

Predicting_bangalore_house_price

Loading data

```
Train_bangalore <- read.csv("Train_banglore_machinehack.csv")
str(Train_bangalore)
```

```
## 'data.frame':    13320 obs. of  9 variables:
## $ area_type      : Factor w/ 4 levels "Built-up Area",...: 4 3 1 4 4 4 4 4 4 3 ..
## $ availability: Factor w/ 81 levels "14-Jul","14-Nov",...: 41 81 81 81 81 81 35
81 81 81 ...
## $ location      : Factor w/ 1306 levels "", "Anekal", "Banaswadi",...: 432 322 121
3 776 729 1288 917 988 815 444 ...
## $ size          : Factor w/ 32 levels "", "1 Bedroom",...: 16 20 19 19 16 16 21 21
19 25 ...
## $ society       : Factor w/ 2689 levels "", "3Codeli", "7 ise P",...: 463 2440 1 214
9 1 609 929 354 1 1 ...
## $ total_sqft    : Factor w/ 2117 levels "1","1.25Acres",...: 71 1289 515 603 240 2
06 1320 1467 370 32 ...
## $ bath          : int  2 5 2 3 2 2 4 4 3 6 ...
## $ balcony       : int  1 3 3 1 1 1 NA NA 1 NA ...
## $ price         : num  39.1 120 62 95 51 ...
```

Data Exploration

```
#we will start by looking at price skewness of all houses in bangalore as per data given
summary(Train_bangalore$price)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      8.0    50.0    72.0   112.6   120.0   3600.0
```

So the price starts from 8lakh and goes upto 3600lakhs.the data is right skewed,where 75% of the houses are less than 1.2crore in price range.

```
#how is the price distribution among diffferent areas in bangalore? without outliers
plot_ly(x=Train_bangalore$price[!Train_bangalore$price %in% boxplot.stats(Train_bangalore$price)$out],type="histogram",color = Train_bangalore$area_type[!Train_bangalore$price %in% boxplot.stats(Train_bangalore$price)$out])
```

```
## Warning: package 'bindrcpp' was built under R version 3.4.2
```

```
#only outliers
```

```
plot_ly(x=Train_bangalore$price[Train_bangalore$price %in% boxplot.stats(Train_bangalore$price)$out], type="histogram", color = Train_bangalore$area_type[Train_bangalore$price %in% boxplot.stats(Train_bangalore$price)$out])
```

We created two plot one with outliers and and other without. From the first(without outlier) we see that most of the houses are in range 50-100lakhs and number is high for type ,superbuilt-up Area

Since Bangalore is an area where lots of working professional live with a better income, Super Built up area provide amenities like covered community centres/ clubs, other covered common facilities.

Dealing Missing values

```
#whether belongs to other society
Train_bangalore$Othersociety<-ifelse(Train_bangalore$society=="", "1", "0")
Train_bangalore[Train_bangalore==""]<-NA
miss_val<-sapply(Train_bangalore, function(x) sum(is.na(x)))
miss_val
```

```
##      area_type availability      location      size      society
##           0           0           1          16          5502
## total_sqft      bath      balcony      price Othersociety
##           0          73          609           0           0
```

Percentage of Missing values

```
percent_mis<-as.data.frame(round((miss_val/nrow(Train_bangalore))*100,1))
names(percent_mis)<-c("Percentage")
percent_mis
```

```
##      Percentage
## area_type      0.0
## availability    0.0
## location        0.0
## size            0.1
## society         41.3
## total_sqft      0.0
## bath            0.5
## balcony         4.6
## price           0.0
## Othersociety    0.0
```

41 percent of missing values are in Society, assuming missing values here as the following property are not in any society but some other place. Replacing Na's in society with other

```
Train_bangalore$society<-as.character(Train_bangalore$society)
Train_bangalore$society[is.na(Train_bangalore$society)]<-"Other"
Train_bangalore$society<-as.factor(Train_bangalore$society)
length(Train_bangalore$society[is.na(Train_bangalore$society)])
```

```
## [1] 0
```

Balconies?

From the plot, we see that number of balconies are one or two in maximum properties, with two slightly higher.

Since missing values are less than .05 percent, replacing with mode

```
#imputing missing value for balcony with mode

Train_bangalore$balcony<-as.character (Train_bangalore$balcony)

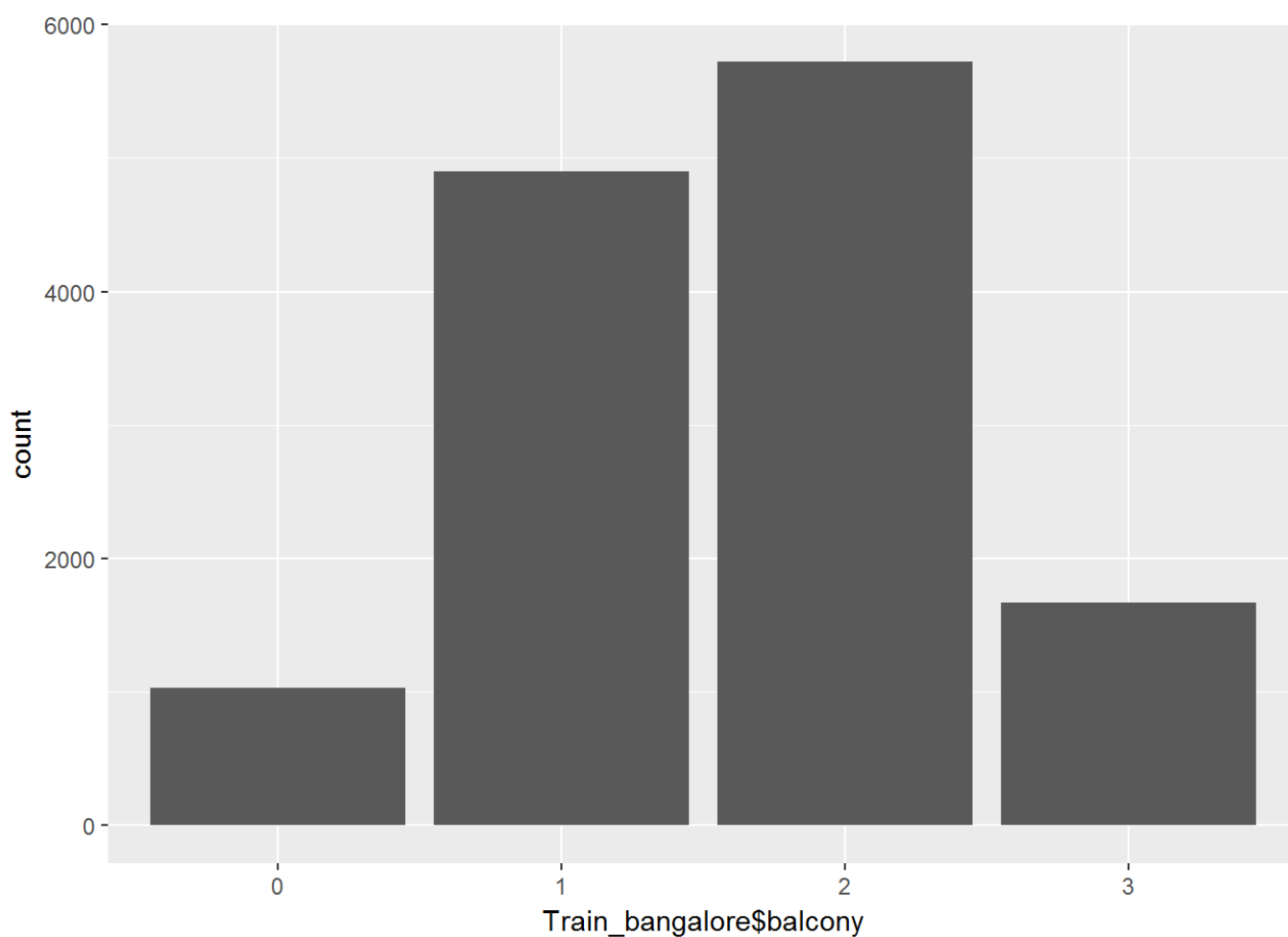
Train_bangalore$balcony[is.na (Train_bangalore$balcony) ]<-"2"

Train_bangalore$balcony<-as.factor (Train_bangalore$balcony)

table (Train_bangalore$balcony)
```

```
##
##    0    1    2    3
## 1029 4897 5722 1672
```

```
ggplot (Train_bangalore,aes (Train_bangalore$balcony) )+geom_bar ()
```



```
Train_bangalore$bath[Train_bangalore$bath>=6]<-6

Train_bangalore$bath<-as.character (Train_bangalore$bath)

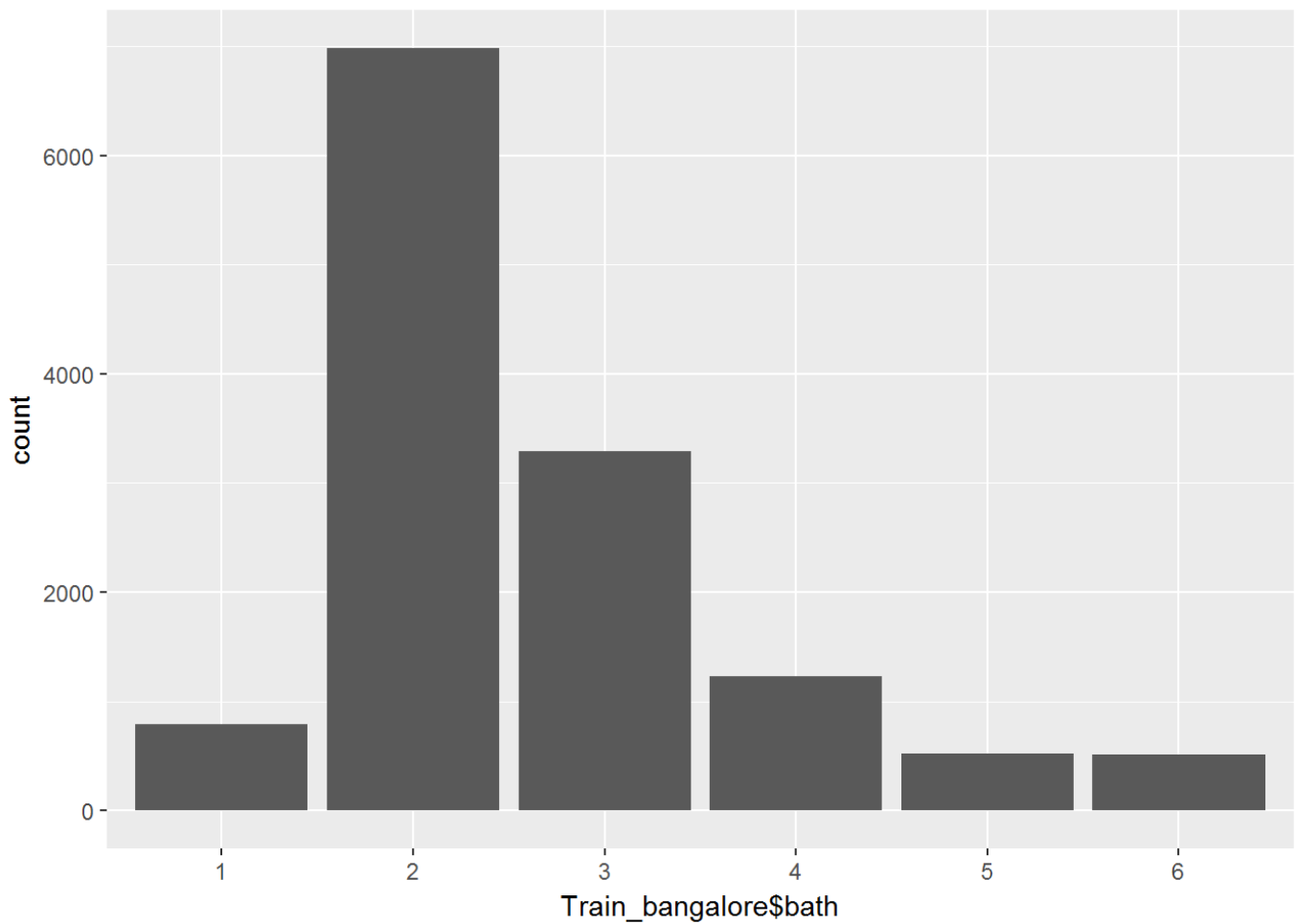
#replacing NA with mode value
Train_bangalore$bath[is.na (Train_bangalore$bath) ]<-"2"

Train_bangalore$bath<-as.factor (Train_bangalore$bath)

table (Train_bangalore$bath)
```

```
##
##      1      2      3      4      5      6
## 788 6981 3286 1226  524  515
```

```
ggplot (Train_bangalore,aes (Train_bangalore$bath)) +geom_bar ()
```



Area Size and price

```
names (Train_bangalore[,1:4])
```

```
## [1] "area_type" "availability" "location" "size"
```

I made a new variable which separate society and not in society,price seems same for each area type ,some of them which are not in societies have prices when area is super-built area.

```
#converting size variable into character variable
Train_bangalore$size<-as.character (Train_bangalore$size)
```

Which locations has highest mean price ?

```
mean_price_per_location<-as.data.frame(arrange(aggregate(Train_bangalore$price,by=location,mean),desc(x)))

names(mean_price_per_location)<-c("location","Price_Mean")
#grouping based on pricing
mean_price_per_location$Price_Range_group<-cut(mean_price_per_location$Price_Mean,c(0,100,200,500,1000,2000),labels = c("lowprice","averagepriced","HighPrice","veryHigh","extreme"))
#places where price is higher
head(mean_price_per_location[mean_price_per_location$Price_Range_group=="extreme"||mean_price_per_location$Price_Range_group=="veryHigh"],20)
```

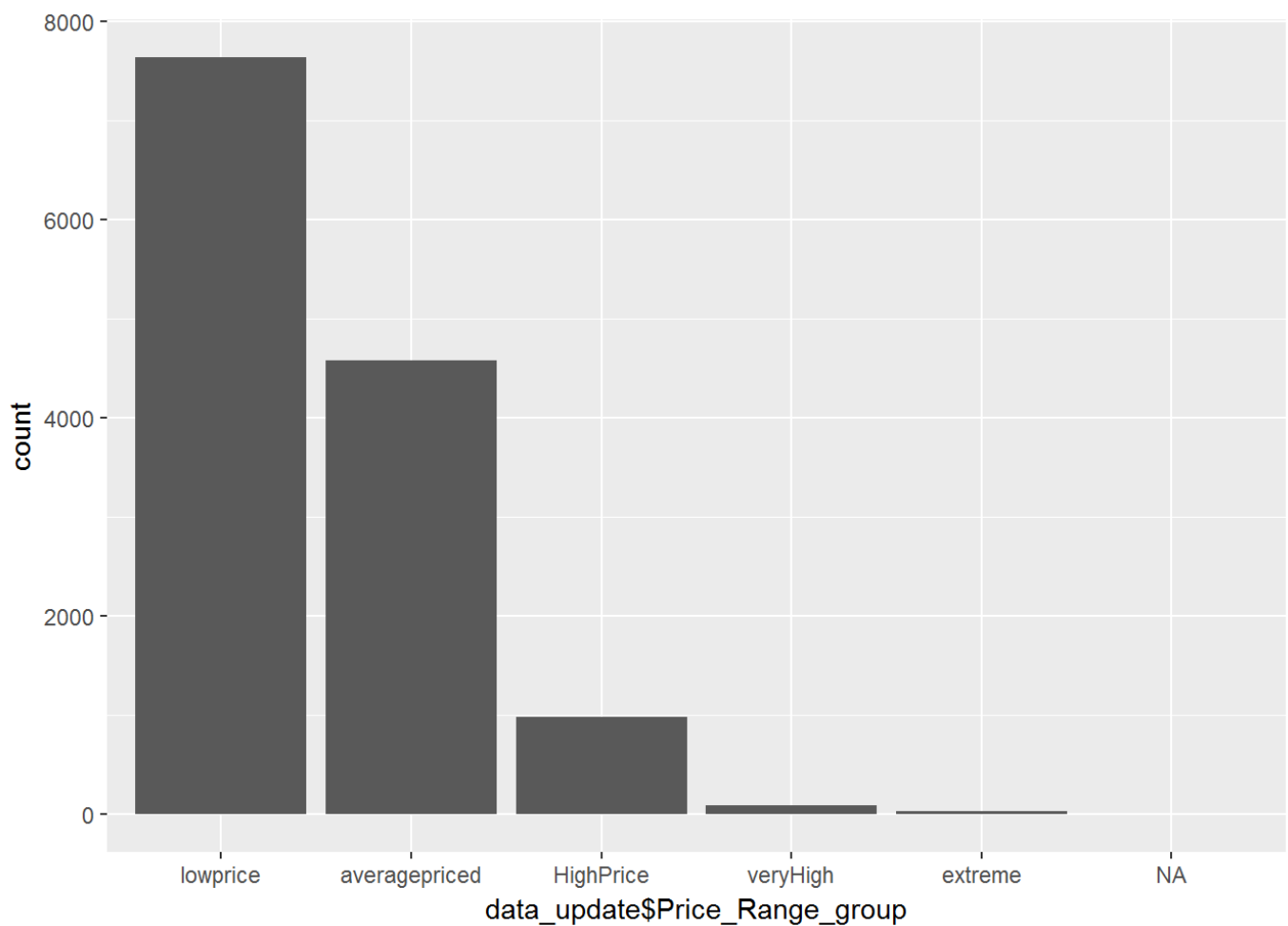
##	location	Price_Mean	Price_Range_group
## 1	Cubbon Road	1900.0000	extreme
## 2	Ashok Nagar	1486.0000	extreme
## 3	Defence Colony	1167.7143	extreme
## 4	Yemlur	1093.3889	extreme
## 5	Church Street	1068.0000	extreme
## 6	D Souza Layout	1015.0000	extreme
## 7	Sadashiva Nagar	1011.1000	extreme
## 8	Sindhi Colony	988.0000	veryHigh
## 9	Srinivas Colony	922.0000	veryHigh
## 10	5th Block Jayanagar	905.0000	veryHigh
## 11	Binnamangala	900.0000	veryHigh
## 12	Cunningham Road	824.3846	veryHigh
## 13	Hunasamaranahalli	787.5000	veryHigh
## 14	2nd Block Koramangala	761.0000	veryHigh
## 15	Shanthala Nagar	754.0000	veryHigh
## 16	Dollars Colony	744.3750	veryHigh
## 17	Kathreguppe	725.0000	veryHigh
## 18	Sector 4 HSR Layout	700.0000	veryHigh
## 19	Rest House Road	690.0000	veryHigh
## 20	Ramakrishnappa Layout	662.8571	veryHigh

merging the new coloumn.

```
data_update<-left_join(Train_bangalore,mean_price_per_location,by="location")
Train_bangalore<-left_join(Train_bangalore,mean_price_per_location,by="location")
```

Distribution of houses based on prices in bangalore

```
ggplot(data_update,aes(data_update$Price_Range_group))+geom_bar()
```

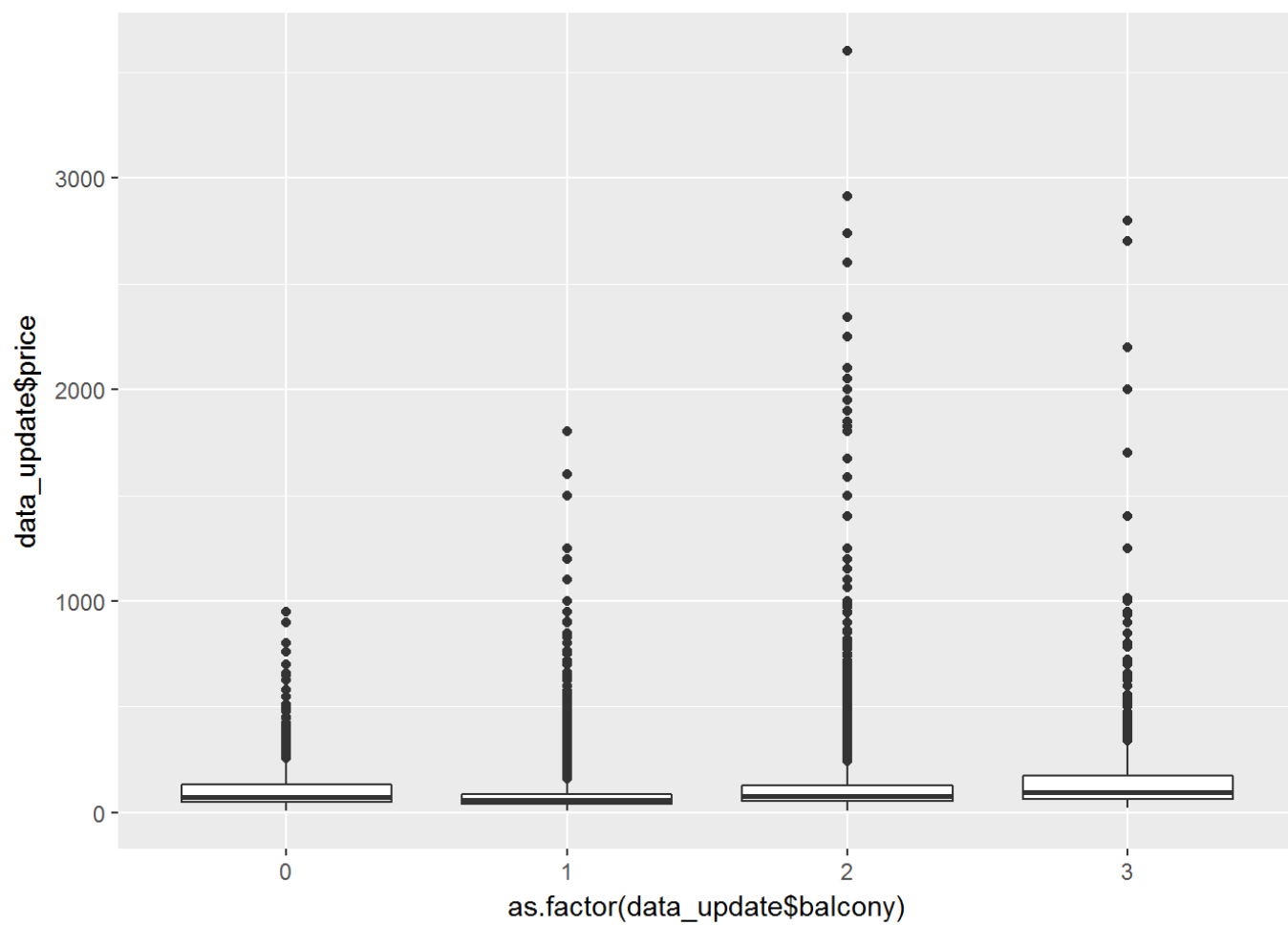


Around 80% of houses are in range less than a crore.

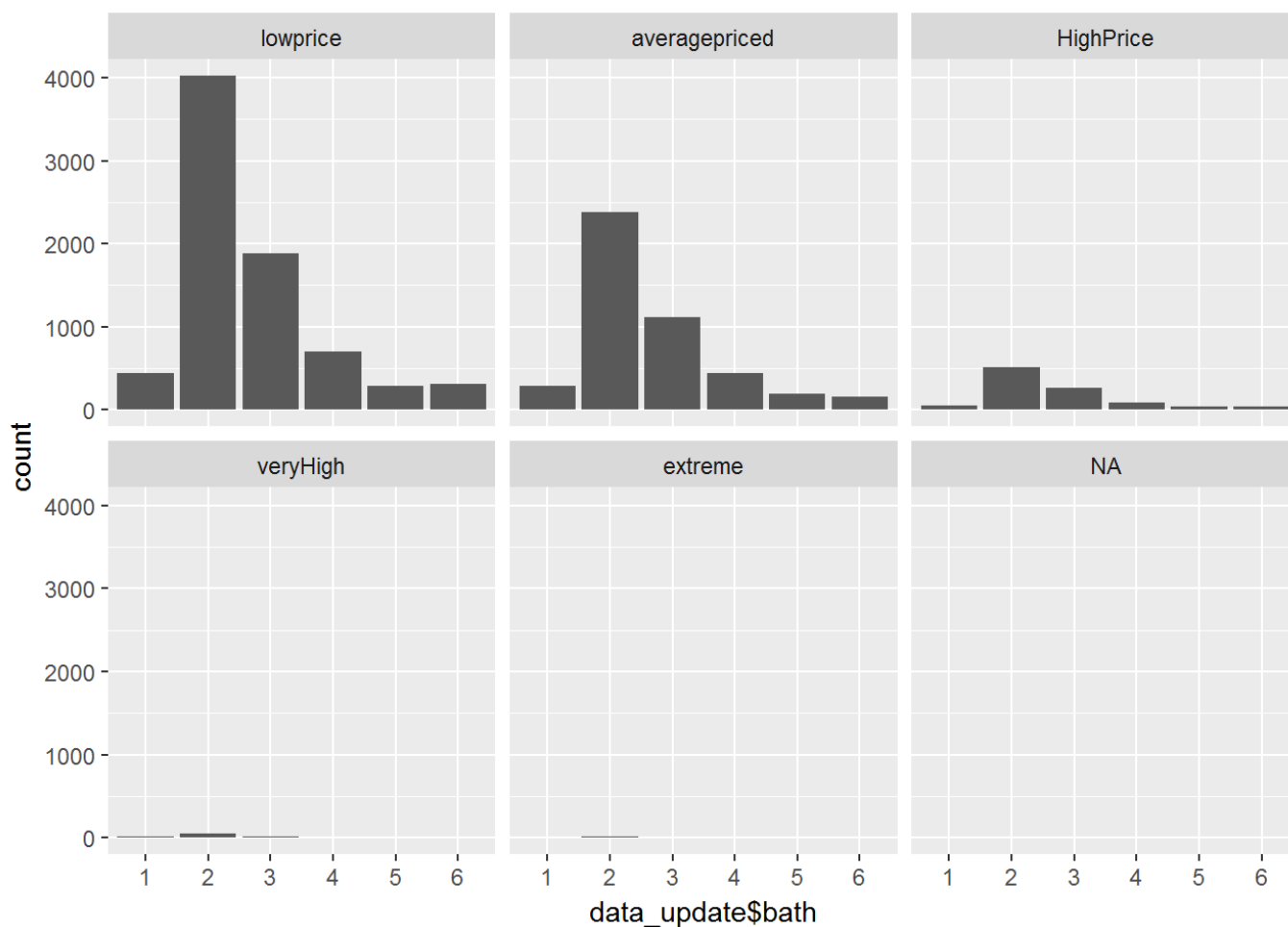
statistical balcony,bath and price relationship

```
data_update$balcony<-as.factor(data_update$balcony)
data_update$bath<-as.factor(data_update$bath)

#is mean price different for all
ggplot(data_update,aes(as.factor(data_update$balcony),data_update$price))+geom_boxplot()
```



```
ggplot(data_update, aes(data_update$bath), fill=data_update$Price_Range_group)+geom_bar()+facet_wrap(~data_update$Price_Range_group)
```

Cleanig Variable Total_sqft

How much is the area of a house or plot. For any analysis we need to clean the variable in a single format i.e. Squarefeet

```
ll<-Train_bangalore$total_sqft

ll<-as.data.frame(ll)

ll$ll<-as.character(ll$ll)
c<-grep("Acres",ll$ll,ignore.case = T)
b<-grep("Sq. Yards",ll$ll,ignore.case = T)
d<-grep("Sq. Meter",ll$ll,ignore.case = T)
e<-grep("Cents",ll$ll,ignore.case = T)
f<-grep("-",ll$ll,ignore.case = T)
g<-grep("Perch",ll$ll,ignore.case = T)

ll[c,<-gsub("Acres", "",ll[c,],ignore.case = T)
ll[b,<-gsub("Sq. Yards", "",ll[b,],ignore.case = T)
ll[d,<-gsub("Sq. Meter", "",ll[d,],ignore.case = T)
ll[e,<-gsub("Cents", "",ll[e,],ignore.case = T)
ll[g,<-gsub("Perch", "",ll[g,],ignore.case = T)

# for values noted as for example 1600-1700, are replaced with their mean values
rr<-as.data.frame(ll[f,])
names(rr)<-"name"
head(rr)
```

```
##          name
## 1 2100 - 2850
## 2 3010 - 3410
## 3 2957 - 3450
## 4 3067 - 8156
## 5 1042 - 1105
## 6 1145 - 1340
```

```
rr$name<-as.character(rr$name)
pp<-separate(rr,name,into = c("firstnumber","secondnumber"),sep = " - ")
head(pp)
```

```
##   firstnumber secondnumber
## 1         2100         2850
## 2         3010         3410
## 3         2957         3450
## 4         3067         8156
## 5         1042         1105
## 6         1145         1340
```

```
pp$firstnumber<-as.numeric(pp$firstnumber)
pp$secondnumber<-as.numeric(pp$secondnumber)
final_pp<-(pp$firstnumber+pp$secondnumber)/2
#final conversions
ll[f,<]<-final_pp
ll$ll<-as.numeric(ll$ll)
```

```
## Warning: NAs introduced by coercion
```

```

ll[which(is.na(ll)),]<-0

# 1 Acres=43560 sqft

ll[c,]<-ll[c,]*4350

#1sq yard = 9 sqft

ll[b,]<-ll[b,]*9

# 1sq.meter=10.76 sqft

ll[d,]<-ll[d,]*10.76

#1 cents =435.61 sqft

ll[e,]<-ll[e,]*435.61

#1 perch =272.75 sqft

ll[g,]<-ll[g,]*435.61

#replacing for values like 1600 -1900 with mean of there values

ll[f,]<-final_pp

Train_bangalore$total_sqft<-ll$ll
Train_bangalore$total_sqft<-as.numeric(Train_bangalore$total_sqft)

# histogram without outliers

cc<-as.data.frame(Train_bangalore$total_sqft[Train_bangalore$total_sqft<=2000])
names(cc)<-"var1"

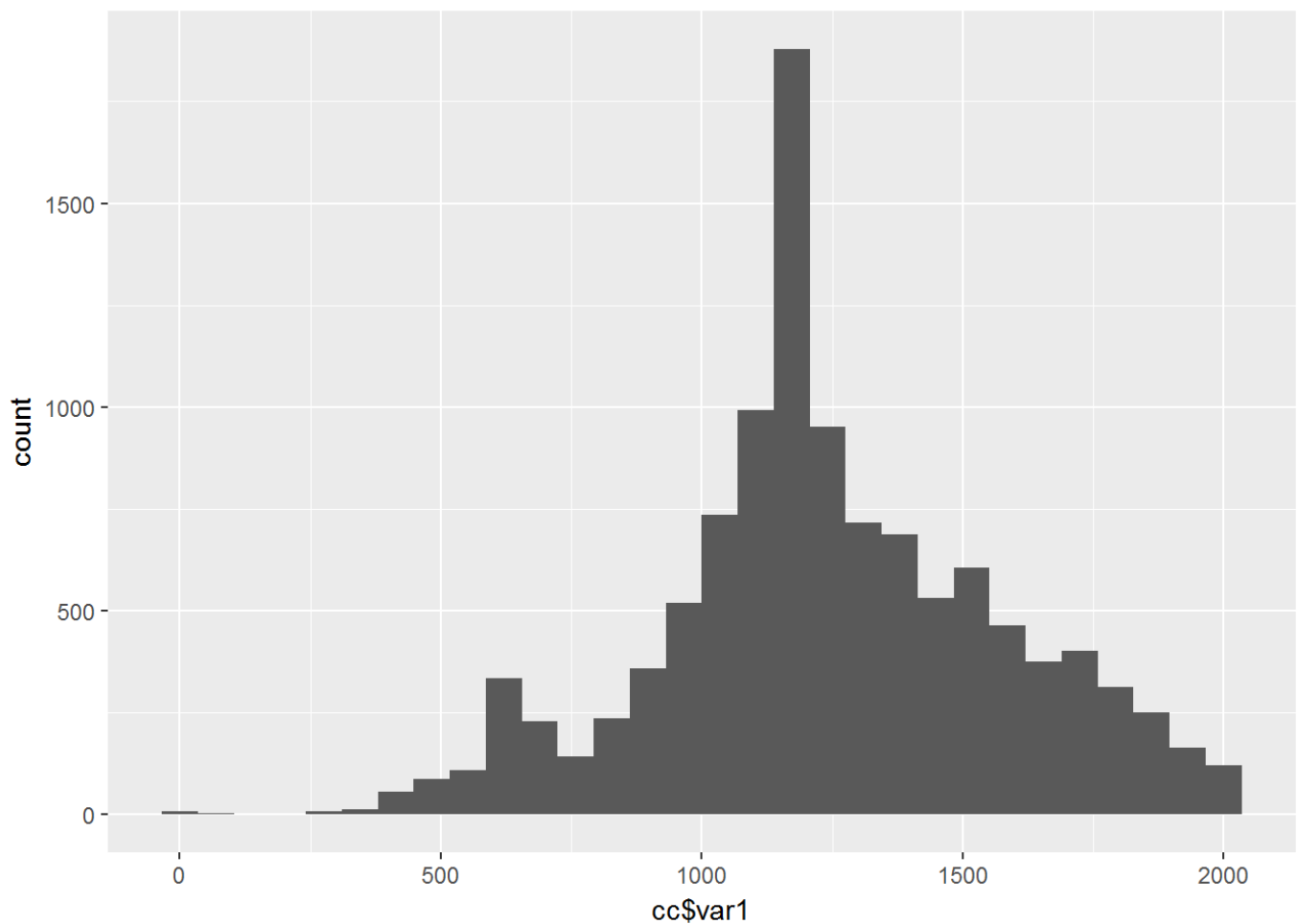
ggplot(cc,aes(cc$var1,))+geom_histogram()

```

```

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

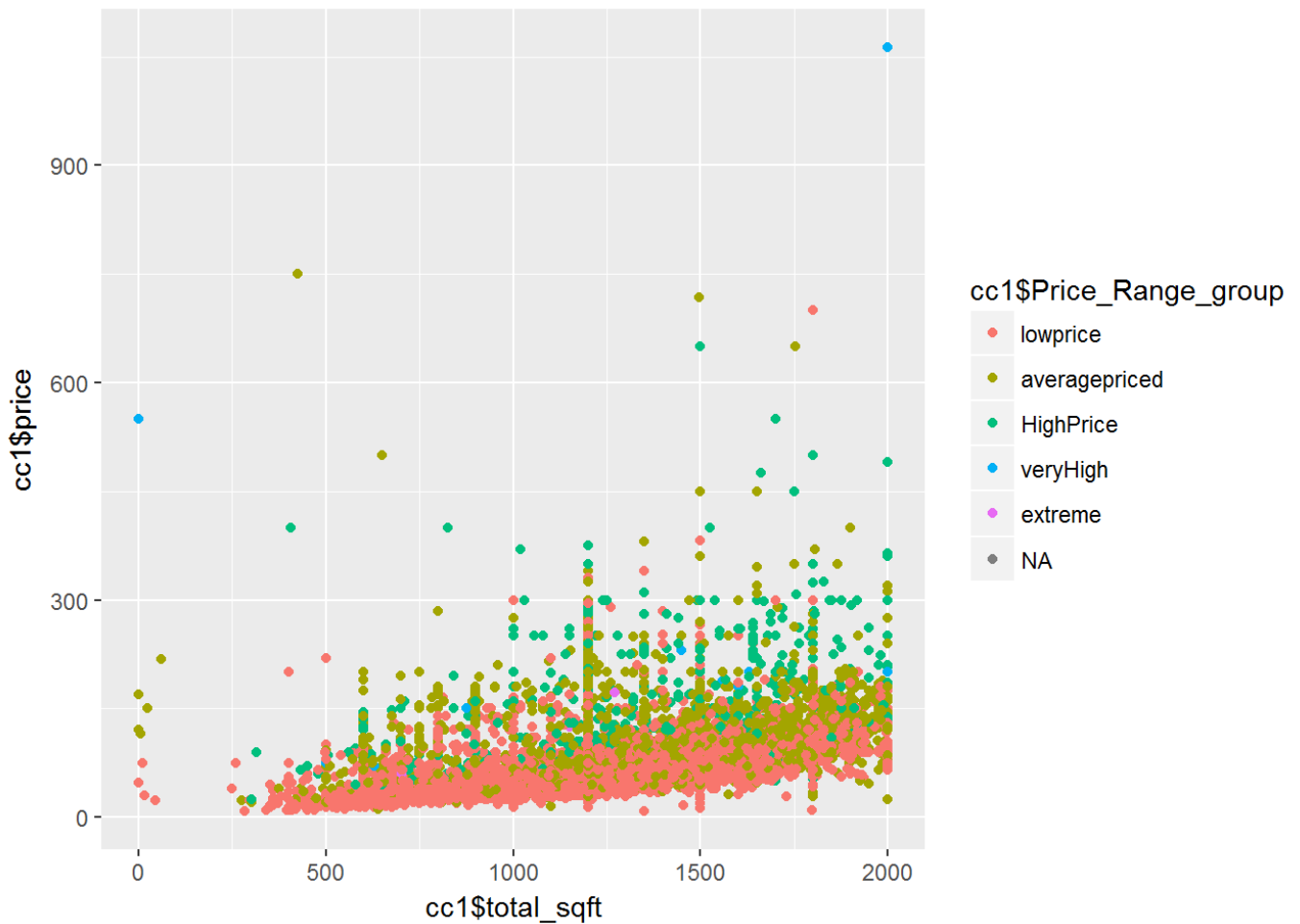
```



most of the plots are in the range from 1000 to 1500 square foot.

Price and and size relationship

```
cc1<-Train_bangalore[Train_bangalore$total_sqft<=2000,]  
ggplot(cc1,aes(cc1$total_sqft,cc1$price,color=cc1$Price_Range_group))+geom_point()
```

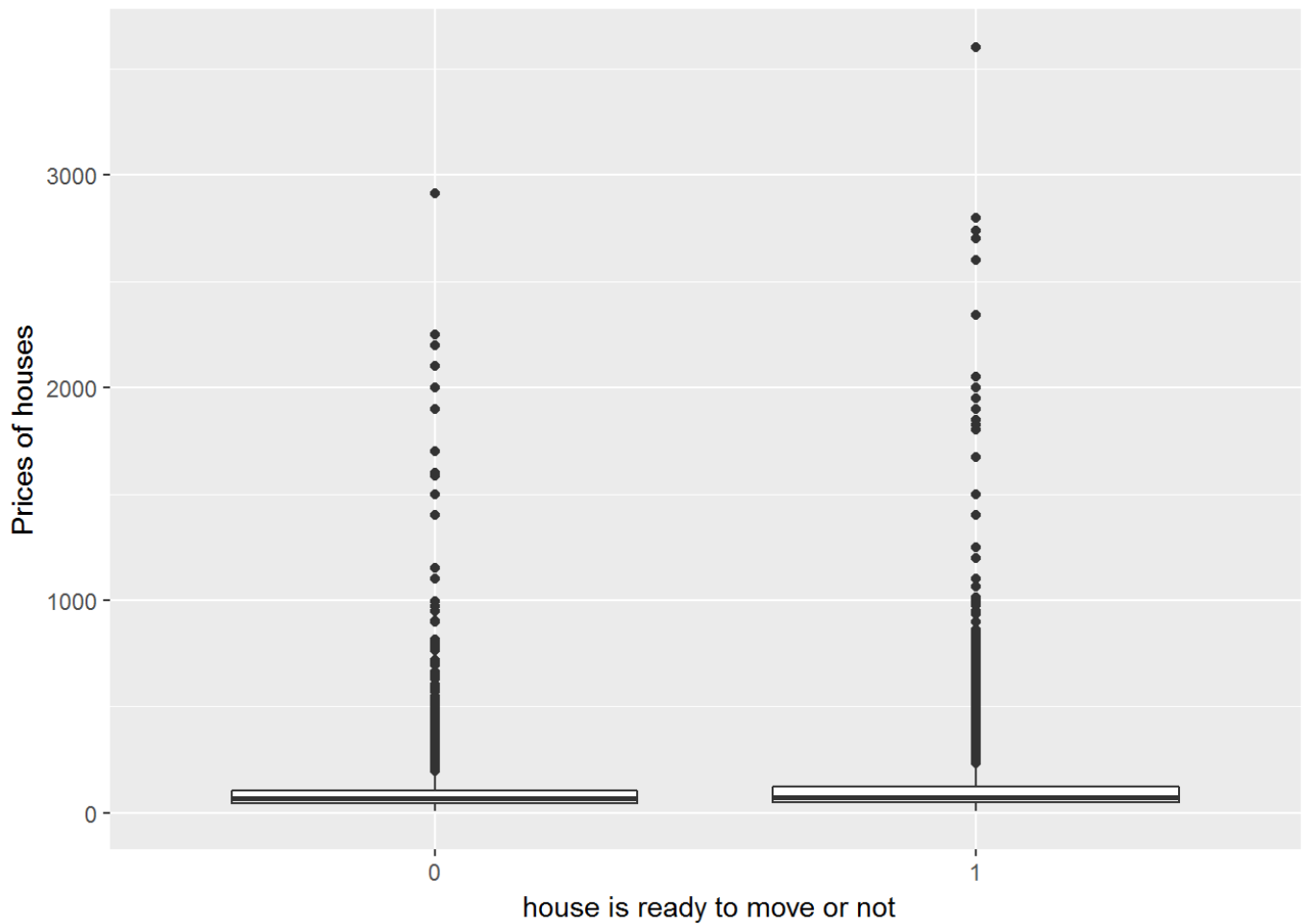


Now we would like to check for the availability and price change. Hypothesis here is that nearer the house available to move in, higher should be the price. Let's analyse the variable Availability

```
head(Train_bangalore$availability)
```

```
## [1] 19-Dec      Ready To Move Ready To Move Ready To Move Ready To Move
## [6] Ready To Move
## 81 Levels: 14-Jul 14-Nov 15-Aug 15-Dec 15-Jun 15-Nov 15-Oct ... Ready To Move
```

```
Train_bangalore$AvailibityCheck<-ifelse(Train_bangalore$availability=="Ready To Move",1,0)
Train_bangalore$AvailibityCheck<-as.factor(Train_bangalore$AvailibityCheck)
ggplot(Train_bangalore,aes(Train_bangalore$AvailibityCheck,Train_bangalore$price))+
  geom_boxplot()+ylab("Prices of houses")+xlab(" house is ready to move or not")
```



```
#Cleaning size variable
ss<-Train_bangalore$size
y<-sapply(strsplit(ss, " "),head,1)
y<-as.numeric(y)
Train_bangalore$size<-y
Train_bangalore$Othersociety<-as.factor(Train_bangalore$Othersociety)
```

Model Builing

Partitioning data into Train and test set for validation

```
Train_bangalore<-na.omit(Train_bangalore)
str(Train_bangalore)
```

```
## 'data.frame':    13303 obs. of  13 variables:
## $ area_type      : Factor w/ 4 levels "Built-up Area",...: 4 3 1 4 4 4 4 4 4 3 ...
## $ availability    : Factor w/ 81 levels "14-Jul","14-Nov",...: 41 81 81 81 81 8 1 35 81 81 81 ...
## $ location        : Factor w/ 1306 levels "", "Anekal", "Banaswadi",...: 432 32 2 1213 776 729 1288 917 988 815 444 ...
## $ size            : num  2 4 3 3 2 2 4 4 3 6 ...
## $ society          : Factor w/ 2689 levels "3Codeli","7 ise P",...: 462 2440 145 1 2149 1451 608 928 353 1451 1451 ...
## $ total_sqft       : num  1056 2600 1440 1521 1200 ...
## $ bath             : Factor w/ 6 levels "1","2","3","4",...: 2 5 2 3 2 2 4 4 3 6 ...
## $ balcony          : Factor w/ 4 levels "0","1","2","3": 2 4 4 2 2 2 3 3 2 3 .
..
## $ price            : num  39.1 120 62 95 51 ...
## $ Othersociety      : Factor w/ 2 levels "0","1": 1 1 2 1 2 1 1 1 2 2 ...
## $ Price_Mean        : num  48.3 115 61.3 115.3 95.6 ...
## $ Price_Range_group: Factor w/ 5 levels "lowprice","averagepriced",...: 1 2 1 2 1 2 2 3 2 3 ...
## $ AvailabilityCheck : Factor w/ 2 levels "0","1": 1 2 2 2 2 2 1 2 2 2 ...
## - attr(*, "na.action")=Class 'omit' Named int [1:17] 569 580 1776 2265 2810 28 63 5334 6424 6637 6720 ...
## .. ..- attr(*, "names")= chr [1:17] "569" "580" "1776" "2265" ...
```

```
Train_bangalore1<-Train_bangalore[,-c(2,3,5)]
Train_bangalore1$size<-as.factor(Train_bangalore1$size)
index<-sample(2,nrow(Train_bangalore1),prob = c(0.7,0.3),replace = TRUE)
bangalore_train<-Train_bangalore1[index==1,]
bangalore_test<-Train_bangalore1[index==2,]
```

Parameter tuning for RandomForest

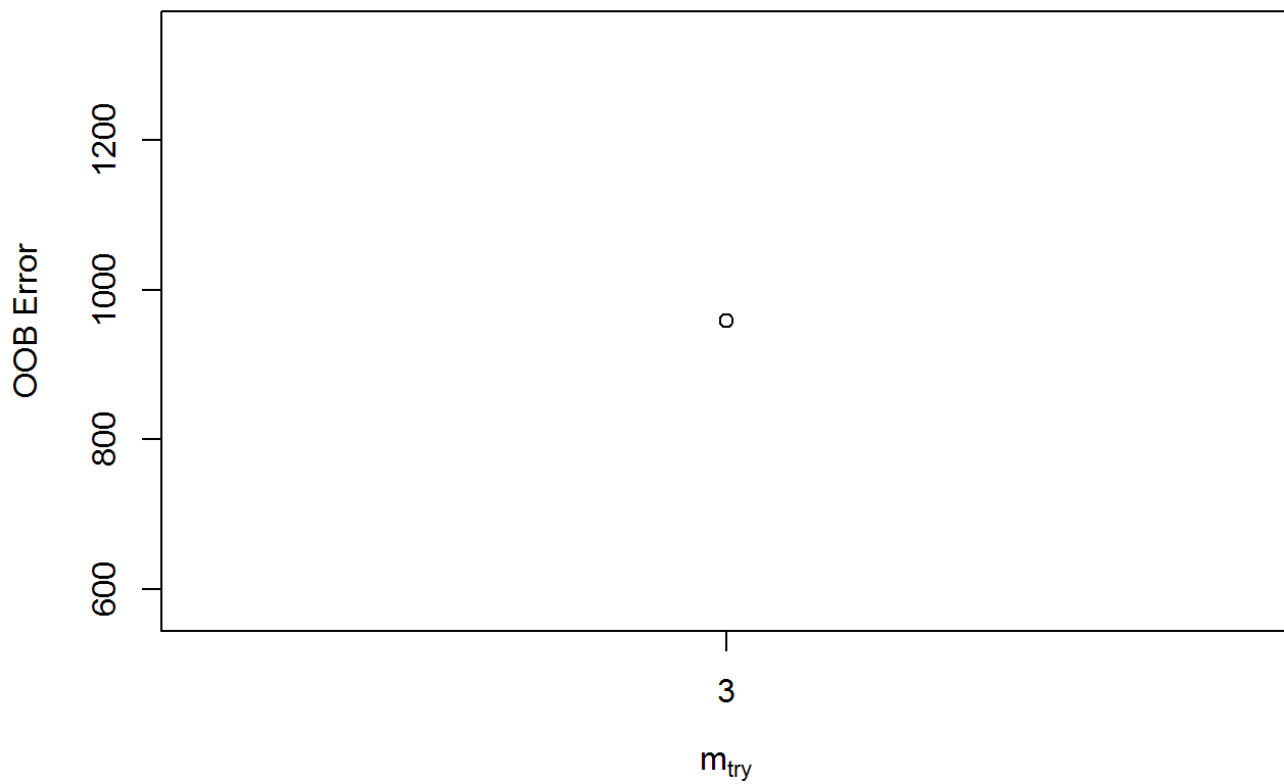
```
#optimised value of mtry
```

```
optimalvalue<-tuneRF(bangalore_train,bangalore_train$price,stepFactor = 1.2,improve
= 0.1,trace = T,plot = T)
```

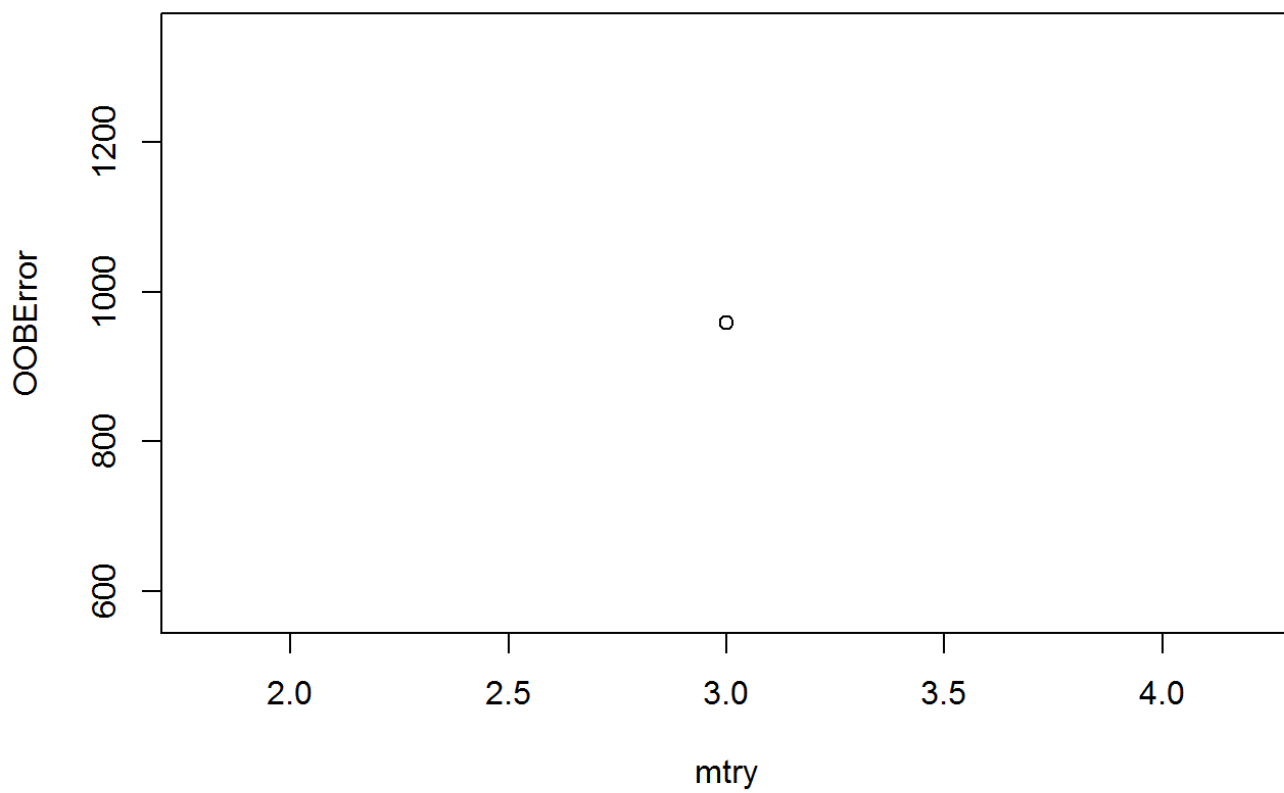
```
## mtry = 3 OOB error = 958.2856
```

```
## Searching left ...
```

```
## Searching right ...
```



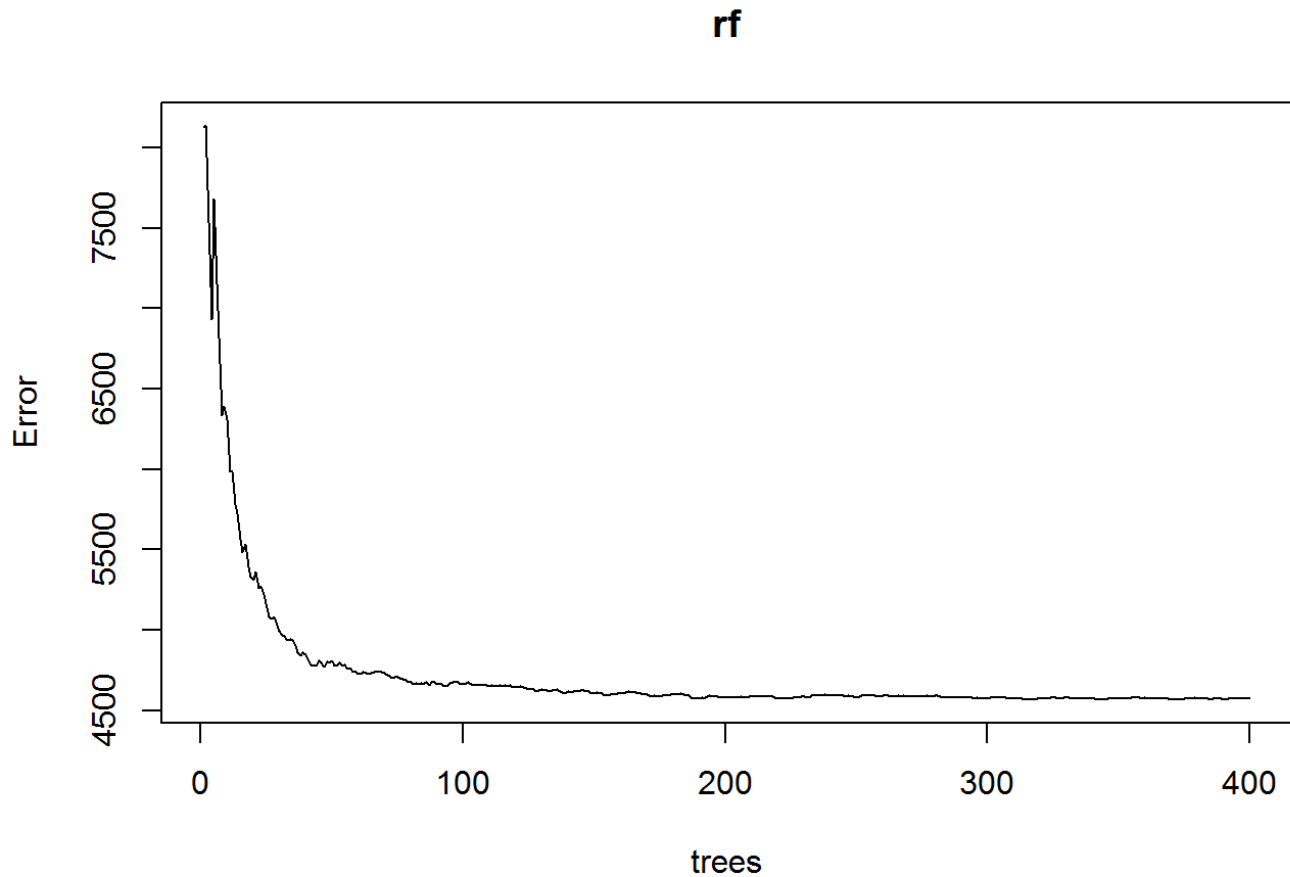
```
plot(optimalvalue)
```



Training Random Forest model

```
rf<-randomForest(bangalore_train$price~.,data = bangalore_train,ntree=400,mtry=3,importance=TRUE)

plot(rf)
```



making prediction

```
#making predictions for test data
pred<-predict(rf,bangalore_test)

temp<-as.data.frame(cbind(pred,bangalore_test$price))

names(temp)<-c("predicted","actual")

temp$difference<-temp$actual-temp$predicted

#function to calculate RMSE
rmse <- function(error)
{
  sqrt(mean(error^2))
}
rmse(temp$difference) #evaluation
```

```
## [1] 85.76328
```

Using Linear regression

```
linearModel<-glm(bangalore_train$price~.,data = bangalore_train)
summary(linearModel)
```

```
##
## Call:
## glm(formula = bangalore_train$price ~ ., data = bangalore_train)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1067.68    -22.86     -1.07     17.05    1964.01
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -1.363e+01  6.624e+00  -2.058 0.039655 *
## area_typeCarpet Area    -8.822e+00  1.216e+01  -0.725 0.468189
## area_typePlot Area      3.051e+01  3.764e+00   8.106 5.88e-16 ***
## area_typeSuper built-up Area -3.916e-01  2.733e+00  -0.143 0.886054
## size2           5.767e+00  8.411e+00   0.686 0.492986
## size3           2.355e+01  8.915e+00   2.642 0.008261 **
## size4           9.333e+01  1.010e+01   9.243 < 2e-16 ***
## size5           8.921e+01  1.186e+01   7.521 5.95e-14 ***
## size6          -9.146e+00  1.352e+01  -0.677 0.498602
## size7           3.083e+01  1.602e+01   1.924 0.054366 .
## size8          -4.320e+01  1.711e+01  -2.526 0.011567 *
## size9          -5.033e+01  2.168e+01  -2.322 0.020263 *
## size10          4.932e+01  4.424e+01   1.115 0.265012
## size11         -8.345e+01  4.890e+01  -1.706 0.087952 .
## size13         -4.686e+01  9.537e+01  -0.491 0.623222
## size14         -7.292e+01  9.533e+01  -0.765 0.444329
## size16         -2.163e+01  9.611e+01  -0.225 0.821919
## size18         -1.152e+02  9.537e+01  -1.208 0.226931
## size27         -6.495e+01  9.540e+01  -0.681 0.495952
## size43          4.329e+02  9.536e+01   4.540 5.69e-06 ***
## total_sqft      9.518e-04  1.410e-04   6.751 1.55e-11 ***
## bath2           1.956e+00  7.746e+00   0.252 0.800673
## bath3           1.873e+01  8.444e+00   2.218 0.026550 *
## bath4           5.450e+01  9.448e+00   5.769 8.24e-09 ***
## bath5           9.953e+01  1.058e+01   9.411 < 2e-16 ***
## bath6           1.559e+02  1.201e+01  12.980 < 2e-16 ***
## balcony1        1.247e+01  4.085e+00   3.052 0.002276 **
## balcony2        1.591e+01  4.126e+00   3.857 0.000116 ***
## balcony3        1.674e+01  4.779e+00   3.503 0.000462 ***
## Othersocety1     -1.554e+01  2.225e+00  -6.986 3.03e-12 ***
## Price_Mean       6.523e-01  3.586e-02  18.188 < 2e-16 ***
## Price_Range_groupaveragepriced -5.058e+00  2.978e+00  -1.698 0.089503 .
## Price_Range_groupHighPrice    1.707e+01  8.442e+00   2.021 0.043273 *
## Price_Range_groupveryHigh     5.090e+01  2.535e+01   2.008 0.044668 *
## Price_Range_groupextreme     1.732e+02  4.388e+01   3.948 7.94e-05 ***
## AvailibilityCheck1          -6.100e+00  2.506e+00  -2.434 0.014950 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for gaussian family taken to be 10015.104)
```

```
## (Dispersion parameter for gaussian family taken to be 8915.124)
##
##      Null deviance: 177313268  on 9299  degrees of freedom
## Residual deviance:  82589705  on 9264  degrees of freedom
## AIC: 111018
##
## Number of Fisher Scoring iterations: 2
```

XGBOOST model

```
#One -Hot coding for factor variables
train1<-sparse.model.matrix(price ~ .,data = bangalore_train)

#train label
train_label<-bangalore_train[, "price"]
train_matrix<-xgb.DMatrix(data = as.matrix(train1),label= train_label)

#repeating steps for test
test1<-sparse.model.matrix(price ~ .,data = bangalore_test)
test_label<-bangalore_test[, "price"]
test_matrix<-xgb.DMatrix(data = as.matrix(test1),label= test_label)

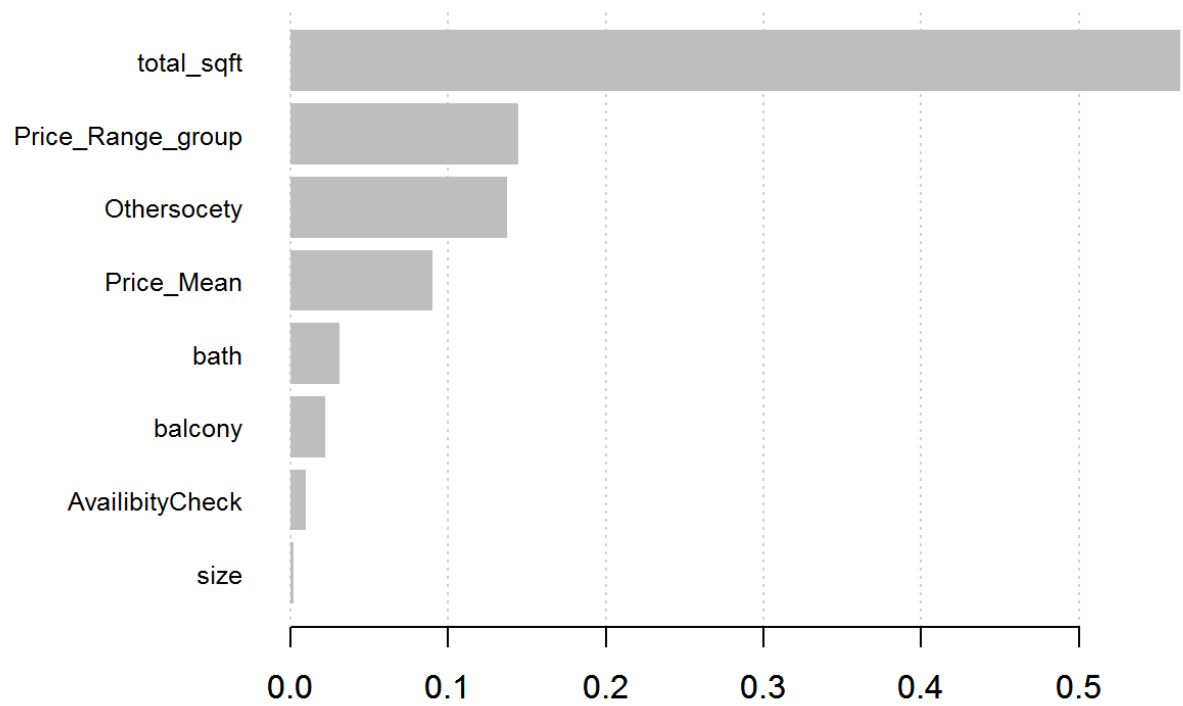
#Parameters

xgb_params<-list(objective="reg:linear",eval_metric="rmse",eta=0.1,max_depth=5,lamb
da=2)
watchlist<-list(train=train_matrix,test=test_matrix)
modelXGB<-xgb.train(params = xgb_params,data = train_matrix,nrounds = 500,watchlist
= watchlist,maximize = FALSE,early_stopping_rounds = 40,print_every_n = 5)
```

```
## [1] train-rmse:163.034042 test-rmse:191.354691
## Multiple eval metrics are present. Will use test_rmse for early stopping.
## Will train until test_rmse hasn't improved in 40 rounds.
##
## [6] train-rmse:111.890579 test-rmse:139.970642
## [11] train-rmse:84.500854 test-rmse:110.927734
## [16] train-rmse:69.635323 test-rmse:96.514130
## [21] train-rmse:61.843849 test-rmse:89.625183
## [26] train-rmse:56.752552 test-rmse:86.063339
## [31] train-rmse:53.623726 test-rmse:84.296738
## [36] train-rmse:51.643616 test-rmse:82.784615
## [41] train-rmse:50.196320 test-rmse:81.158119
## [46] train-rmse:49.092514 test-rmse:80.100334
## [51] train-rmse:48.401028 test-rmse:80.161987
## [56] train-rmse:47.862728 test-rmse:80.031670
## [61] train-rmse:47.390701 test-rmse:79.610474
## [66] train-rmse:46.699856 test-rmse:79.098457
## [71] train-rmse:46.372822 test-rmse:78.818298
## [76] train-rmse:45.838303 test-rmse:78.370789
## [81] train-rmse:45.599636 test-rmse:78.324188
## [86] train-rmse:45.358986 test-rmse:78.219849
## [91] train-rmse:44.973942 test-rmse:78.392731
## [96] train-rmse:44.648945 test-rmse:78.466042
## [101] train-rmse:44.298374 test-rmse:78.538651
## [106] train-rmse:43.959389 test-rmse:78.504784
## [111] train-rmse:43.453743 test-rmse:78.488358
## [116] train-rmse:42.806496 test-rmse:78.343857
## [121] train-rmse:42.465309 test-rmse:78.268715
## Stopping. Best iteration:
## [85] train-rmse:45.426285 test-rmse:78.179047
```

Feature important from XGBoost

```
predictors<-names(bangalore_train)[!names(bangalore_train) %in% c("price")]
gg<-xgb.plot.importance(xgb.importance(model = modelXGB,feature_names = predictors)
)
```



We get same important variable from XGBoost also.