

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*

- Load the data into a dataframe called `ca_pa`.
- How many rows and columns does the dataframe have?
- Run this command, and explain, in words, what this does:
`colSums(apply(ca_pa,c(1,2),is.na))`
- 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.
- How many rows did this eliminate?
- Are your answers in (c) and (e) compatible? Explain.

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.
- 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.

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.

- Add a new column to the dataframe which contains the vacancy rate. What are the minimum, maximum, mean, and median vacancy rates?
- Plot the vacancy rate against median house value.
- Plot vacancy rate against median house value separately for California and for Pennsylvania. Is there a difference?

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).

- Explain what the block of code at the end of this question is supposed to accomplish, and how it does it.
- 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.
- For Alameda, Santa Clara and Allegheny Counties, what were the average percentages of housing built since 2005?
- 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?
- 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.)

```
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)
```