# Exploring the housing crisis
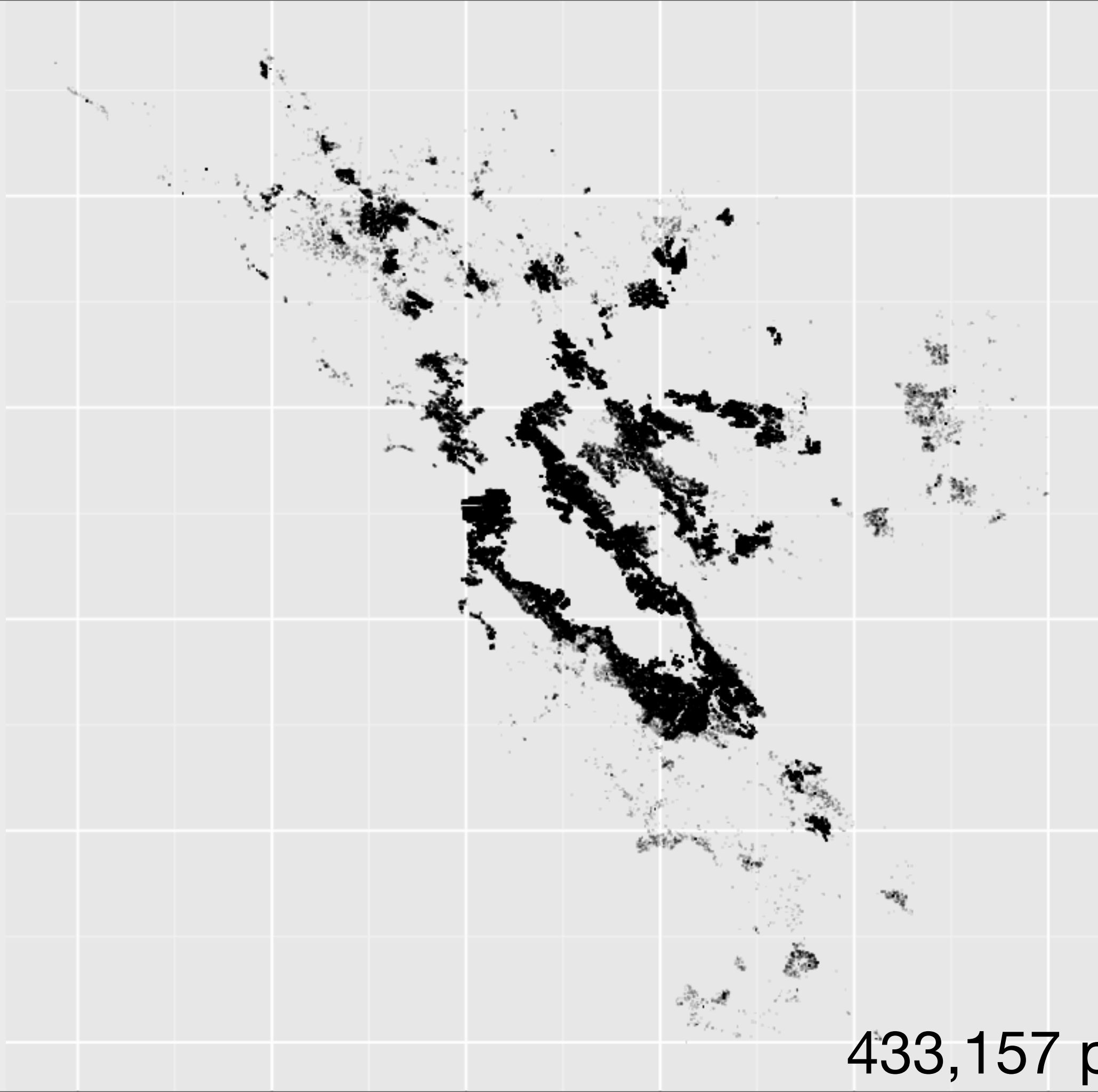
## with ggplot2 and plyr

Hadley Wickham

# About the data

521,495 house sales scraped from SF Chronicle, March 2003—November 2008.

Addresses geocoded using USC WebGIS. (83% to interpolated city block or better)

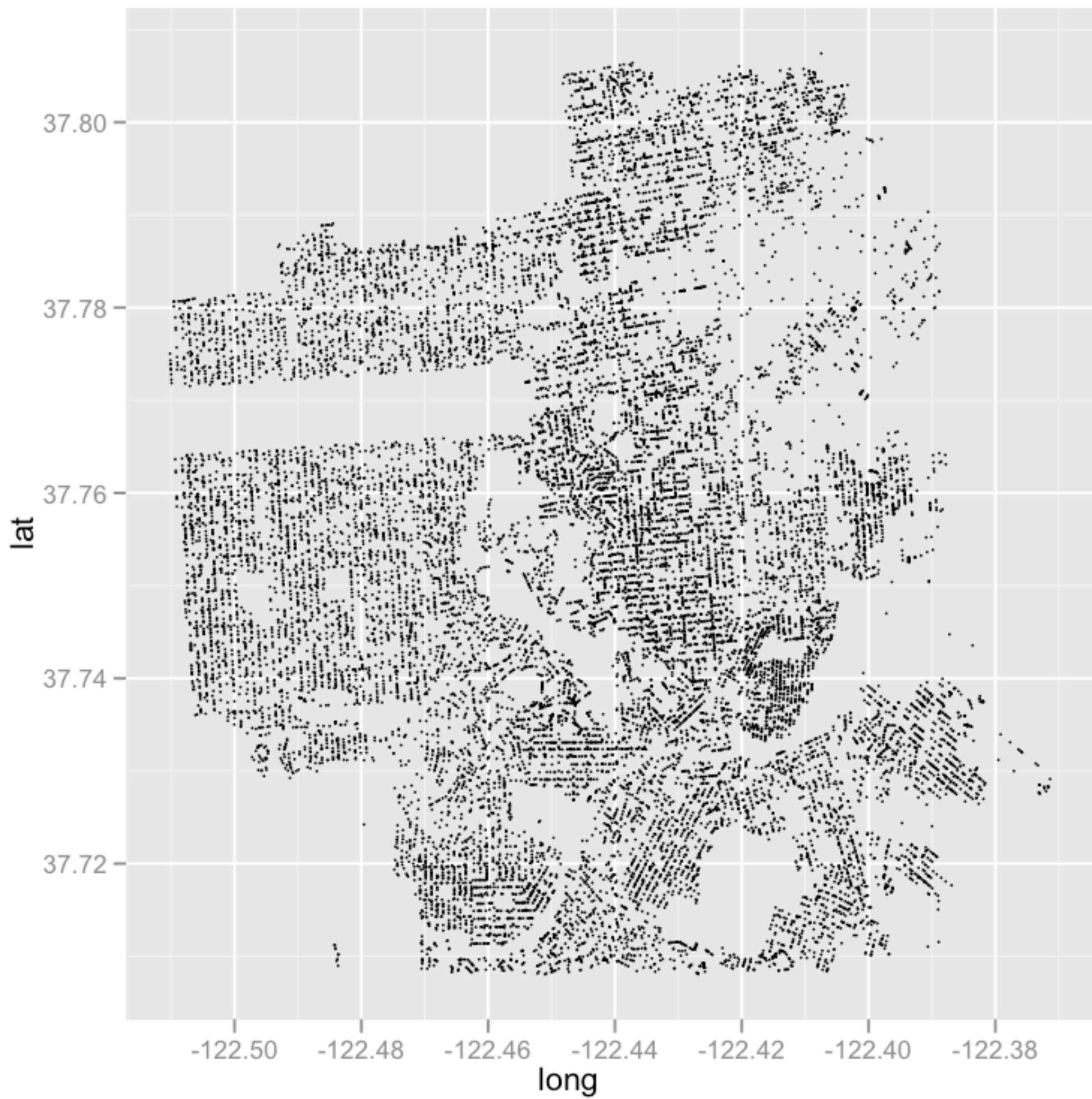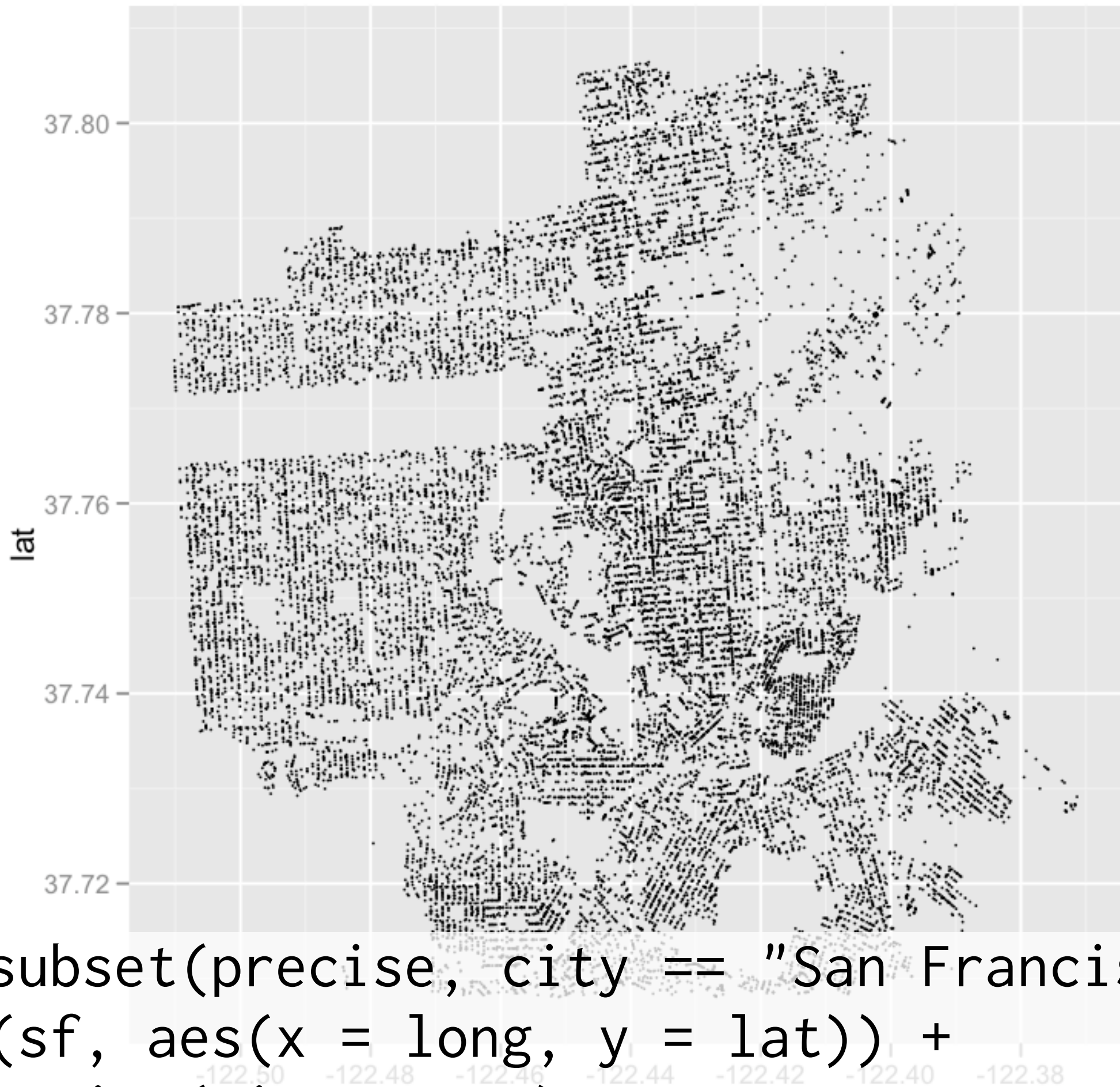**Variables**: longitude, latitude, date, price, bedrooms, area

433,157 points!

# Basic graphics:

Scatterplots, bar charts and histograms
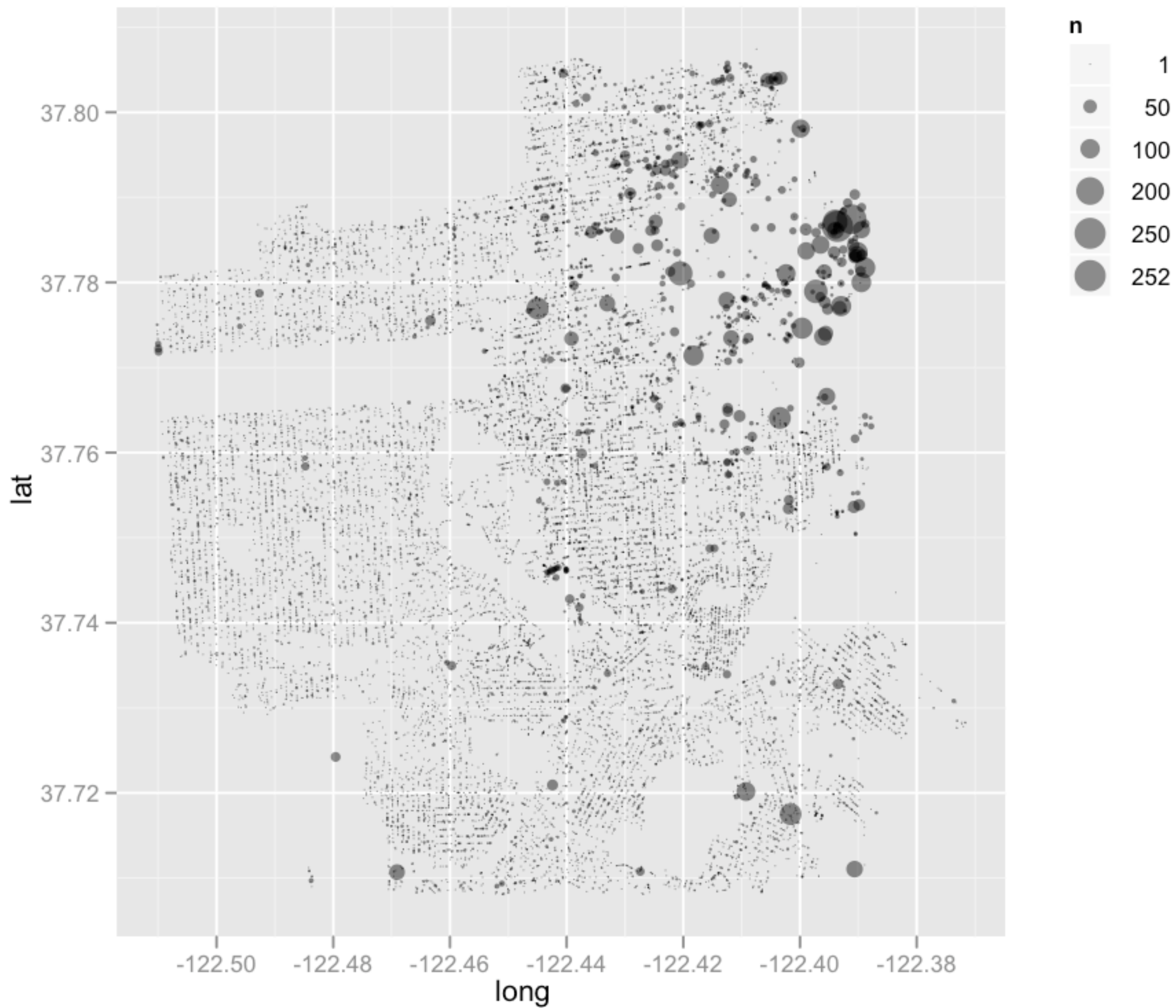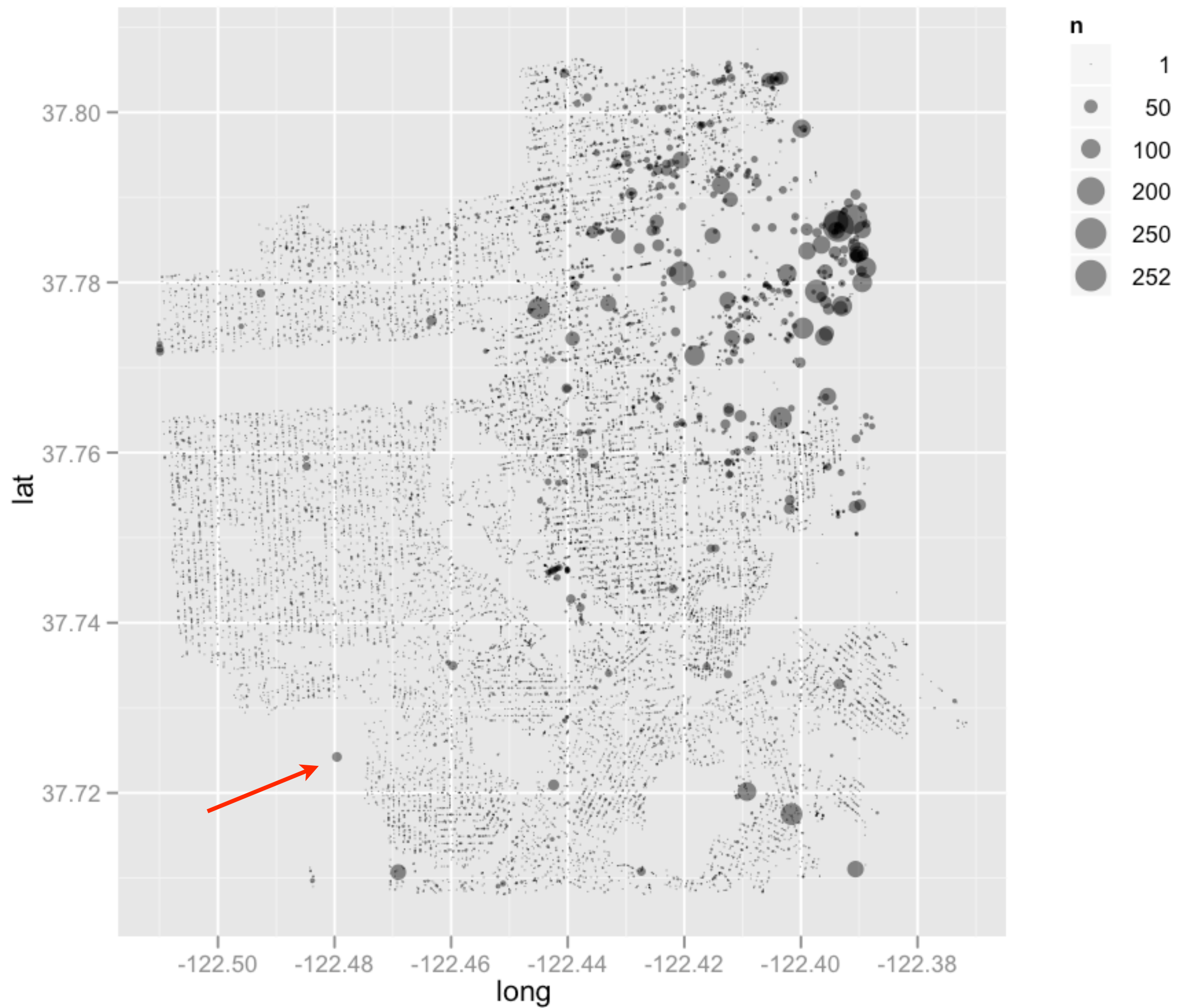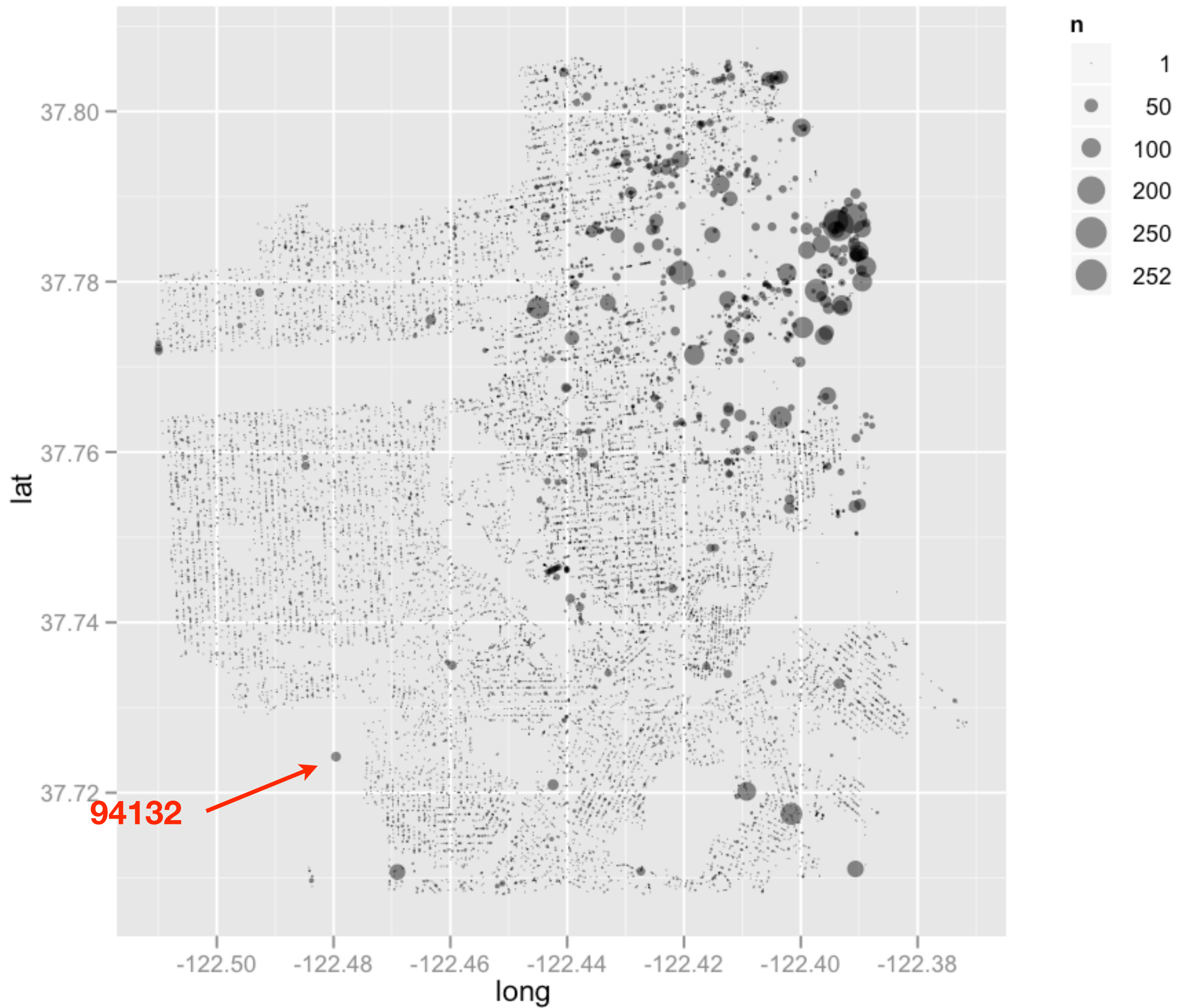
# Where are the houses?

```
sf <- subset(precise, city == "San Francisco")
ggplot(sf, aes(x = long, y = lat)) +
  geom_point(size = 0.5)
```

14,781 points!

```r
# In principle:
ggplot(sf, aes(long, lat)) +
  geom_point(stat = "sum")

# In practice:  (takes 42 seconds)
sfsum <- ddply(sf, c("lat", "long"), summarise,
  n = length(lat),
  avg_year = mean(year, na.rm = TRUE),
  .progress = "text"
)

ggplot(sfsum, aes(long, lat, size = n)) +
  geom_point(alpha = 1/2) +
  scale_area(to = c(0.3, 6), breaks = c(1, 50, 100, 200, 252))
```
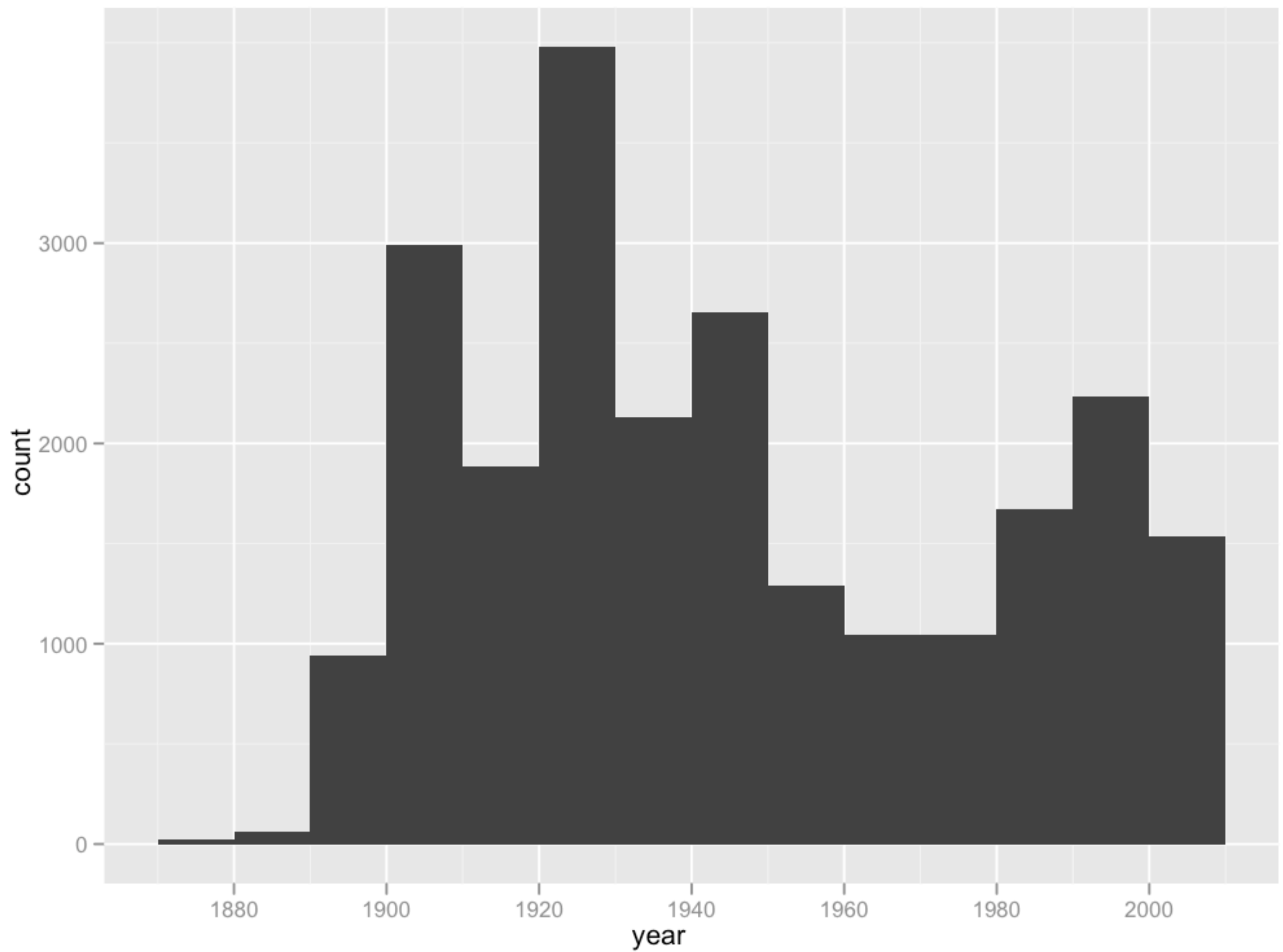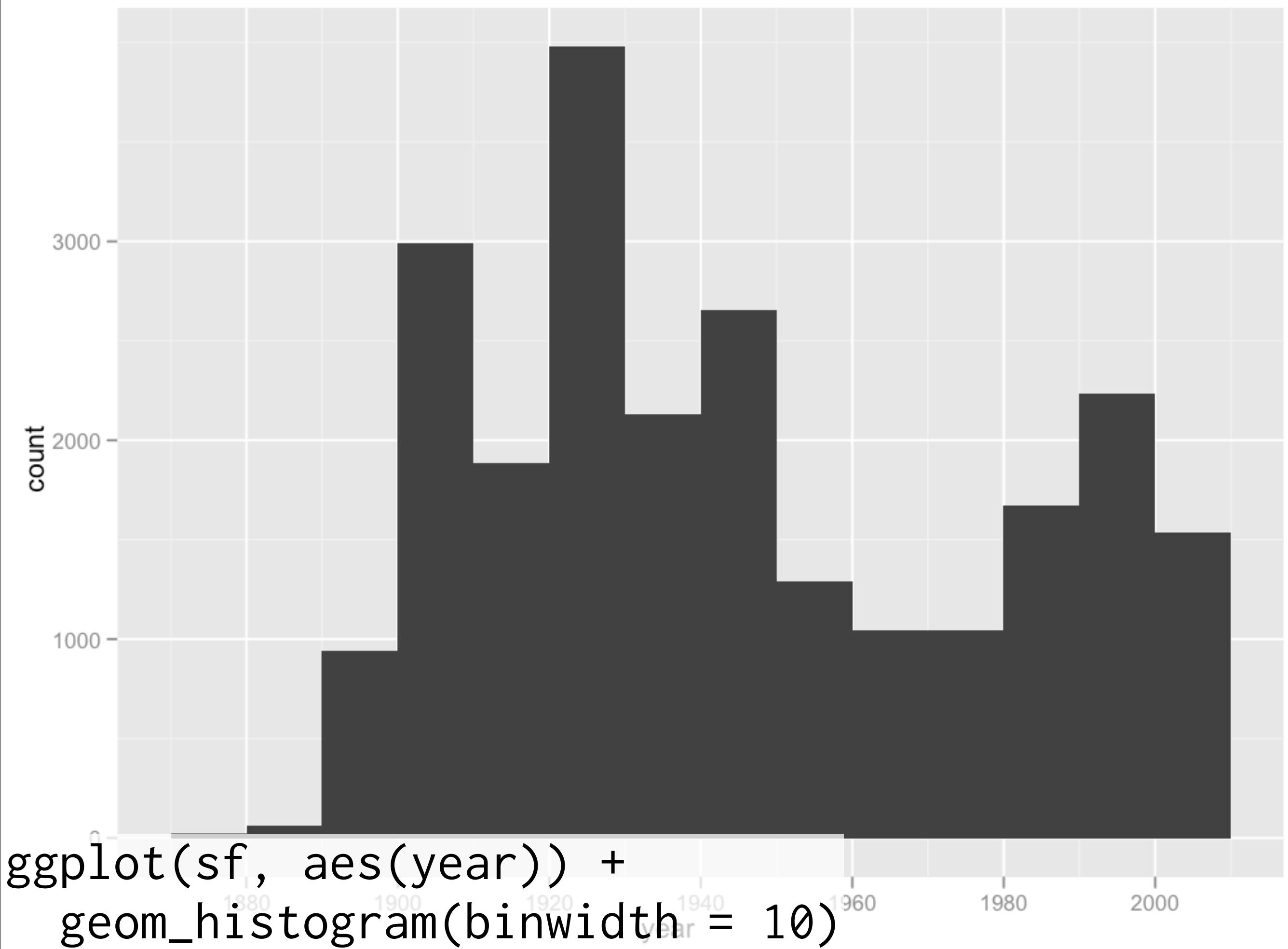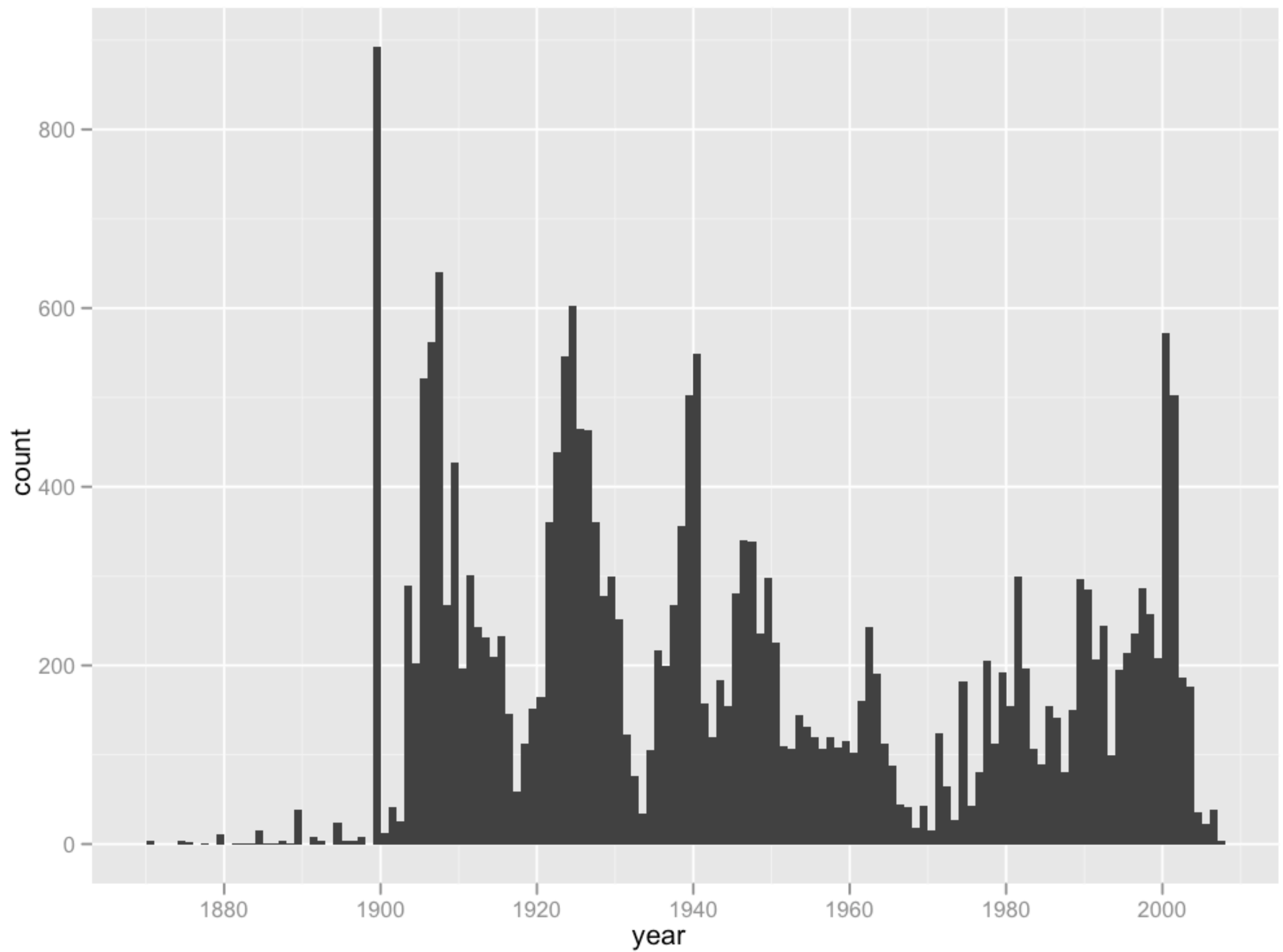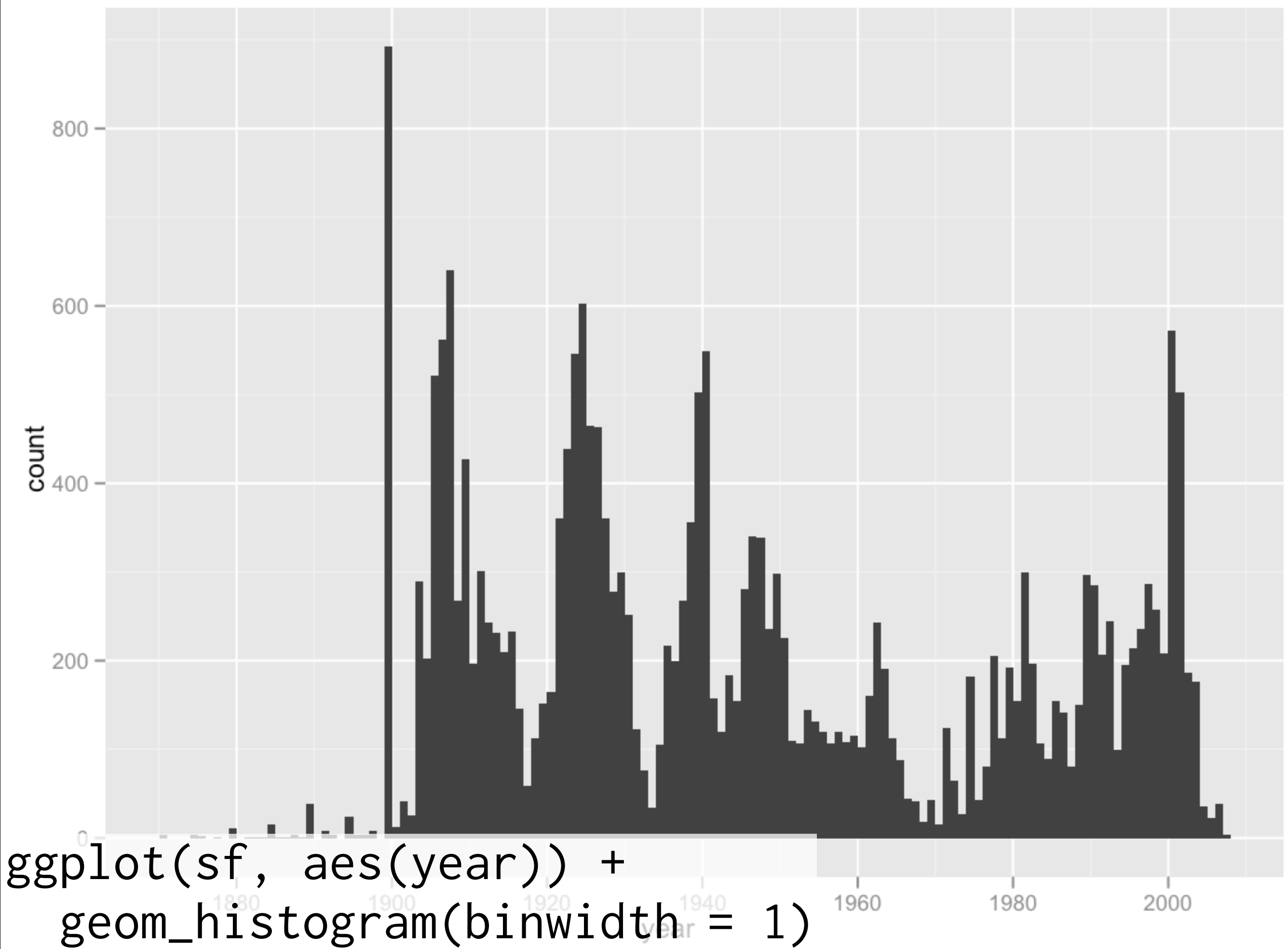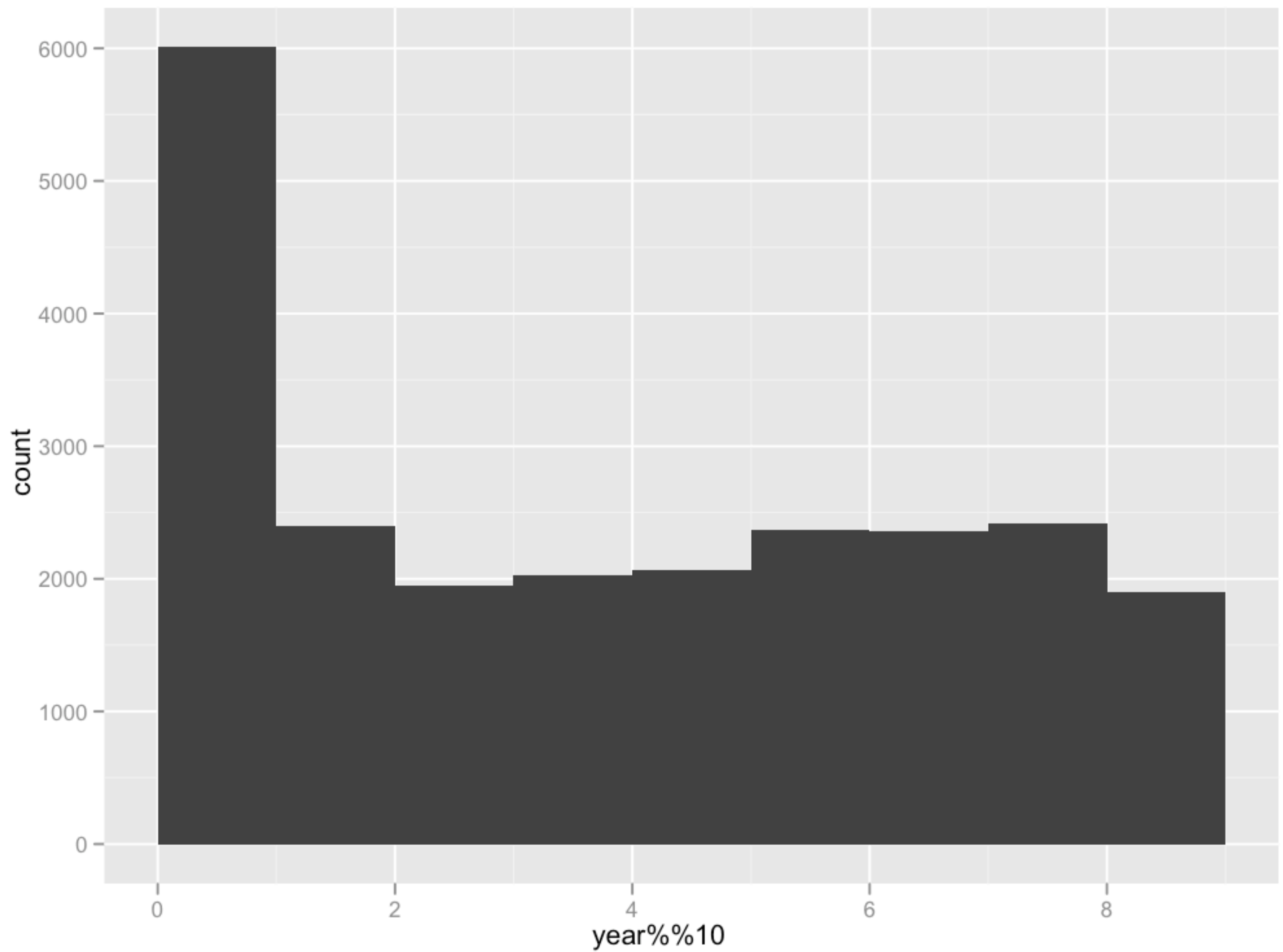
# When were they build?

```
ggplot(sf, aes(year)) +
  geom_histogram(binwidth = 10)
```

```
ggplot(sf, aes(year)) +
  geom_histogram(binwidth = 1)
```
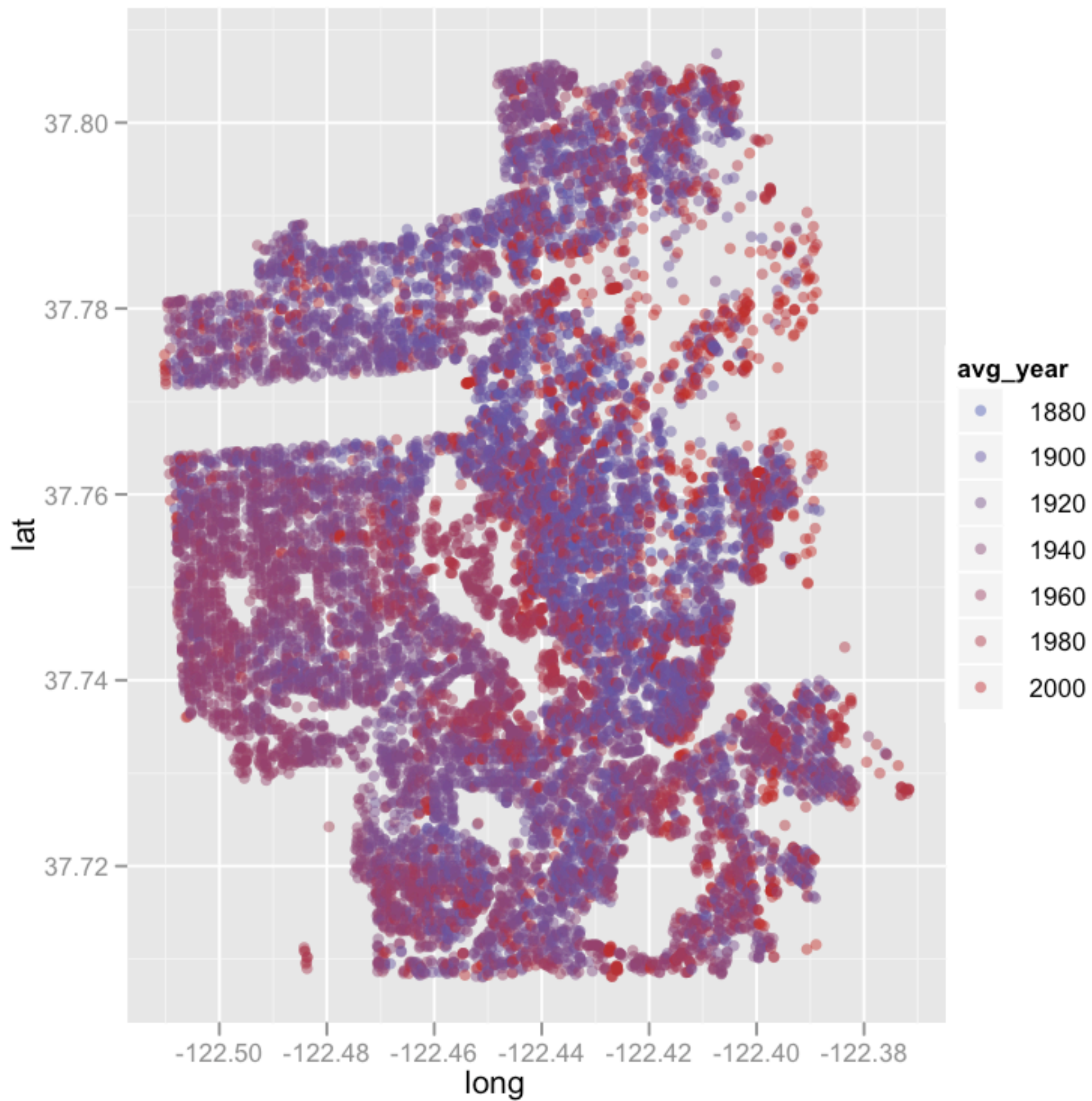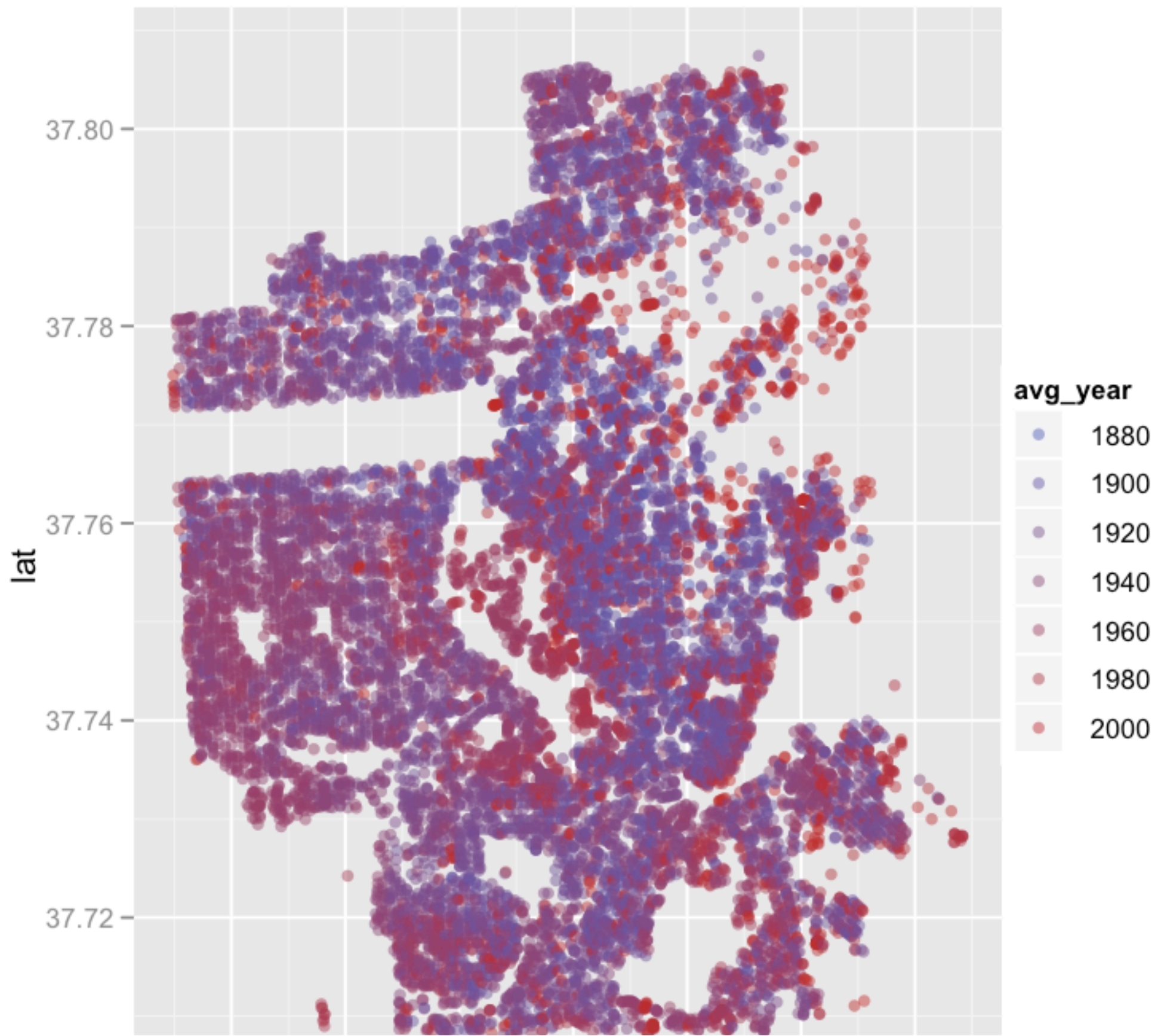
# How does location vary with age?

# Additional variables

Can supplement any existing plot with additional variables in two ways: adding **aesthetics**, or creating **facets**.

Saw example of using size.  Other aesthetics are: colour, fill, shape, linetype, alpha.

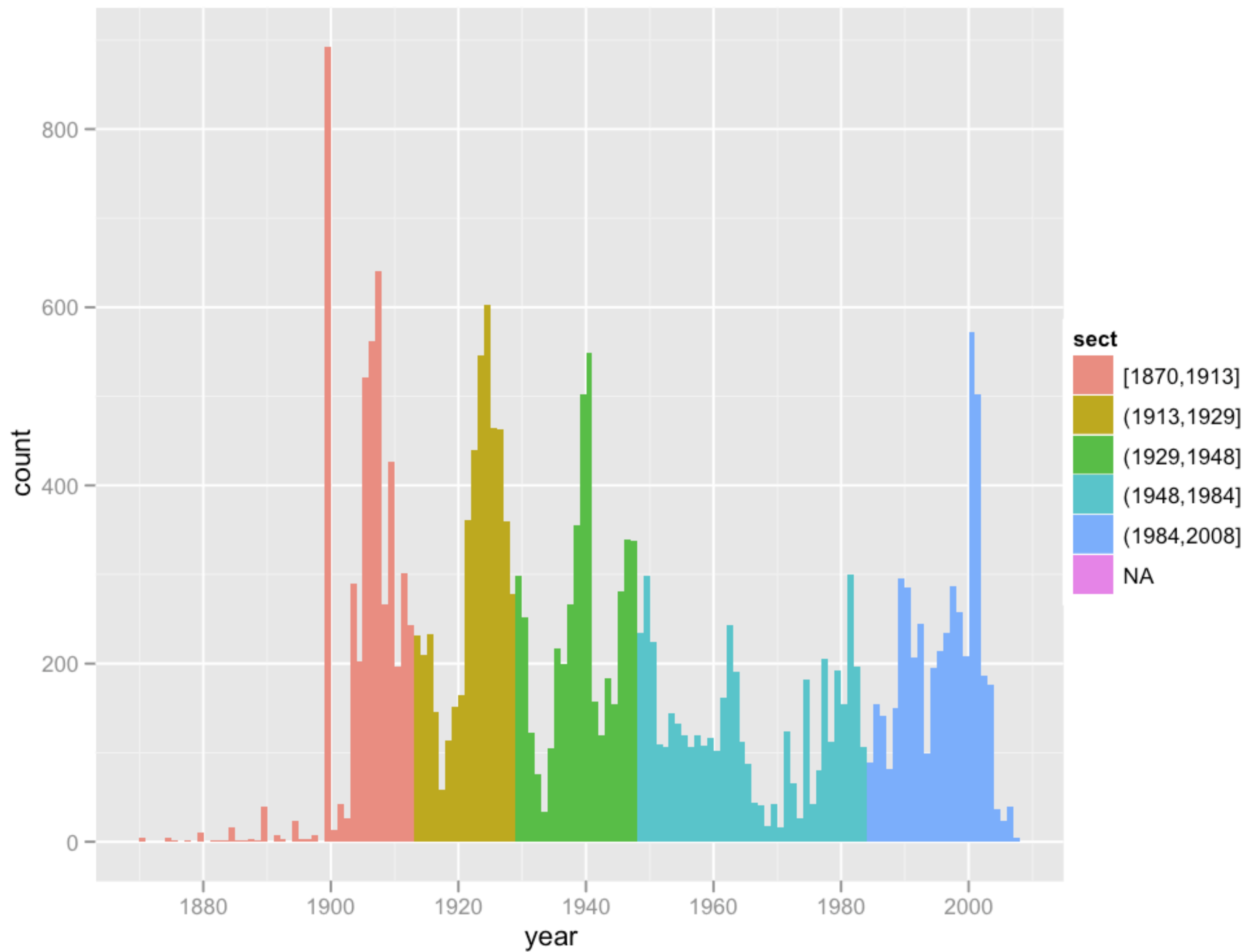Facetting creates small multiples of subsets of the data.

```
ggplot(sfsum, aes(long, lat, colour = avg_year)) +
    geom_point(alpha = 1/2, size = 2)
```
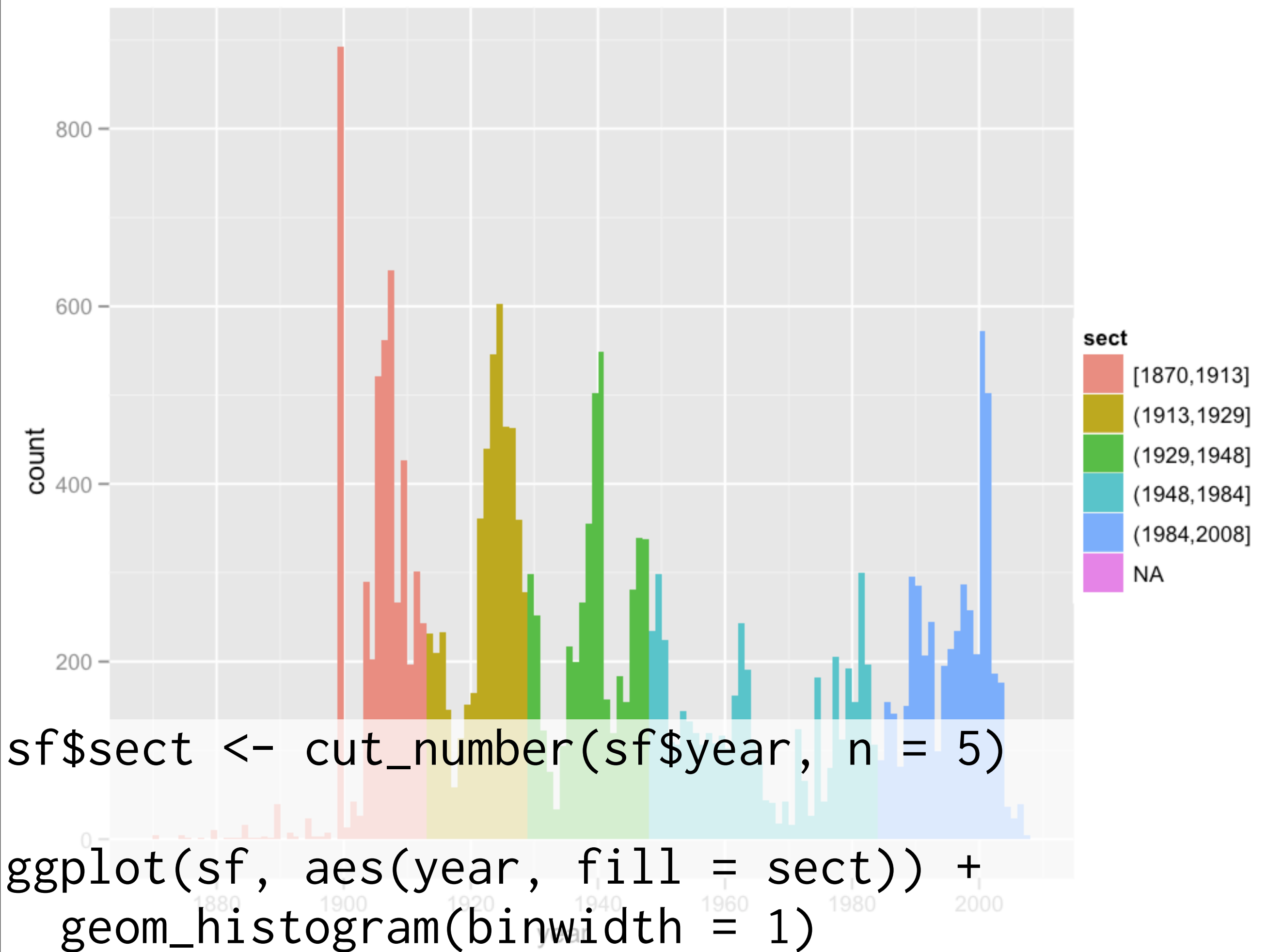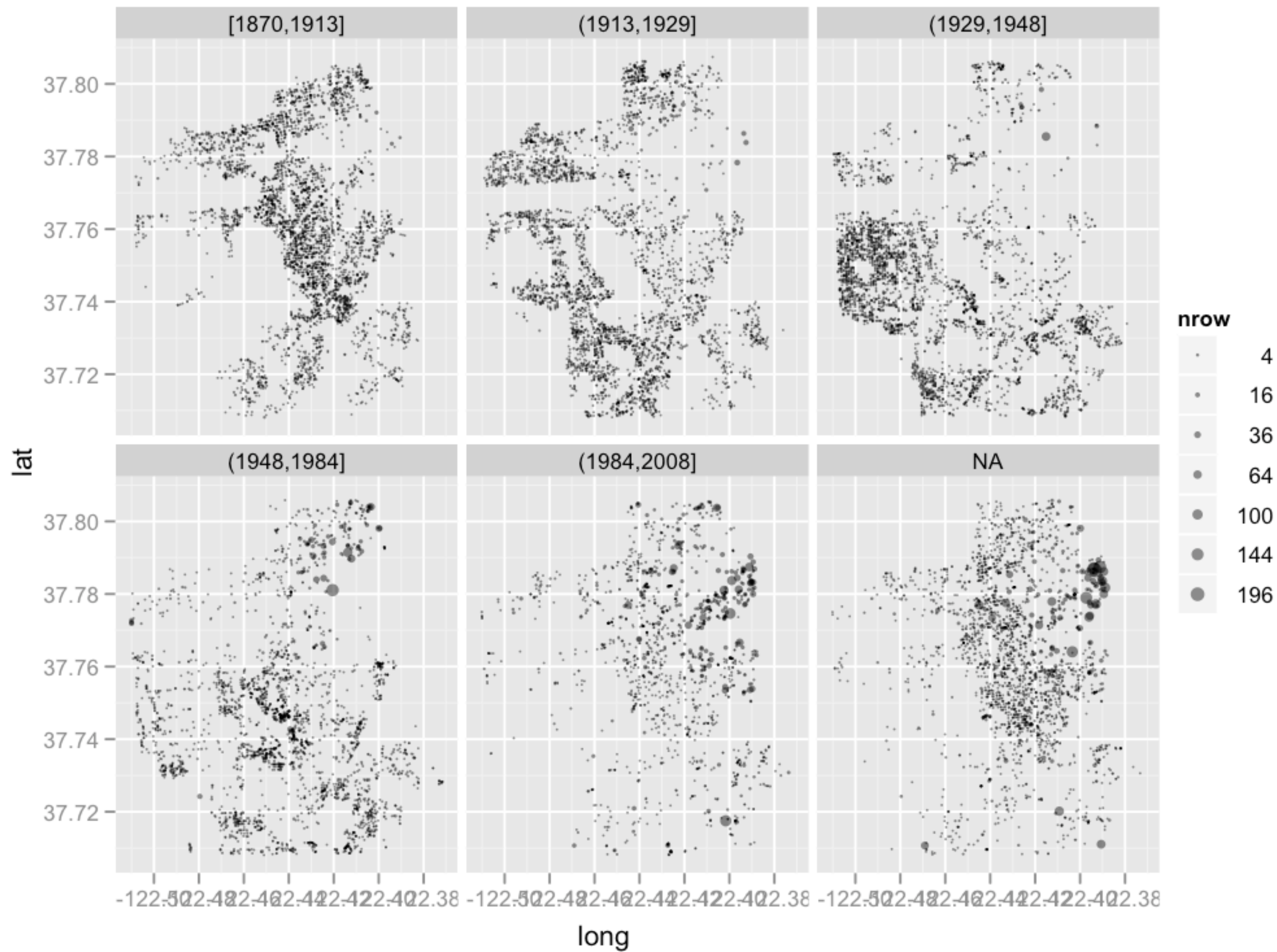
# Interactive demo

(a taste of future directions)

# Let's try facetting

First, need a categorical variable to facet by

```
sf$sect <- cut_number(sf$year, n = 5)

ggplot(sf, aes(year, fill = sect)) +
    geom_histogram(binwidth = 1)
```

```r
# Recount number at each location
sfsumsect <- ddply(sf, c("lat", "long" , "sect"),
  "nrow", .progress = "text")

# Display spatial distribution of each cut
ggplot(sfsumsect, aes(long, lat, size = nrow)) +
  geom_point(alpha = 1/2) +
  scale_area(to = c(0.5, 3)) +
  facet_wrap(~ sect)
```
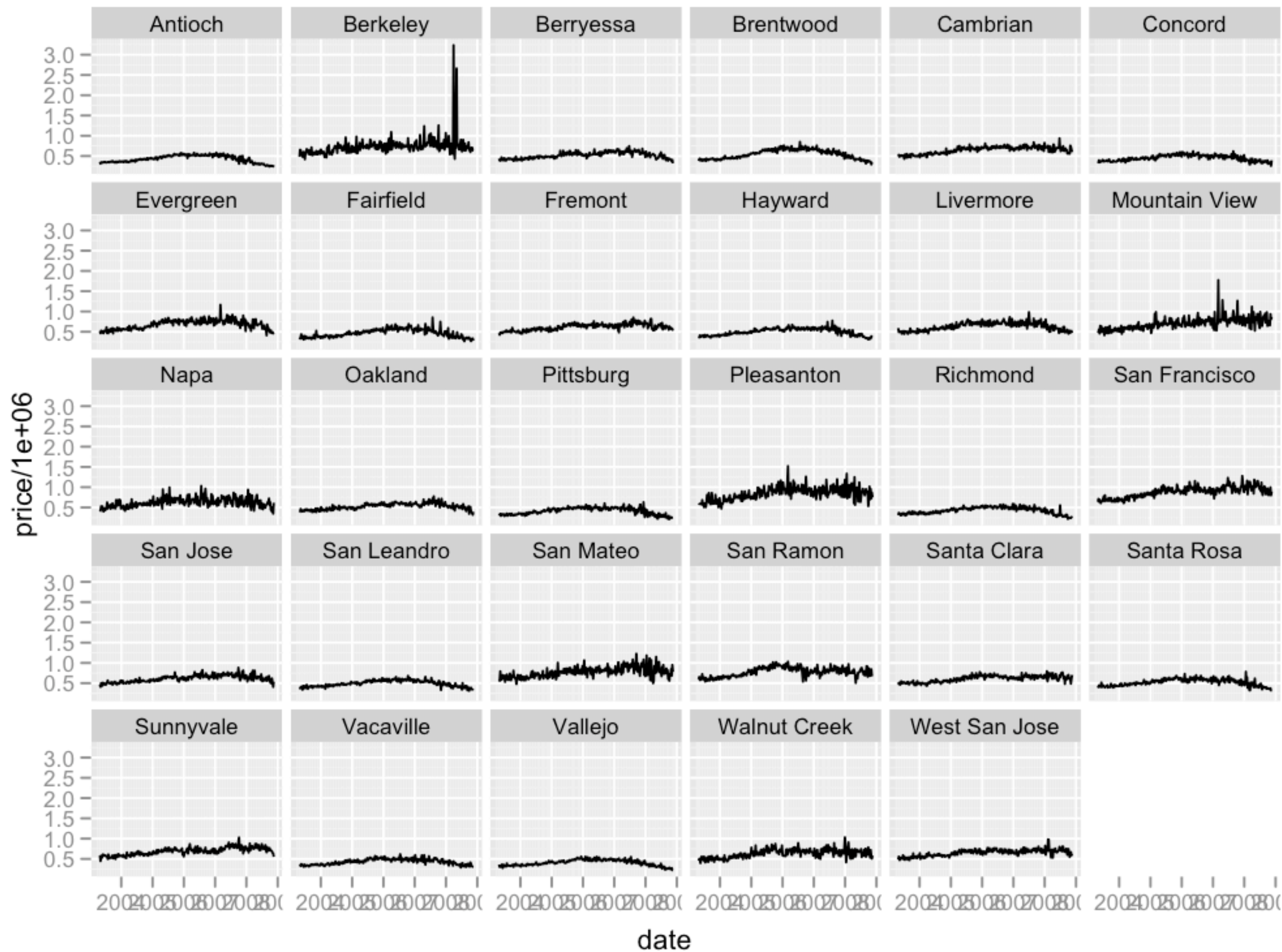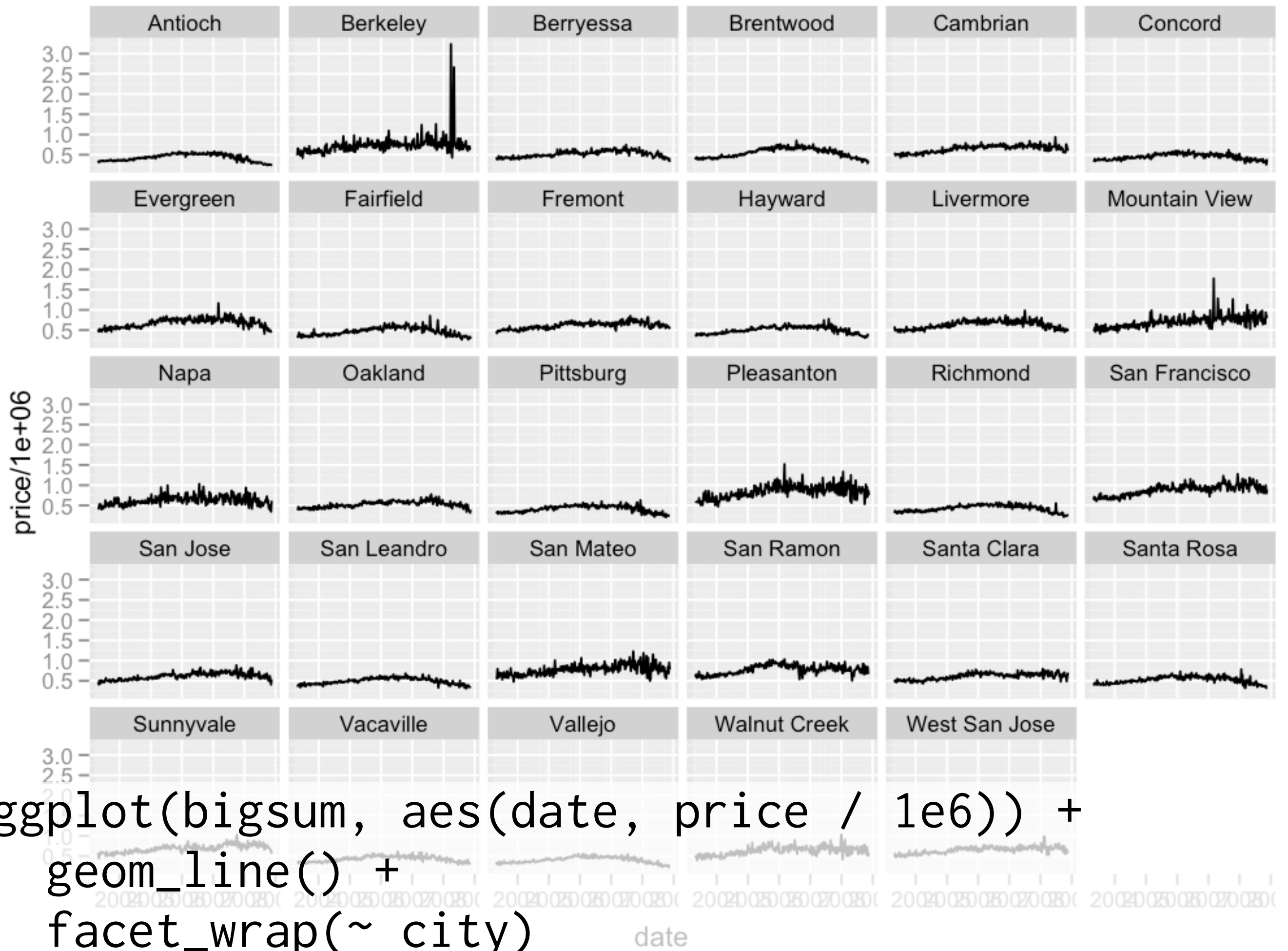
# Time series

## using plyr to manipulate data

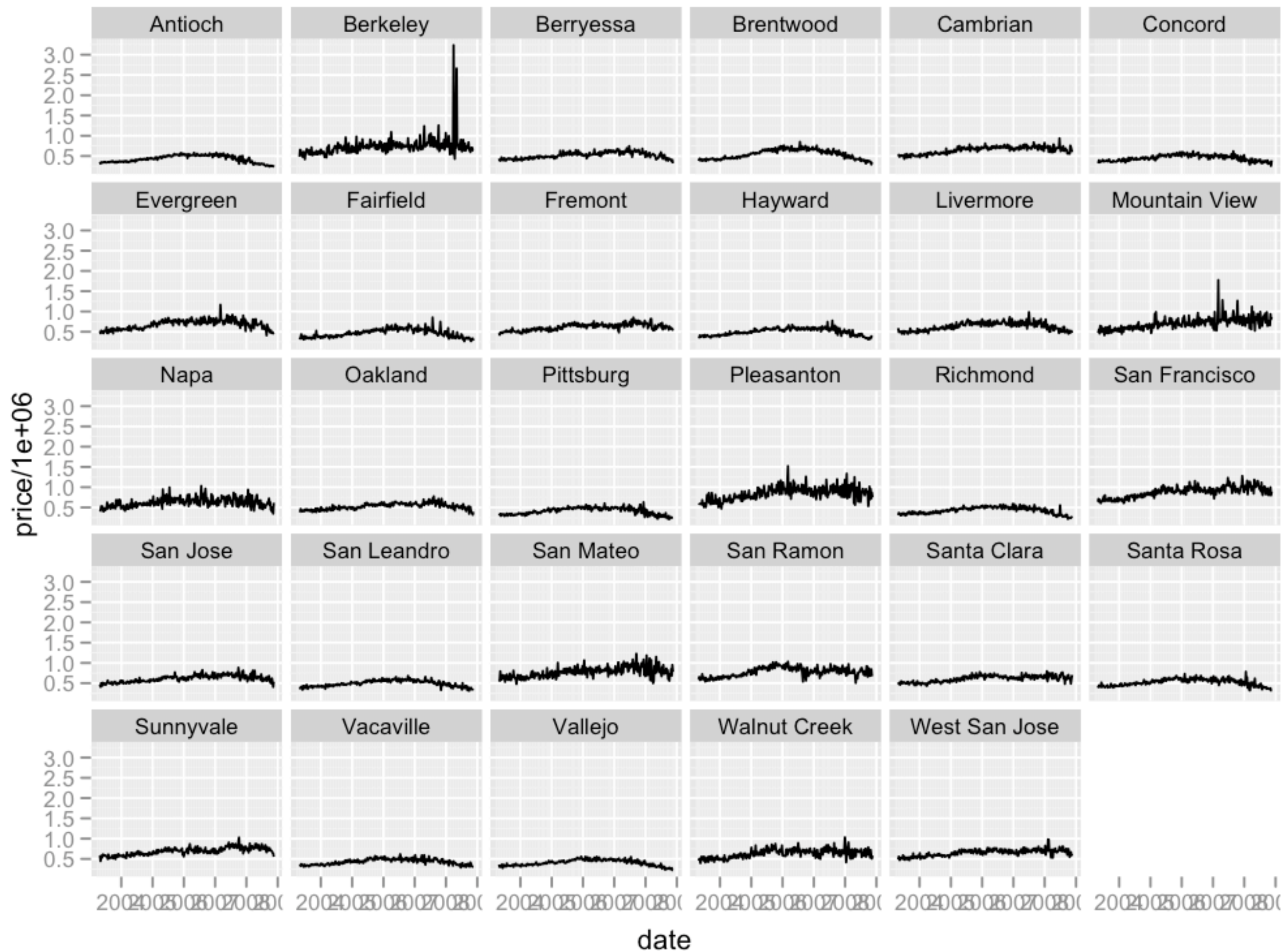Zoom out to the Bay area.
What's going on in different cities?

Focus on average sale price,
**summarising** patterns, then
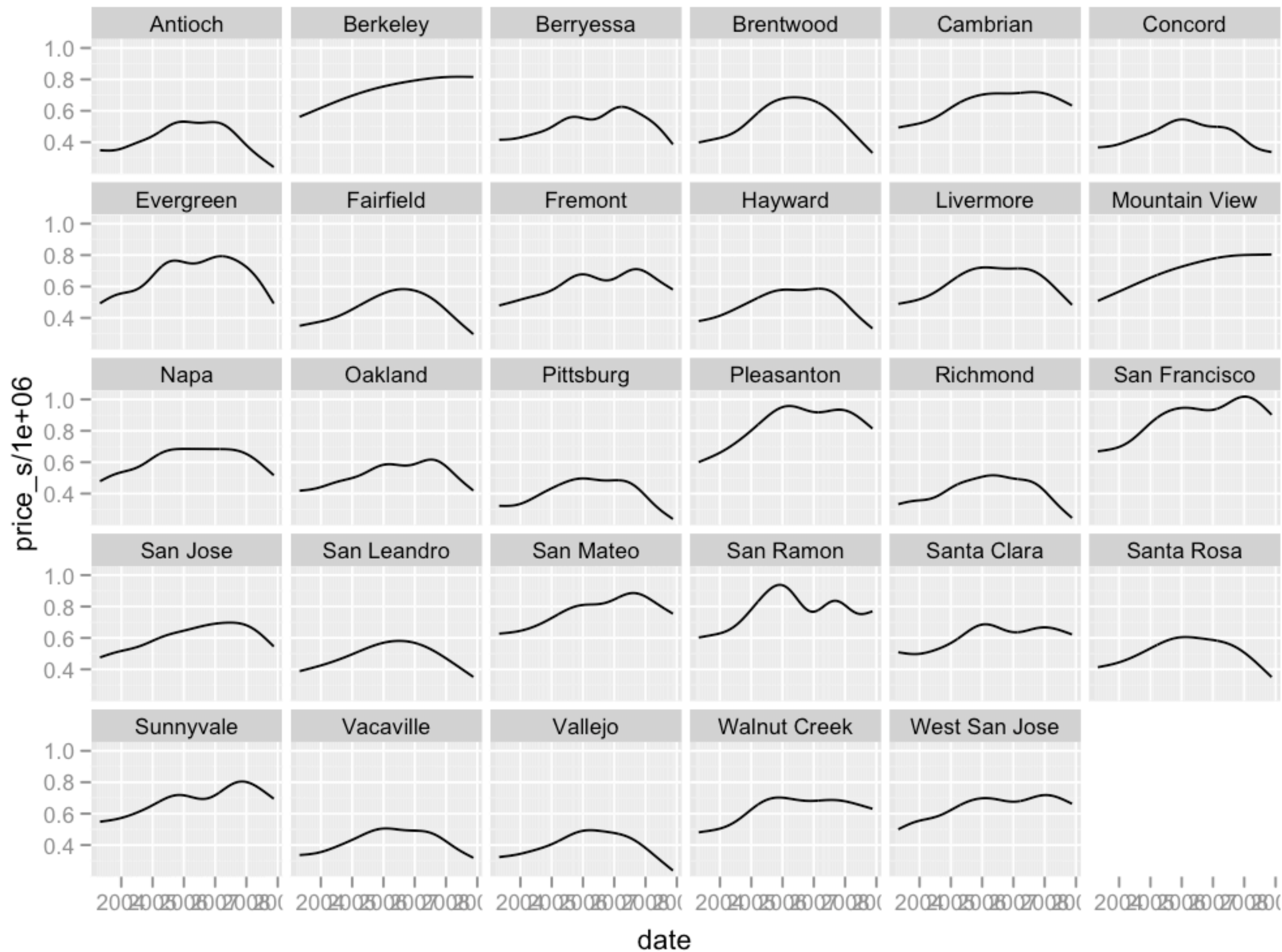**investigating** possible explanatory
variables.

```
ggplot(bigsum, aes(date, price / 1e6)) +
  geom_line() +
  facet_wrap(~ city)
```

```r
library(mgcv)
smooth <- function(y, x) {
  as.numeric(predict(gam(y ~ s(x))))
}

bigsum <- ddply(bigsum, "city", transform,
  price_s = smooth(price, as.numeric(date)))
```

```
index <- function(y, x) {
  y / y[order(x)[1]]
}

bigsum <- ddply(bigsum, "city", transform,
  price_si = index(price_s, date))
```
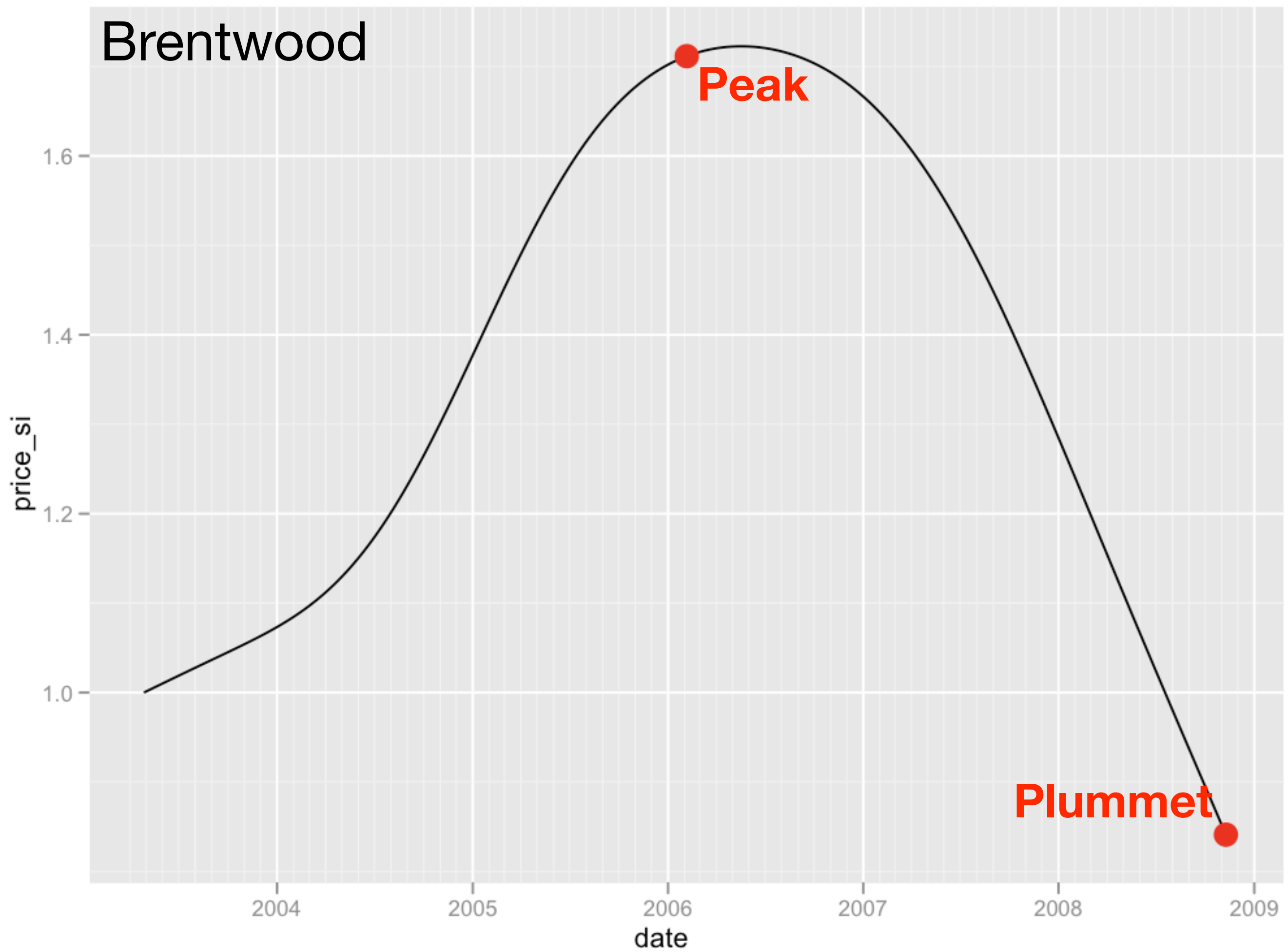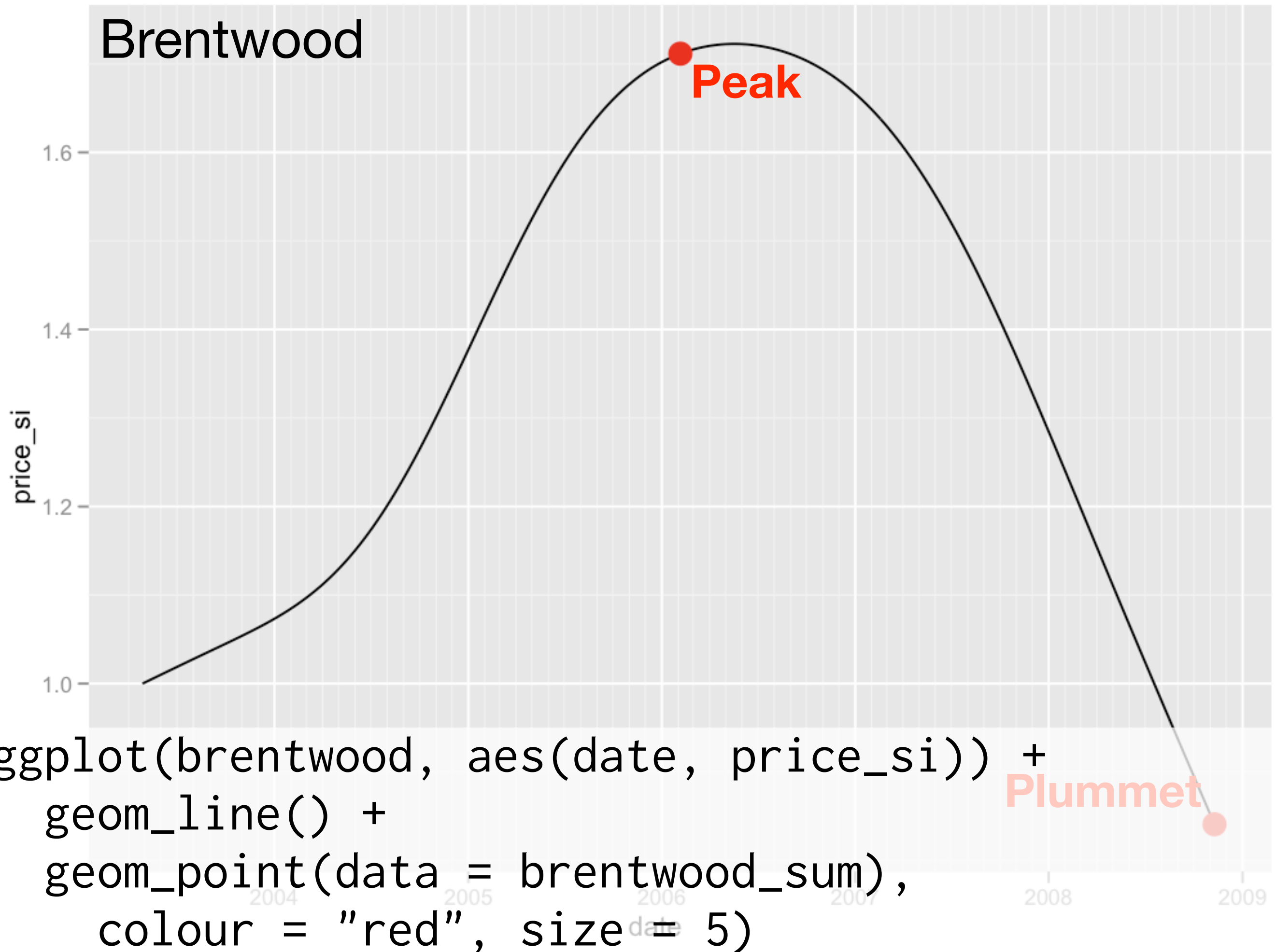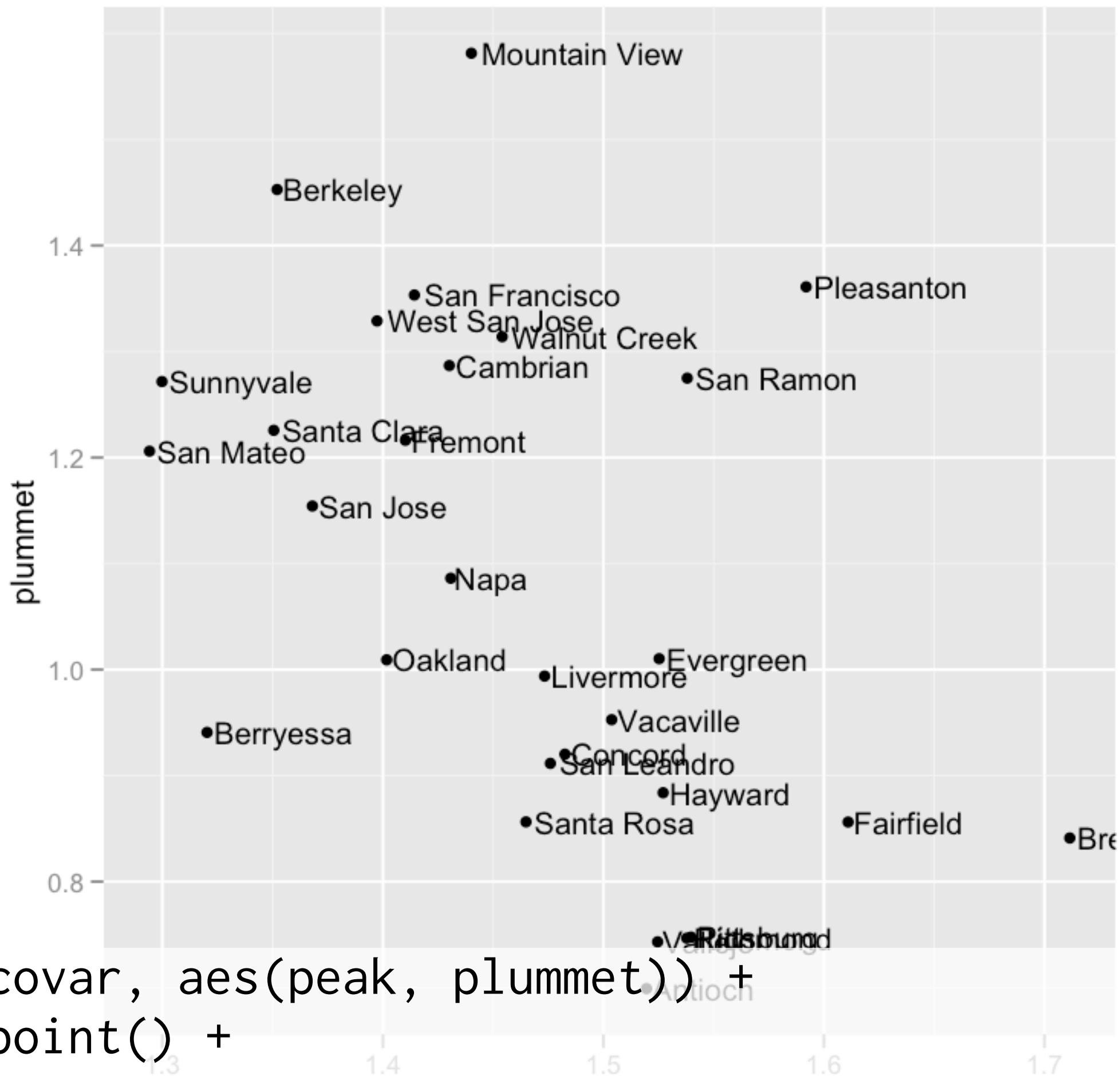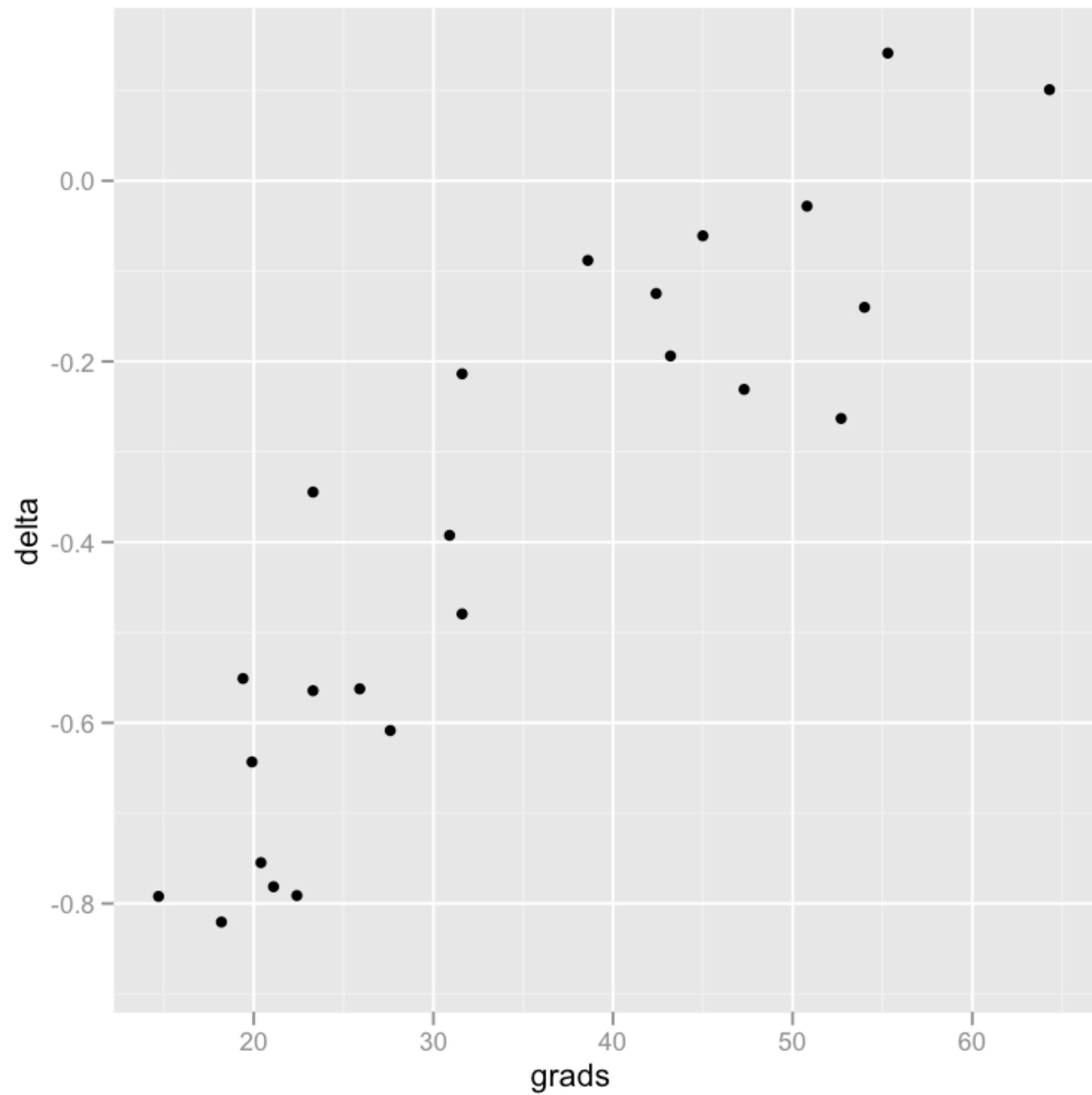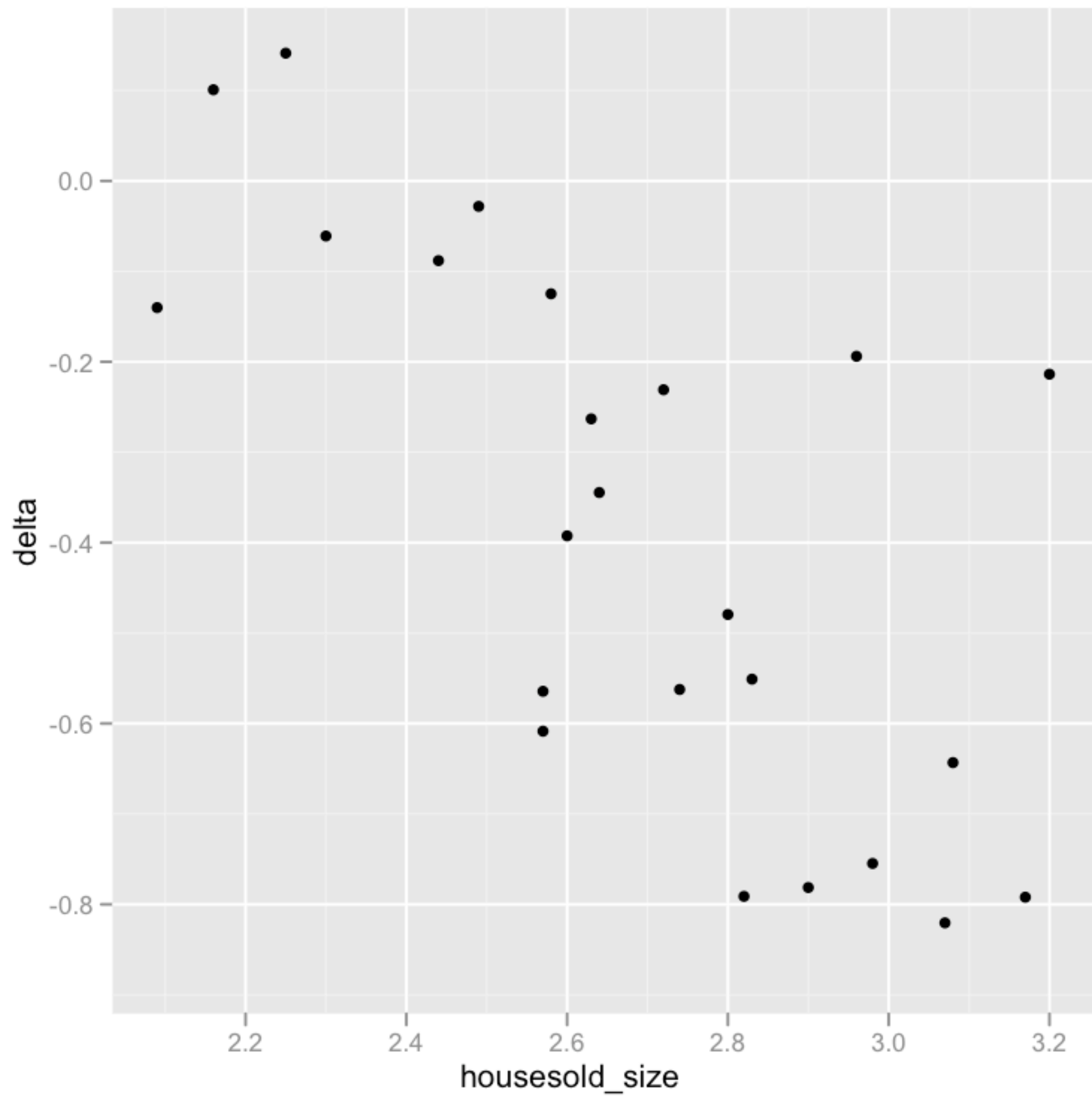
```
ggplot(covar, aes(peak, plummet)) +
    geom_point() +
    geom_text(aes(label = city), size = 4, hjust = -0.05)
```
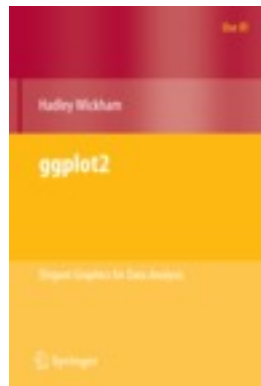
```r
covar$delta <- with(covar, plummet - peak)

census <- read.csv("census-city.csv")
covar <- join(covar, census, by = "city")
```
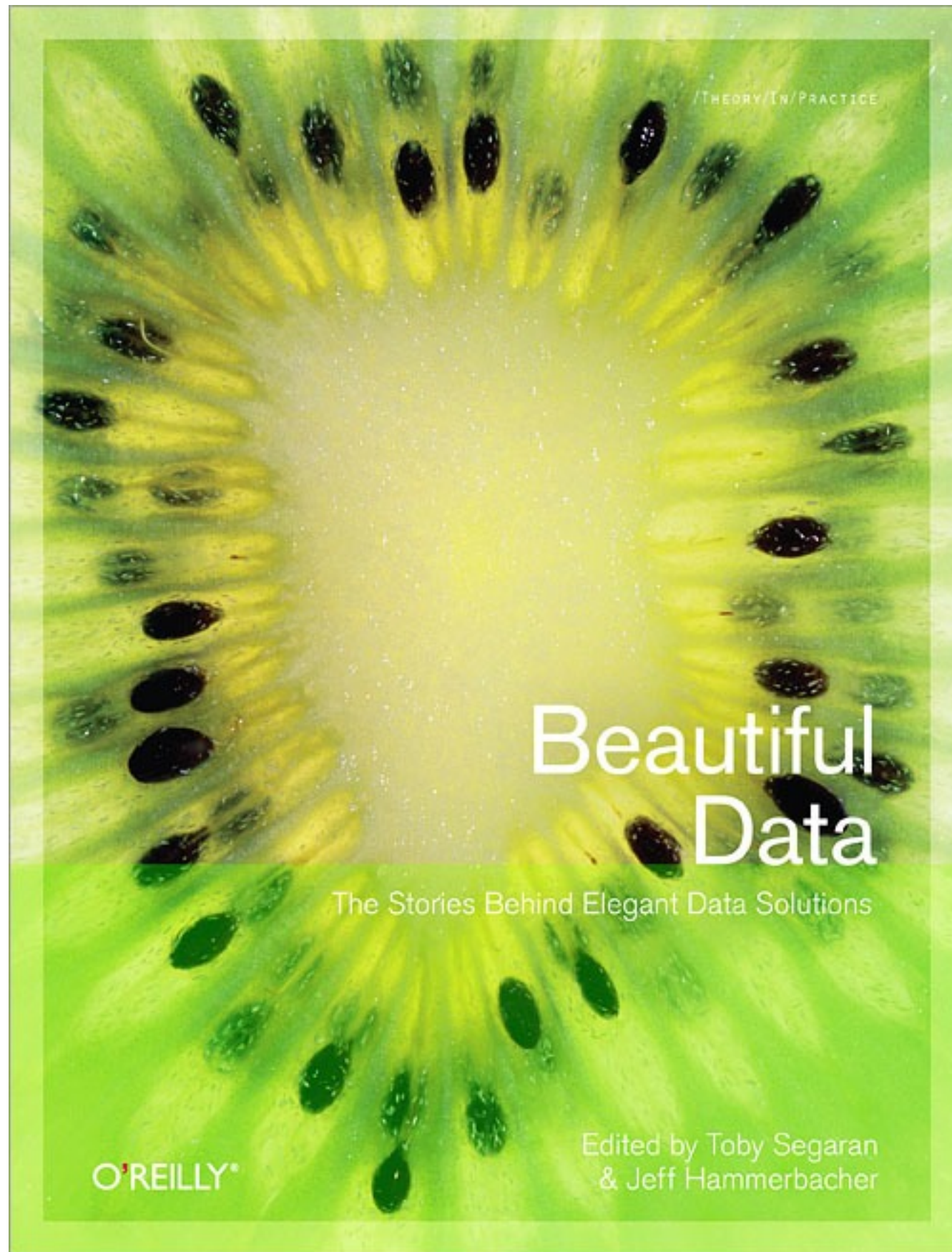
# Learn more

**ggplot2**
http://had.co.nz/ggplot2
http://groups.google.com/group/ggplot2



**plyr**
http://had.co.nz/plyr
http://groups.google.com/group/manipulatr

Chapter 18

All code and data openly licensed: http://github.com/hadley/sfhousing

Beautiful data is **reproducible**!

# Thank you!