

STA 141A HW2

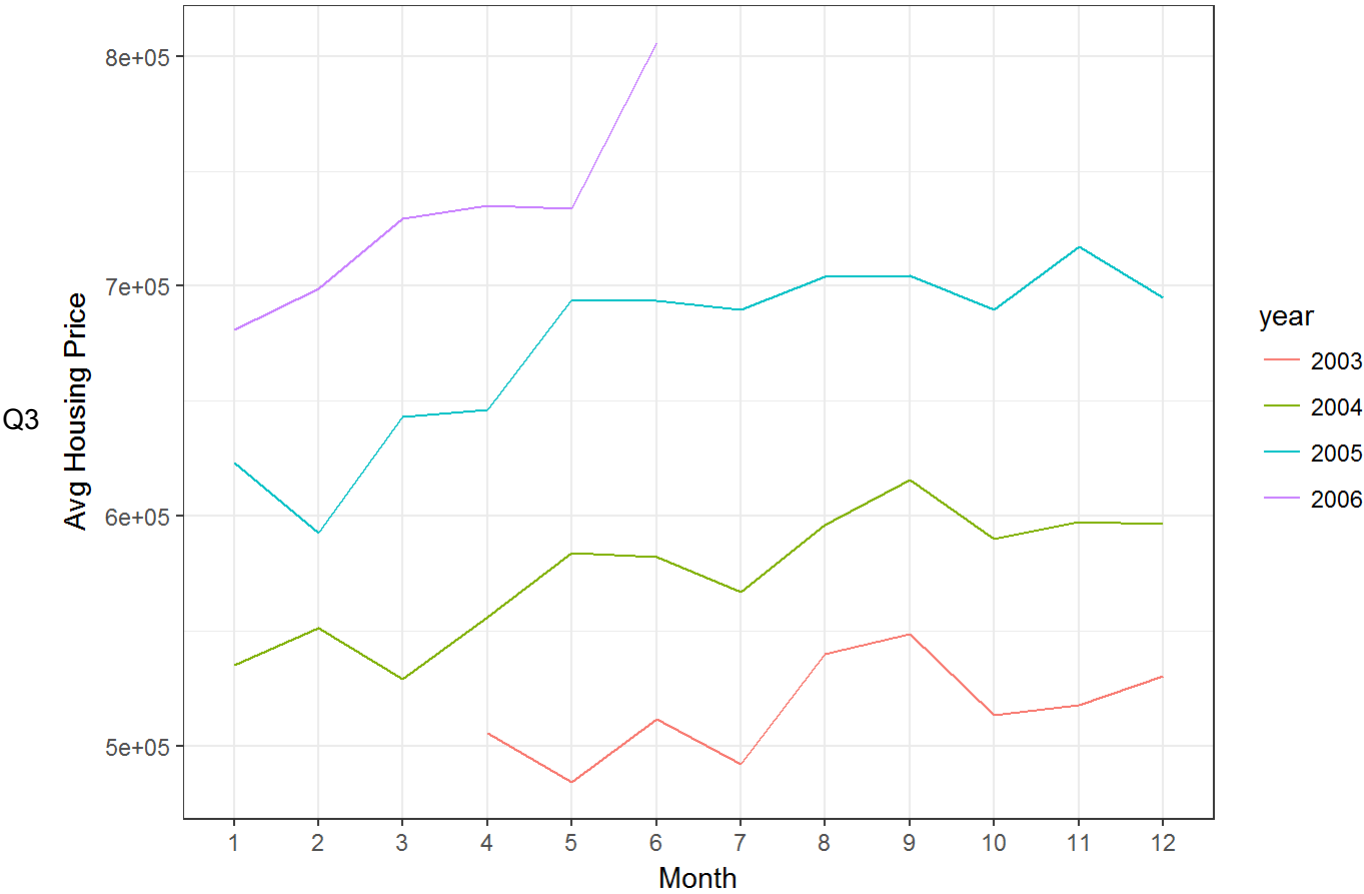
Zihan Mo 914998952

April 29, 2018

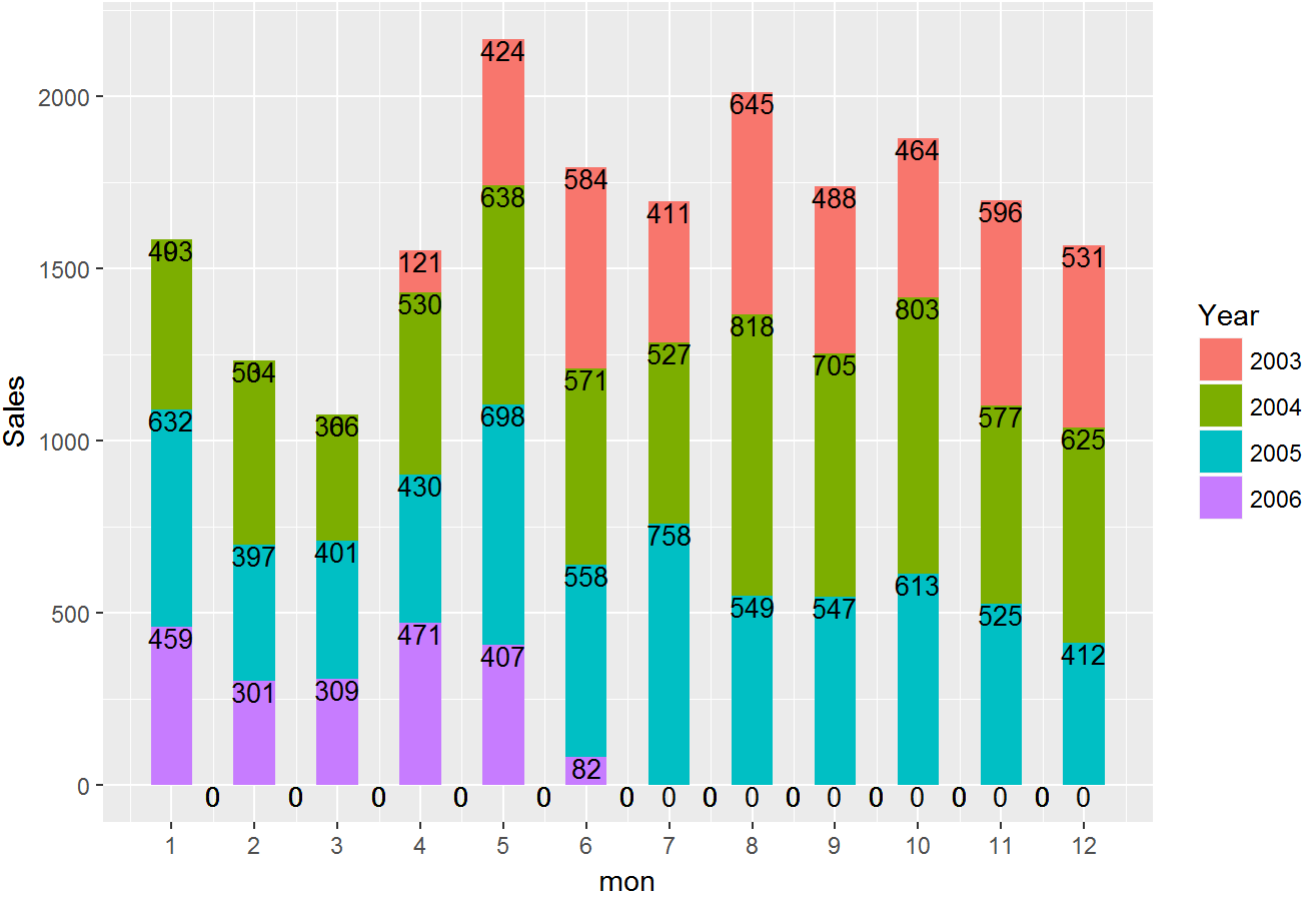
Q1 Justified by code in the appendix

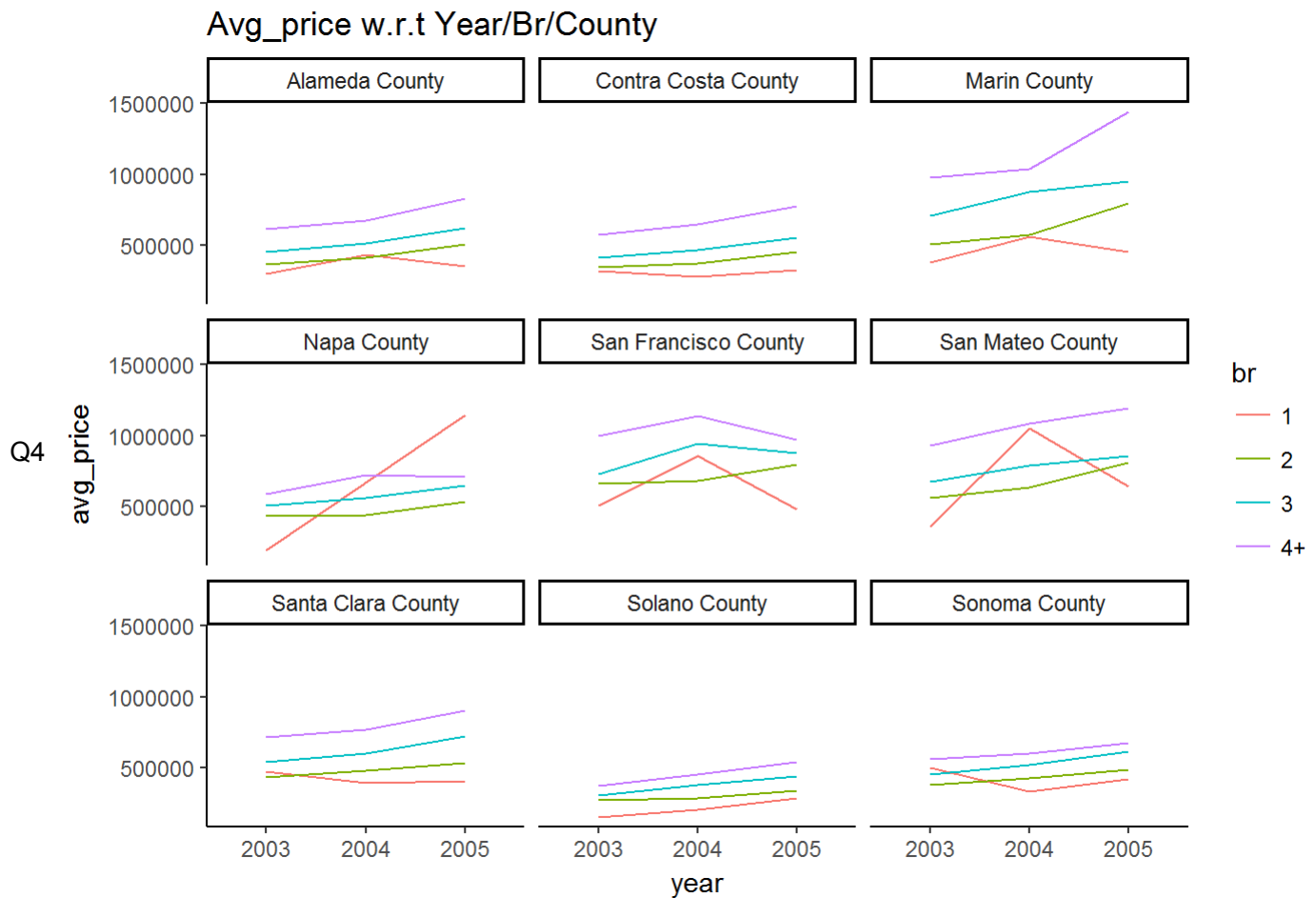
Q2 The timespan of the housing sales cover from 2003-04-27 to 2006-06-04; The timespan of the construction year of houses cover from 1885 to 2005.

Avg Price Over Time



Sales Over Time



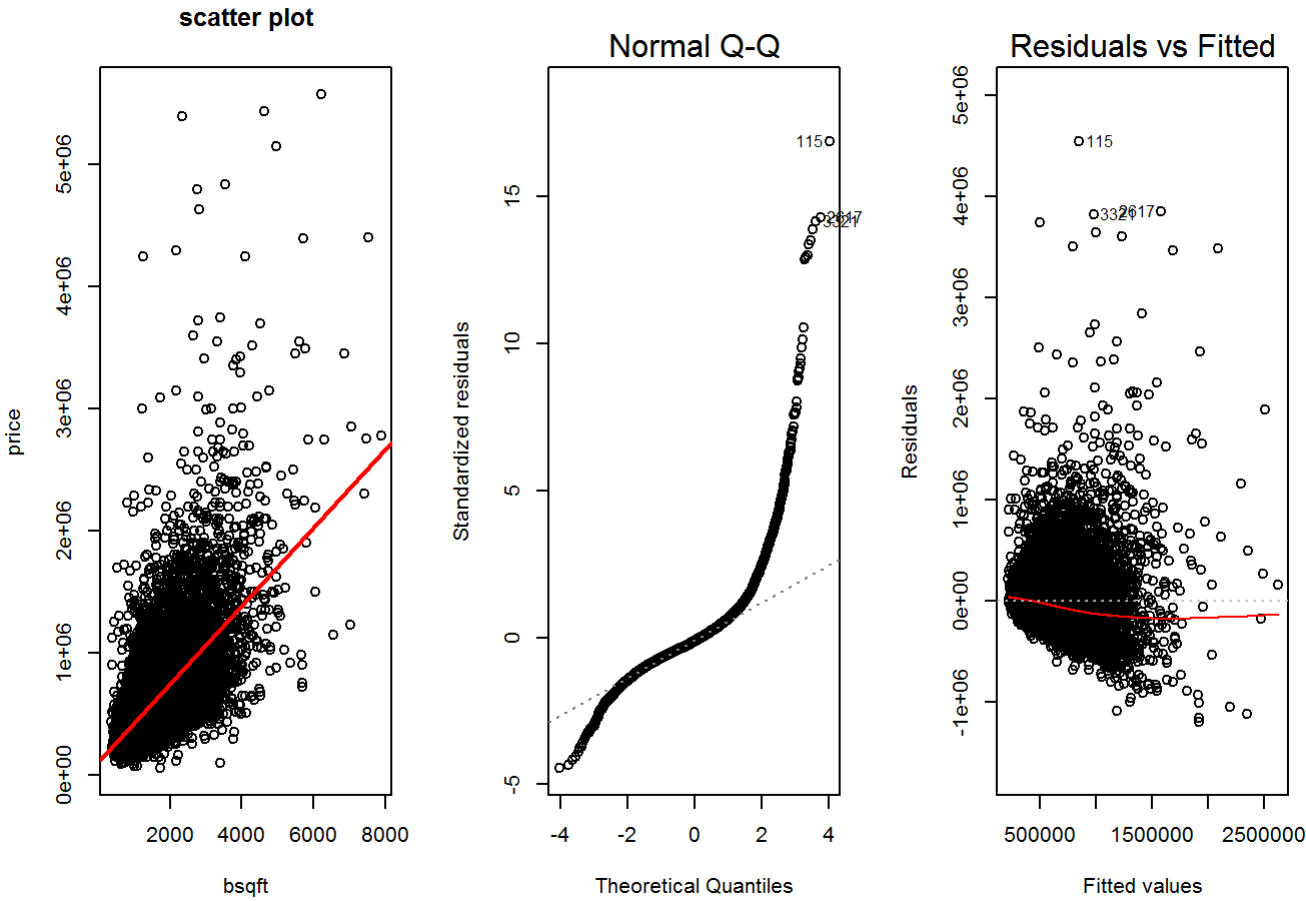


Q5

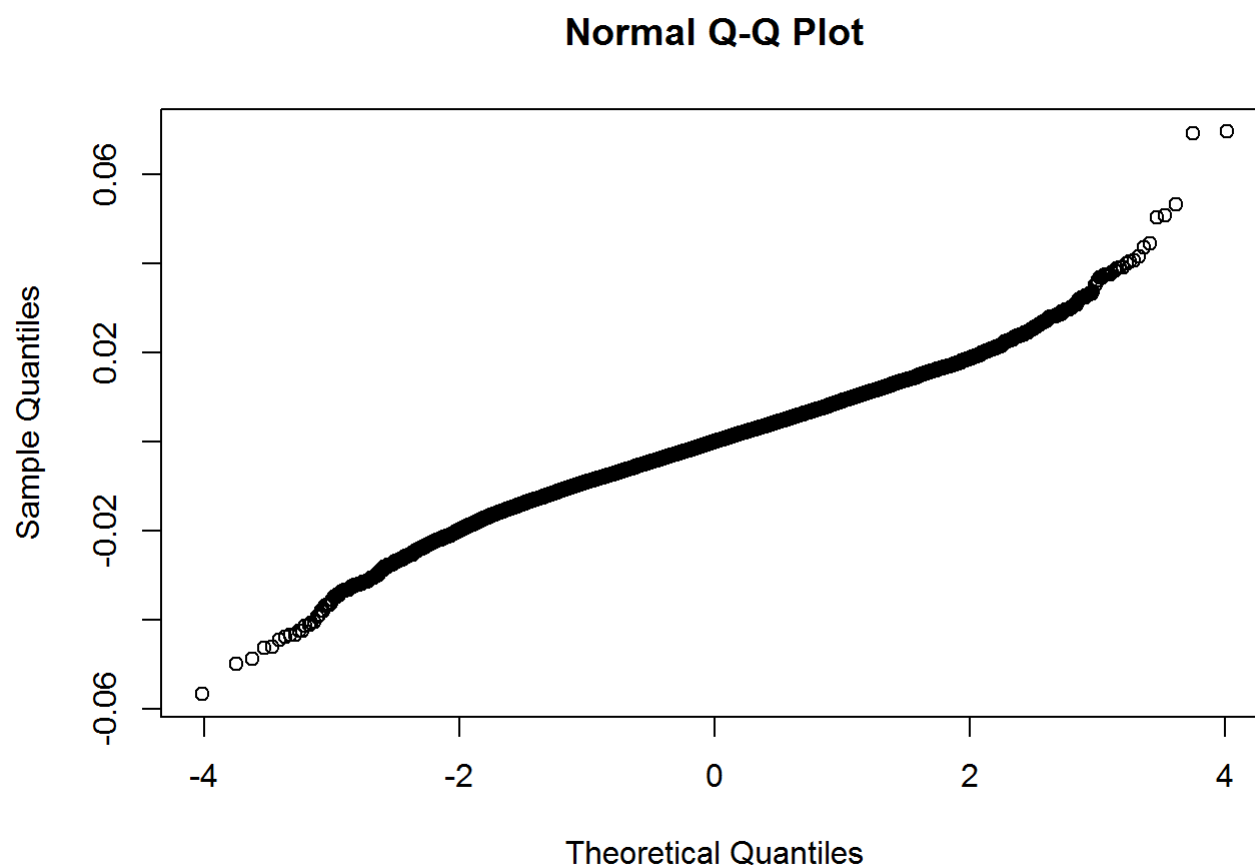
Only Vallejo city has sales in more than one county, which are Napa County and Solano County.

Q6

After taking out extreme outliers, like price equals zero and the building size larger than 10000, we have plots. Based on the plots below, the linear regression model seems provide a good fit to the data. But the QQ-plot and the residuals plot imply the normality assumption isn't met and the variance of residual is not constant.



Based on the QQ-plot below, after BoxCox transformation, the QQ-plot roughly follows a straight line.



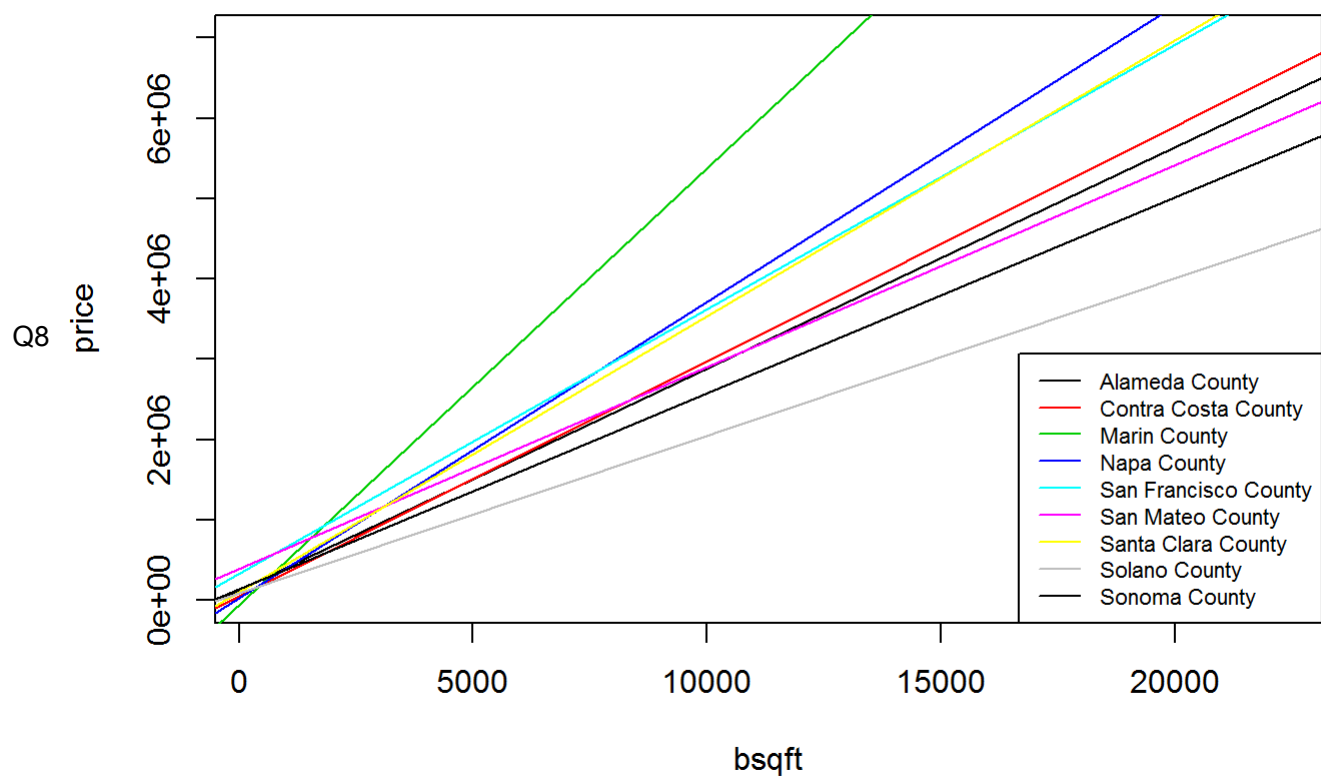
Q7

```
model3<-lm(data$price~data$bsqft+data$lsqft)
model3_est<-coef(summary(model3))[,1]
model3_se<-coef(summary(model3))[,2]
t_star<-(model3_est[2]-model3_est[3])/(model3_se[2])
critical_val<-qt(0.99,19997)
t_star>=critical_val
```

```
## data$bsqft
##      TRUE
```

H Null:; H1:; Using t test, because t_{star} is greater than critical_val , conclude null hypothesis.

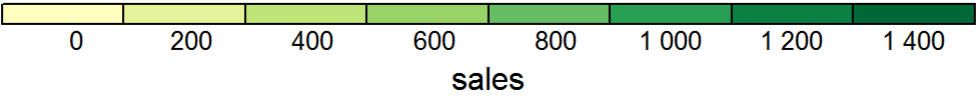
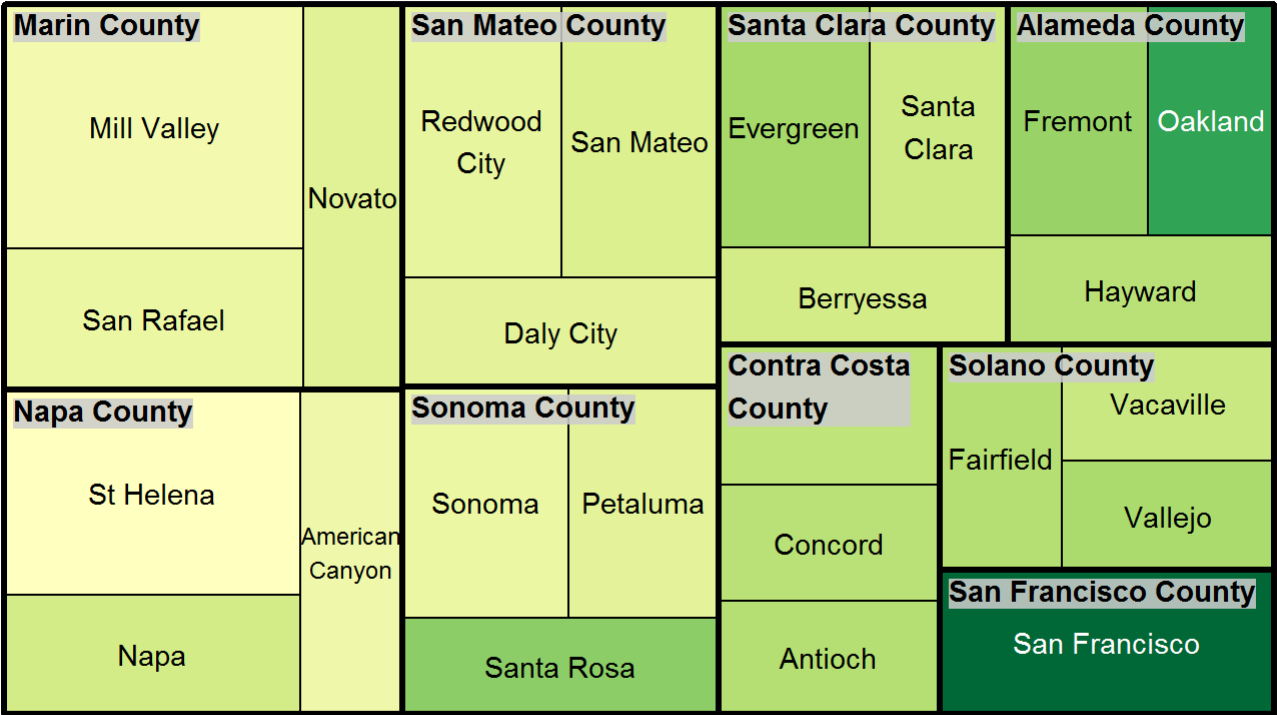
regression models

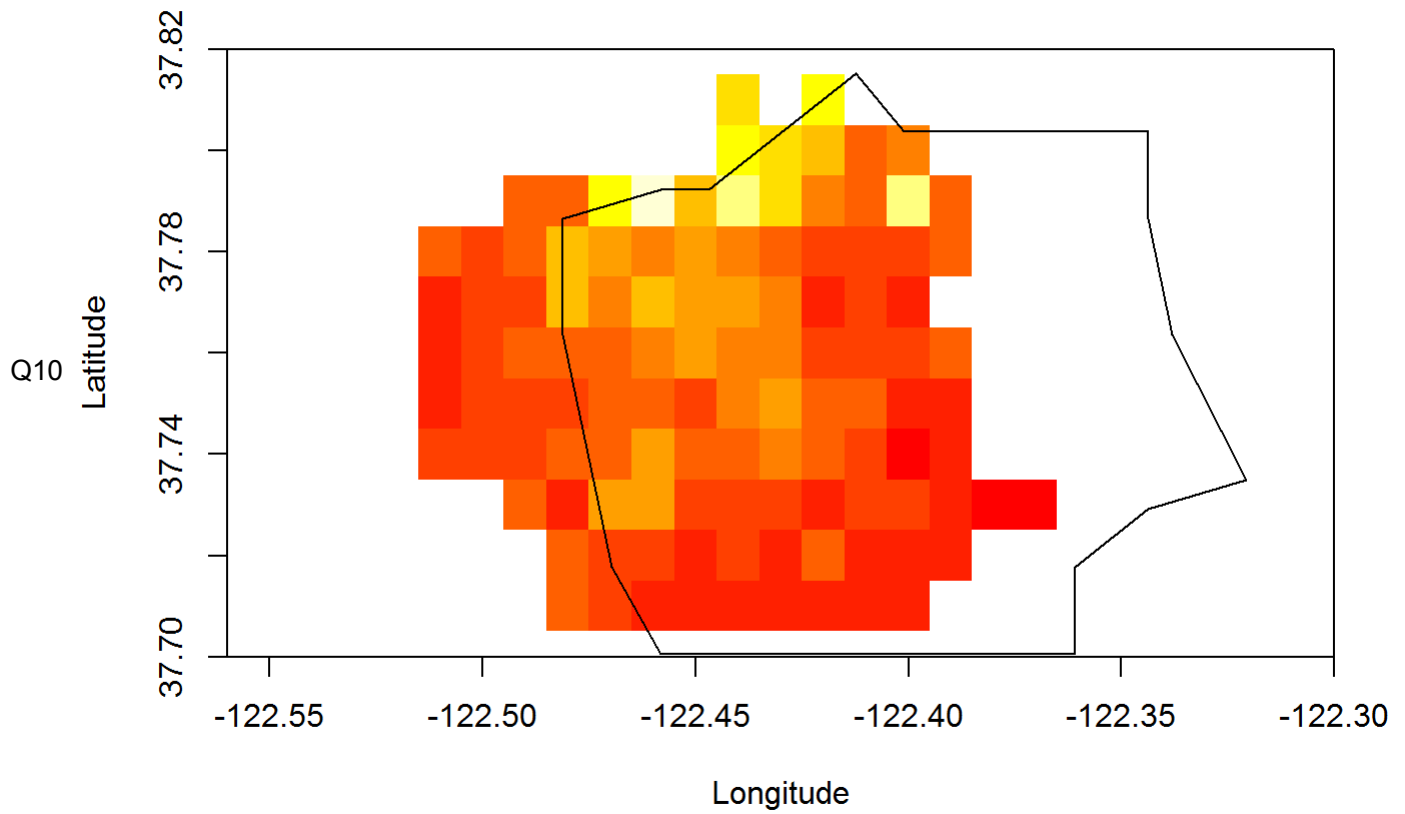


It's obvious that the simple regression lines with respect to different county are not parallel. Therefore, we can conclude the regression lines depend on county.

top3 sales with avg_price in differetn county

Q9



SF Heatmap of Average Housing Prices**SF Heatmap of Housing Sales**