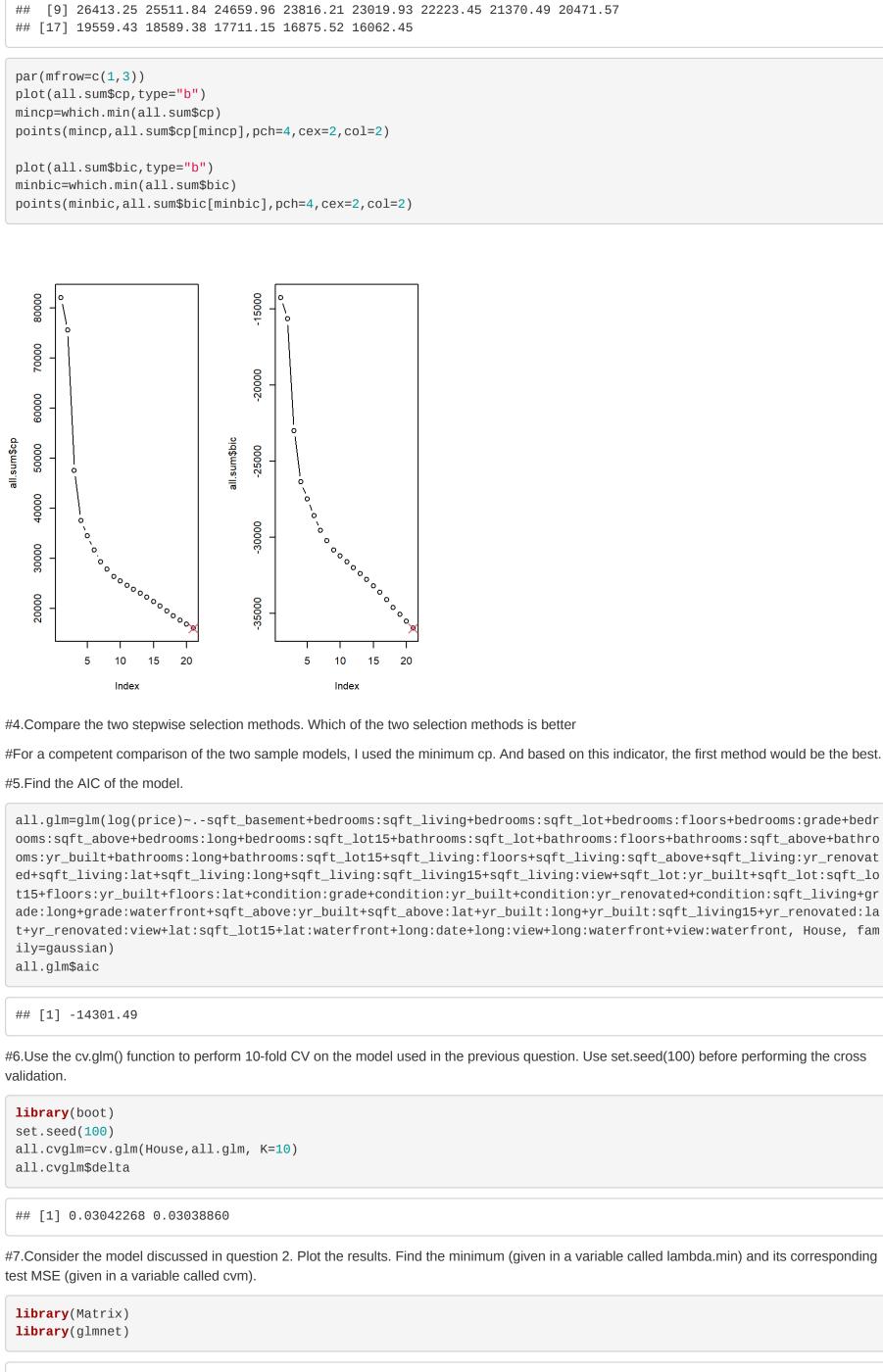
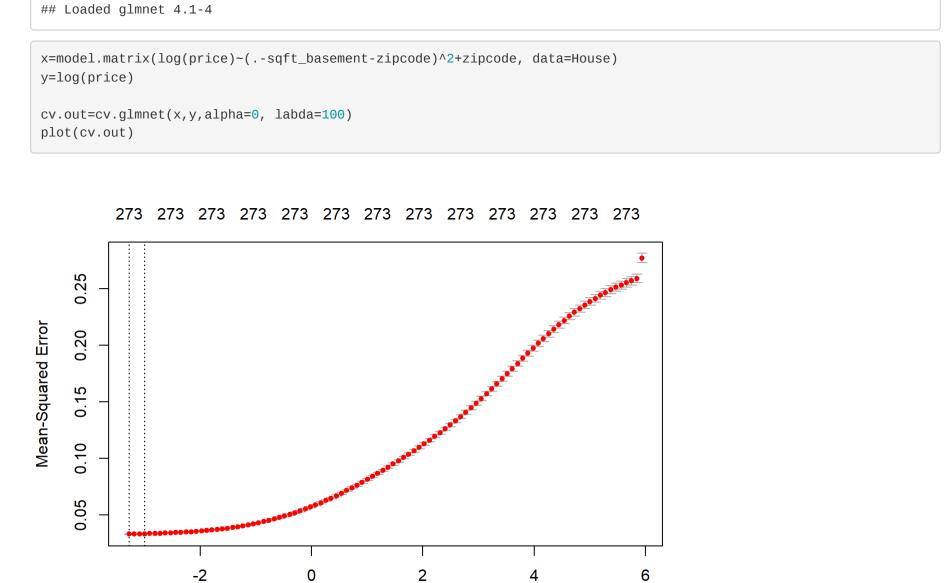
## Methods and Models Comparison Elena Kargopoltsev 01 03 2021 #1.Read and revise the data set House=read.csv("kc house sales.csv") House\$date=as.Date(House\$date, format="%Y%m%d") House\$view=as.factor(House\$view) House\$waterfront=as.factor(House\$waterfront) House\$zipcode=as.factor(House\$zipcode) House[15871,3]=3 attach(House) #2. Consider a multiple regression with log(price) as the response and all other variables and their two-way interactions as the predictors. Among the variables, exclude sqft basement since it is linearly dependent on other variables. Also, for zipcode consider the main effect only (do not consider its interaction with any other variables). Find the minimum Cp and the number of predictors in the best model. Plot BIC and mark the minimum point on the plot. Find the minimum BIC and the number of predictors in the best model. library(leaps) all.fit=regsubsets(log(price)~(.-sqft\_basement-zipcode)^2+zipcode, data=House, nvmax=20, method="forward") ## Reordering variables and trying again: all.sum=summary(all.fit) all.sum\$cp ## [1] 75964.07 43873.04 38396.44 34814.25 31298.66 29002.95 27127.30 25096.23 ## [9] 23667.52 22330.60 21131.72 20179.34 19257.24 18404.42 17551.73 16684.65 ## [17] 15769.09 14764.00 13980.81 13183.67 12466.26 par(mfrow=c(1,3))plot(all.sum\$cp,type="b") mincp=which.min(all.sum\$cp) points(mincp, all.sum\$cp[mincp], pch=4, cex=2, col=2) plot(all.sum\$bic,type="b") minbic=which.min(all.sum\$bic) points(minbic, all.sum\$bic[minbic], pch=4, cex=2, col=2) -20000 50000 15 10 15 Index Index #3.Repeat question 2 for the backward stepwise selection. all.fit=regsubsets(log(price)~(.-sqft\_basement-zipcode)^2+zipcode, data=House, nvmax=20, method="backward") ## Warning in leaps.setup(x, y, wt = wt, nbest = nbest, nvmax = nvmax, force.in = ## force.in, : 1 linear dependencies found ## Reordering variables and trying again: all.sum=summary(all.fit) all.sum\$cp ## [1] 82022.33 75564.04 47581.52 37565.25 34522.03 31717.20 29372.43 27806.55 [9] 26413.25 25511.84 24659.96 23816.21 23019.93 22223.45 21370.49 20471.57 ## [17] 19559.43 18589.38 17711.15 16875.52 16062.45 par(mfrow=c(1,3))plot(all.sum\$cp, type="b") mincp=which.min(all.sum\$cp) points(mincp, all.sum\$cp[mincp], pch=4, cex=2, col=2) plot(all.sum\$bic,type="b") minbic=which.min(all.sum\$bic) points(minbic, all.sum\$bic[minbic], pch=4, cex=2, col=2) all.sum\$cp all.sum\$bic 20000 -25000 15 15 10





 $Log(\lambda)$ 

cv.out\$lambda.min

## [1] 0.03777521

cv.out\$cvm[min.r]

## [1] 0.03278221

cv.out\$lambda.min

## [1] 8.726563e-05

cv.out\$cvm[min.r]

## [1] 0.03178324

0.1430

0.1426

1422

##

##

## AIC: 19176

iteration is 2.

15

4

<del>1</del>3

0

2000

4000

6000

#14.Repeat question 13 using lat as the predictor. There may be warnings due to non-unique values.

sqft\_living

8000

10000

12000

14000

log(price)

## Coefficients:

## (Intercept)

## -1.22606 -0.28482 0.01267 0.26152 1.27019

##  $ns(sqft_living, df = 12)1 0.58139$ 

##  $ns(sqft_living, df = 12)2$  0.63208

## ns(sqft\_living, df = 12)3 0.74925

##  $ns(sqft_living, df = 12)4 0.75321$ 

##  $ns(sqft_living, df = 12)5 0.80095$ 

## ns(sqft\_living, df = 12)6 0.91725

##  $ns(sqft_living, df = 12)7$  1.00108

##  $ns(sqft_living, df = 12)8$  1.10500 ## ns(sqft\_living, df = 12)9 1.27341

##  $ns(sqft_living, df = 12)10 2.68240$ 

##  $ns(sqft_living, df = 12)11 3.40100$ 

## ns(sqft\_living, df = 12)12 3.01637

## Number of Fisher Scoring iterations: 2

log(price). Add a line displaying the natural cubic spline.

12.14237

## Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

## (Dispersion parameter for gaussian family taken to be 0.1420859)

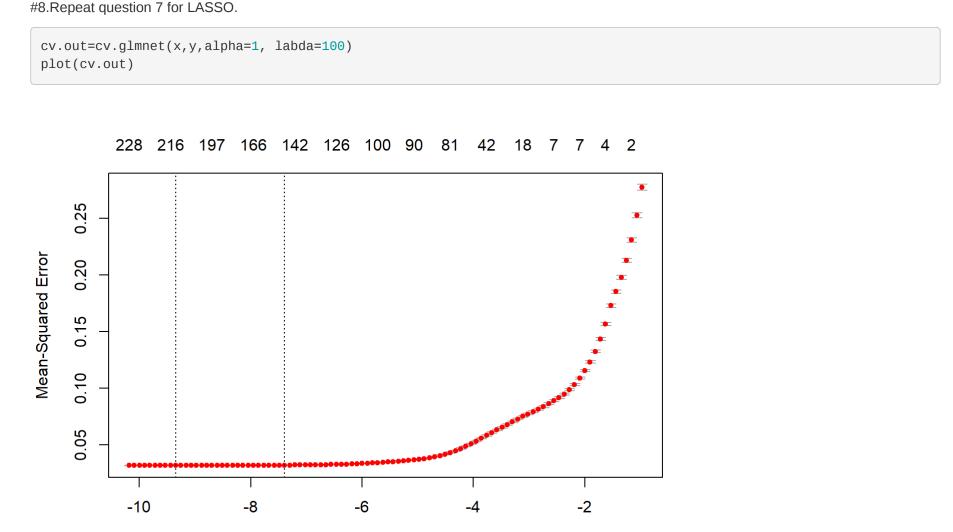
Null deviance: 5998.2 on 21612 degrees of freedom

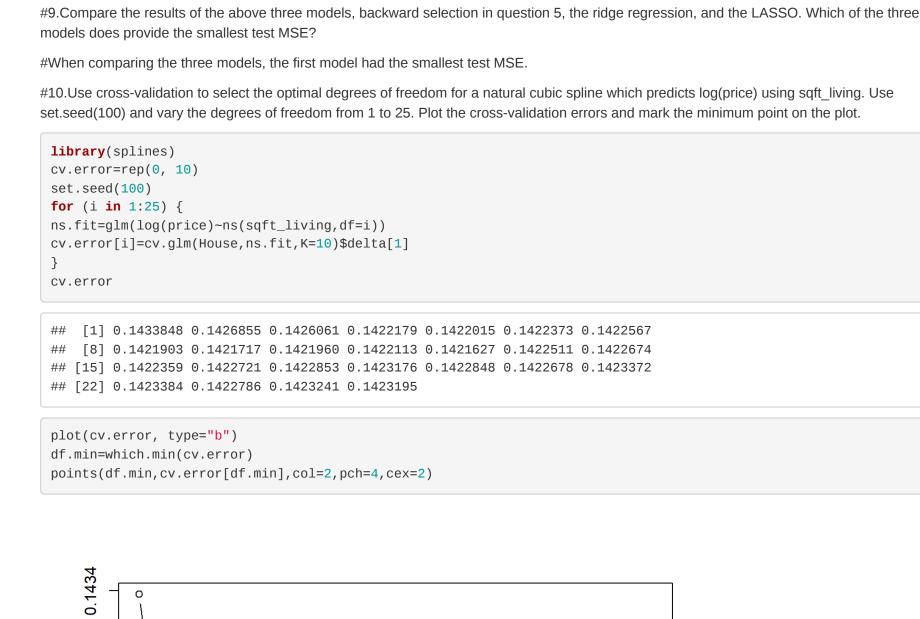
## Residual deviance: 3069.1 on 21600 degrees of freedom

min.r=which.min(cv.out\$cvm)

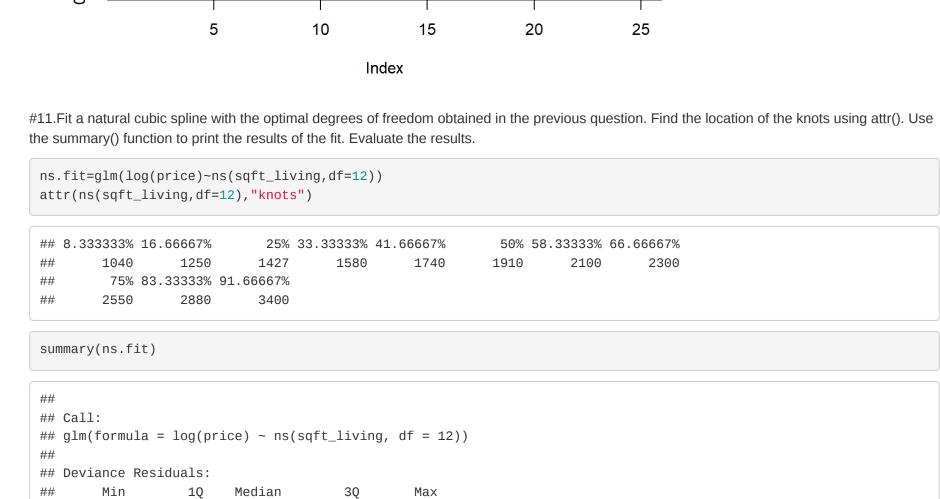
min.r=which.min(cv.out\$cvm)

## Warning: package 'glmnet' was built under R version 4.0.5





 $Log(\lambda)$ 



Estimate Std. Error t value Pr(>|t|)

0.06201 9.376

0.07331 8.622

0.06780 11.050

0.06950 11.524

0.05876 45.651

0.17397 19.549

0.06735 180.287 <2e-16 \*\*\*

0.07106 10.599 <2e-16 \*\*\*

0.07012 13.081 <2e-16 \*\*\*

0.06951 14.403 <2e-16 \*\*\* 0.06928 15.951 <2e-16 \*\*\*

0.06823 18.663 <2e-16 \*\*\*

0.21019 14.351 <2e-16 \*\*\*

#For this model the null deviance is 5998 what is good and also residual deviance is 3069 that also good, aic 19176 and number of fisher scoring

#12.Create a grid of sqft.living using the seq() function varying from the minimum to the maximum value in the data set with an increment of one.

Predict log(price) for the grid using the natural cubic spline model in the previous question. Construct a scatter plot between sqft\_living and

<2e-16 \*\*\*

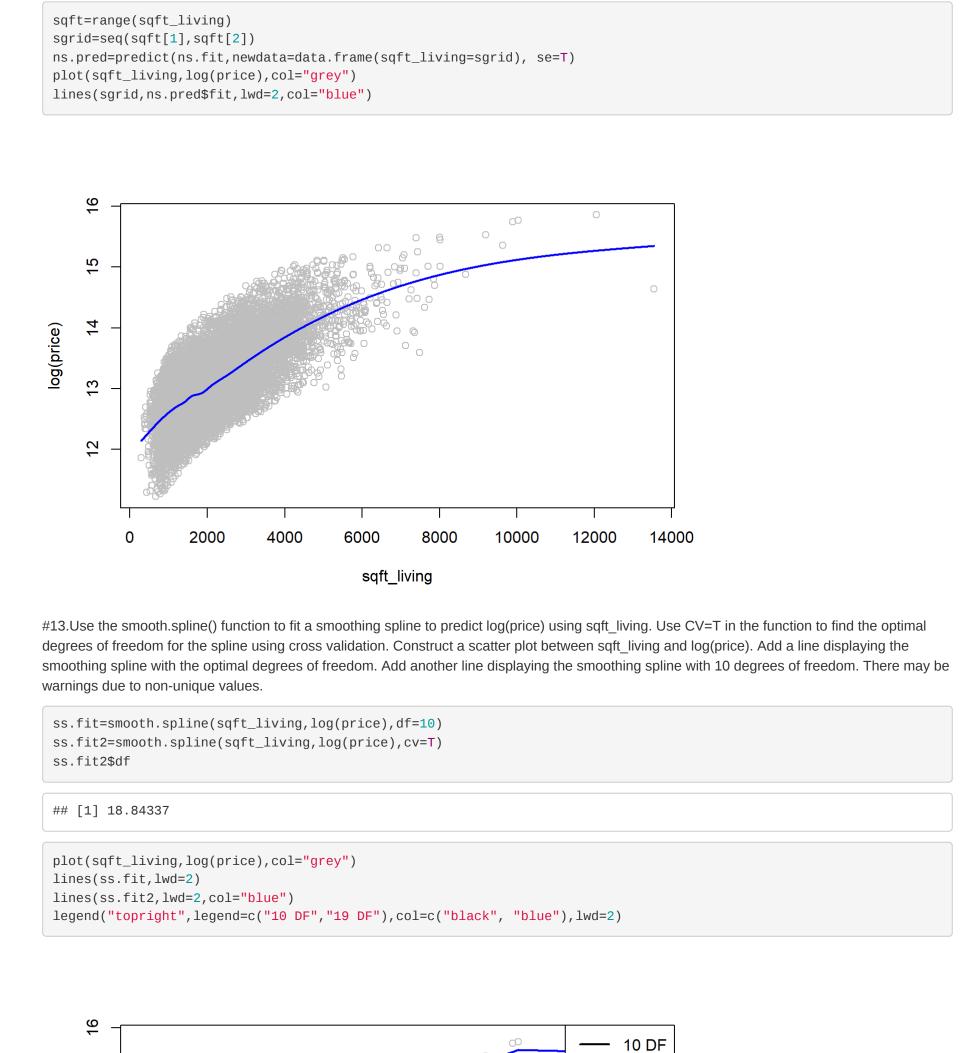
<2e-16 \*\*\*

<2e-16 \*\*\*

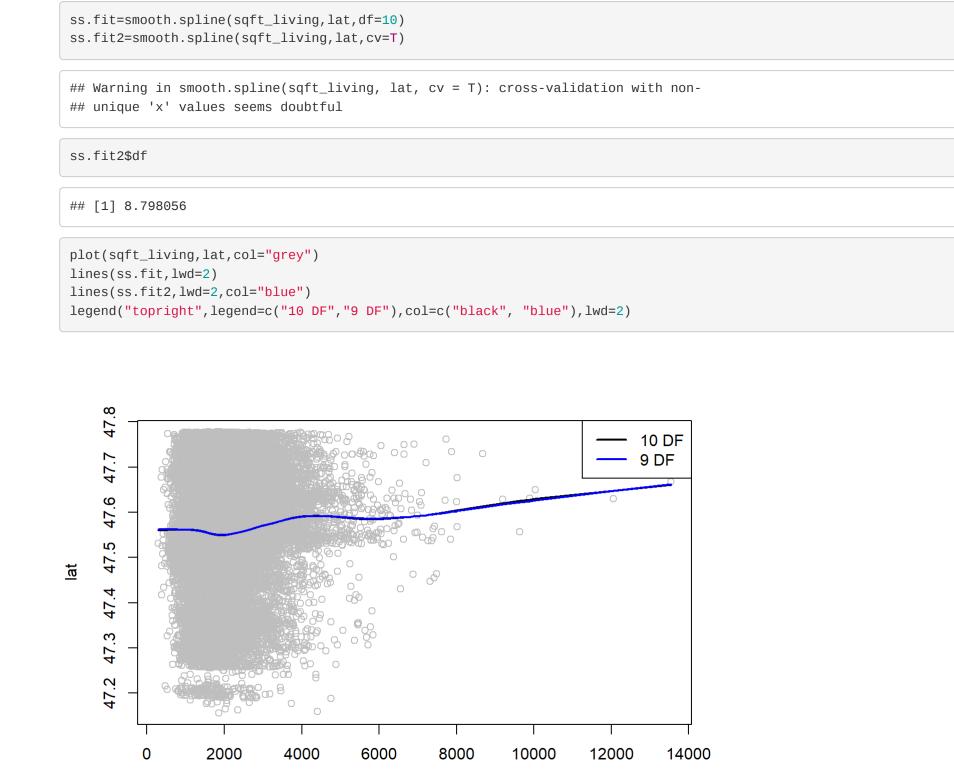
<2e-16 \*\*\*

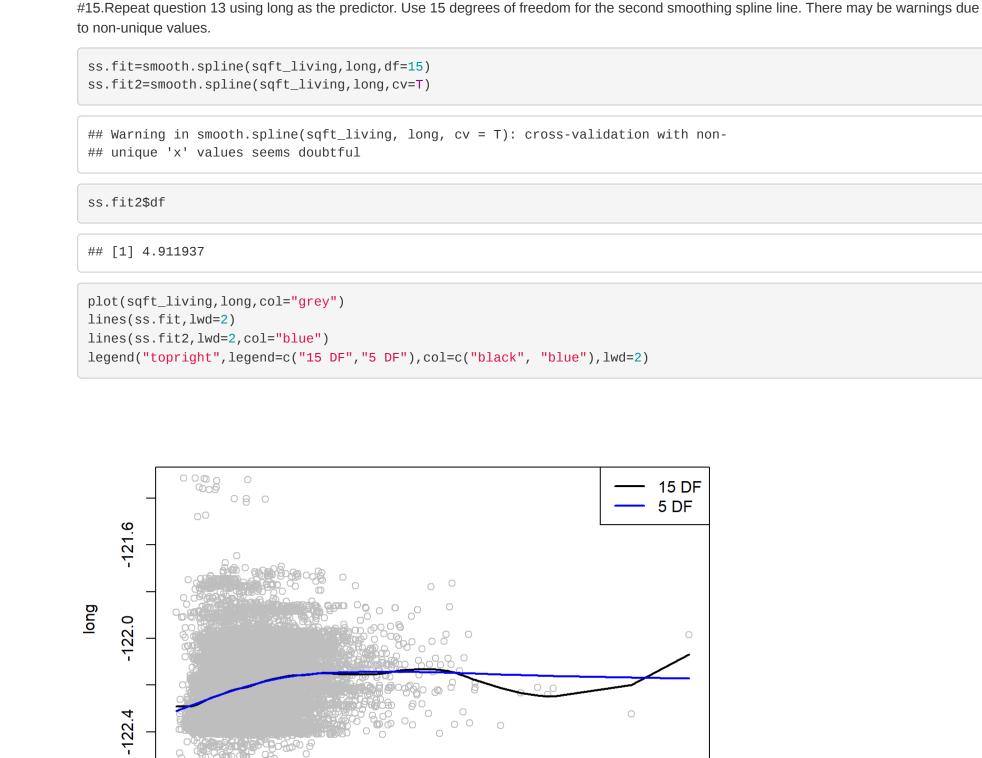
<2e-16 \*\*\*

<2e-16 \*\*\*

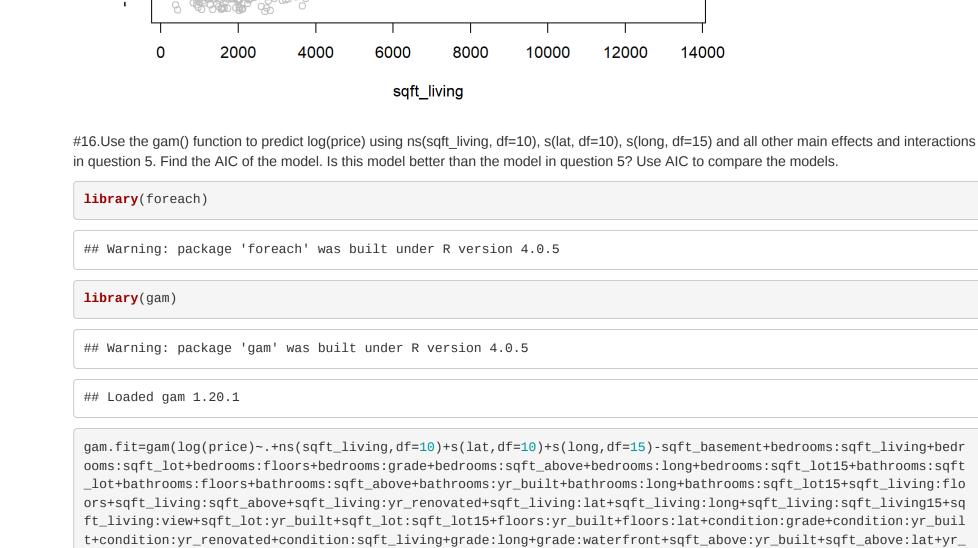


19 DF





sqft\_living



cv.error=rep(0, 10)

for(i in 1:k){

built:long+yr\_built:sqft\_living15+yr\_renovated:lat+yr\_renovated:view+lat:sqft\_lot15+lat:waterfront+long:date+lon g:view+long:waterfront+view:waterfront, data=House) gam.fit\$aic ## [1] -15629.64 #This model is best because has a smaller aic. #Extra set.seed(10)

 $house.fit=glm(log(price) - . + ns(sqft\_living, df=10) + s(lat, df=10) + s(long, df=15) - sqft\_basement + bedrooms: sqft\_living + bedrooms: sqft\_livi$  $drooms: sqft\_lot+bedrooms: floors+bedrooms: grade+bedrooms: sqft\_above+bedrooms: long+bedrooms: sqft\_lot15+bathrooms: sqft\_lot15+b$  $ft\_lot+bathrooms:floors+bathrooms:sqft\_above+bathrooms:yr\_built+bathrooms:long+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_above+bathrooms:yr\_built+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_lot15+sqft\_living:floors+bathrooms:sqft\_living:sqft\_li$ loors+sqft\_living:sqft\_above+sqft\_living:yr\_renovated+sqft\_living:lat+sqft\_living:long+sqft\_living:sqf  $\verb|sqft_living:view+sqft_lot:yr_built+sqft_lot:$  $ilt+condition: yr\_renovated+condition: sqft\_living+grade: long+grade: waterfront+sqft\_above: yr\_built+sqft\_above: lat+yr-built+sqft\_above: lat-yr-built+sqft\_above: lat-y$  $r\_built:long+yr\_built:sqft\_living15+yr\_renovated:lat+yr\_renovated:view+lat:sqft\_lot15+lat:waterfront+long:date+louble.equilibrium:date+louble.equili$ 

ng:view+long:waterfront+view:waterfront, data=House, family=gaussian) cv.error[i]=cv.glm(House, house.fit, K=10)\$delta[1] cv.error ## [1] 0.03031268 0.03026746 0.03030724 0.03025463 0.03029429 0.03028774 ## [7] 0.03035161 0.03032994 0.03034355 0.03028529 #This model is best because has a smaller test MSE.