将 Alameda County 的所有数据读入到 acca，然后将 Alameda County 的每个区的房屋价值中位数读入到 accamhv，再对 accamhv 求中位数。

- 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(subset(ca_pa, STATEFP==6 & COUNTYFP==1)[,10])
```

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

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

```
mean(subset(ca_pa, STATEFP==6 & COUNTYFP==1)[,16])
```

```
## [1] 2.820468
```

```
mean(subset(ca_pa, STATEFP==6 & COUNTYFP==85)[,16])
```

```
## [1] 3.200319
```

```
mean(subset(ca_pa, STATEFP==42 & COUNTYFP==3)[,16])
```

```
## [1] 1.474219
```

- 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?

```
cor(ca_pa[,10], ca_pa[,16])
```

```
## [1] -0.01893186
```

```
cor(subset(ca_pa, STATEFP==6)[,10], subset(ca_pa, STATEFP==6)[,16])
```

```
## [1] -0.1153604
```

```
cor(subset(ca_pa, STATEFP==42)[,10], subset(ca_pa, STATEFP==42)[,16])
```

```
## [1] 0.2681654
```

```
cor(subset(ca_pa, STATEFP==6 & COUNTYFP==1)[,10], subset(ca_pa, STATEFP==6 & COUNTYFP==1)[,16])
```

```
## [1] 0.01303543
```

```
cor(subset(ca_pa, STATEFP==6 & COUNTYFP==85)[,10], subset(ca_pa, STATEFP==6 & COUNTYFP==85)[,16])
```

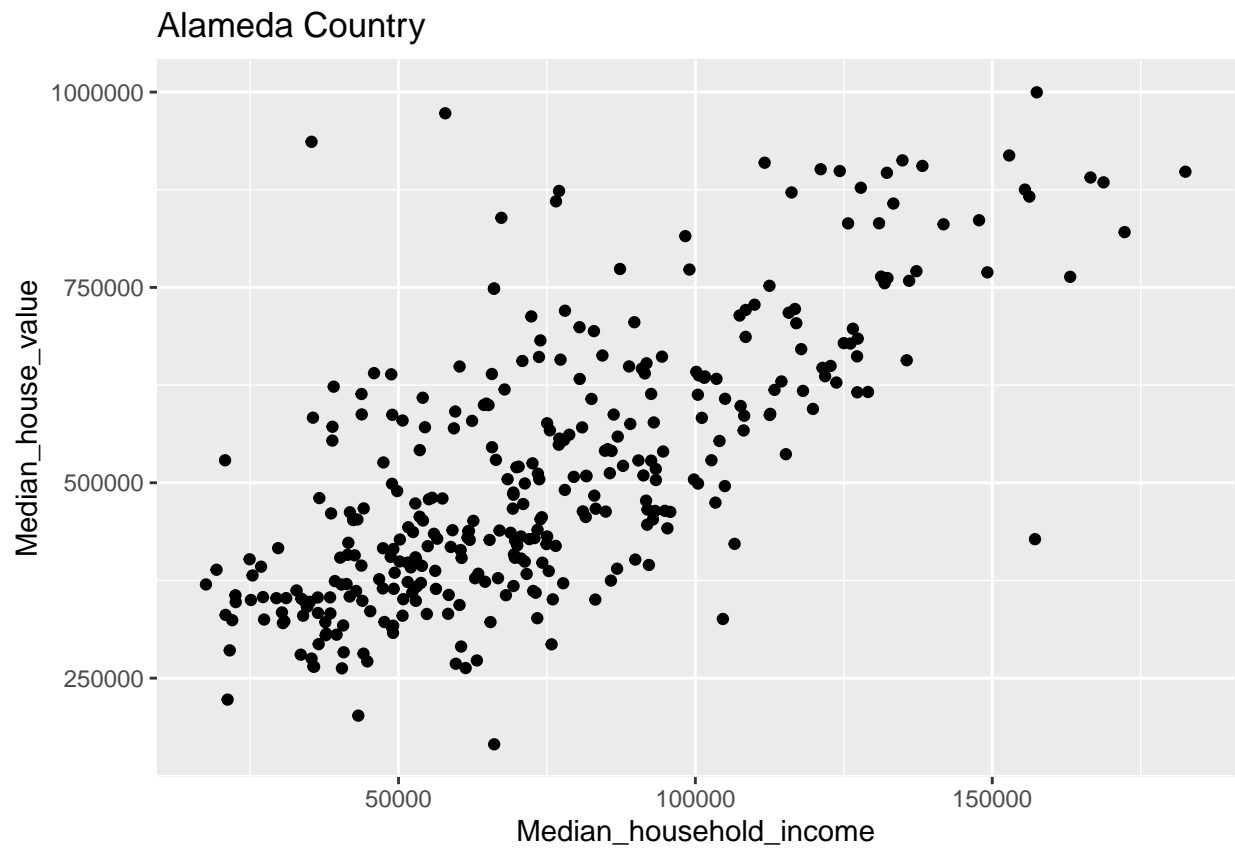
```
## [1] -0.1726203
```

```
cor(subset(ca_pa, STATEFP==42 & COUNTYFP==3)[,10], subset(ca_pa, STATEFP==42 & COUNTYFP==3)[,16])
```

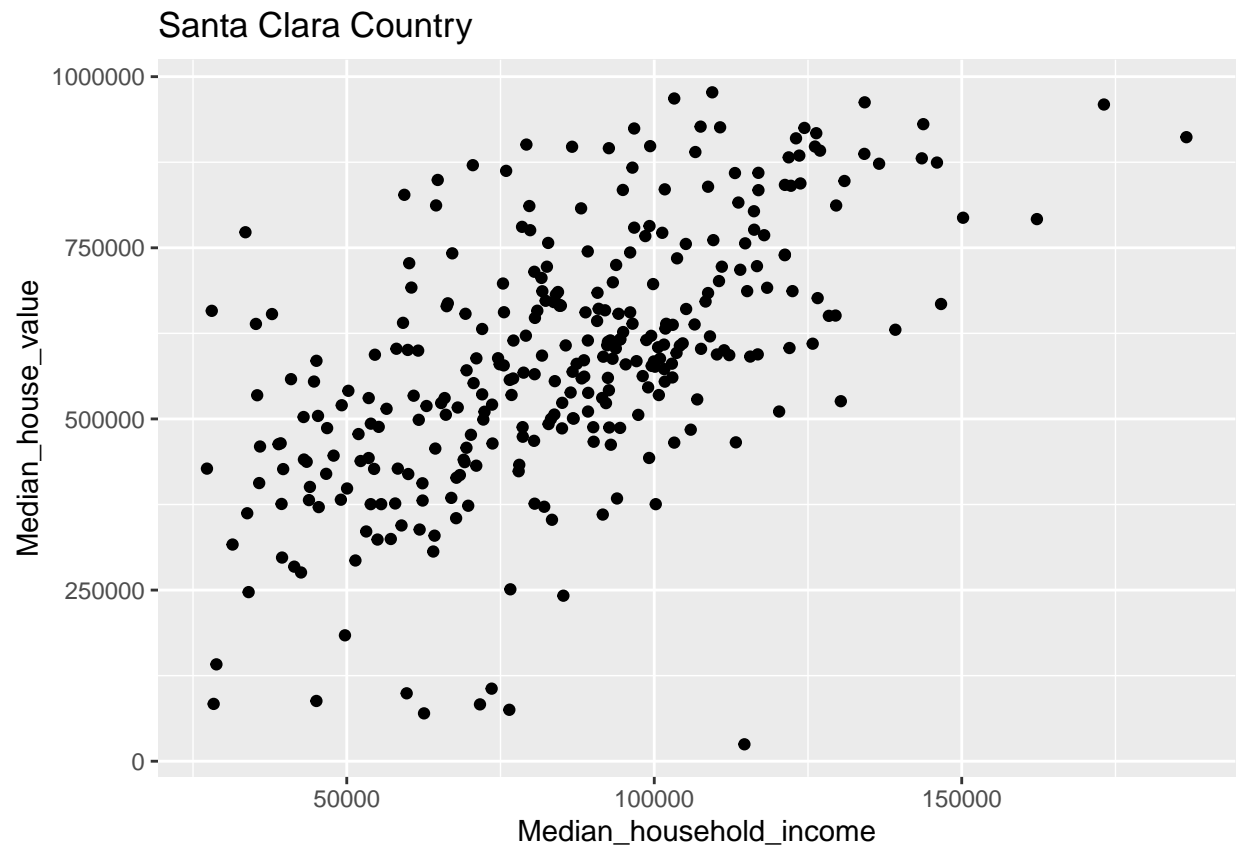
```
## [1] 0.1939652
```

- e. Make three plots, showing median house values against median income, for Alameda, Santa Clara, and Allegheny Counties. (If you can fit the information into one plot, clearly distinguishing the three counties, that's OK too.)

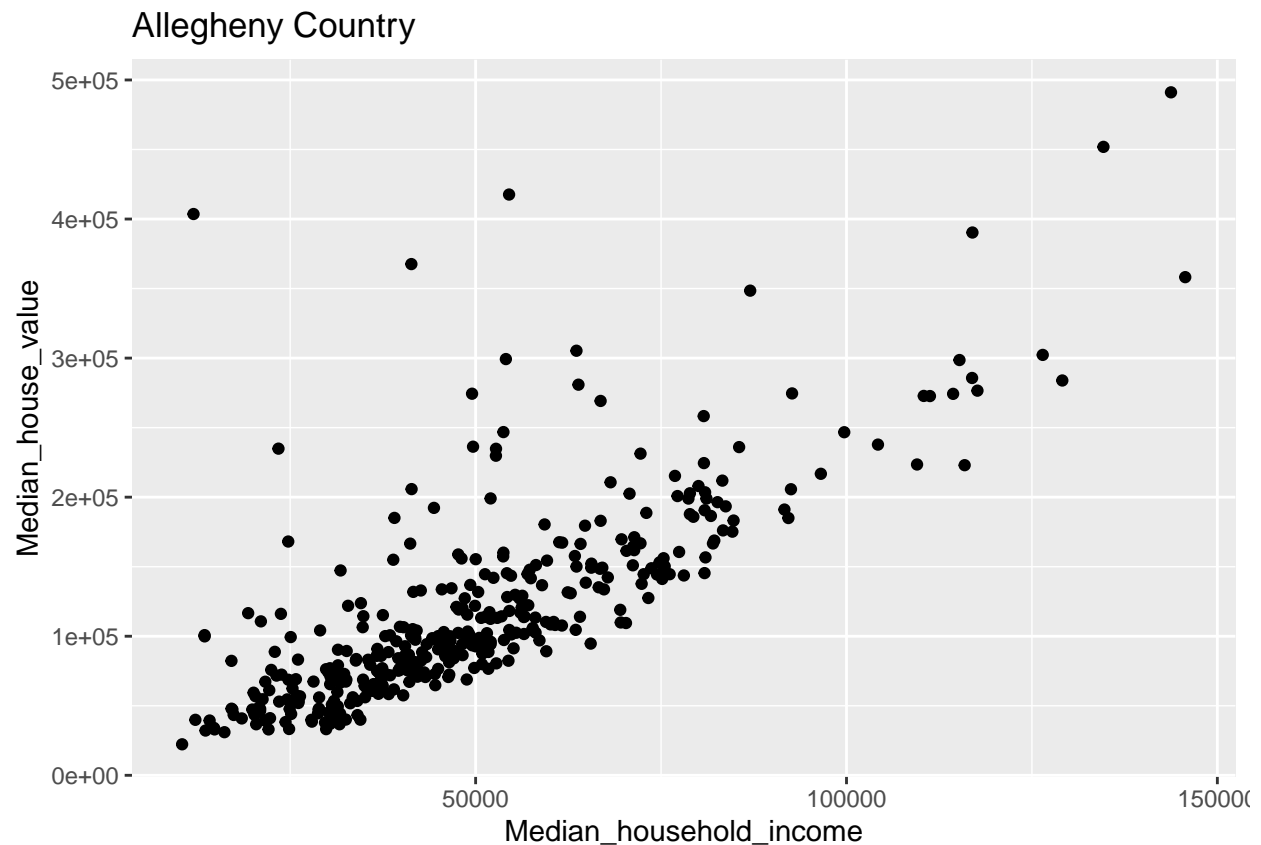
```
p4<-ggplot(data = subset(ca_pa, STATEFP==6 & COUNTYFP==1))+  
  geom_point(aes(x=Median_household_income, y=Median_house_value))+  
  ggtitle("Alameda Country")  
p4
```



```
p5<-ggplot(data = subset(ca_pa,STATEFP==6 & COUNTYFP==85))+  
  geom_point(aes(x=Median_household_income,y=Median_house_value))+  
  ggtitle("Santa Clara Country")  
p5
```

```
p6<-ggplot(data = subset(ca_pa,STATEFP==42 & COUNTYFP==3))+  
  geom_point(aes(x=Median_household_income,y=Median_house_value))+  
  ggtitle("Allegheny Country")  
p6
```



MB.Ch1.11. Run the following code:

```
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
```

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

table() 显示 factor 数据 levels 中每个元素的名称和频数。

第一个 table() 显示了 gender 中 female 和 male 的个数。

第二个 table() 显示了 gender 中 male 和 female 的个数。

第三个 table() 显示了 Male 和 female 的个数。显然 gender 中不包含 Male 类别的数据，所以 Male 的个数为零。

第四个 table() 显示了 Male 和 female 及其他类别数据的个数，“exclude=NULL”表示在表中包含 NA 值。Male 的个数为零，female 的个数为 91，NA 的个数为 92（即 male 的个数）。

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

```
f<-function(x,cutoff){  
  sum(x>cutoff)/length(x)  
}
```

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

```
x1<-c(1:100)  
f(x1,10)
```

```
## [1] 0.9
```

```
f(x1,70)
```

```
## [1] 0.3
```

- (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.

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

....

```
cbind(Treatment = unstack(Rabbit, Treatment ~ Animal)[,1],  
      Dose = unstack(Rabbit, Dose ~ Animal)[,1],  
      unstack(Rabbit, BPchange ~ Animal))
```

```
##      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