

Metrics of Successful websites and companies

Danai Avratoglou

18 February 2017

```
#we upload the dataset
total_500 <- read.csv("~/GitHub/thesis_msc_business_analytics/Python/total_500_new.csv", sep=";", na.st
#we see how many observations and how many variables we have
dim(total_500)

## [1] 500 730

#We create a subset to make some changes to the data
total_500_sub <- total_500
#Change the decimal point for the 4 variables
total_500_sub$Assets.. <- gsub(",",".", total_500_sub$Assets.. )
total_500_sub$Market.value.. <- gsub(",",".", total_500_sub$Market.value.. )
total_500_sub$Revenues.. <- gsub(",",".", total_500_sub$Revenues.. )
total_500_sub$Total.Stockholder.Equity.. <- gsub(",",".", total_500_sub$Total.Stockholder.Equity.. )
#Make the variables numeric
for(i in 1:18){
  total_500_sub[,i] <- as.numeric(total_500_sub[,i])}

## Warning: NAs introduced by coercion

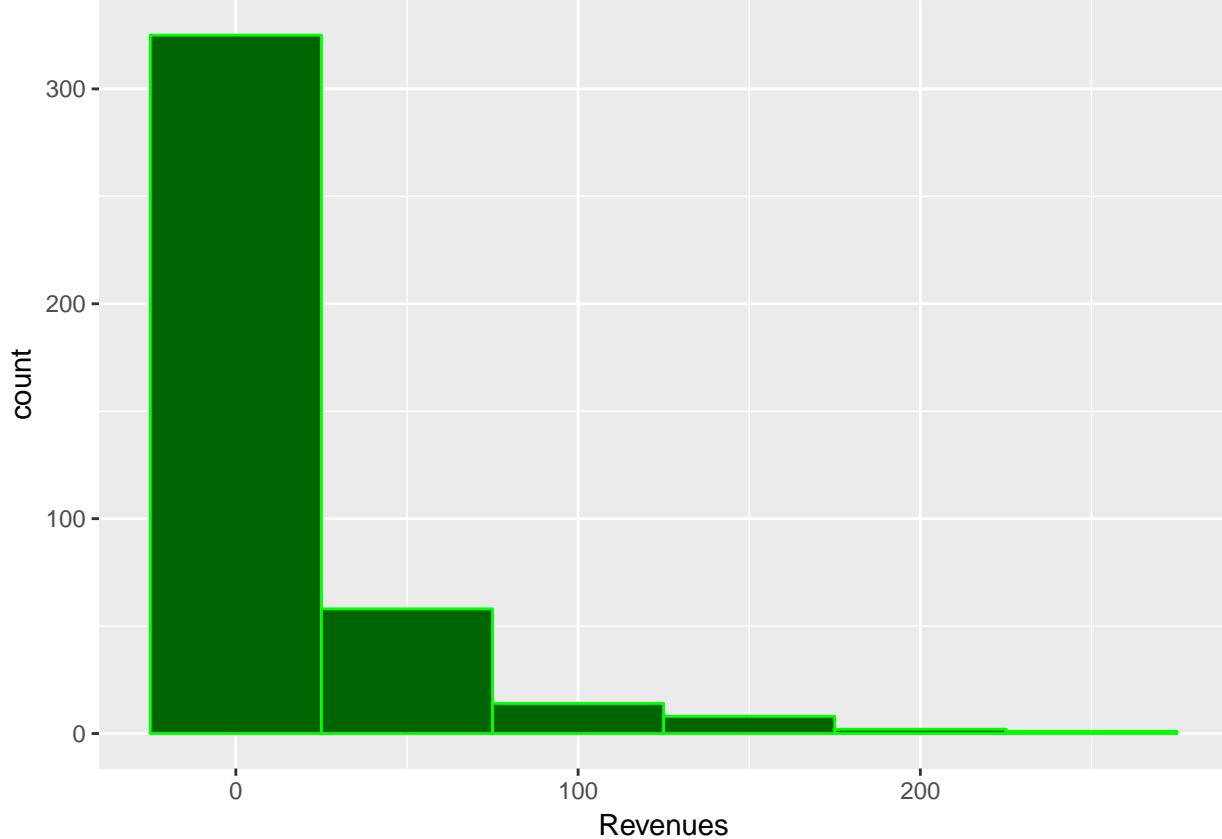
## Warning: NAs introduced by coercion

for(i in 20:730){
  total_500_sub[,i] <- as.numeric(total_500_sub[,i])}
#We omit the nas from the analysis
total_500_final <- na.omit(total_500_sub)
#We rename variable X as Ranking
colnames(total_500_final)[1] <- "Ranking"
#Change the names of some variables to be more easily readable
colnames(total_500_final)[2] <- "Assets"
colnames(total_500_final)[3] <- "Market_Value"
colnames(total_500_final)[4] <- "Revenues"
colnames(total_500_final)[6] <- "Total_SH_Equity"
#Delete the variables we will not need
total_500_final$Revenues...1 <- NULL #Revenues %
total_500_final$company <- NULL #company name
total_500_final$url<- NULL # company url
#we upload the libraries beneath that we will use in the analysis
library(ggplot2)
library(reshape2)
library(DAAG)

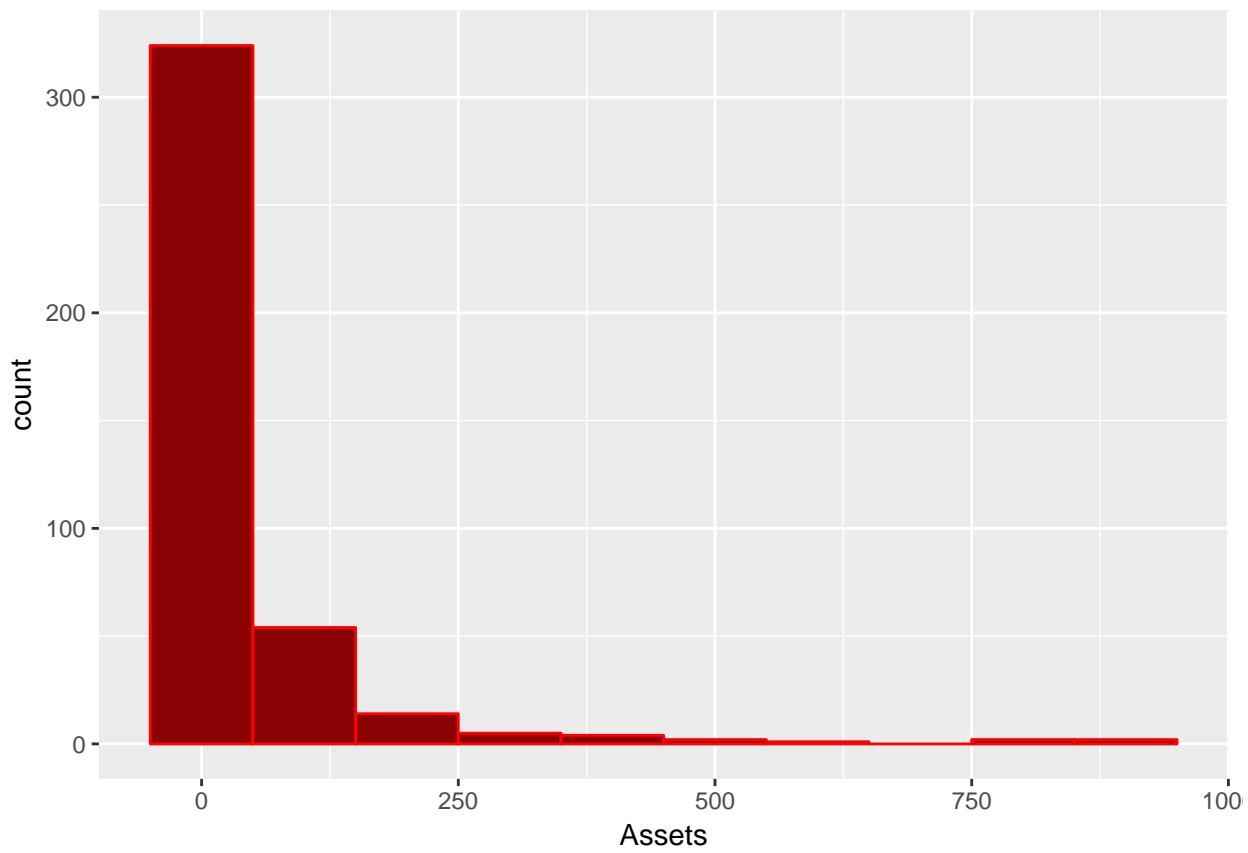
## Loading required package: lattice
#Final number of observation and variables we will use
dim(total_500_final)

## [1] 408 727
```

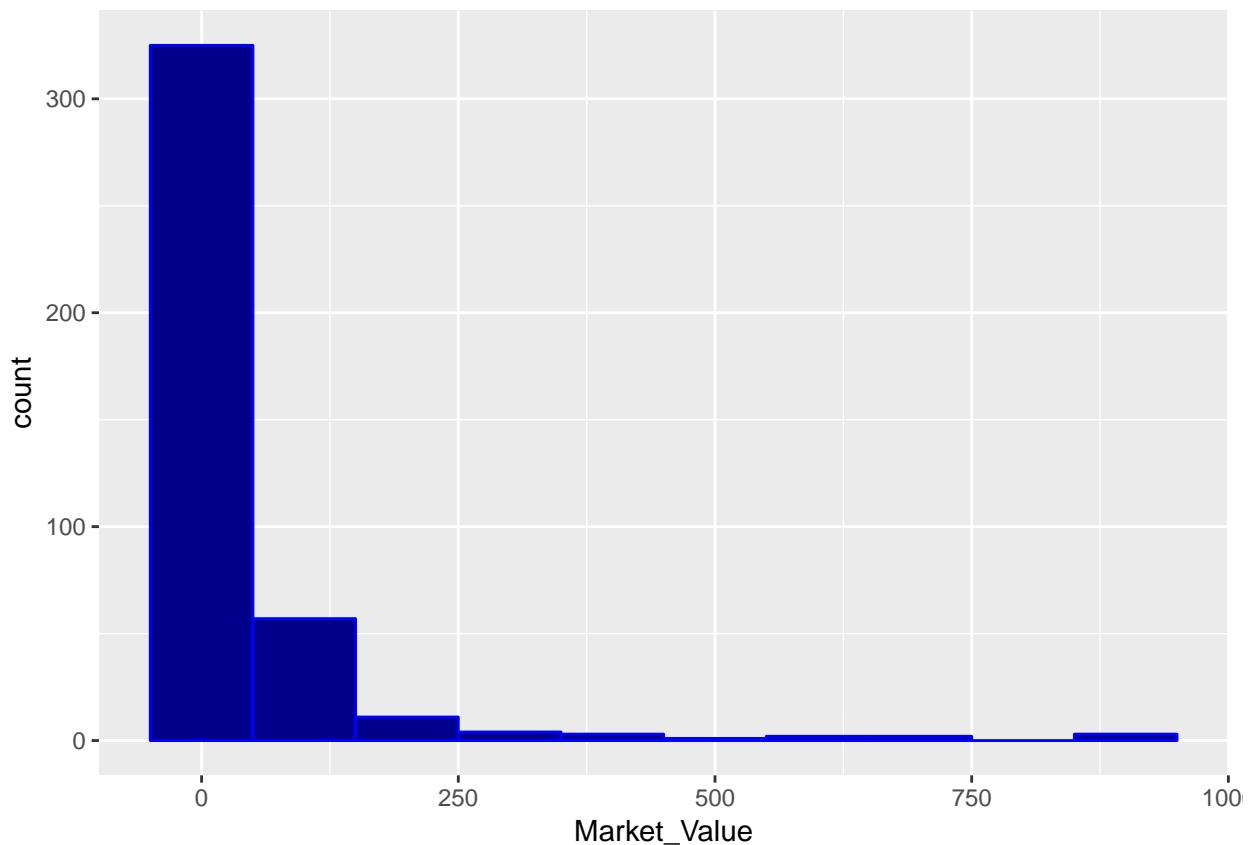
```
#####
#we first see the summary of the Fortune variables and then we create their histogram so as to have a
#good grasp of how they are distributed
ggplot(data=total_500_final,aes(x=Revenues))+geom_histogram(binwidth=50, colour = "green", fill ="darkg
```



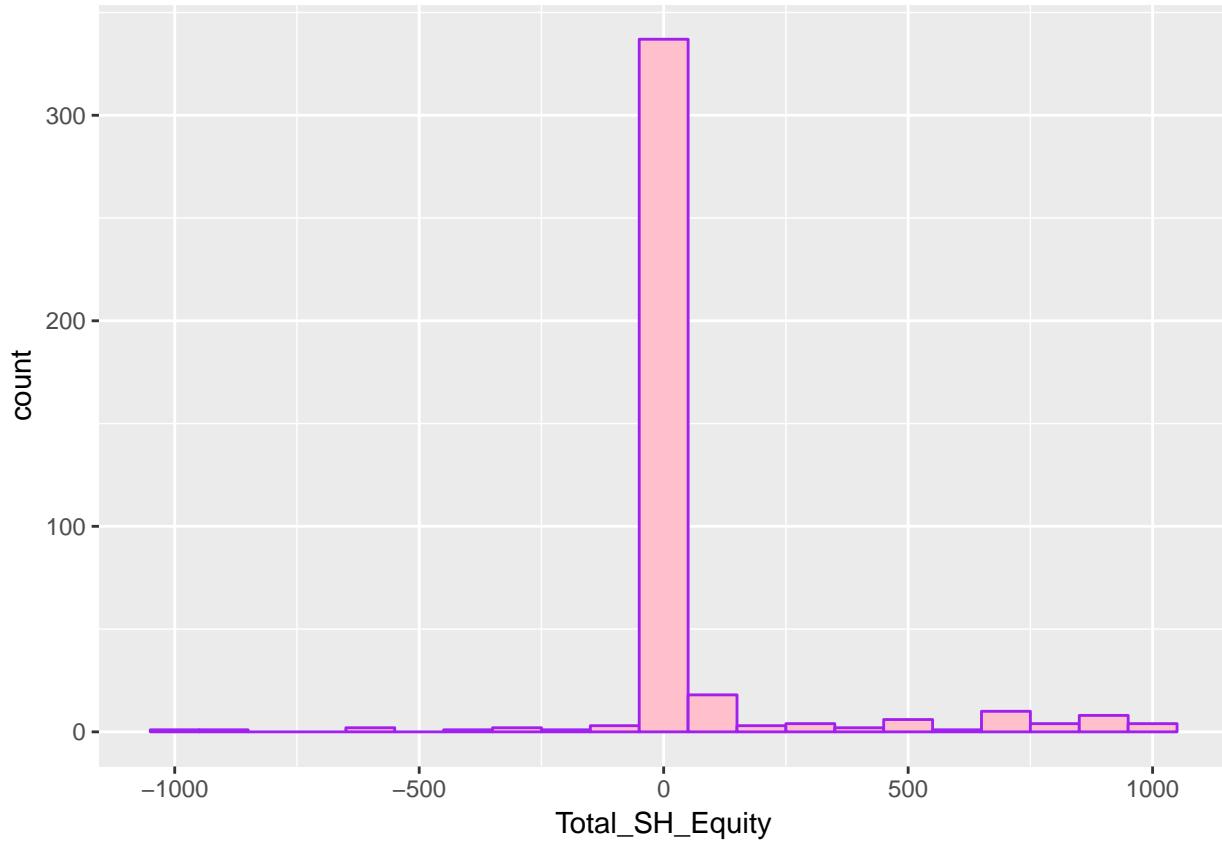
```
ggplot(data=total_500_final,aes(x=Assets))+geom_histogram(binwidth=100, colour = "red", fill ="darkred")
```



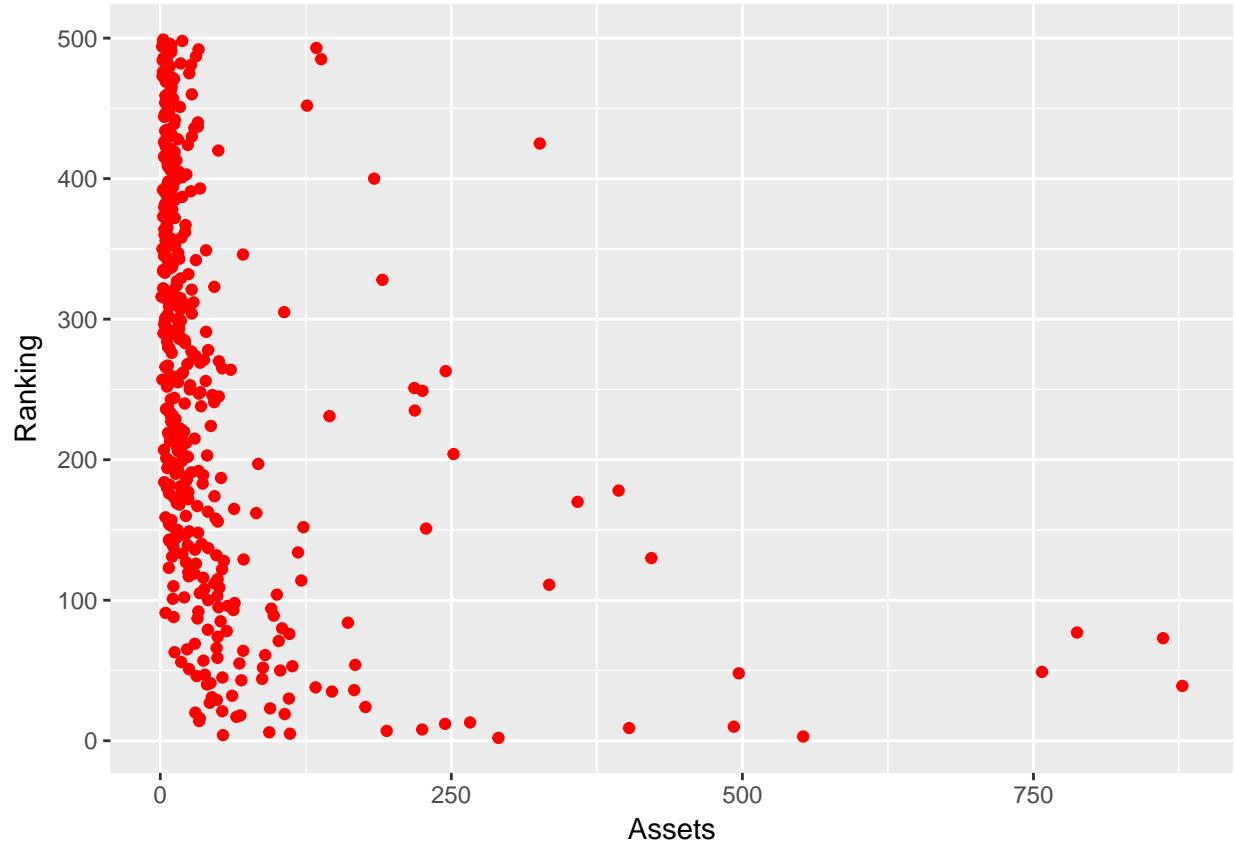
```
ggplot(data=total_500_final,aes(x=Market_Value))+geom_histogram(binwidth=100, colour = "blue", fill = "darkred")
```



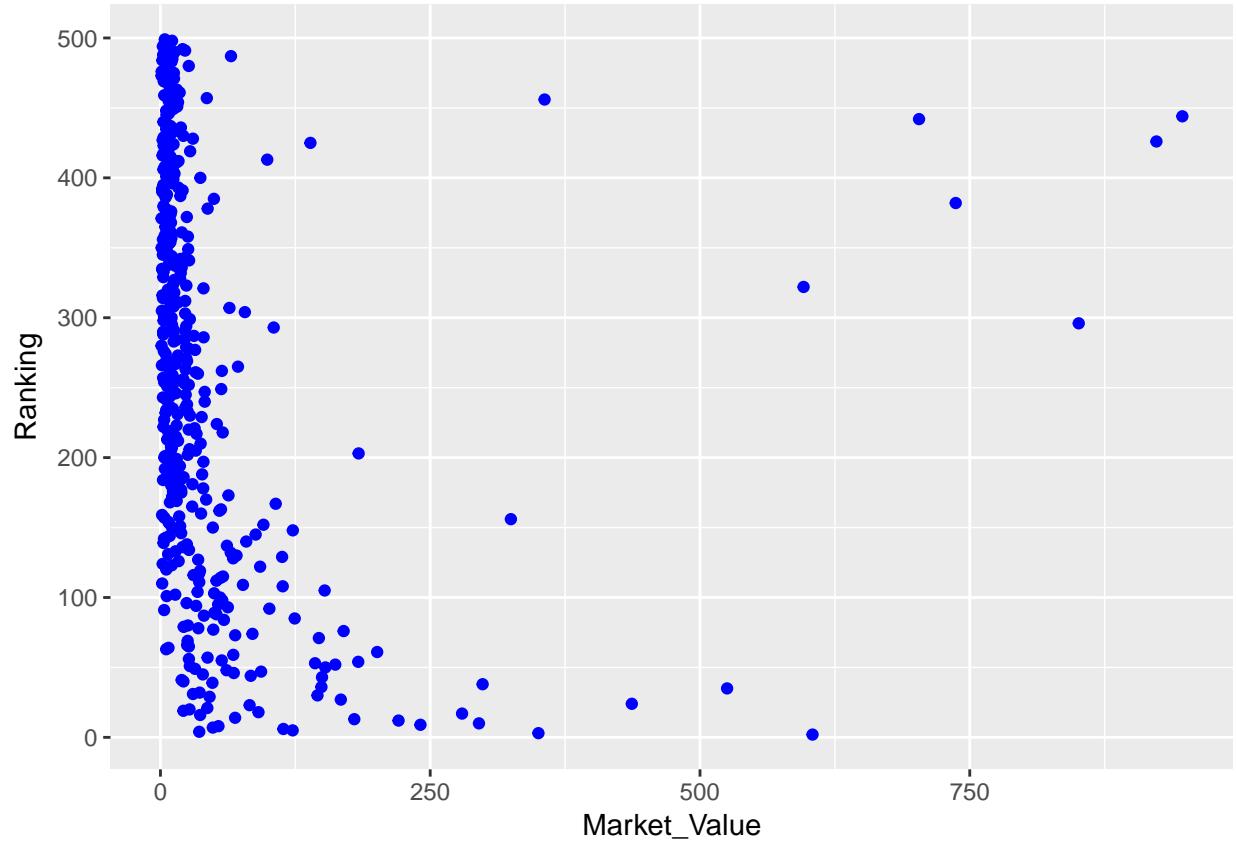
```
ggplot(data=total_500_final,aes(x=Total_SH_Equity))+geom_histogram(binwidth=100, colour = "purple", fill
```

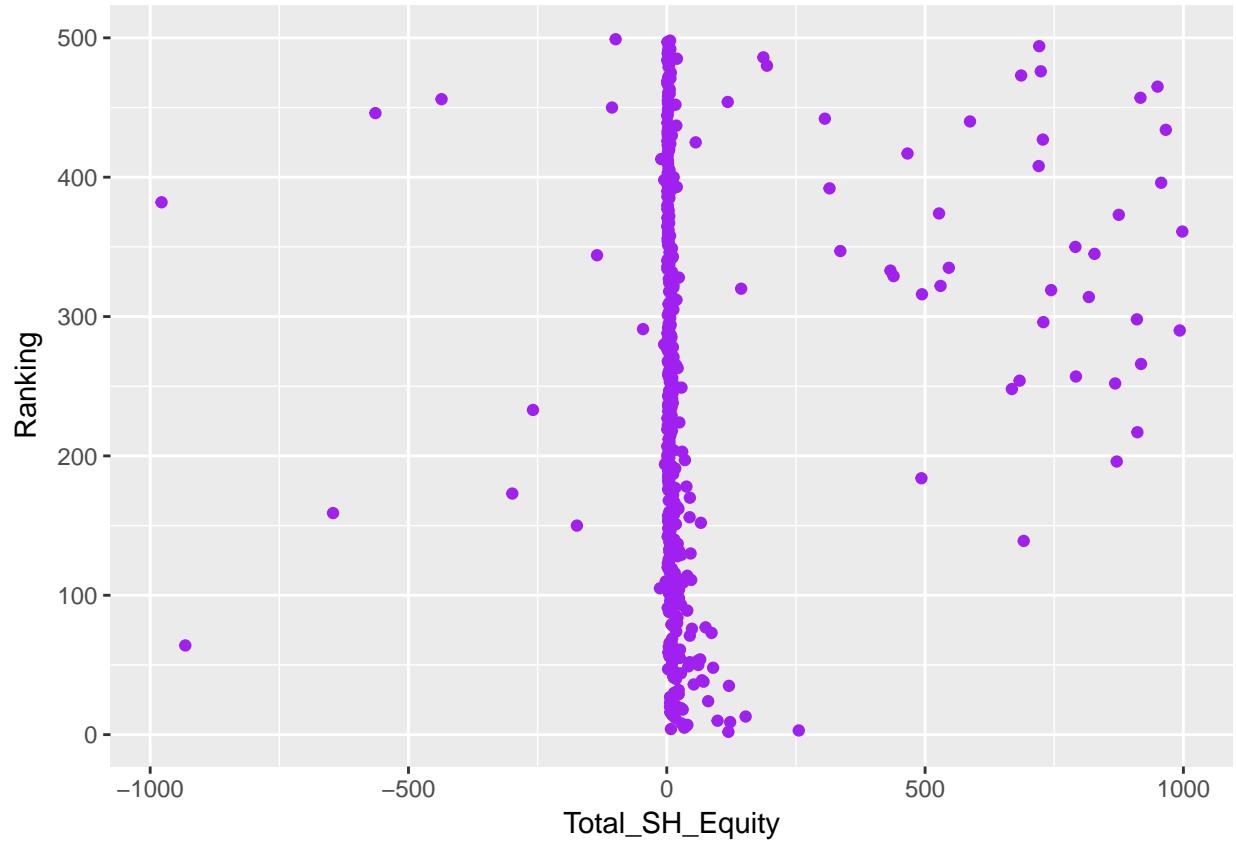


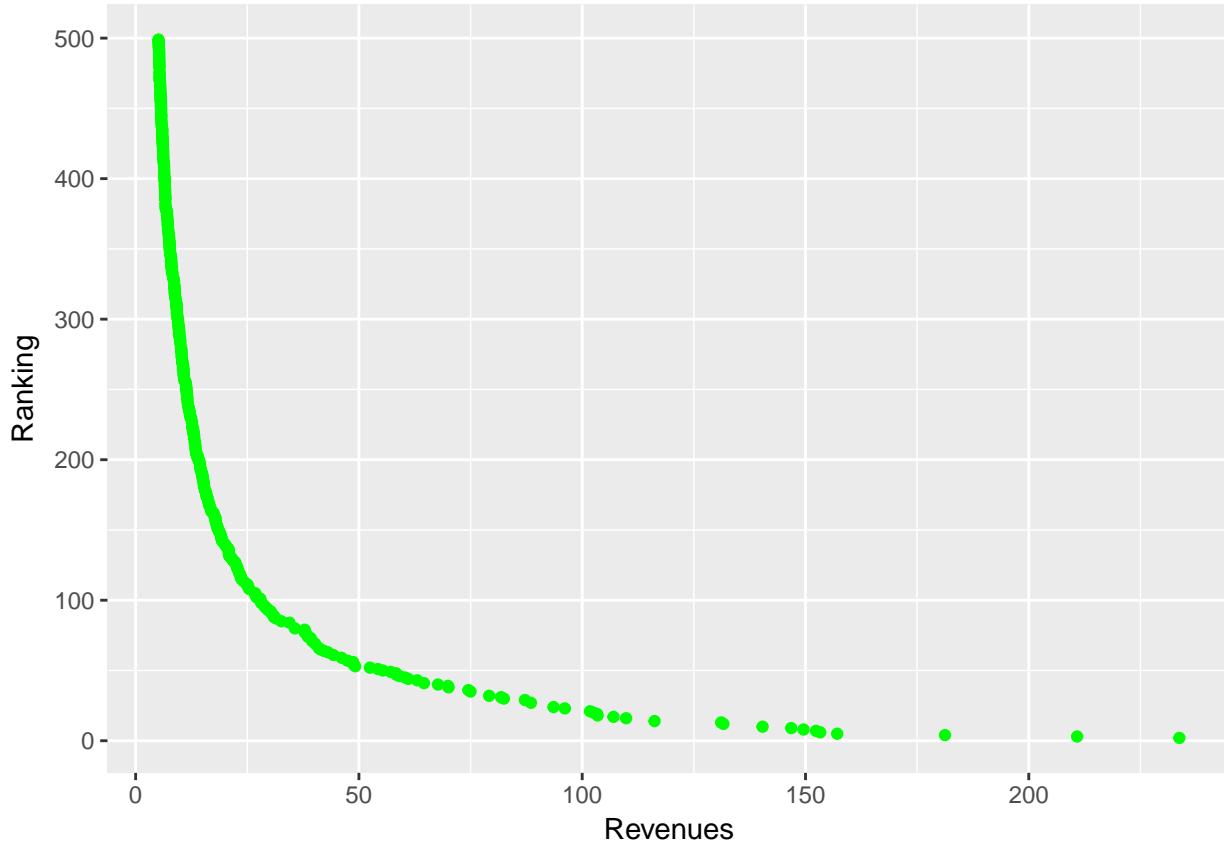
```
#####
#We make plots to see how the variables we got from Fortune 500 are related with the Ranking
ggplot(total_500_final, aes(Assets,Ranking)) + geom_point(colour = "red")
```



```
ggplot(total_500_final, aes(Market_Value, Ranking)) + geom_point(colour = "blue")
```

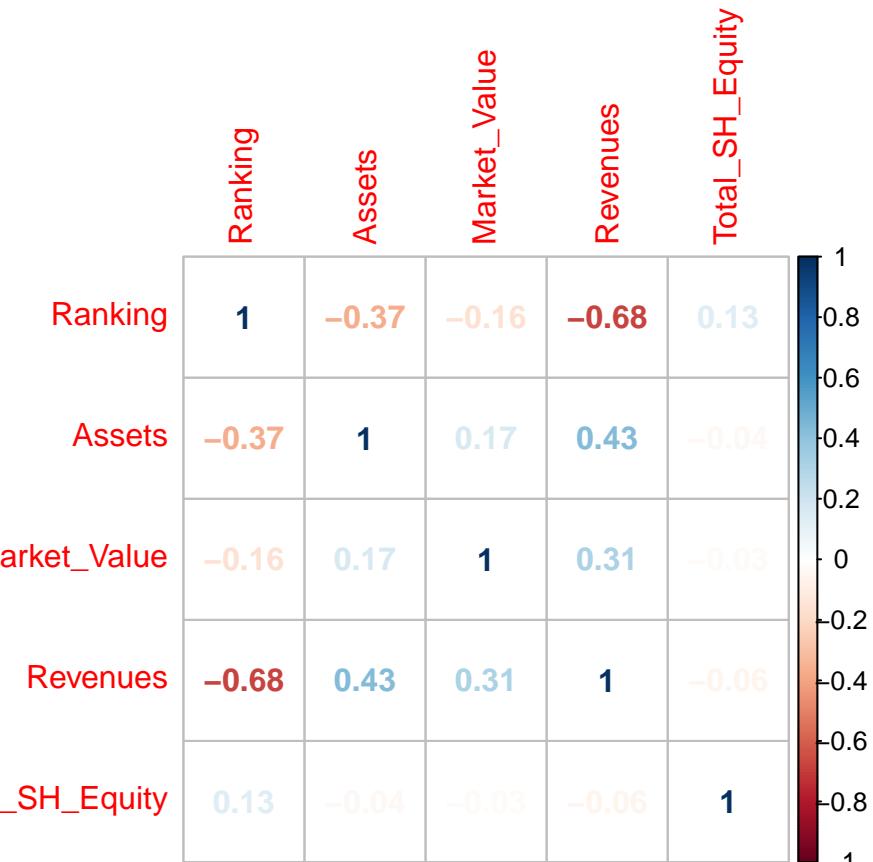






```
#We can see that the Ranking has a linear relationship with the Revenues so we will use one of those 2
#In order to have a more clear look we also create a correlation diagram
total_500_fortune <- total_500_final[,c(1:5)]
library(corrplot)
library(caret)
sm <- cor(total_500_fortune)
sm

##          Ranking      Assets Market_Value   Revenues
## Ranking 1.0000000 -0.36673307 -0.15959008 -0.67511457
## Assets -0.3667331  1.00000000  0.16787320  0.43479882
## Market_Value -0.1595901  0.16787320  1.00000000  0.31085660
## Revenues -0.6751146  0.43479882  0.31085660  1.00000000
## Total_SH_Equity 0.1327272 -0.03638159 -0.02912268 -0.05616772
##          Total_SH_Equity
## Ranking          0.13272724
## Assets          -0.03638159
## Market_Value     -0.02912268
## Revenues         -0.05616772
## Total_SH_Equity  1.00000000
corrplot(cor(total_500_fortune),method="number")
```



#From this plot we understand that the Ranking and the Revenues have very high correlation.

```
#Firstly we will analyze the social media relevance with the sites.
```

```
#We will see how many of the sites have social media and what type of social media
#Facebook
```

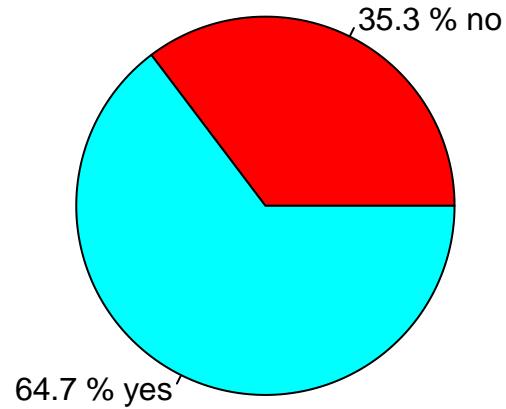
```
social_media_facebook <- round(table(total_500_final$facebook)/408,3)
social_media_facebook
```

```
##
##      0      1
## 0.353 0.647
```

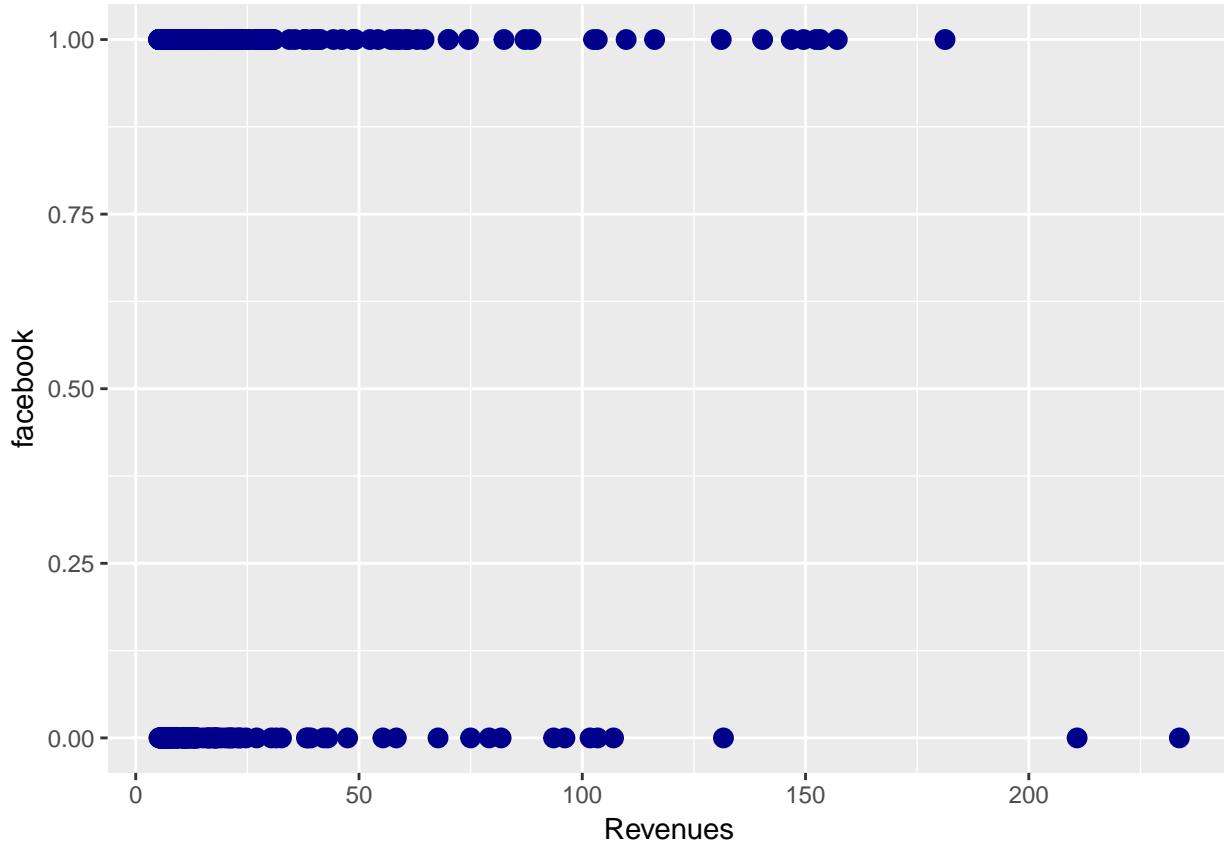
```
slicelable <- c(paste(35.3,"% no"),paste(64.7,"% yes"))
```

```
pie(social_media_facebook,label = slicelable,main="Share of companies with Facebook",col=rainbow(length
```

Share of companies with Facebook



```
ggplot(total_500_final, aes(Revenues, facebook)) + geom_point(size=3, colour = "darkblue")
```



#Twitter

```
social_media_twitter <- round(table(total_500_final$twitter)/408,3)
social_media_twitter
```

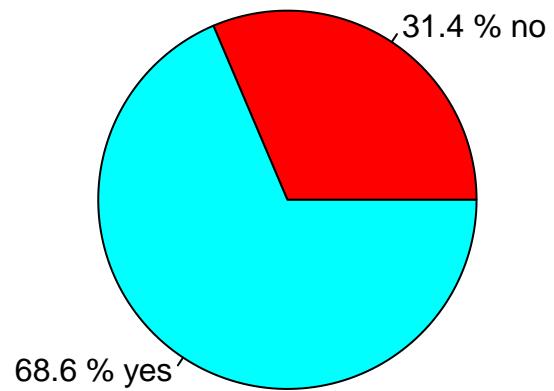
##

```
##      0      1
## 0.314 0.686
```

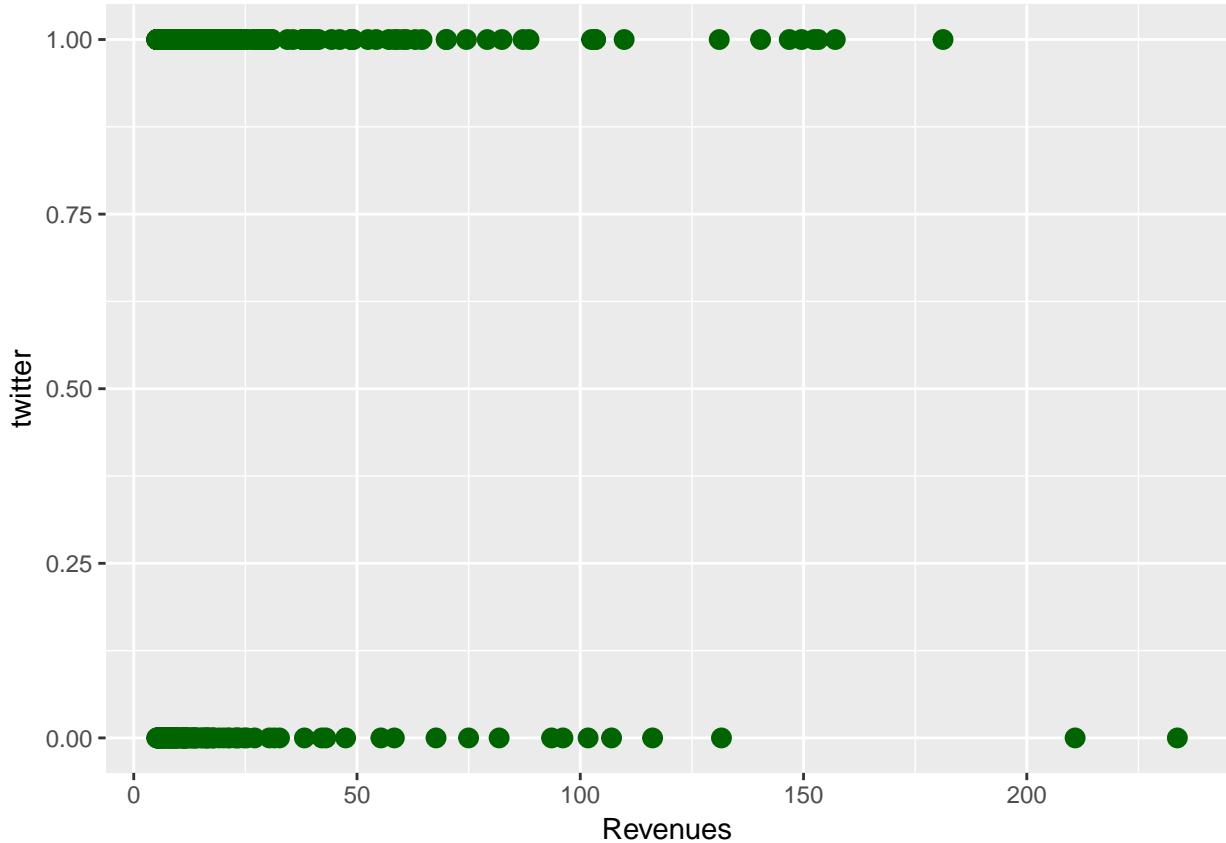
```
slicelable <- c(paste(31.4,"% no"),paste(68.6,"% yes"))
```

```
pie(social_media_twitter,label = slicelable,main="Share of companies with Twitter",col=rainbow(length(slicelable)))
```

Share of companies with Twitter



```
ggplot(total_500_final, aes(Revenues, twitter)) + geom_point(size=3, colour = "darkgreen")
```

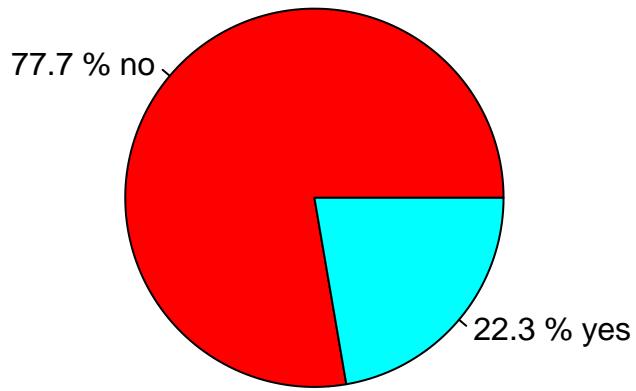


```
#Instagram
social_media_instagram <- round(table(total_500_final$instagram)/408,3)
social_media_instagram

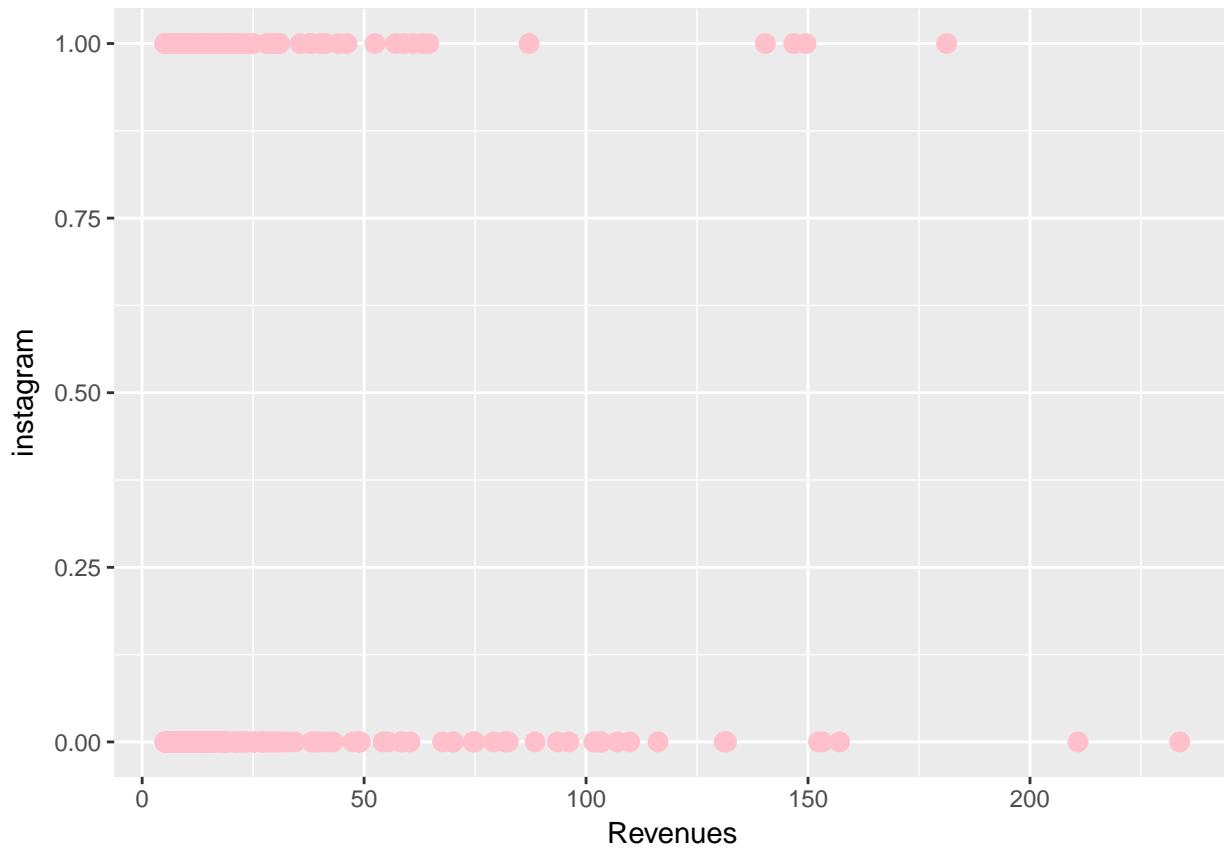
##
##      0      1
## 0.777 0.223

slicelable <- c(paste(77.7,"% no"),paste(22.3,"% yes"))
pie(social_media_instagram,label = slicelable,main="Share of companies with Instagram",col=rainbow(length(slicelable)))
```

Share of companies with Instagram



```
ggplot(total_500_final, aes(Revenues, instagram)) + geom_point(size=3, colour = "pink")
```

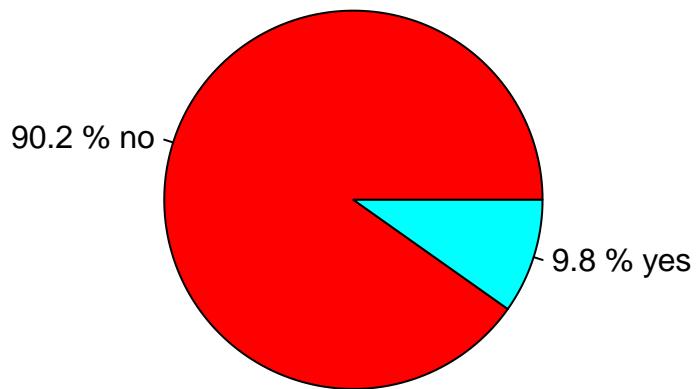


```
#Pinterest
social_media_pinterest <- round(table(total_500_final$pinterest)/408,3)
social_media_pinterest

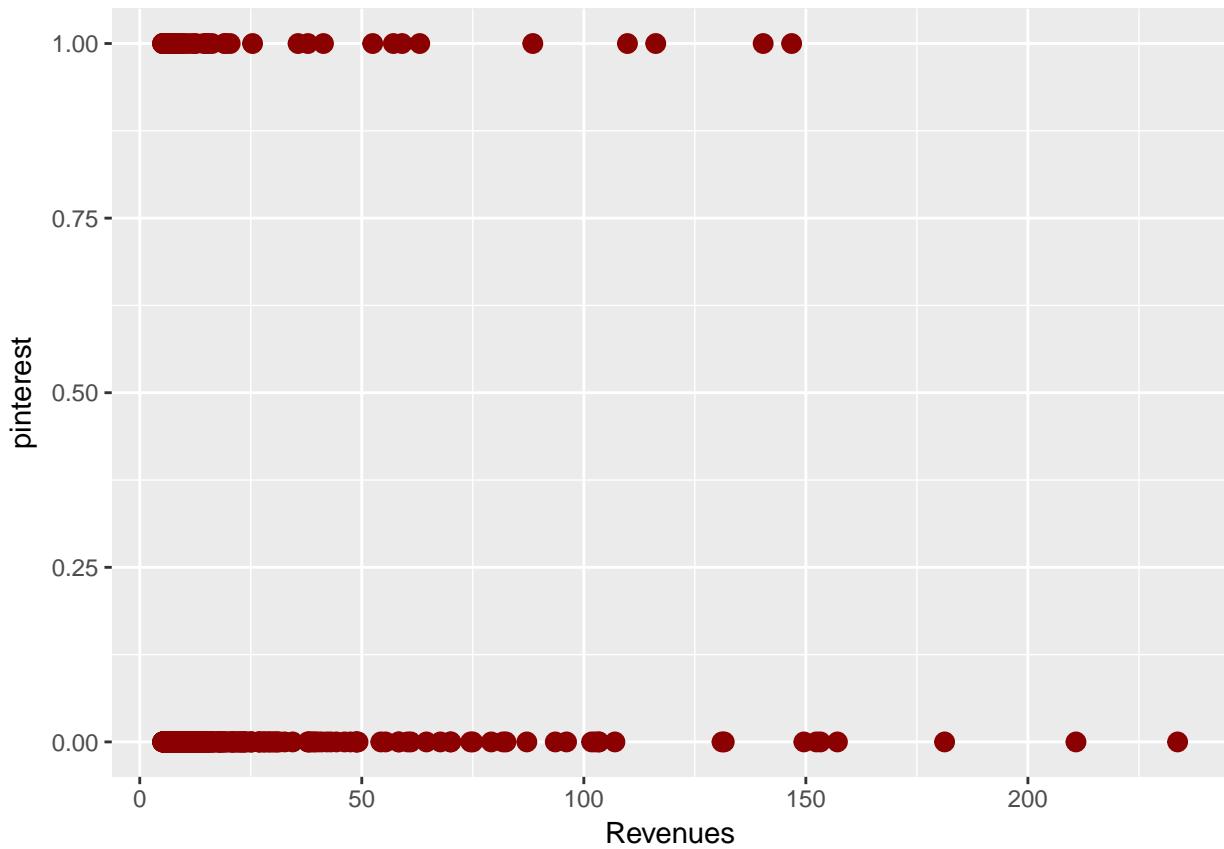
##
##      0      1
## 0.902 0.098

slicelable <- c(paste(90.2,"% no"),paste(9.8,"% yes"))
pie(social_media_pinterest,label = slicelable,main="Share of companies with Pinterest",col=rainbow(leng
```

Share of companies with Pinterest



```
ggplot(total_500_final, aes(Revenues, pinterest)) + geom_point(size=3, colour = "darkred")
```



```
#Youtube
```

```
social_media_youtube <- round(table(total_500_final$youtube)/408,3)
social_media_youtube
```

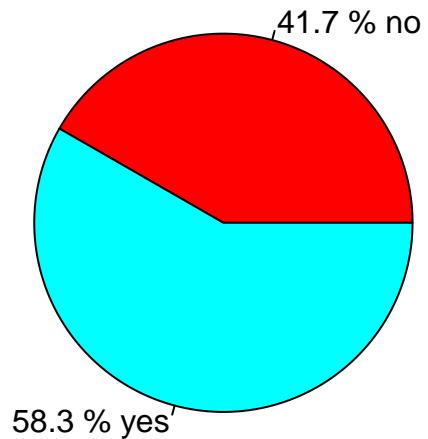
```
##
```

```
##      0      1
## 0.417 0.583
```

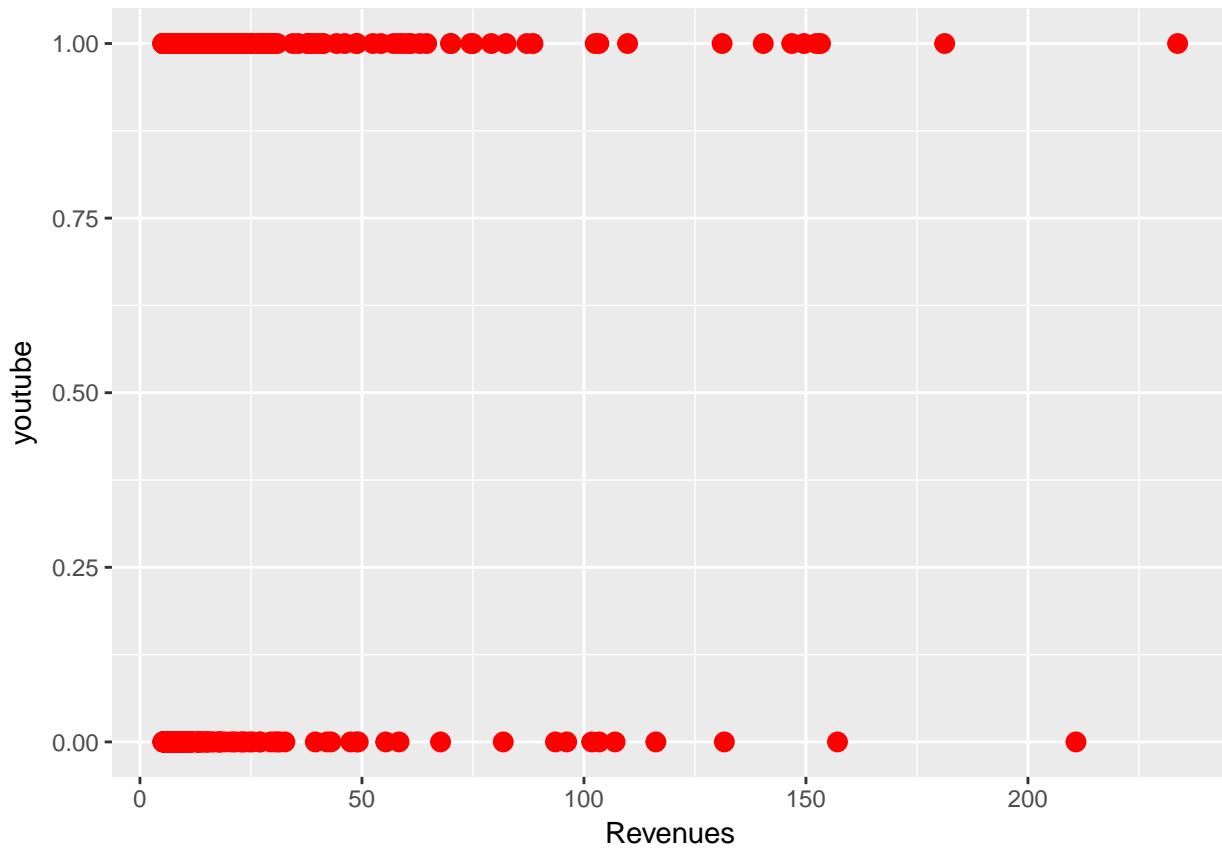
```
slicelable <- c(paste(41.7,"% no"),paste(58.3,"% yes"))
```

```
pie(social_media_youtube,label = slicelable,main="Share of companies with Youtube",col=rainbow(length(slicelable)))
```

Share of companies with Youtube



```
ggplot(total_500_final, aes(Revenues, youtube)) + geom_point(size=3, colour = "red")
```



#LinkedIn

```
social_media_linkedin <- round(table(total_500_final$linkedin)/408,3)
social_media_linkedin
```

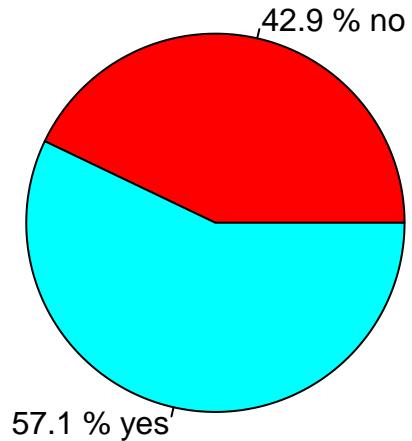
##

```
##      0      1
## 0.429 0.571
```

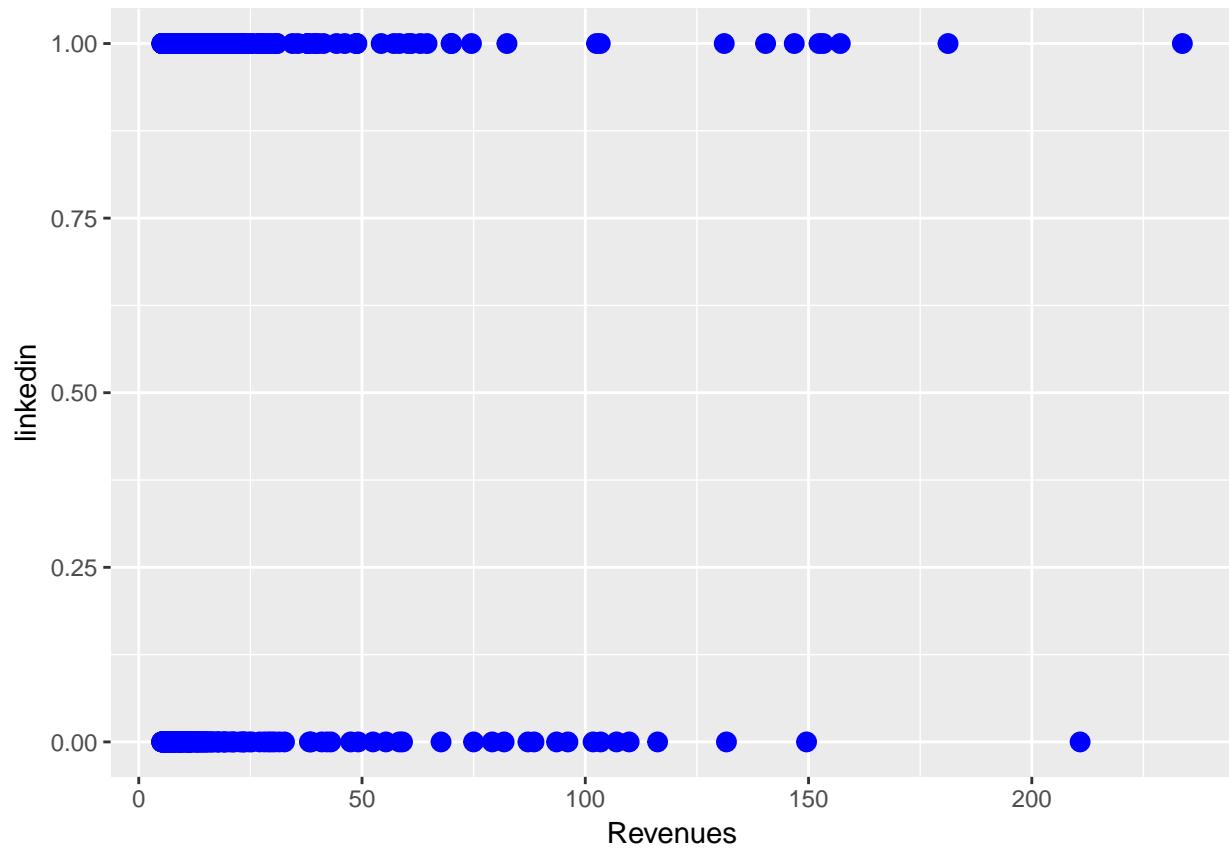
```
slicelable <- c(paste(42.9,"% no"),paste(57.1,"% yes"))
```

```
pie(social_media_linkedin,label = slicelable,main="Share of companies with LinkedIn",col=rainbow(length
```

Share of companies with LinkedIn



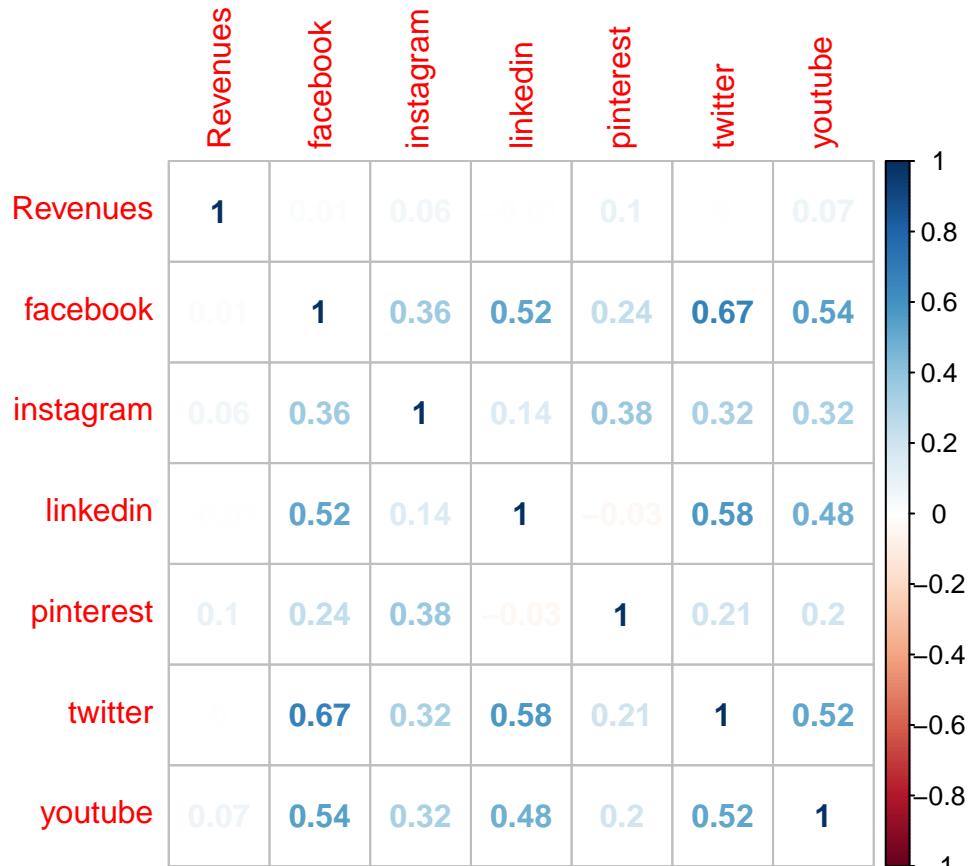
```
ggplot(total_500_final, aes(Revenues, linkedin)) + geom_point(size=3, colour = "blue")
```



```
#And we can also see for correlations
total_500_social_media <- total_500_final[,c(4,10:15)]
library(corrplot)
library(caret)
sm <- cor(total_500_social_media)
sm

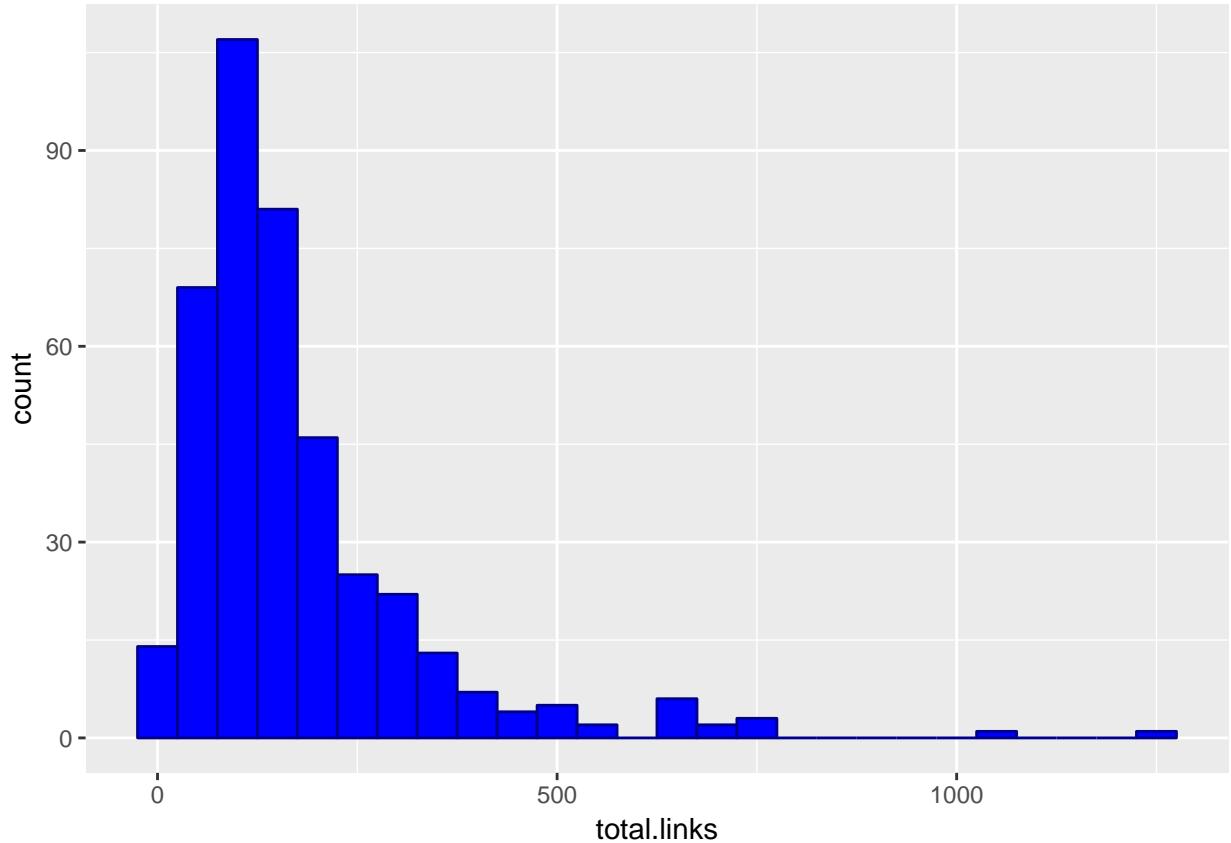
##           Revenues   facebook  instagram   linkedin   pinterest
## Revenues  1.000000000 0.01121852 0.05771665 -0.008311532 0.09686843
## facebook  0.011218524 1.00000000 0.35874256  0.520581725 0.24349238
## instagram 0.057716654 0.35874256 1.00000000  0.143134960 0.37774489
## linkedin -0.008311532 0.52058172 0.14313496  1.000000000 -0.03069495
## pinterest 0.096868426 0.24349238 0.37774489 -0.030694951 1.00000000
## twitter   0.002185367 0.67230226 0.32419034  0.577378625 0.20514804
## youtube   0.074833925 0.54096275 0.32145351  0.482997415 0.19504737
##           twitter   youtube
## Revenues  0.002185367 0.07483393
## facebook  0.672302259 0.54096275
## instagram 0.324190344 0.32145351
## linkedin  0.577378625 0.48299741
## pinterest 0.205148042 0.19504737
## twitter   1.000000000 0.52142857
## youtube   0.521428571 1.00000000

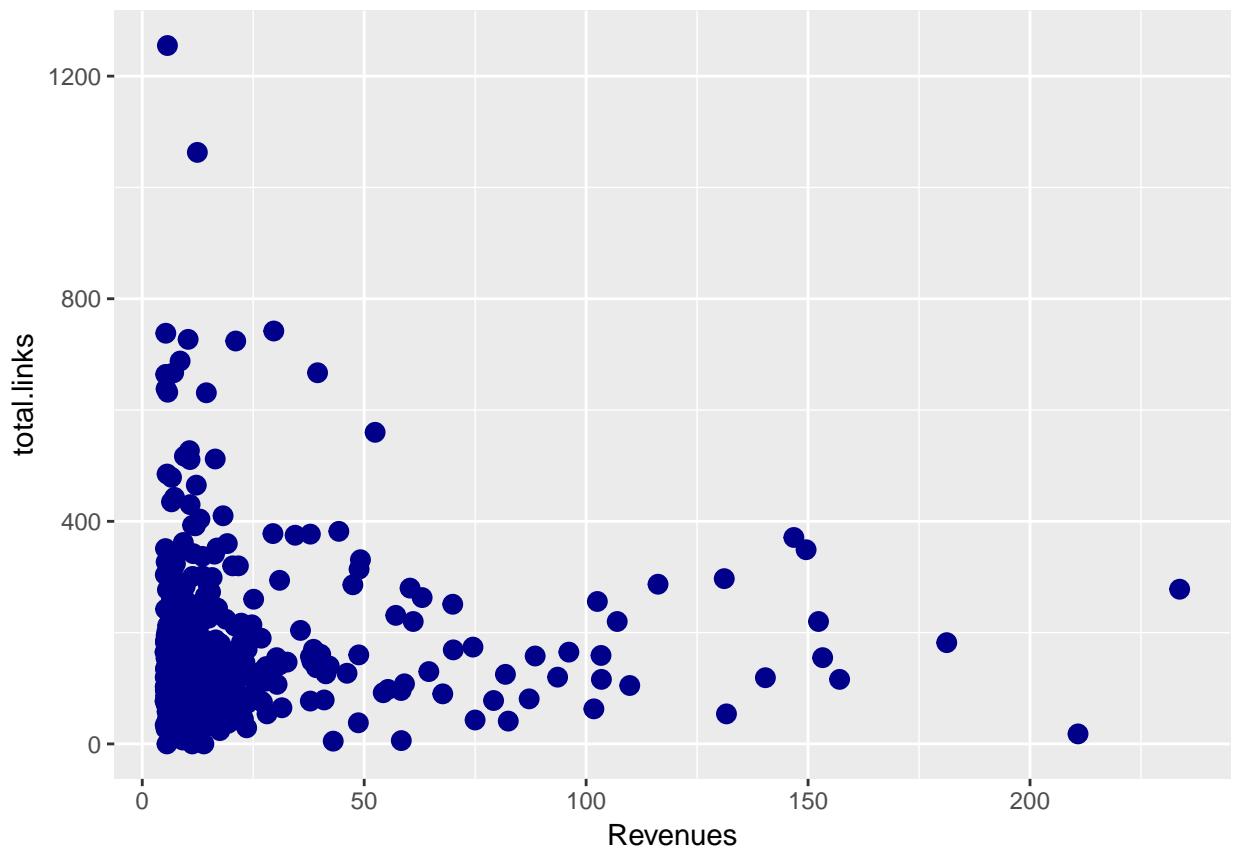
corrplot(cor(total_500_social_media),method="number")
```



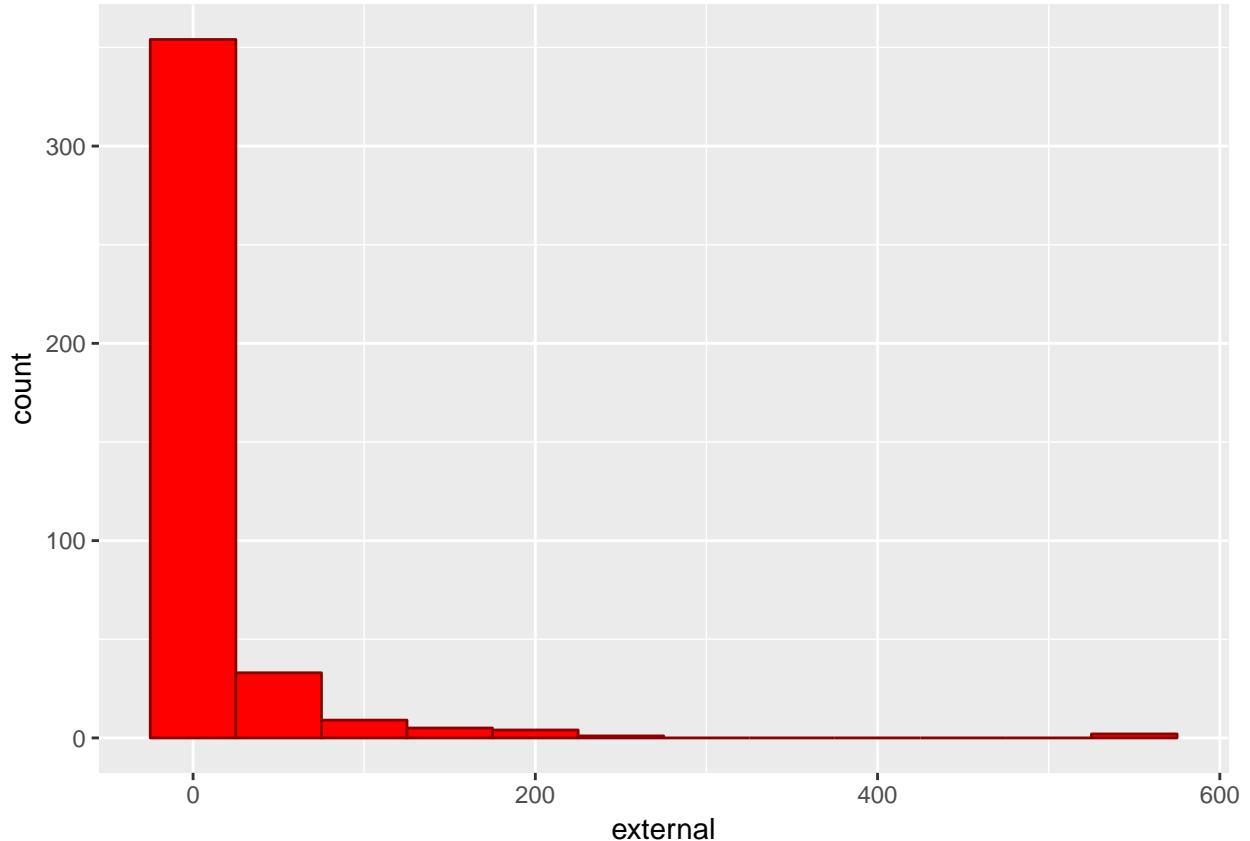
```
#we see that facebook has correlation more than 50% with twitter, youtube and linkedin
#And that the smallest correlations are those of pinterest and instagram
```

```
#We will now check the links by creating an histogram
#Then we create ggplots in order to see in what frequency the links appear
par(mfrow=c(1,1))
library(ggplot2)
ggplot(data=total_500_final,aes(x=total.links))+geom_histogram(binwidth=50, colour = "darkblue", fill =
```

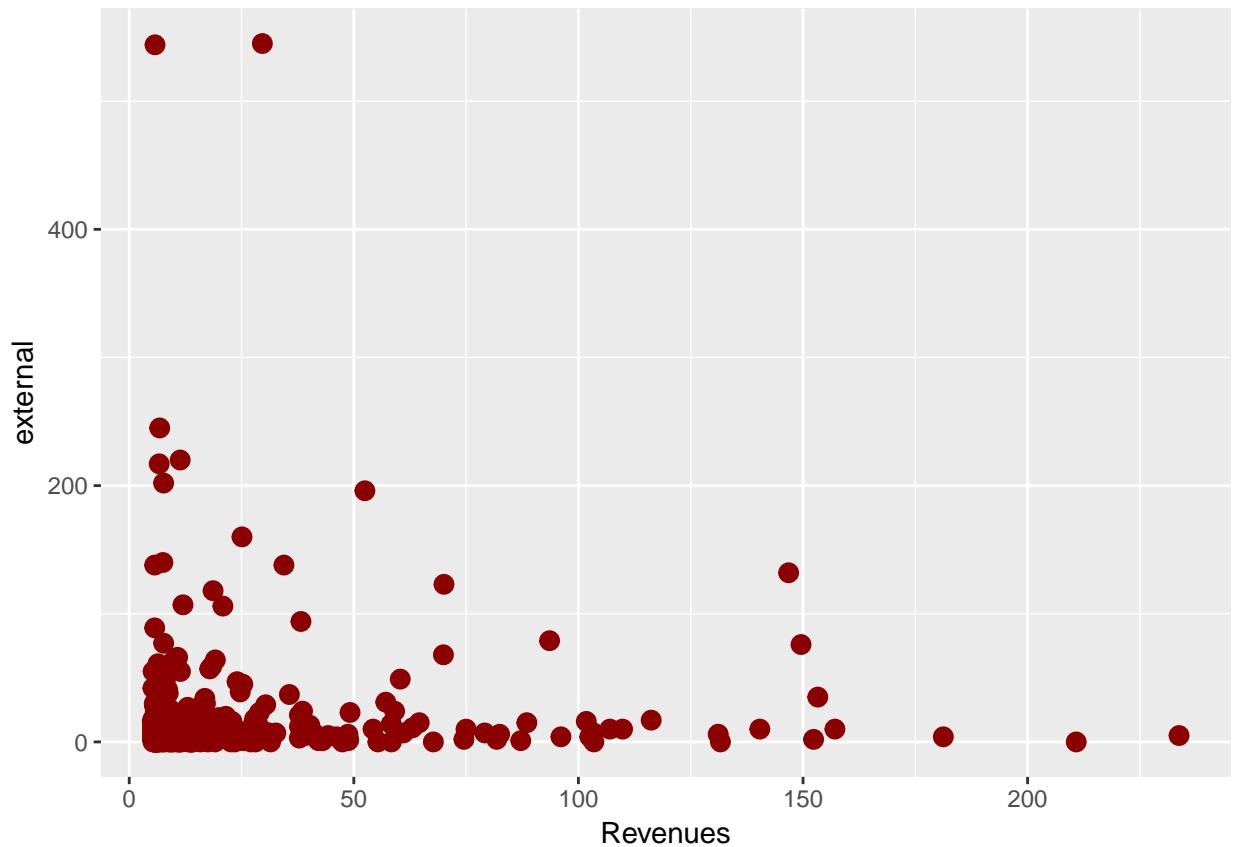


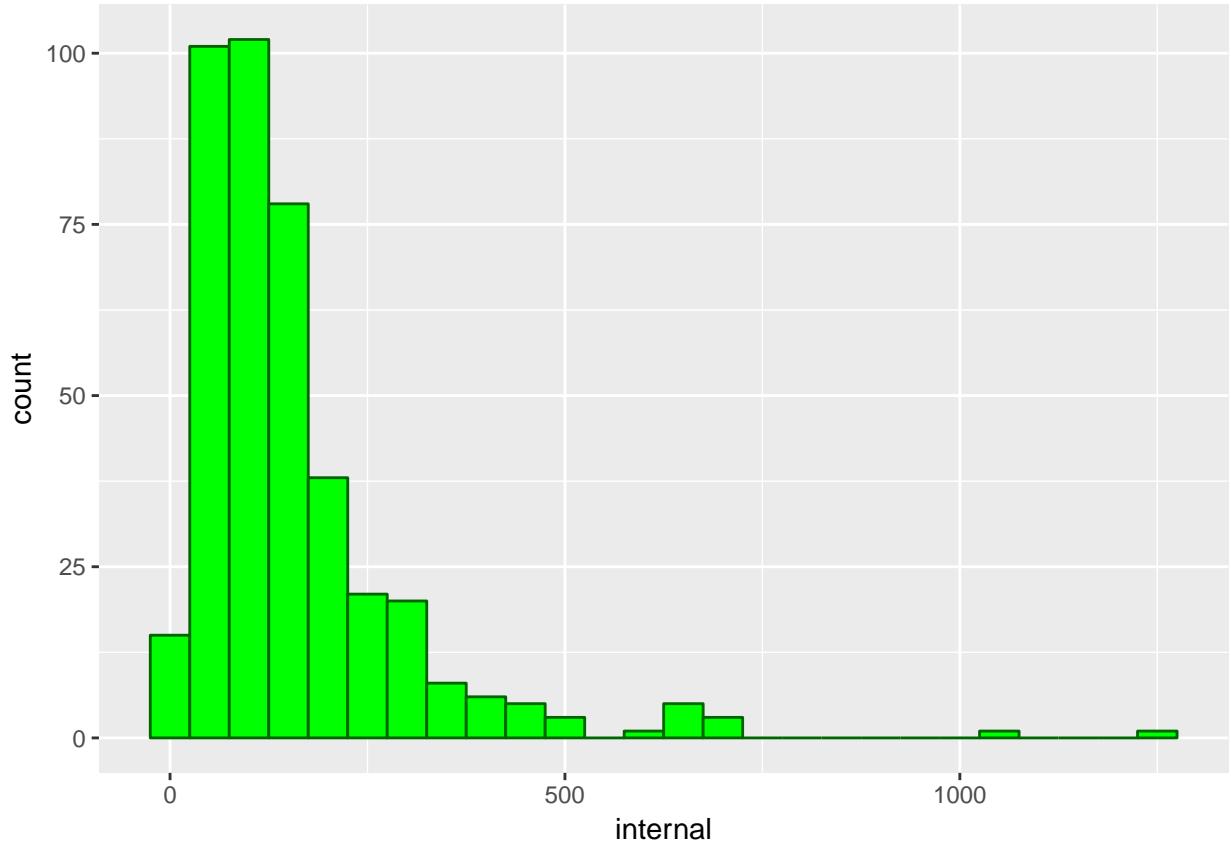


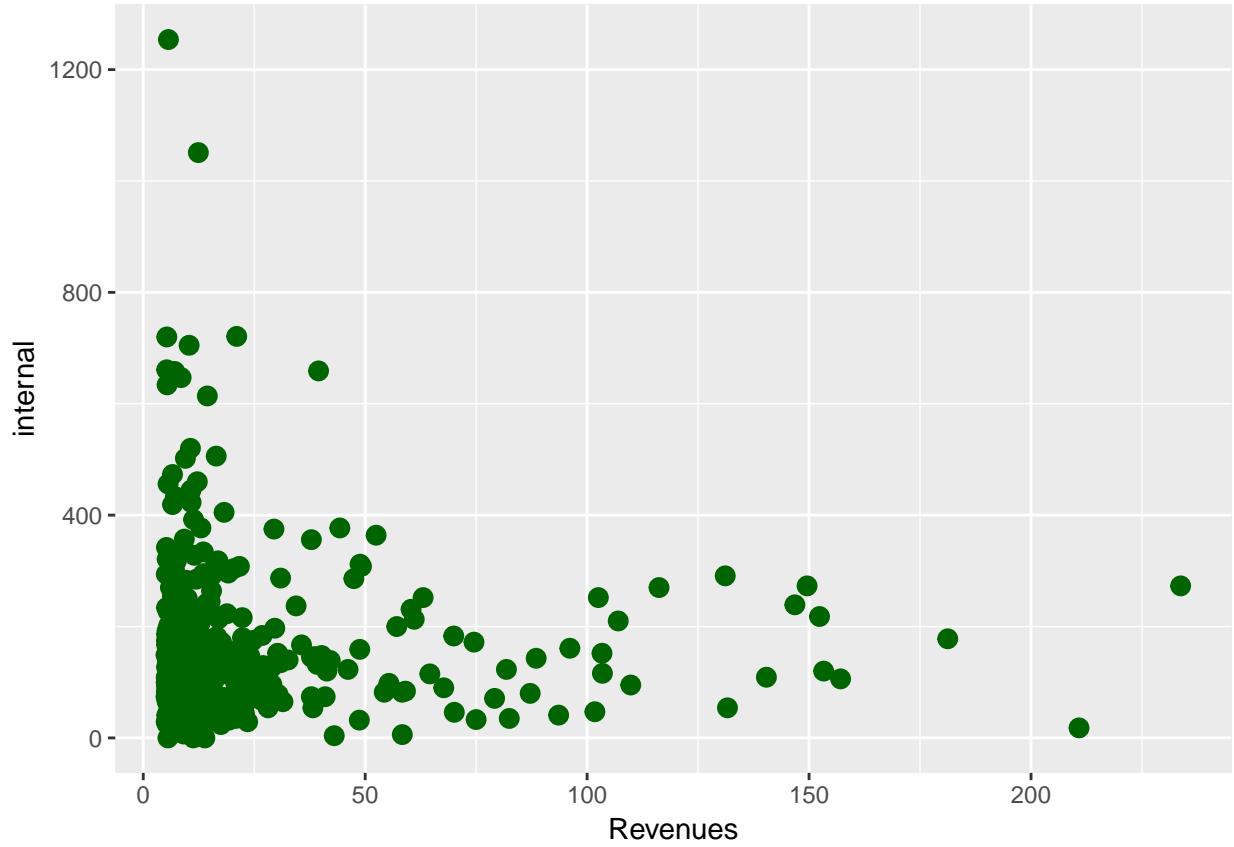
```
ggplot(data=total_500_final,aes(x=external))+geom_histogram(binwidth=50, colour = "darkred", fill ="red")
```



```
ggplot(total_500_final, aes(Revenues, external)) + geom_point(size=3, colour = "darkred")
```

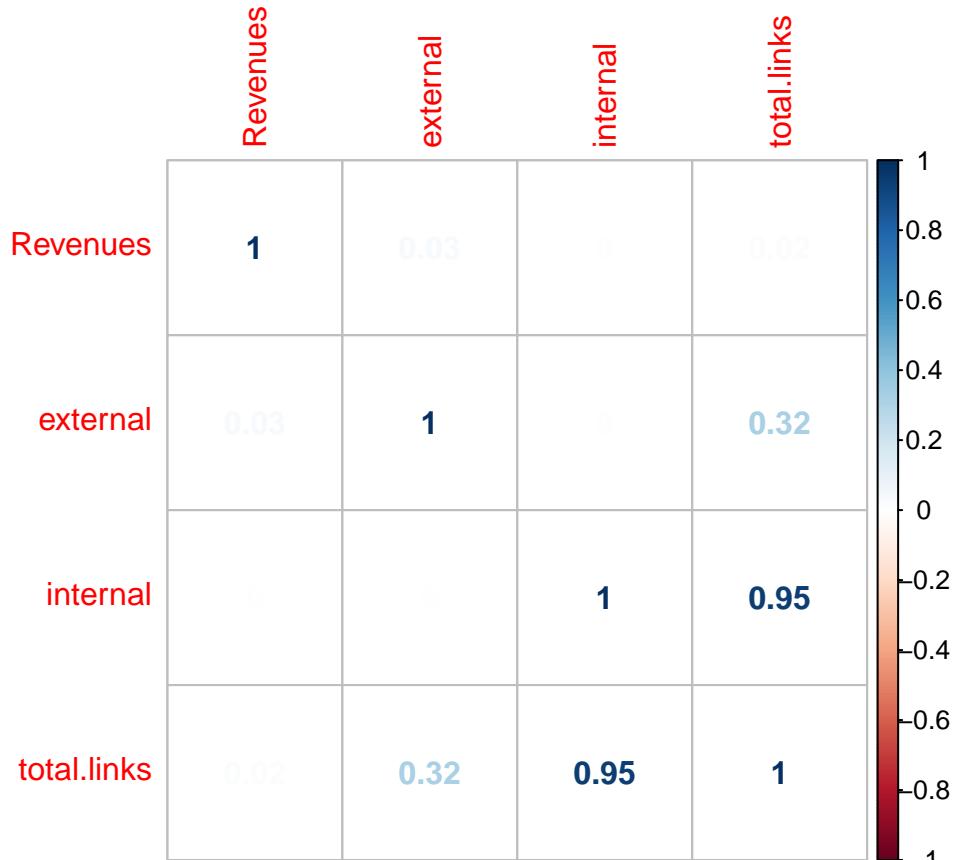






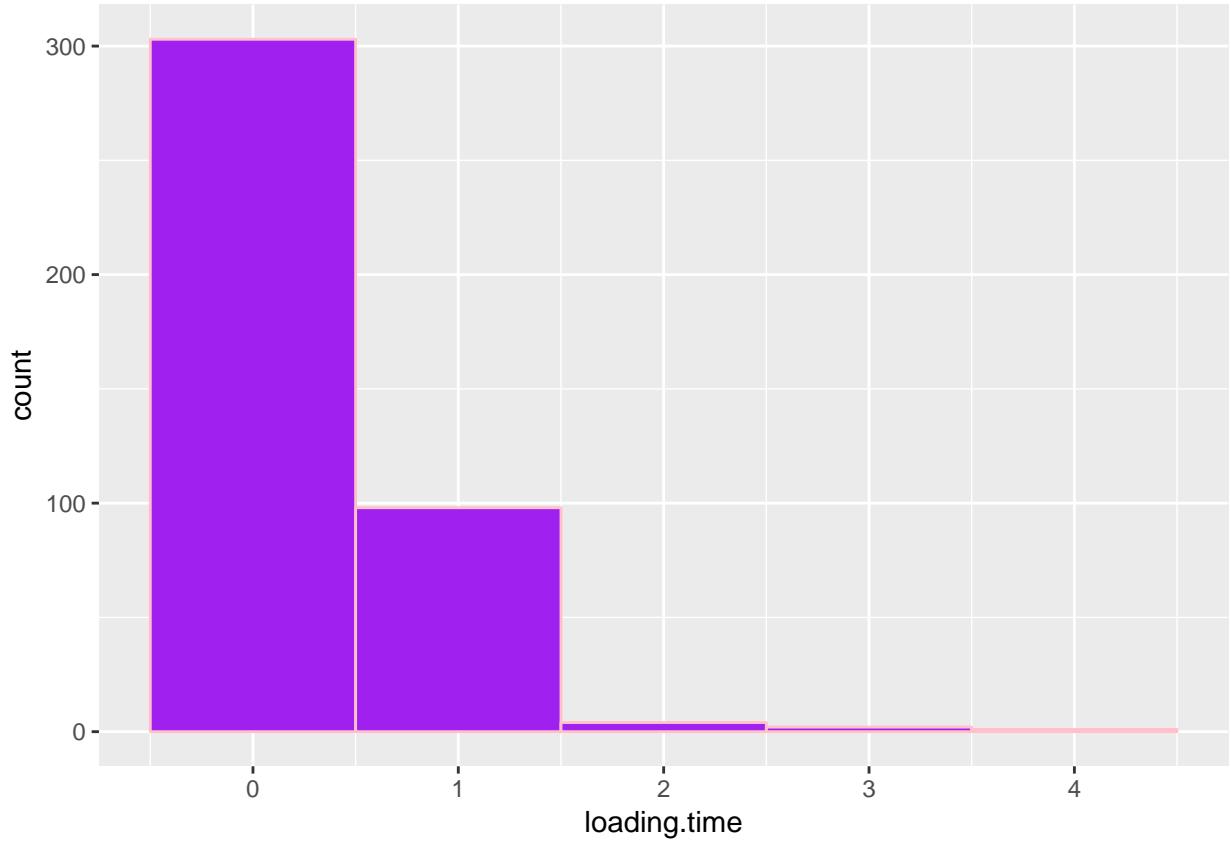
```
#And we can also see for correlations
total_500_links <- total_500_final[,c(4,21:23)]
library(corrplot)
library(caret)
tl <- cor(total_500_links)
tl

##             Revenues      external      internal total.links
## Revenues    1.00000000  0.034100506  0.004559950  0.01538199
## external     0.03410051  1.000000000 -0.002593961  0.32202419
## internal     0.00455995 -0.002593961  1.000000000  0.94589294
## total.links  0.01538199  0.322024191  0.945892937  1.00000000
corrplot(cor(total_500_links),method="number")
```

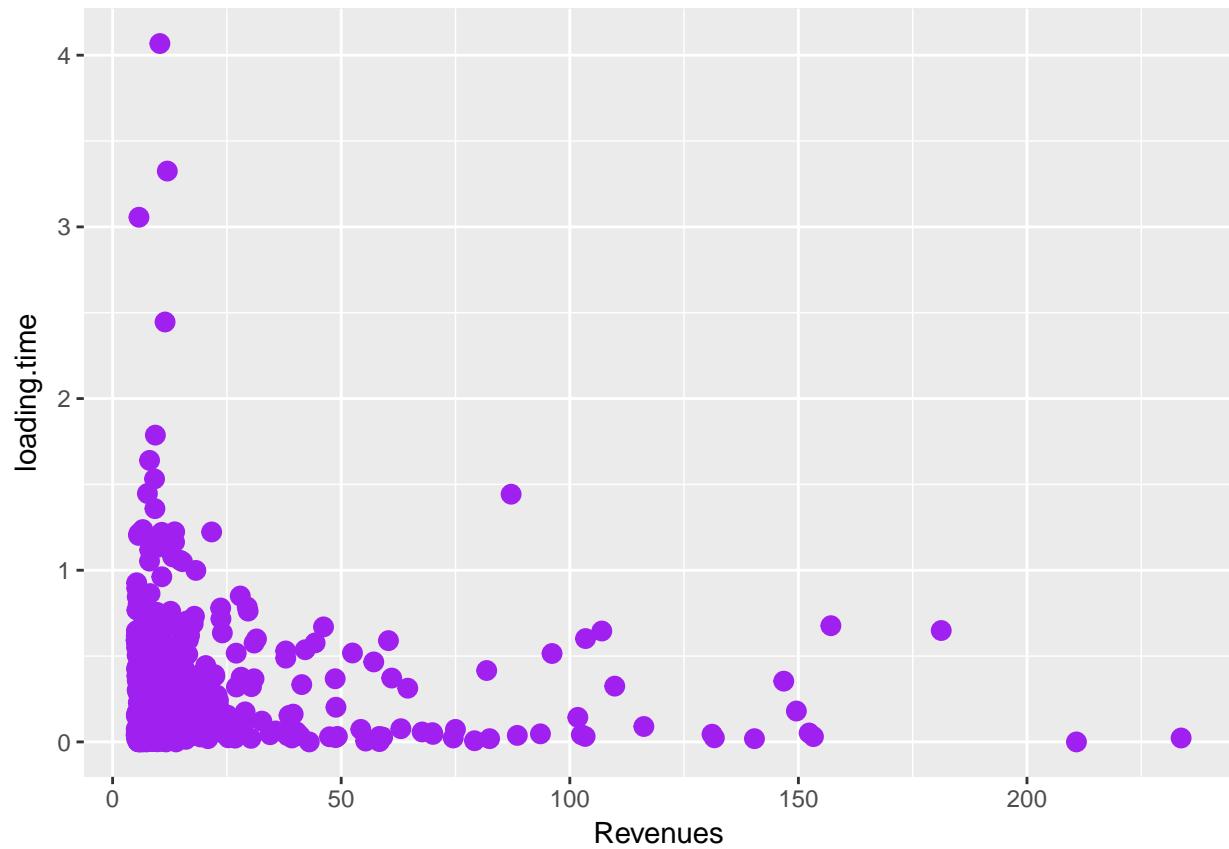


```
#We can see that the total links with the internal links have a correlation almost 95%.
#So we will not include the total links in the regression model
```

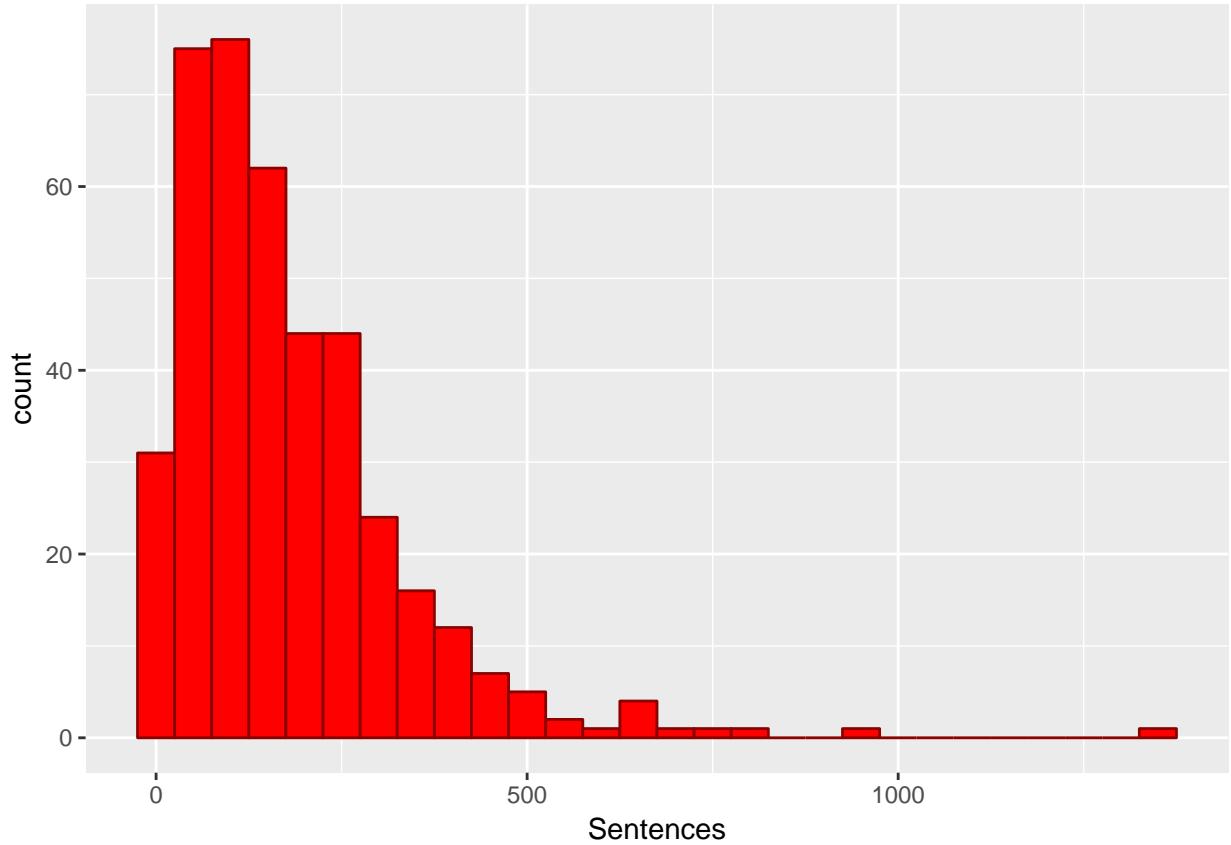
```
#####
#Now we will see the loading time per site
ggplot(data=total_500_final,aes(x=loading.time))+geom_histogram(binwidth=1, colour = "pink", fill ="purple")
```

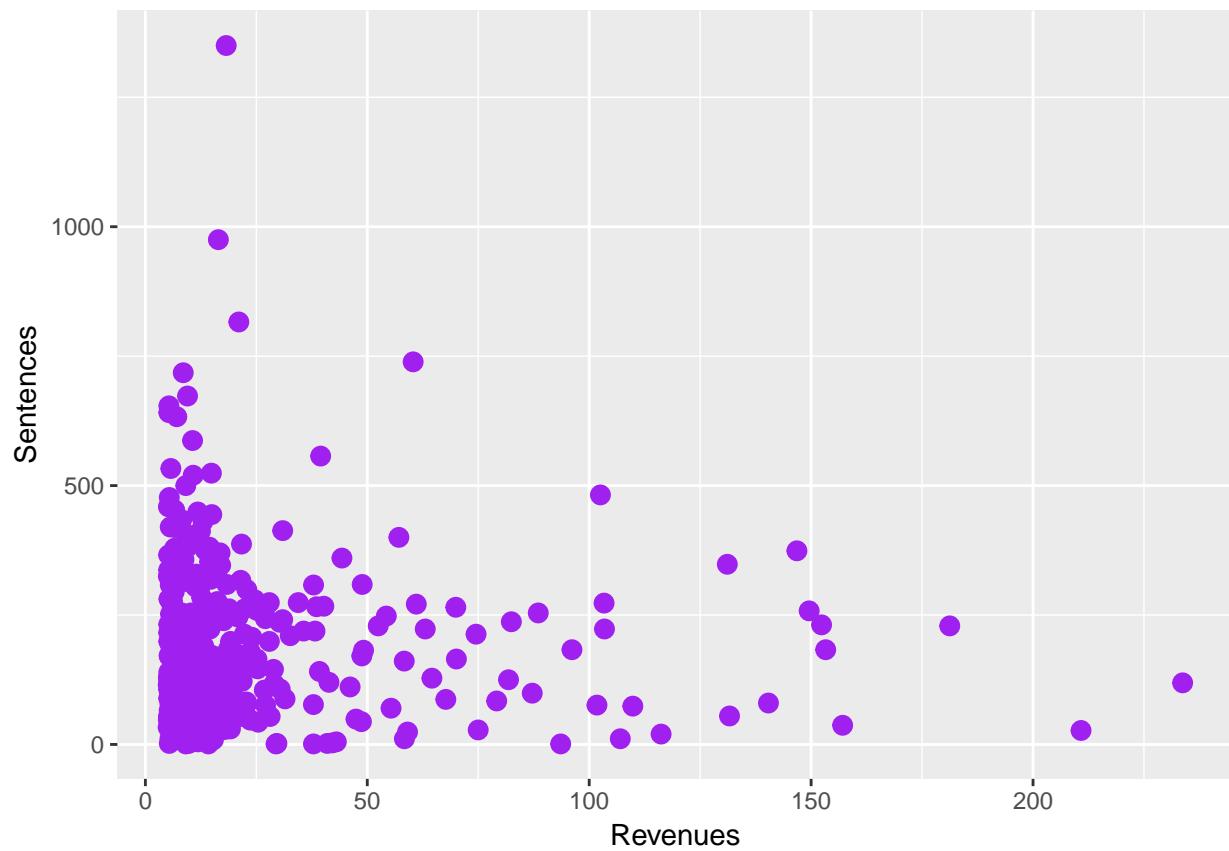


```
ggplot(total_500_final, aes(Revenues, loading.time)) + geom_point(size=3, colour = "purple")
```

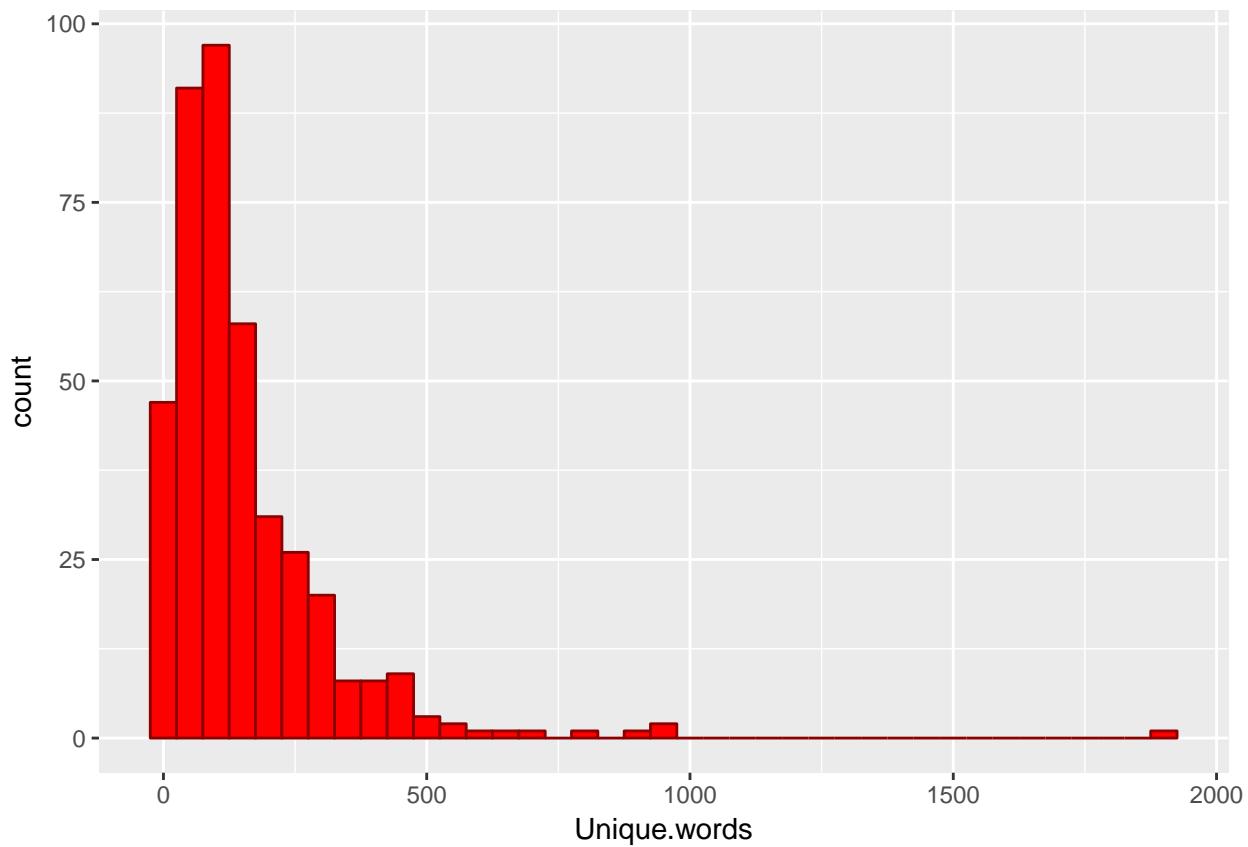


```
#####
#Now we will see the total words, the unique words and the sentences how are distributed alone and in r
ggplot(data=total_500_final,aes(x=Sentences))+geom_histogram(binwidth=50, colour = "darkred", fill = "red")
```

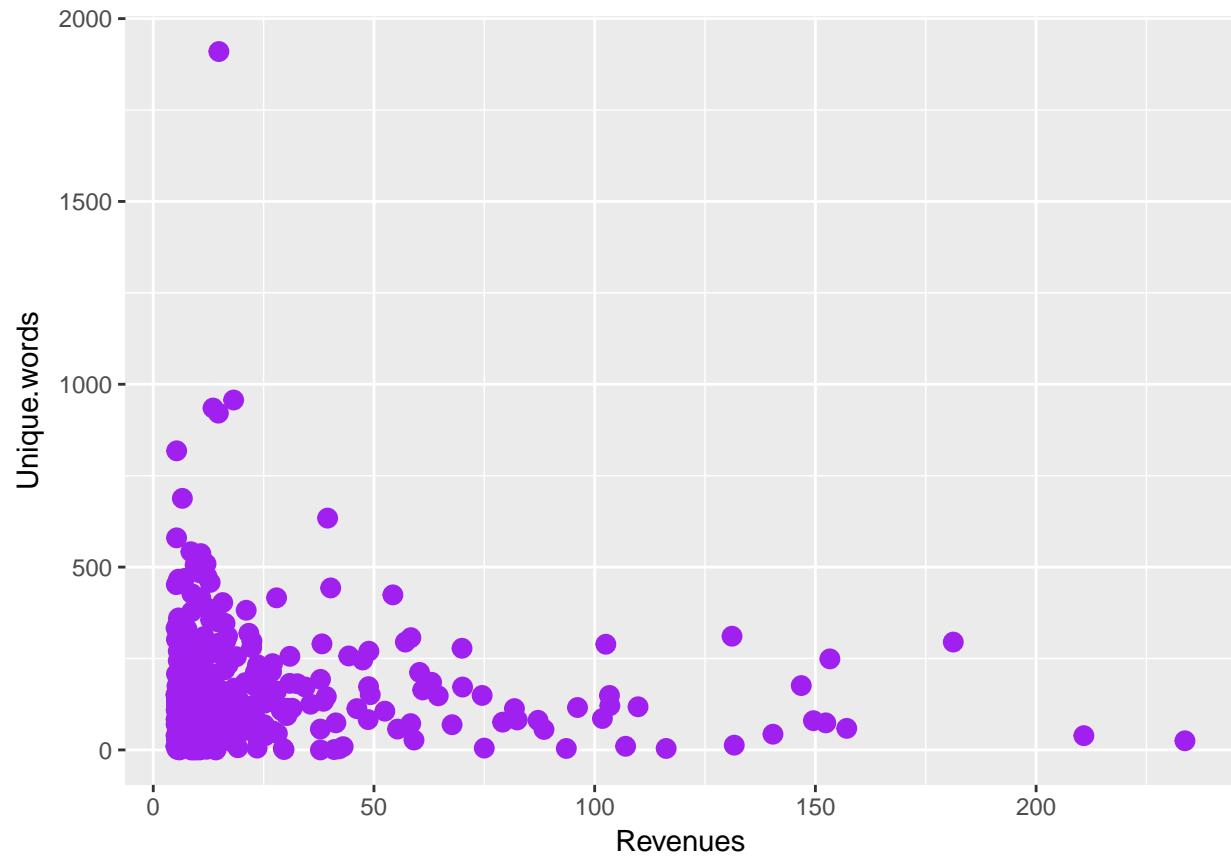




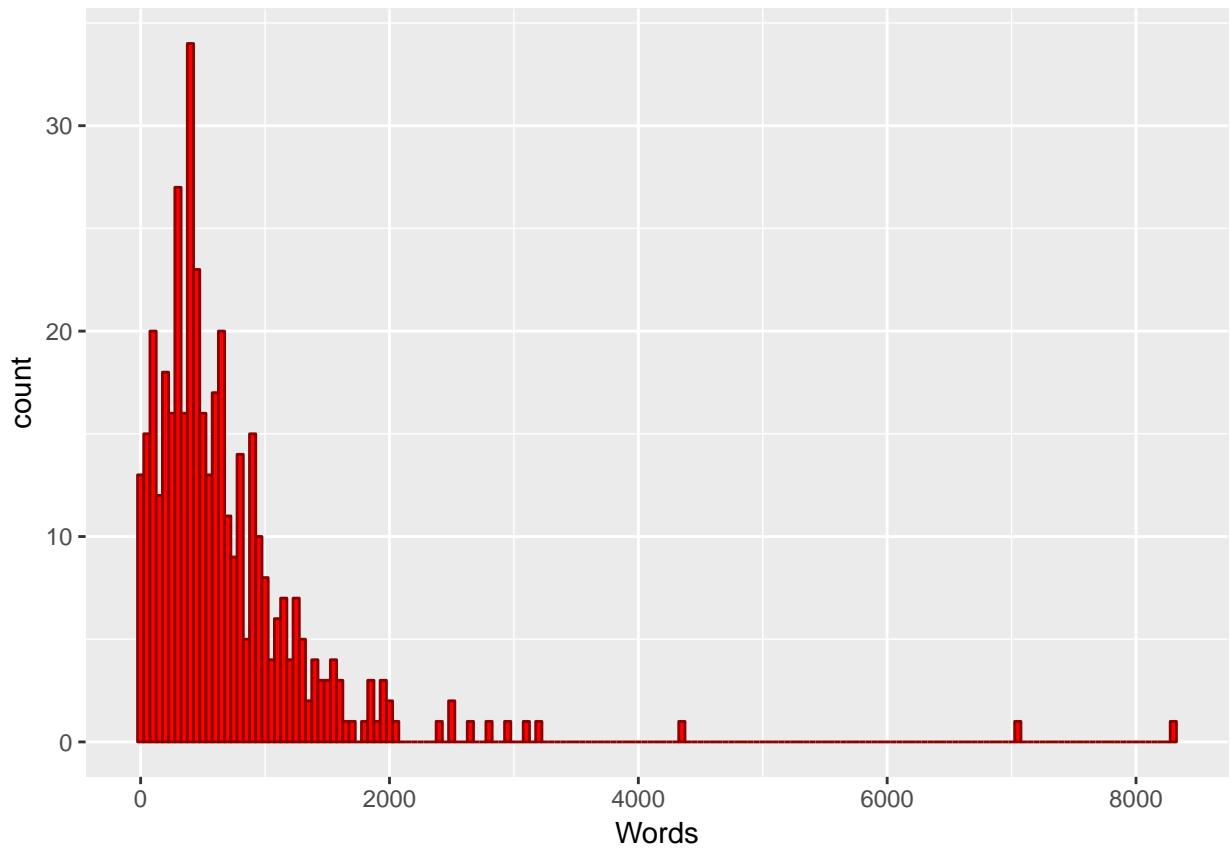
```
#####
ggplot(data=total_500_final,aes(x=Unique.words))+geom_histogram(binwidth=50, colour = "darkred", fill =
```

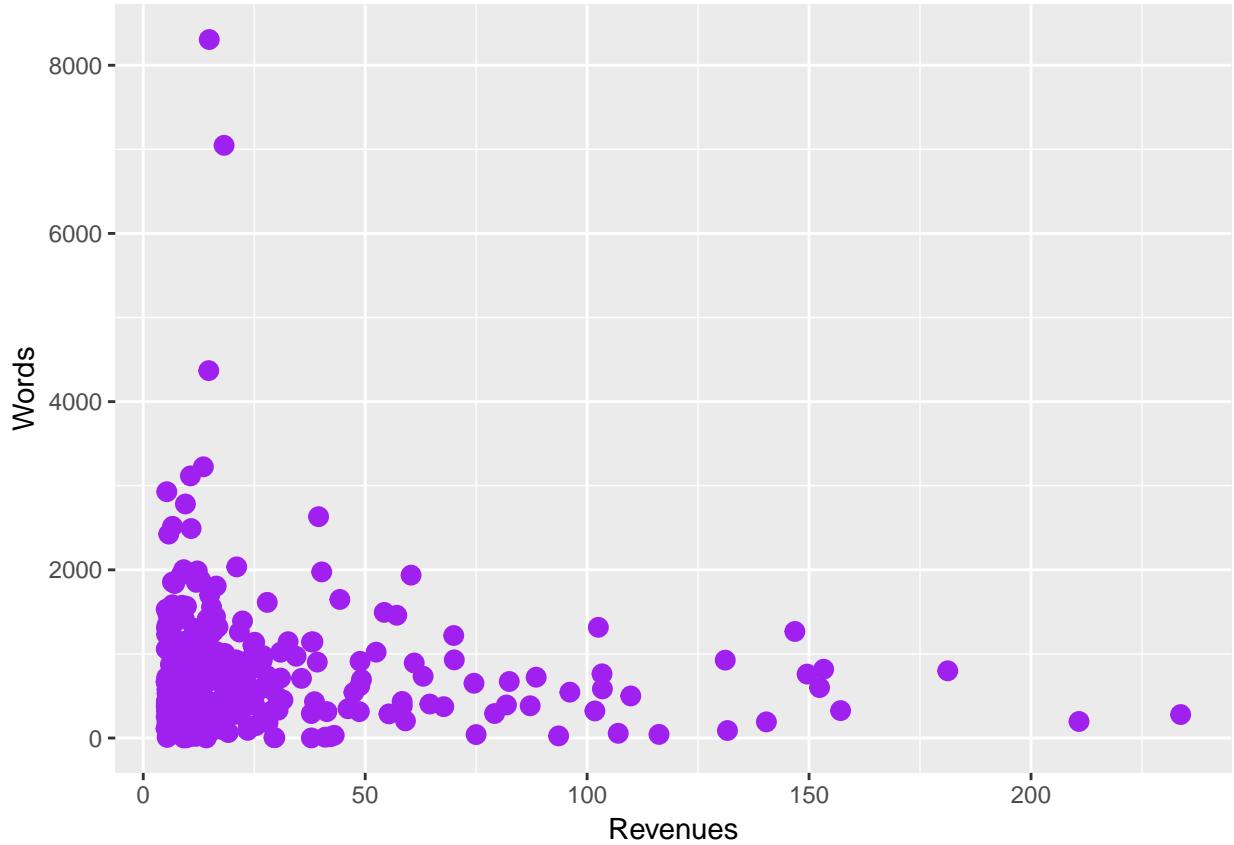


```
ggplot(total_500_final, aes(Revenues, Unique.words)) + geom_point(size=3, colour = "purple")
```



```
#####
ggplot(data=total_500_final,aes(x=Words))+geom_histogram(binwidth=50, colour = "darkred", fill ="red")
```



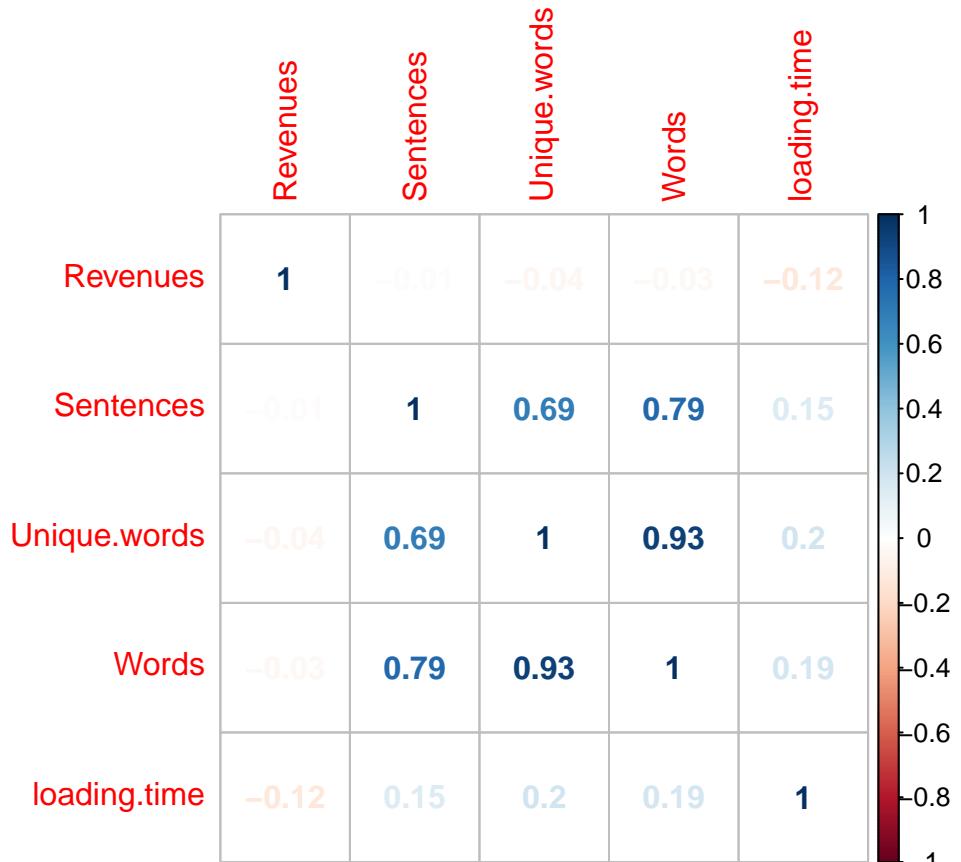


```
#####
#And we can also see for correlations
```

```
total_500_lt_w <- total_500_final[,c(4,18:20,727)]
library(corrplot)
library(caret)
tl <- cor(total_500_lt_w)
tl
```

	Revenues	Sentences	Unique.words	Words	loading.time
## Revenues	1.00000000	-0.01183819	-0.04362118	-0.03479049	-0.1212650
## Sentences	-0.01183819	1.00000000	0.69454327	0.78851979	0.1497520
## Unique.words	-0.04362118	0.69454327	1.00000000	0.93243940	0.1994296
## Words	-0.03479049	0.78851979	0.93243940	1.00000000	0.1857922
## loading.time	-0.12126500	0.14975205	0.19942956	0.18579225	1.0000000

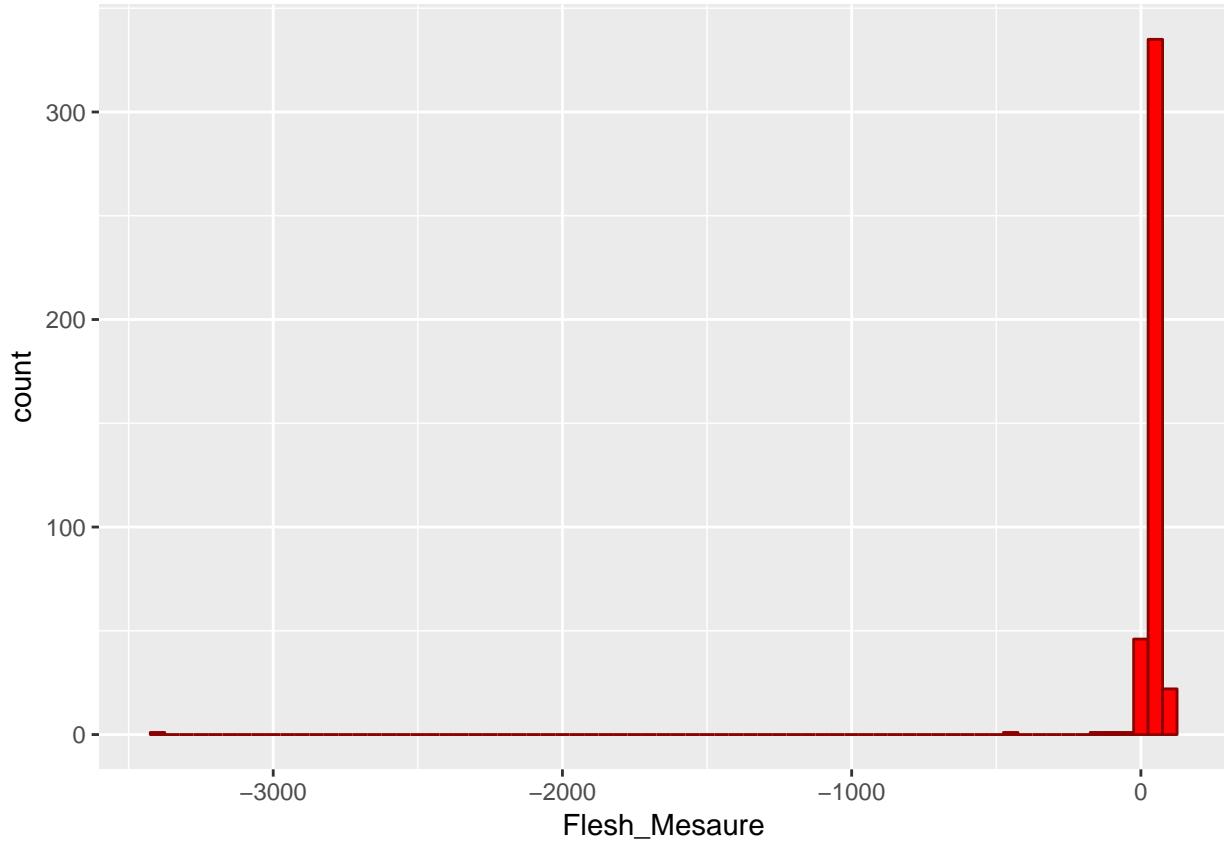
```
corrplot(cor(total_500_lt_w),method="number")
```



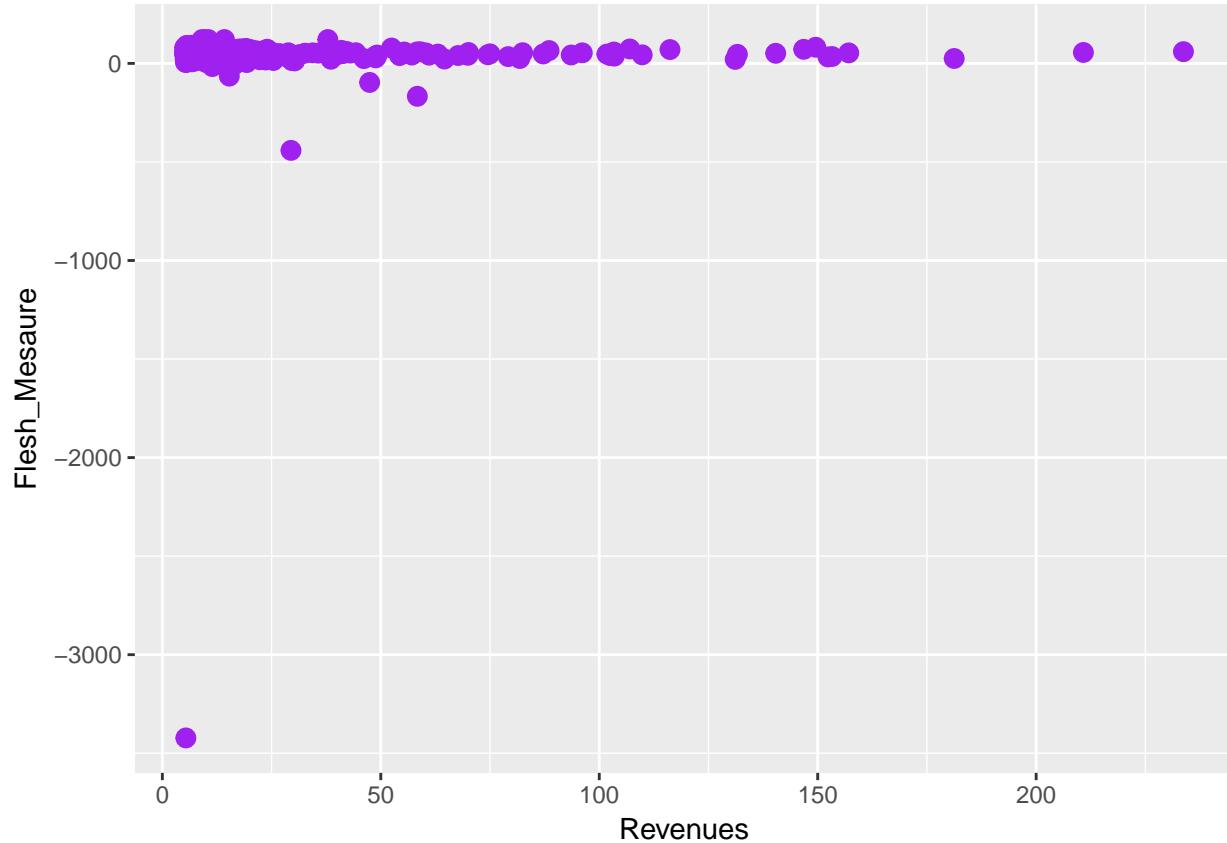
```
#####
#####
```

```
#Next we will check the Flesh Measure alone and in relationship with revenues
```

```
ggplot(data=total_500_final,aes(x=Flesh_Mesaure))+geom_histogram(binwidth=50, colour = "darkred", fill =
```



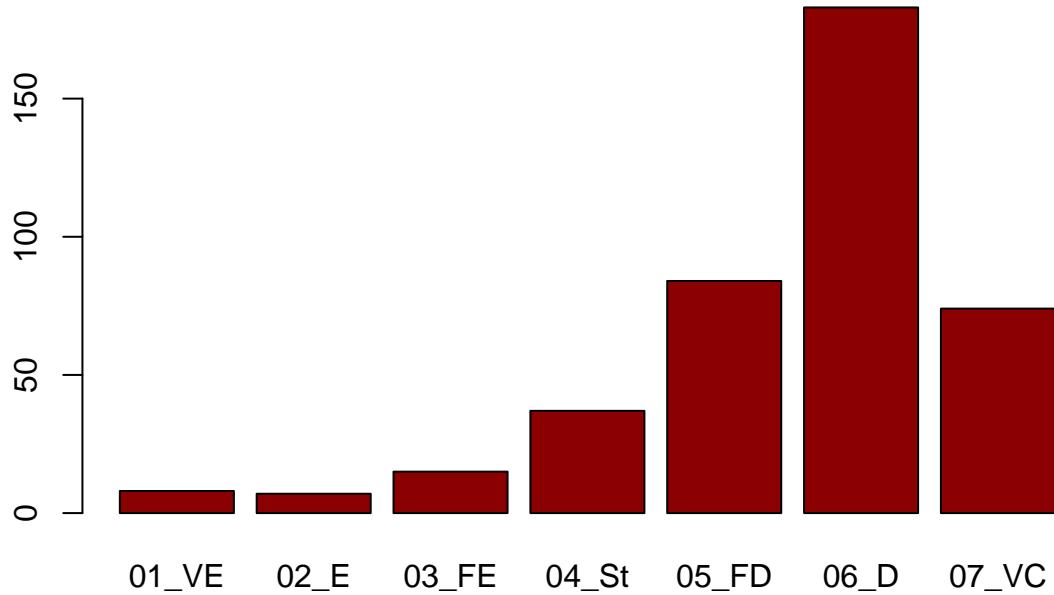
```
ggplot(total_500_final, aes(Revenues, Flesh_Mesaure)) + geom_point(size=3, colour = "purple")
```



```
#####

```

```
total_500_final$Readability <- gsub("Very easy", "01_VE", total_500_final$Readability )
total_500_final$Readability <- gsub("Easy", "02_E", total_500_final$Readability )
total_500_final$Readability <- gsub("Fairly easy", "03_FE", total_500_final$Readability )
total_500_final$Readability <- gsub("Standard", "04_St", total_500_final$Readability )
total_500_final$Readability <- gsub("Fairly difficult", "05_FD", total_500_final$Readability )
total_500_final$Readability <- gsub("Difficult", "06_D", total_500_final$Readability )
total_500_final$Readability <- gsub("Very Confusing", "07_VC", total_500_final$Readability )
barplot(table(total_500_final$Readability), col ="dark red")
```

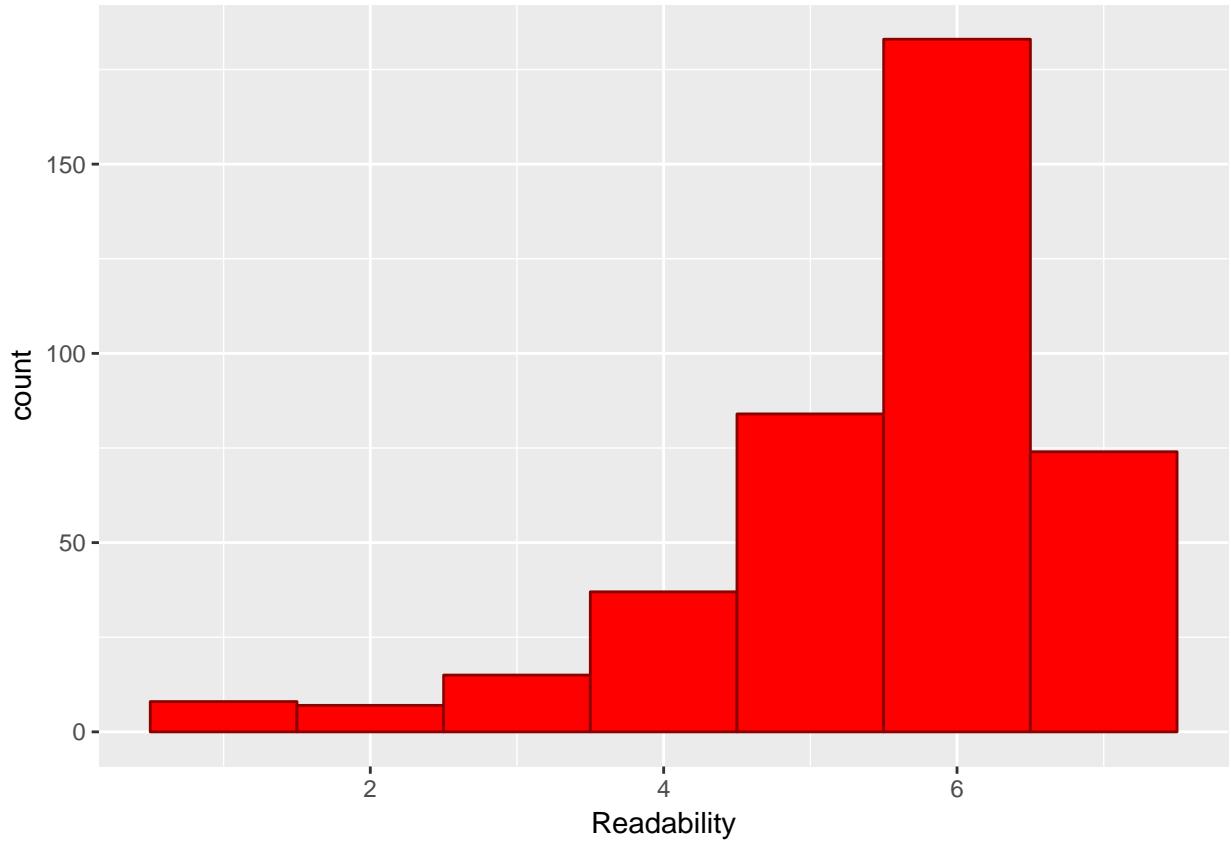


```

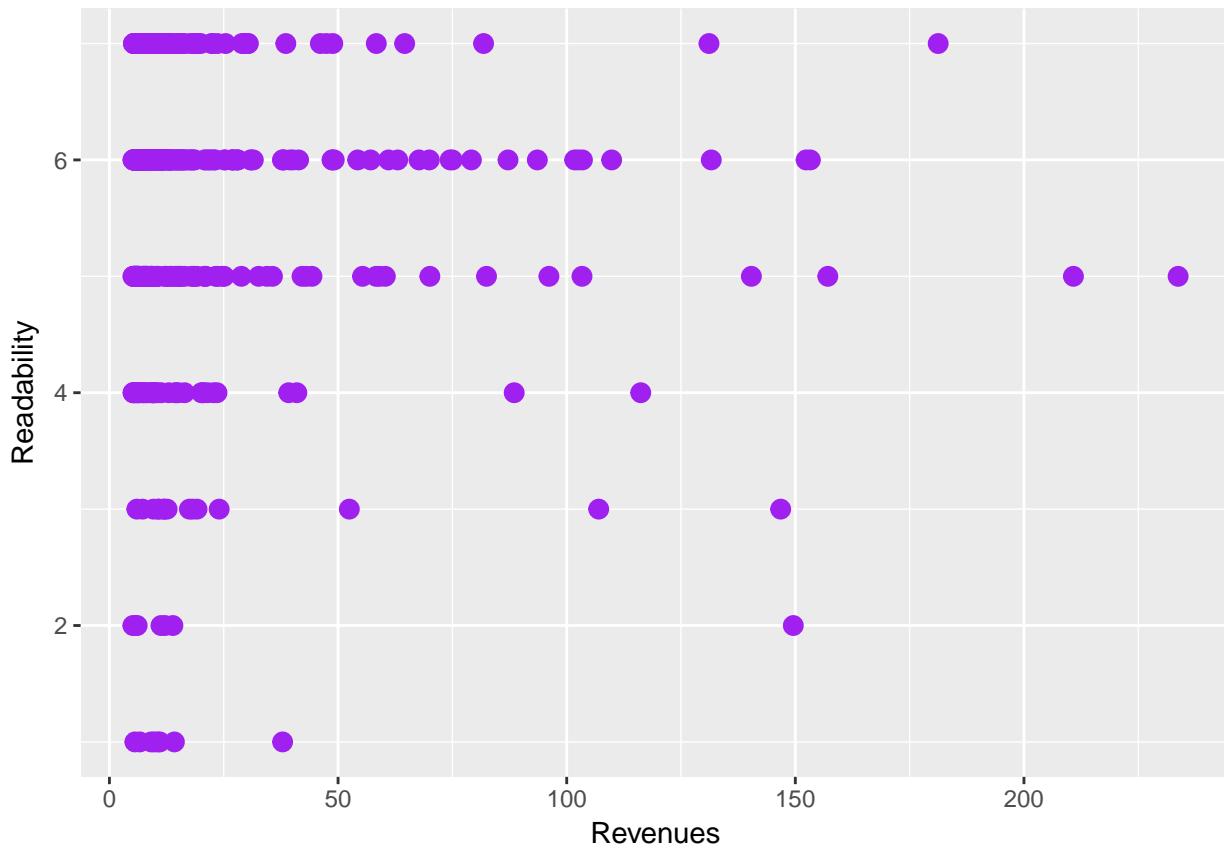
total_500_final$Readability <- gsub("01_VE", "1", total_500_final$Readability )
total_500_final$Readability <- gsub("02_E", "2", total_500_final$Readability )
total_500_final$Readability <- gsub("03_FE", "3", total_500_final$Readability )
total_500_final$Readability <- gsub("04_St", "4", total_500_final$Readability )
total_500_final$Readability <- gsub("05_FD", "5", total_500_final$Readability )
total_500_final$Readability <- gsub("06_D", "6", total_500_final$Readability )
total_500_final$Readability <- gsub("07_VC", "7", total_500_final$Readability )
total_500_final$Readability <- as.numeric(total_500_final$Readability )
ggplot(data=total_500_final,aes(x=Readability))+geom_bar(binwidth=1, colour = "darkred", fill ="red")

## Warning: `geom_bar()` no longer has a `binwidth` parameter. Please use
## `geom_histogram()` instead.

```

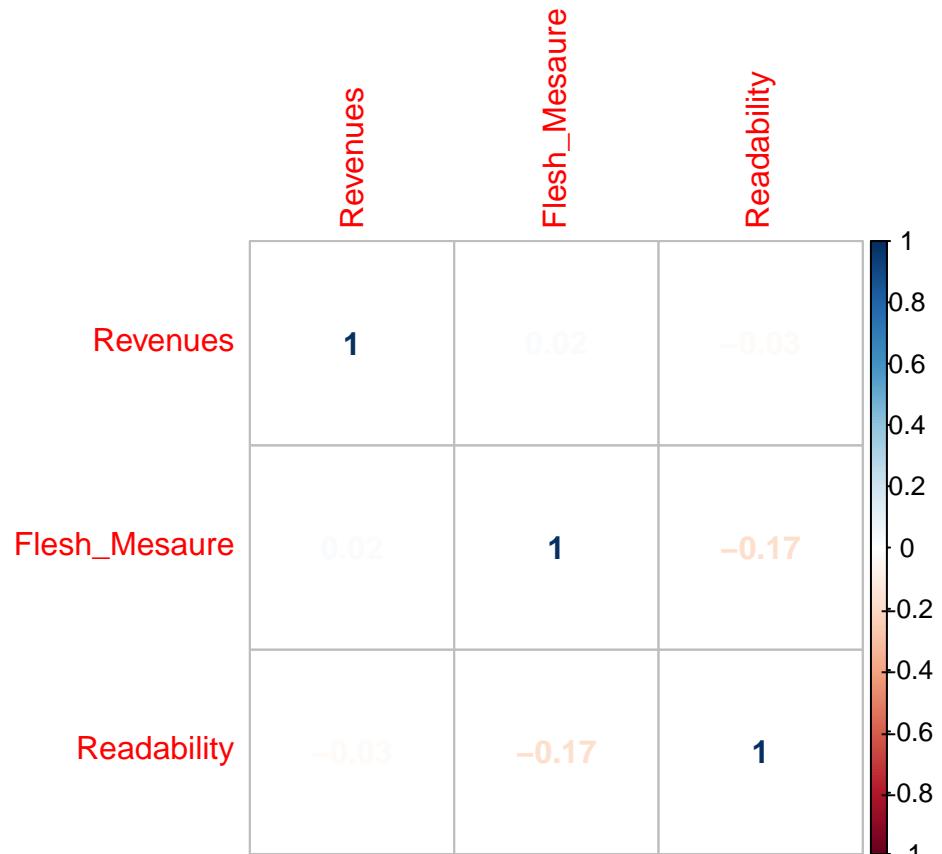


```
ggplot(total_500_final, aes(Revenues, Readability)) + geom_point(size=3, colour = "purple")
```

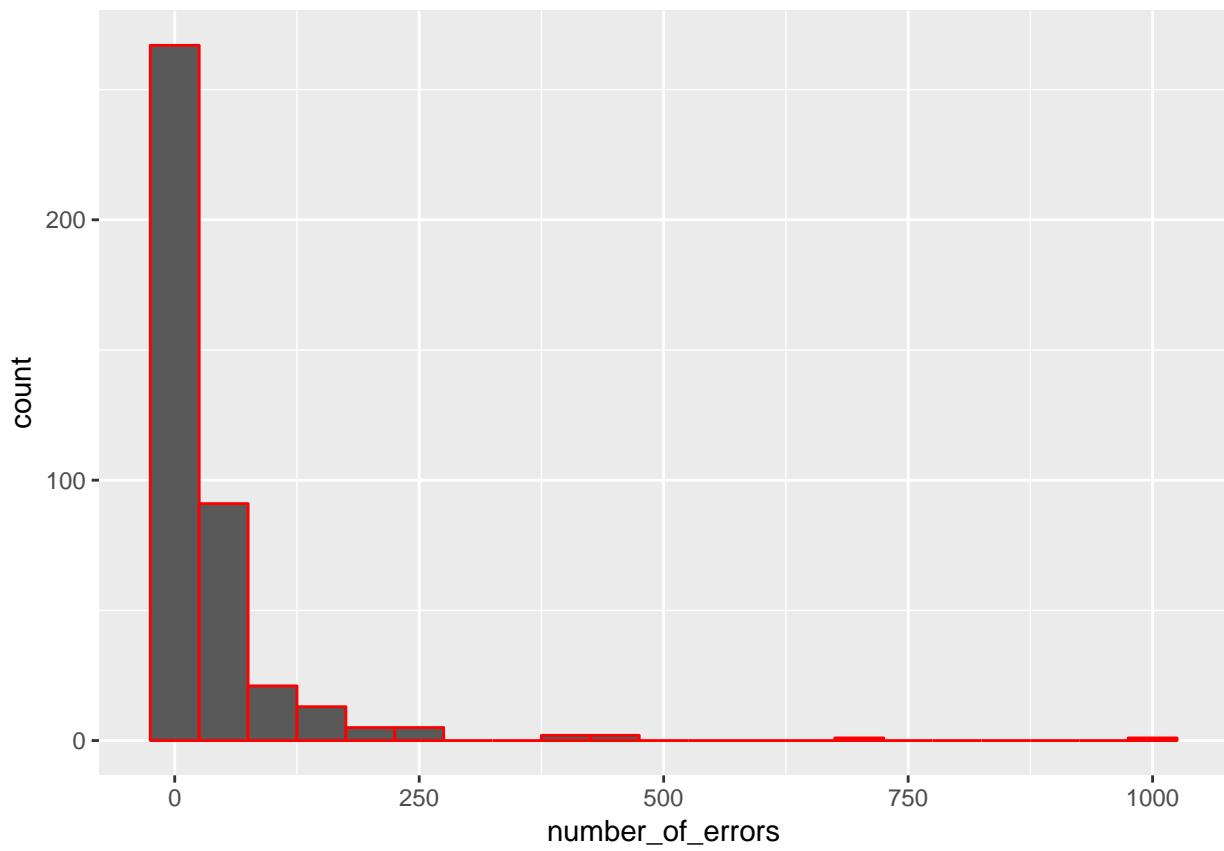


```
#And we can also see for correlations
total_500_r <- total_500_final[,c(4,16,17)]
library(corrplot)
library(caret)
tl <- cor(total_500_r)
tl

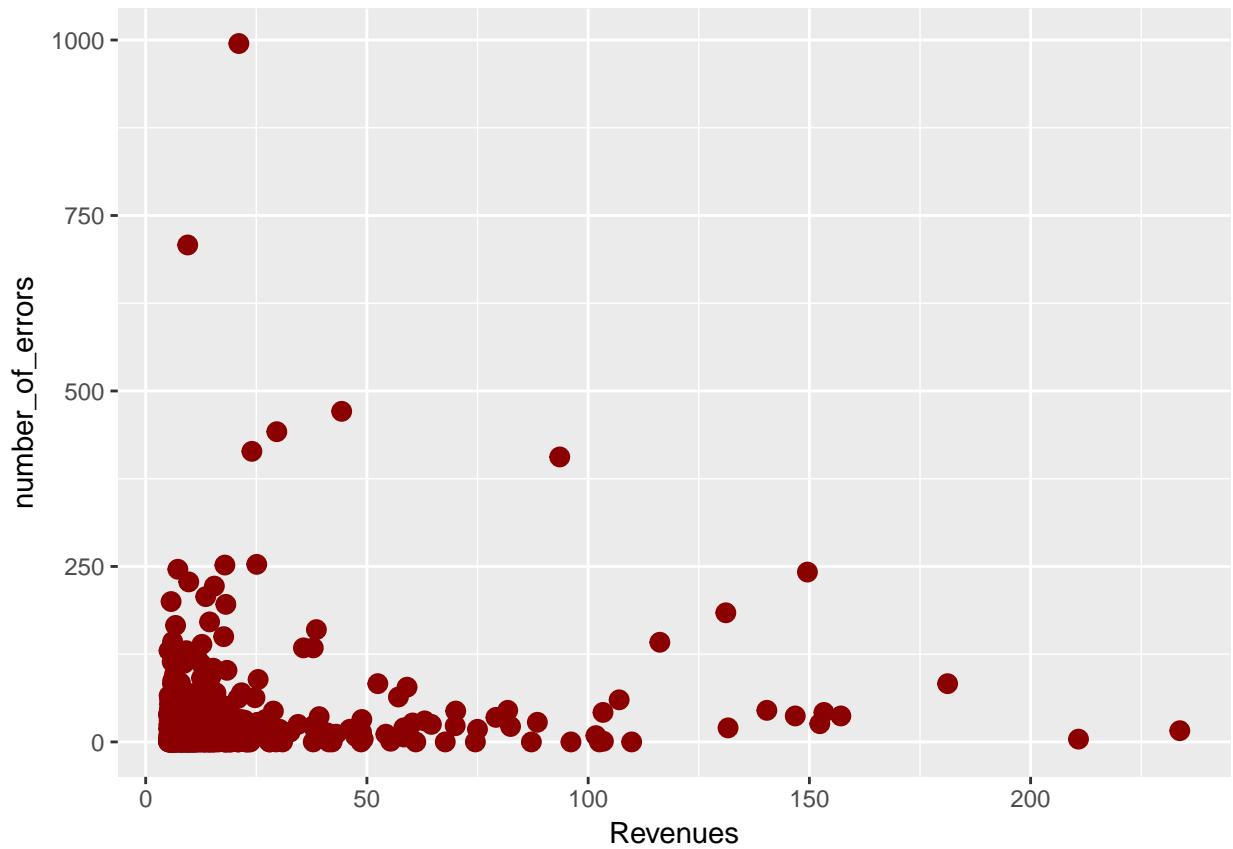
##          Revenues Flesh_Mesaure Readability
## Revenues      1.00000000  0.02476229 -0.02694931
## Flesh_Mesaure  0.02476229  1.00000000 -0.17094994
## Readability   -0.02694931 -0.17094994  1.00000000
corrplot(cor(total_500_r),method="number")
```

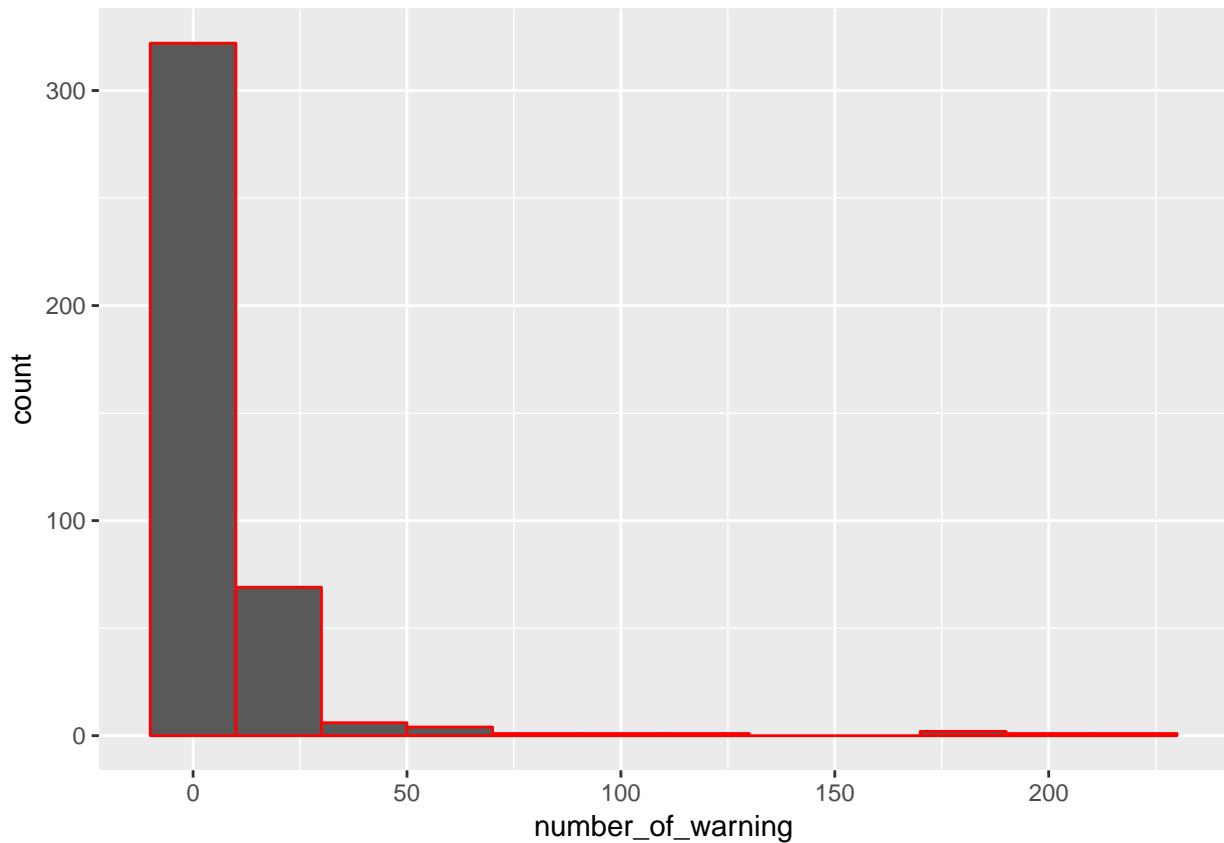


```
#####
#Now we will see the number of errors and warnings alone and in relationship with the Revenues
ggplot(data=total_500_final,aes(x=number_of_errors))+geom_histogram(binwidth=50, colour = "red")
```

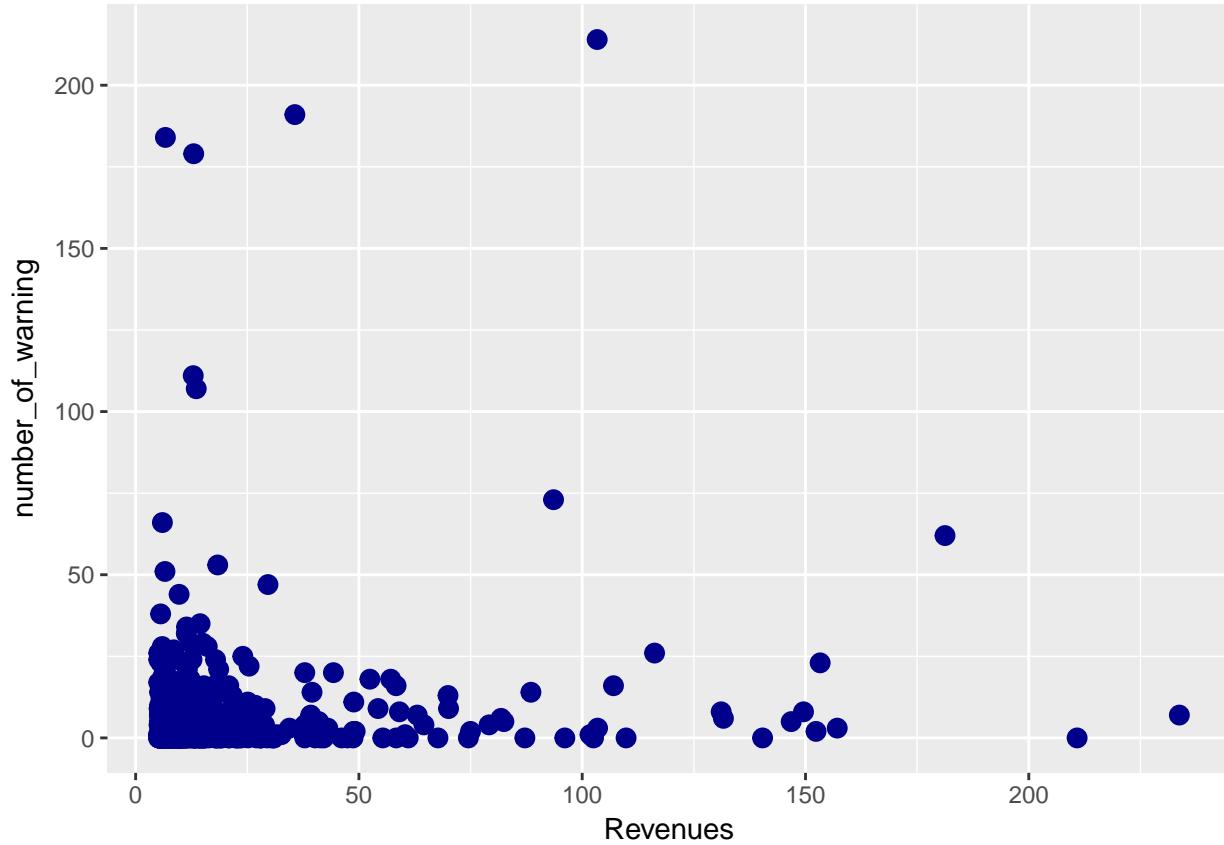


```
ggplot(total_500_final, aes(Revenues, number_of_errors)) + geom_point(size=3, colour = "dark red")
```

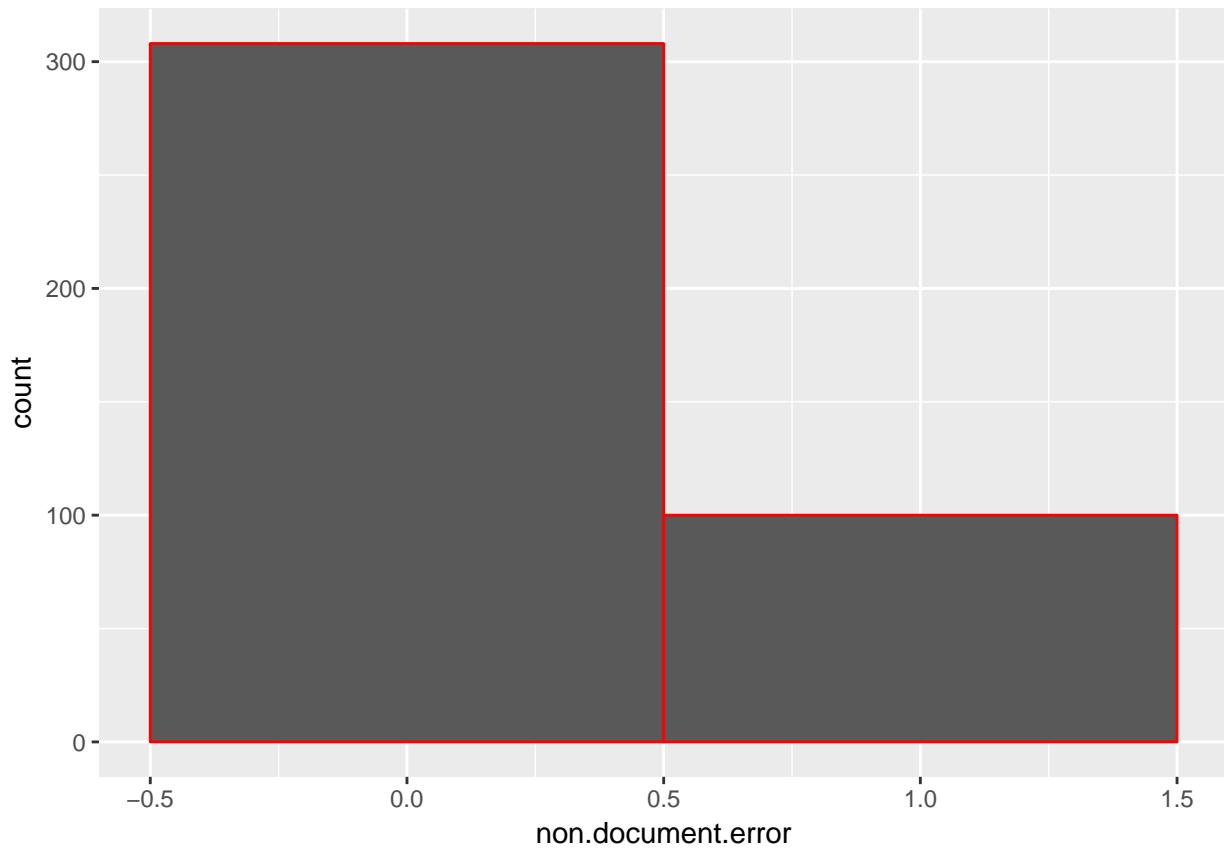


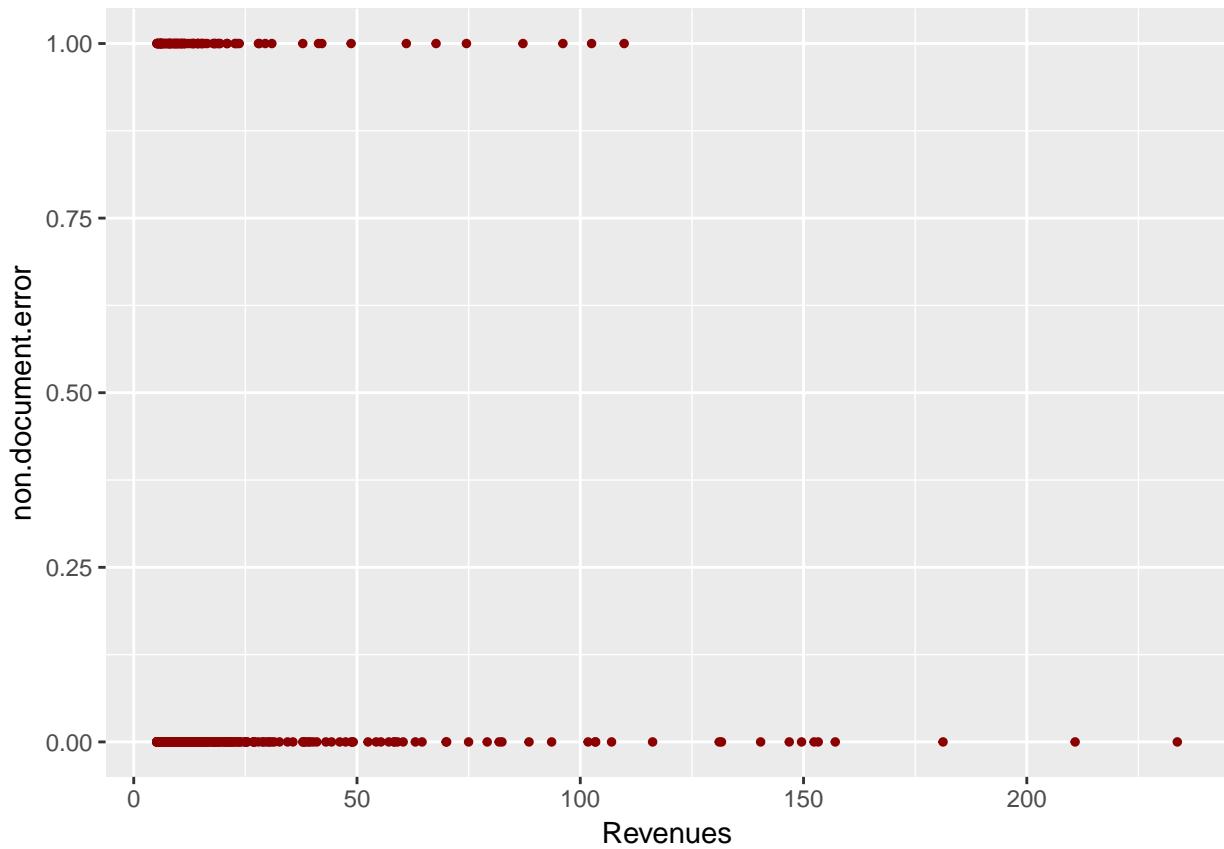


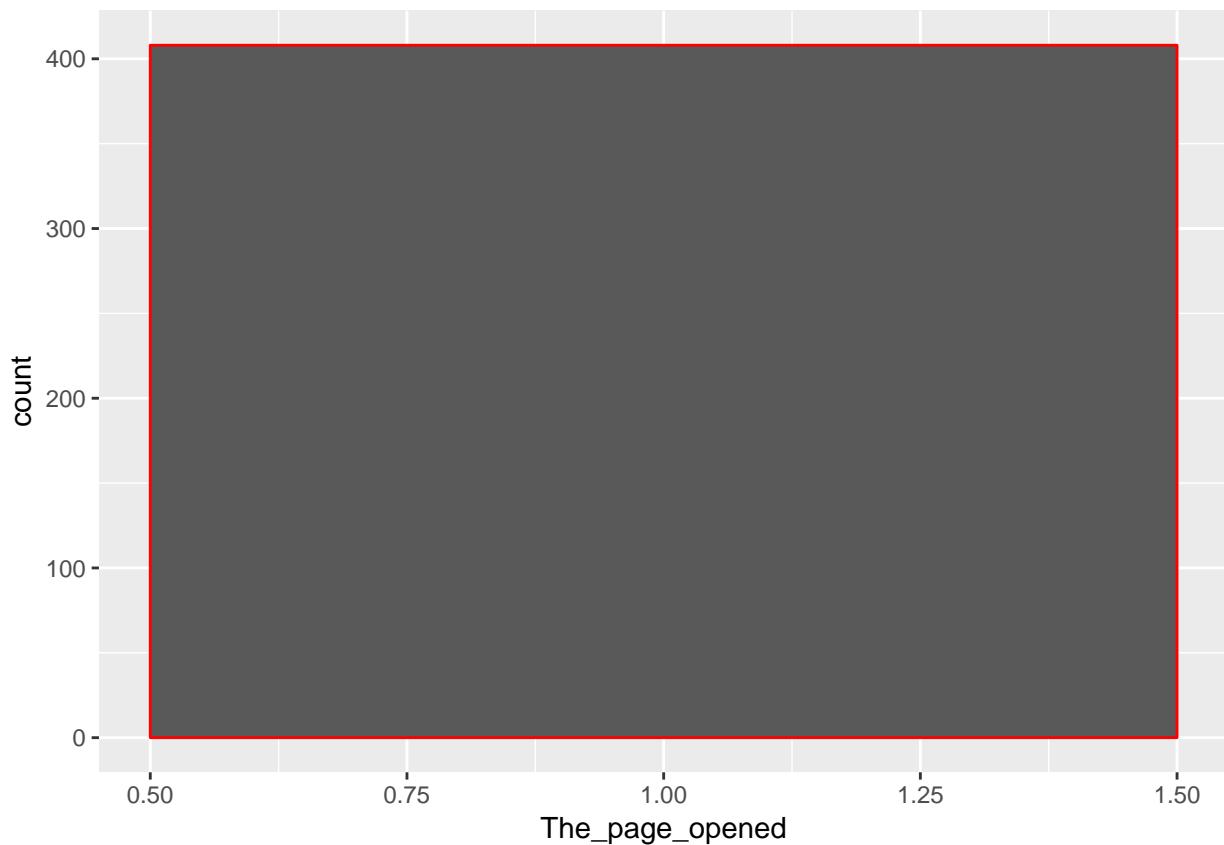
```
ggplot(total_500_final, aes(Revenues, number_of_warning)) + geom_point(size=3, colour = "dark blue")
```



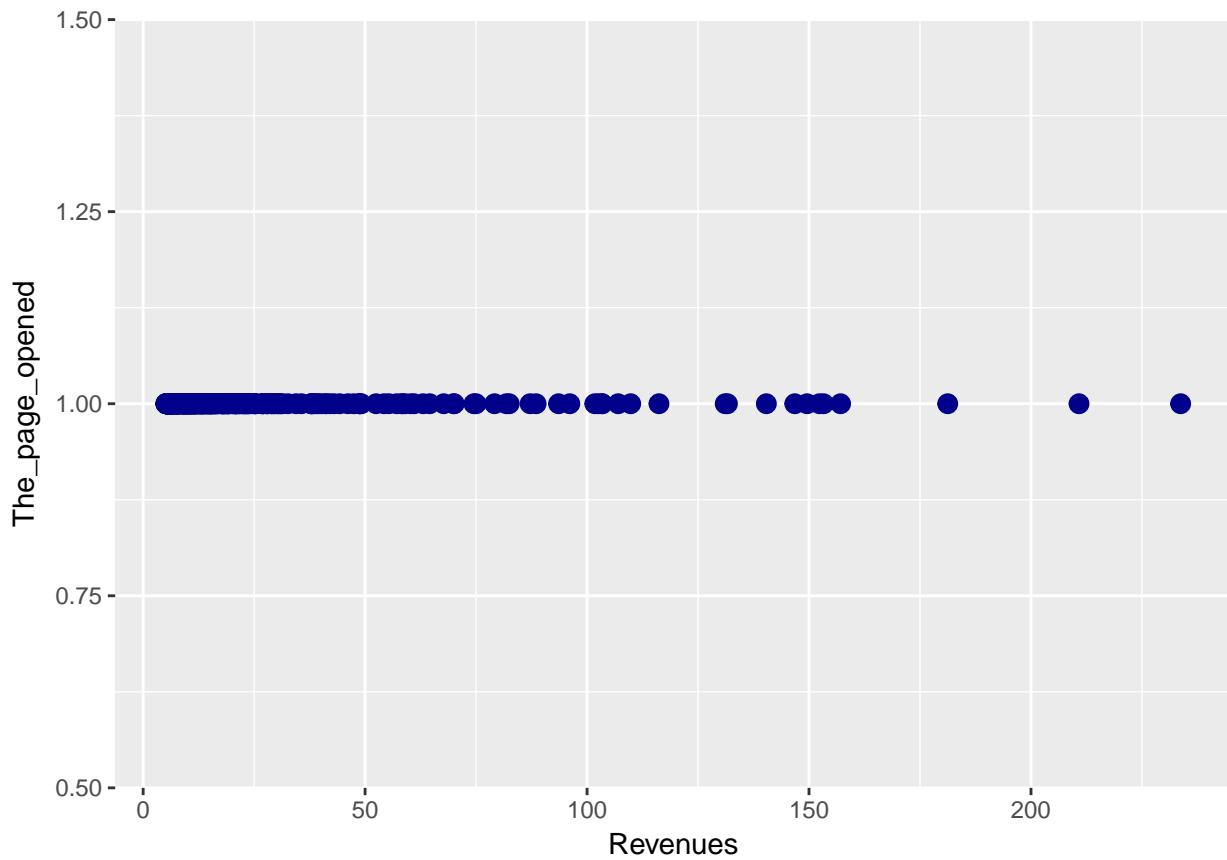
```
#####
#####  
##### Now we will see the non.document.error and the page not opened variables alone and in relationship with  
ggplot(data=total_500_final,aes(x=non.document.error))+geom_histogram(binwidth=1, colour = "red")
```







```
ggplot(total_500_final, aes(Revenues, The_page_opened)) + geom_point(size=3, colour = "dark blue")
```



```

#In the page not opened we can see that the variable has only the price 1 that means that the page open
#####
##### #And we can also see for correlations
total_500_html <- total_500_final[,c(4,7:9)]
library(corrplot)
library(caret)
tl <- cor(total_500_html)
tl

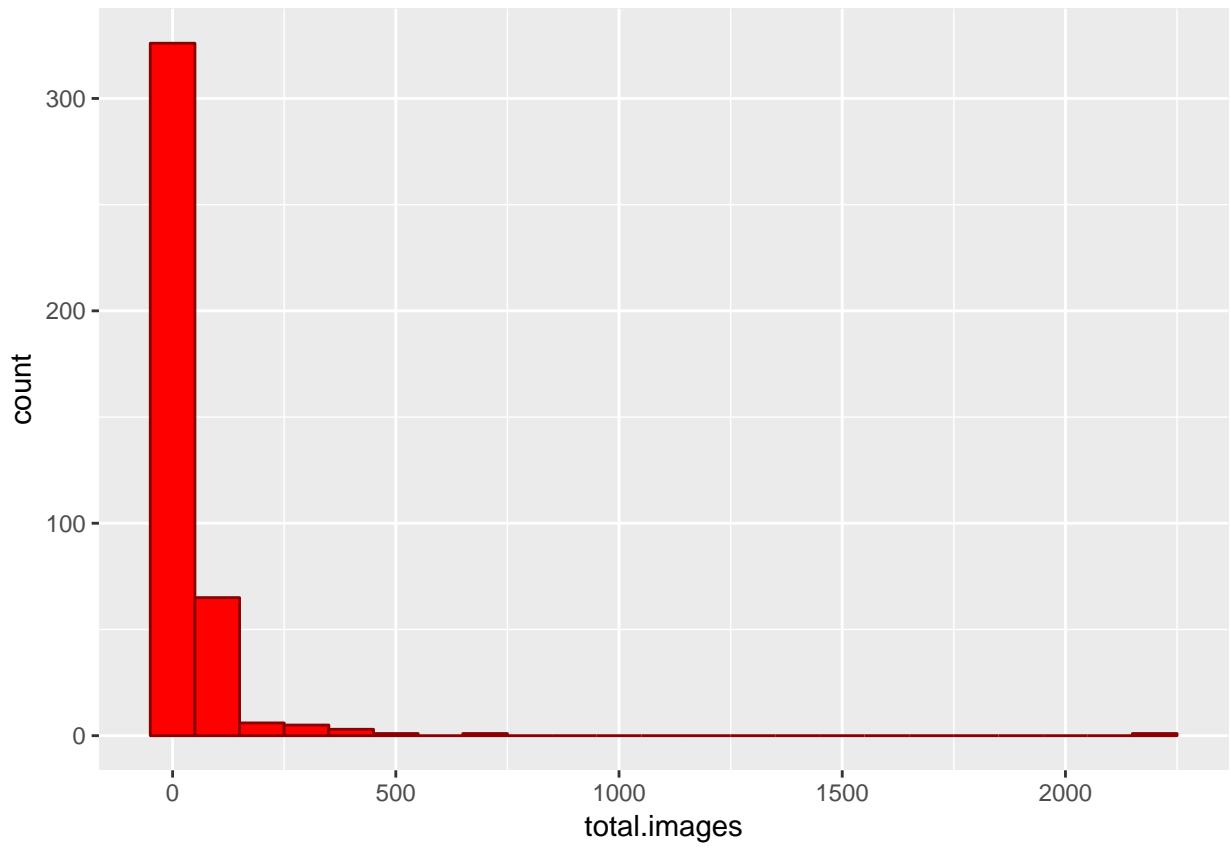
##          Revenues non.document.error number_of_errors
## Revenues      1.000000000 -0.0748407      0.0800205
## non.document.error -0.07484070      1.0000000     -0.2545301
## number_of_errors    0.08002050     -0.2545301      1.0000000
## number_of_warning   0.09505013     -0.2242315      0.2309578
##          number_of_warning
## Revenues           0.09505013
## non.document.error -0.22423152
## number_of_errors     0.23095778
## number_of_warning    1.000000000
corrplot(cor(total_500_html),method="number")

```

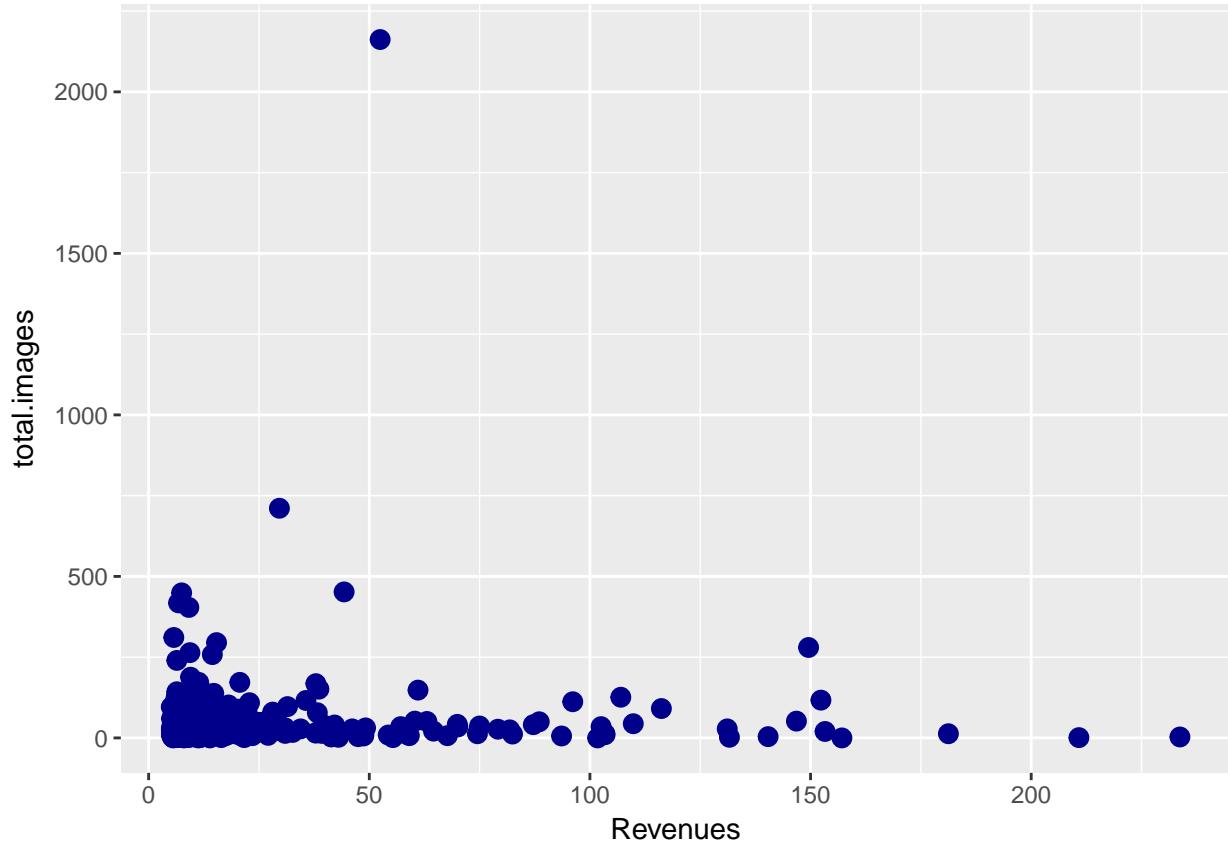


#Now we will see the total images alone and in relationship with the revenues

```
ggplot(data=total_500_final,aes(x=total.images))+geom_histogram(binwidth=100, colour = "darkred", fill = "white")
```

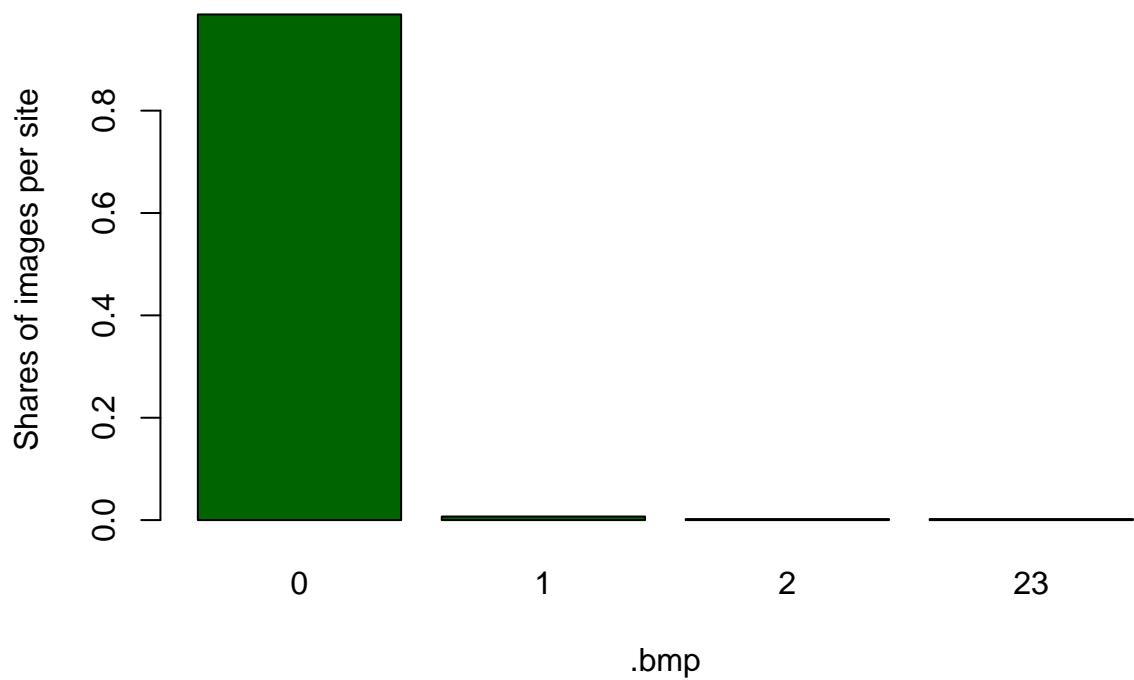


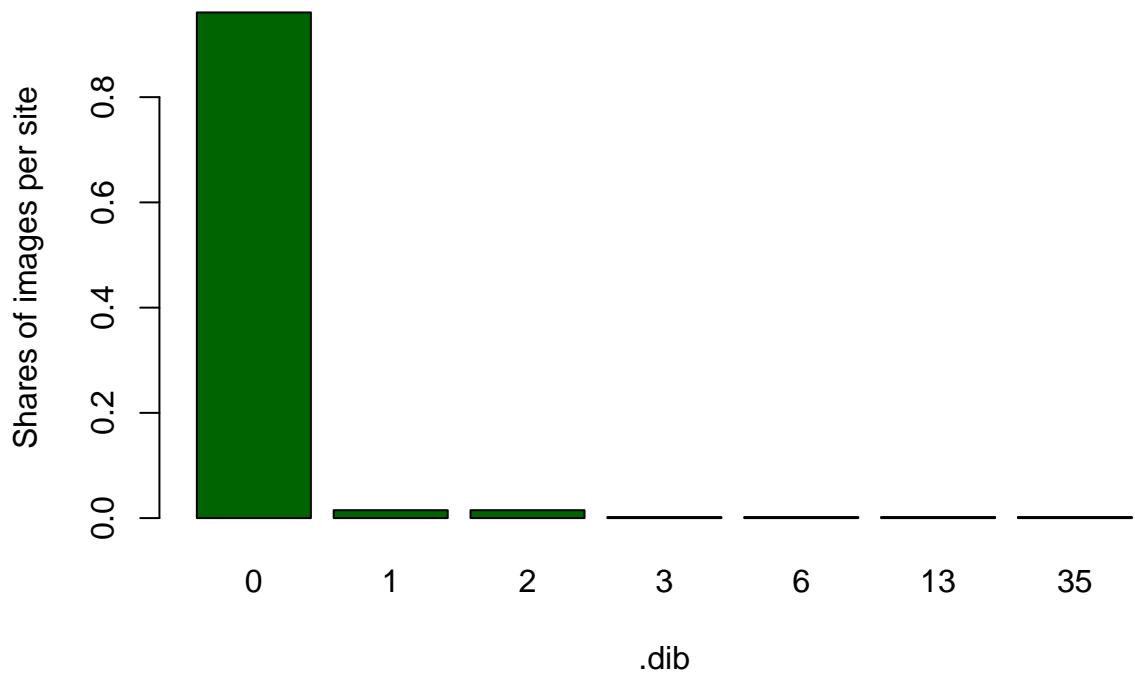
```
ggplot(total_500_final, aes(Revenues, total.images)) + geom_point(size=3, colour = "dark blue")
```

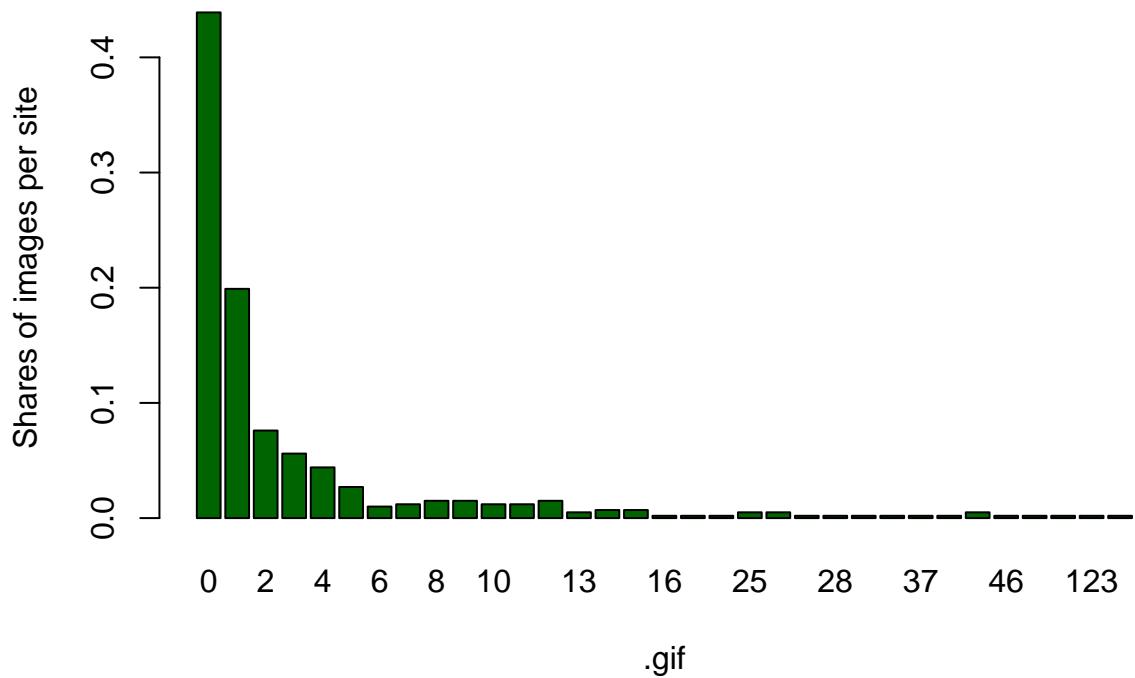


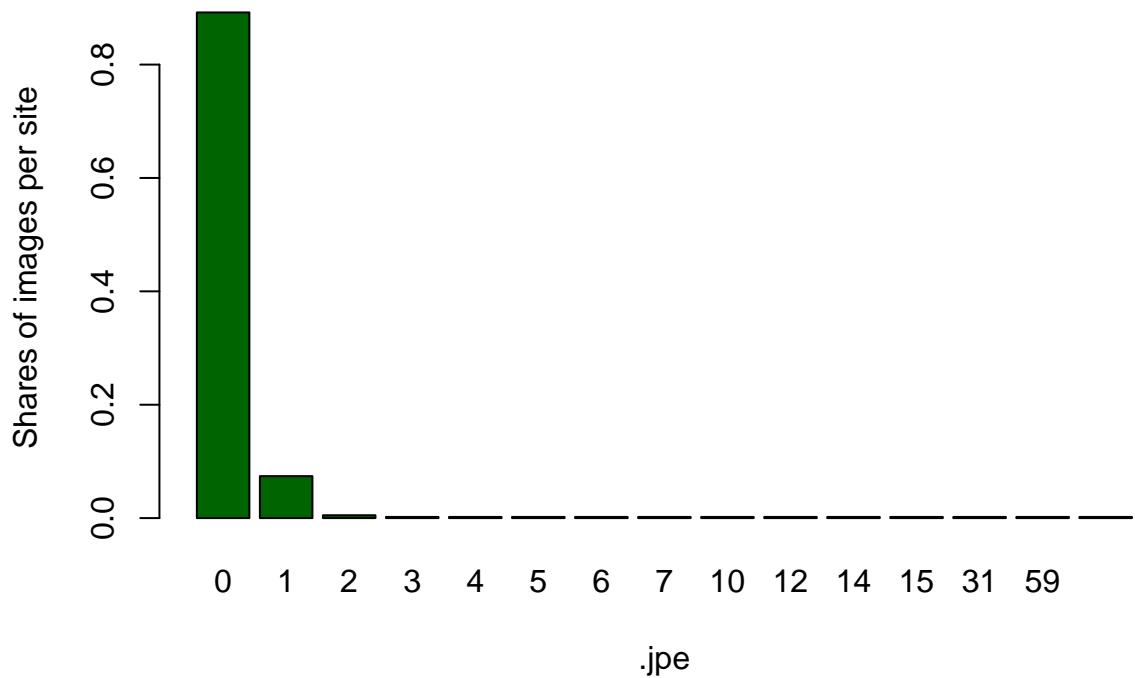
```
#####
#We will see now the frequency of image types that is being used
```

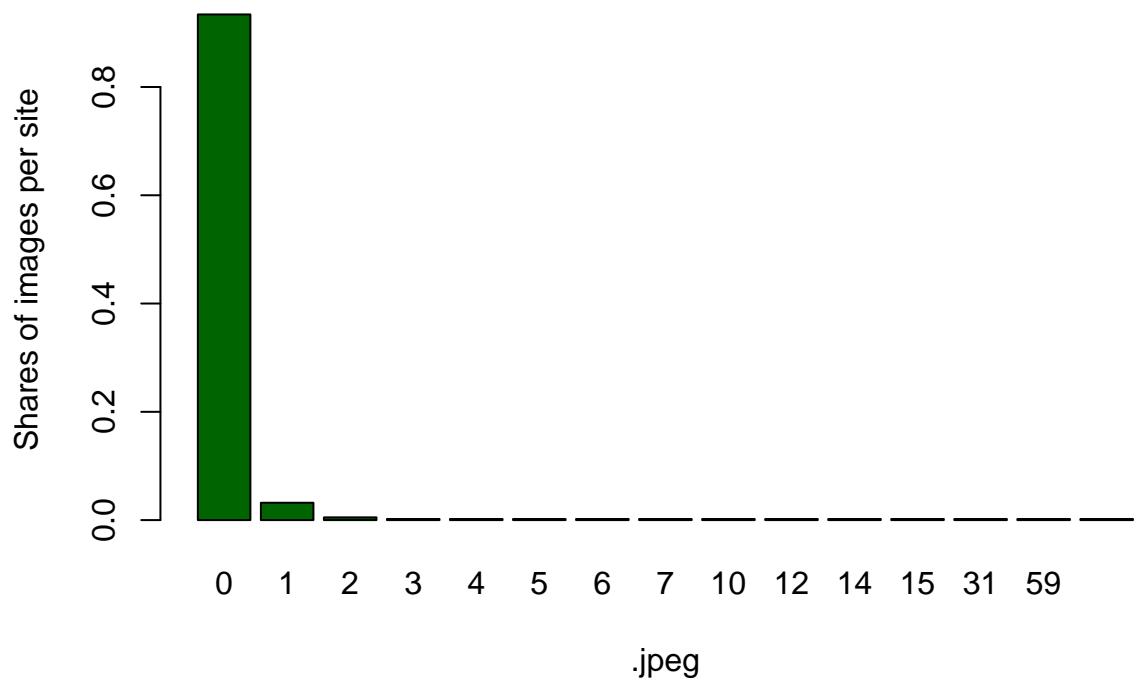
```
par(mfrow=c(1,1))
k = c(717:725)
for(i in 1:9){
  a <- k[i]
  image_type<- round(table(total_500_final[,a])/408,3)
  barplot(image_type,xlab=names(total_500_final)[a],ylab = "Shares of images per site", col = "dark green")
```

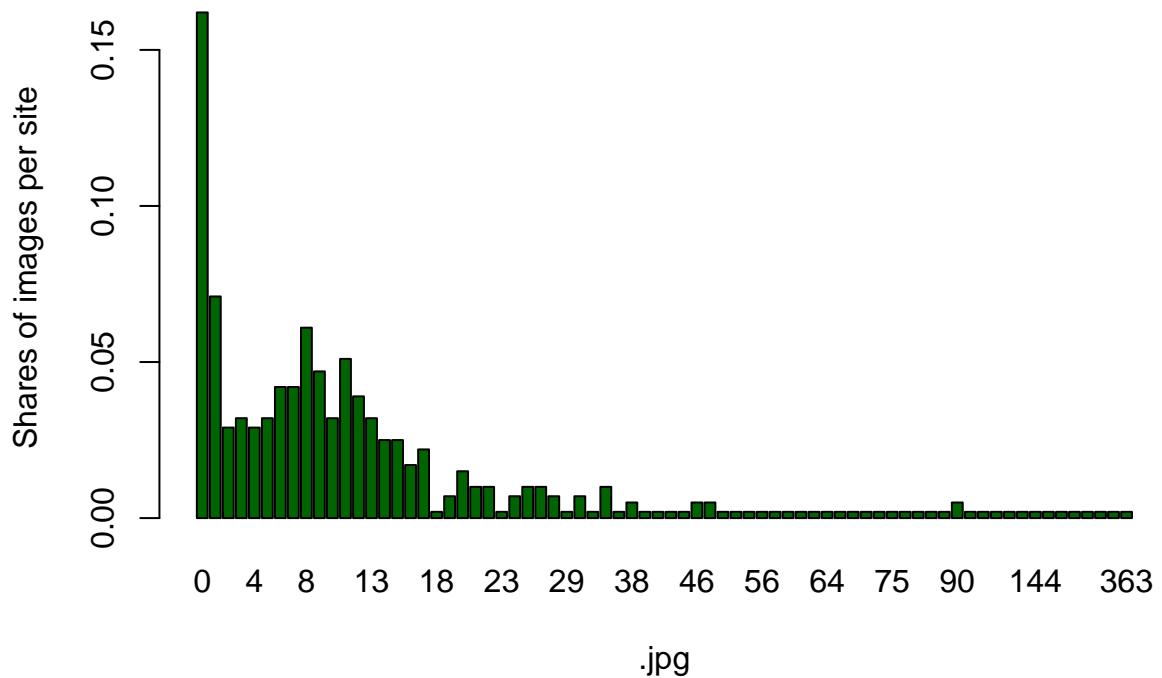


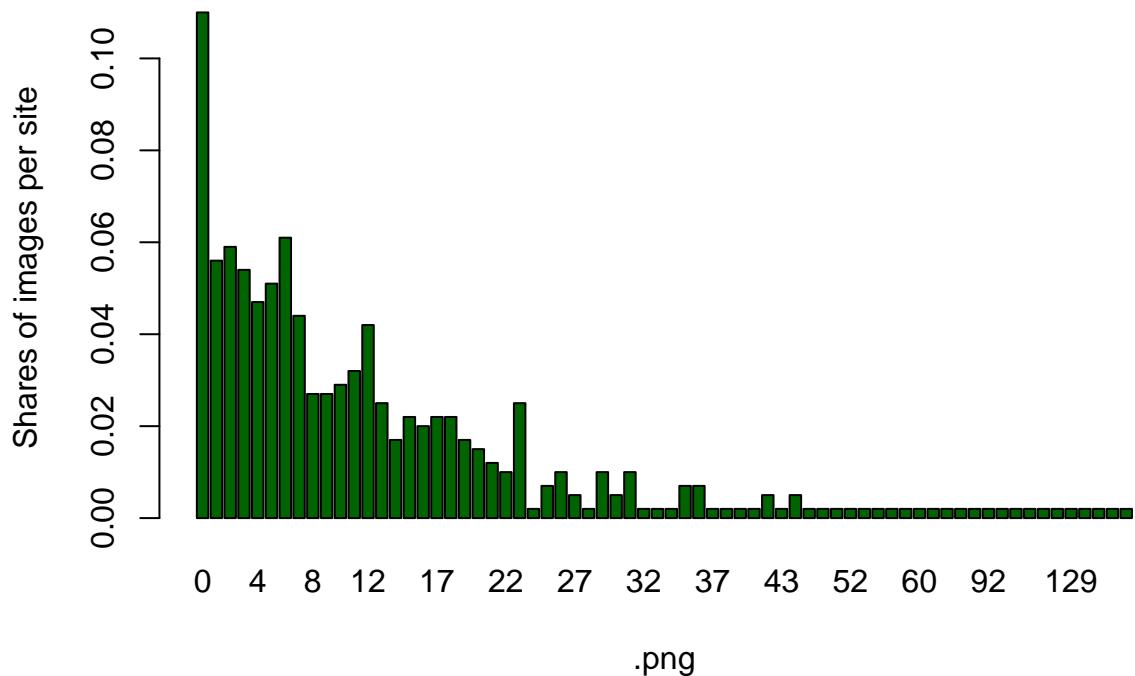


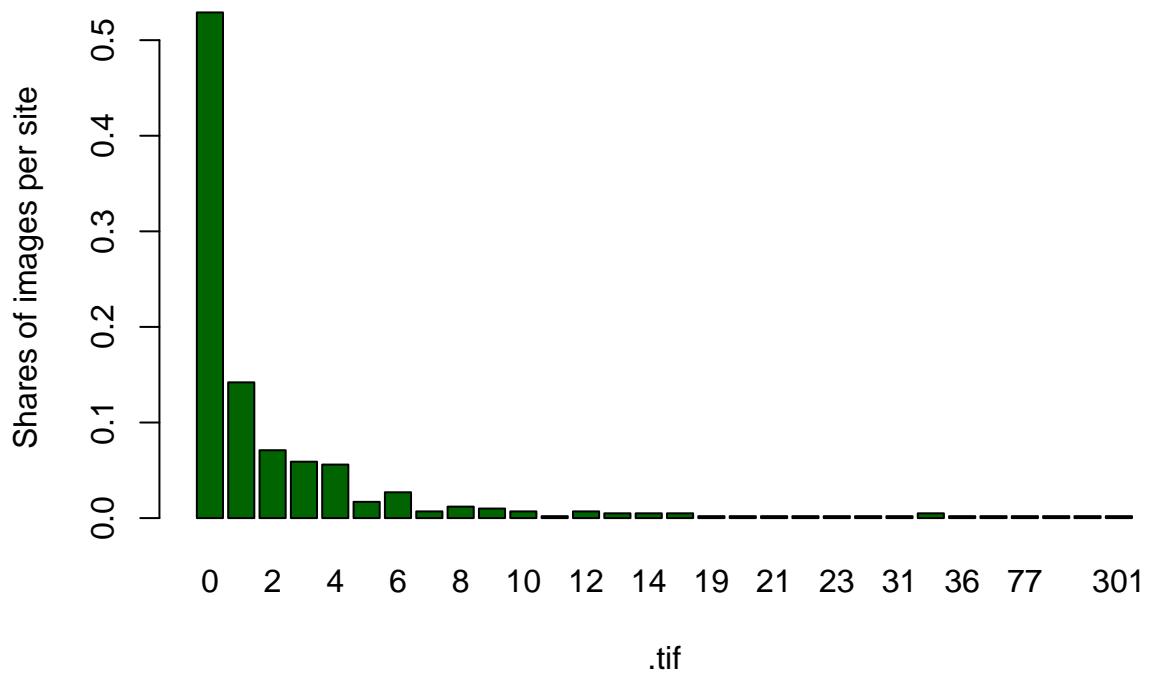


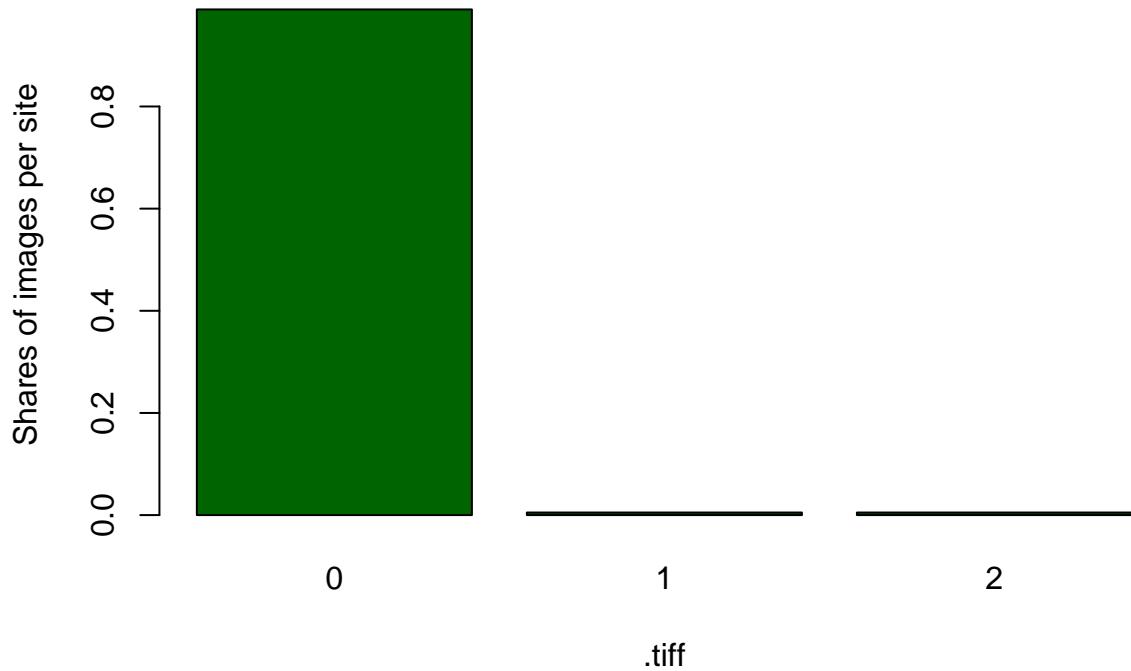




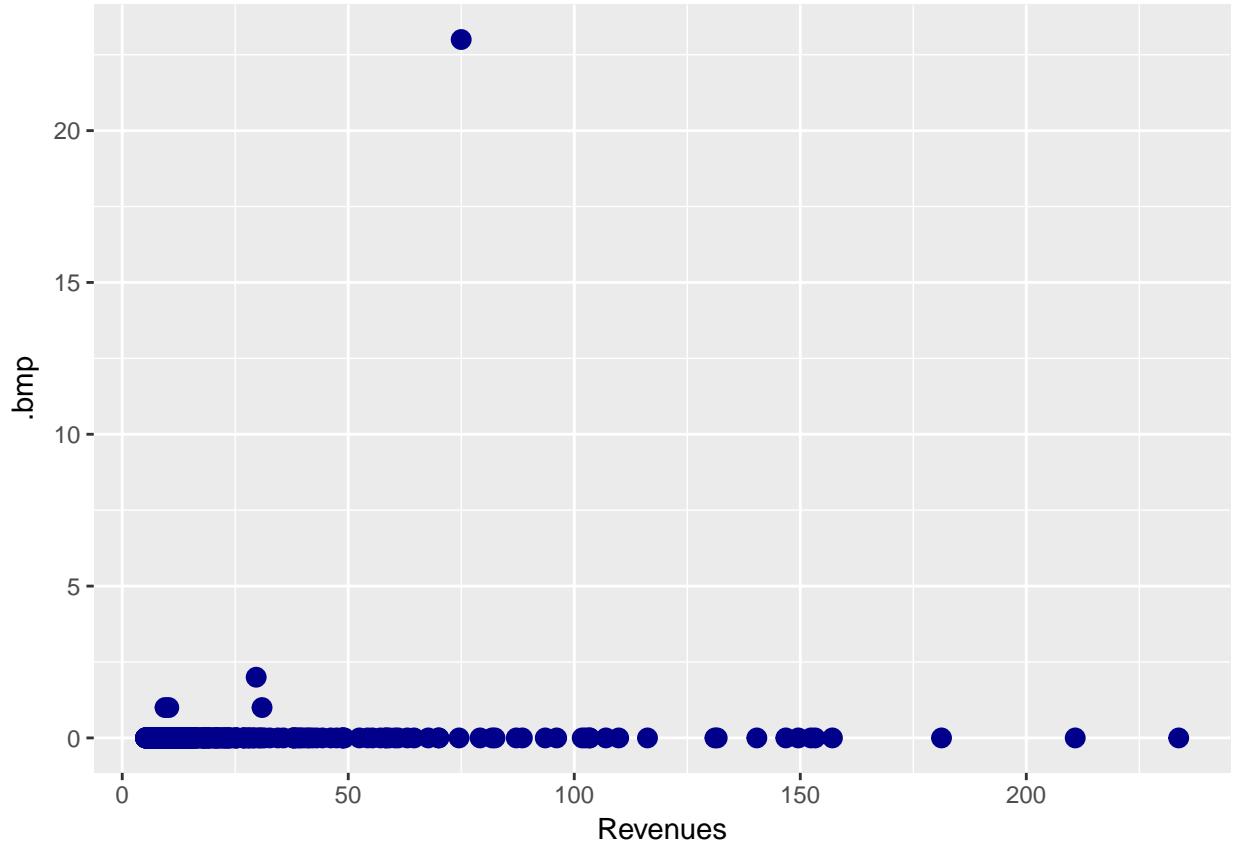


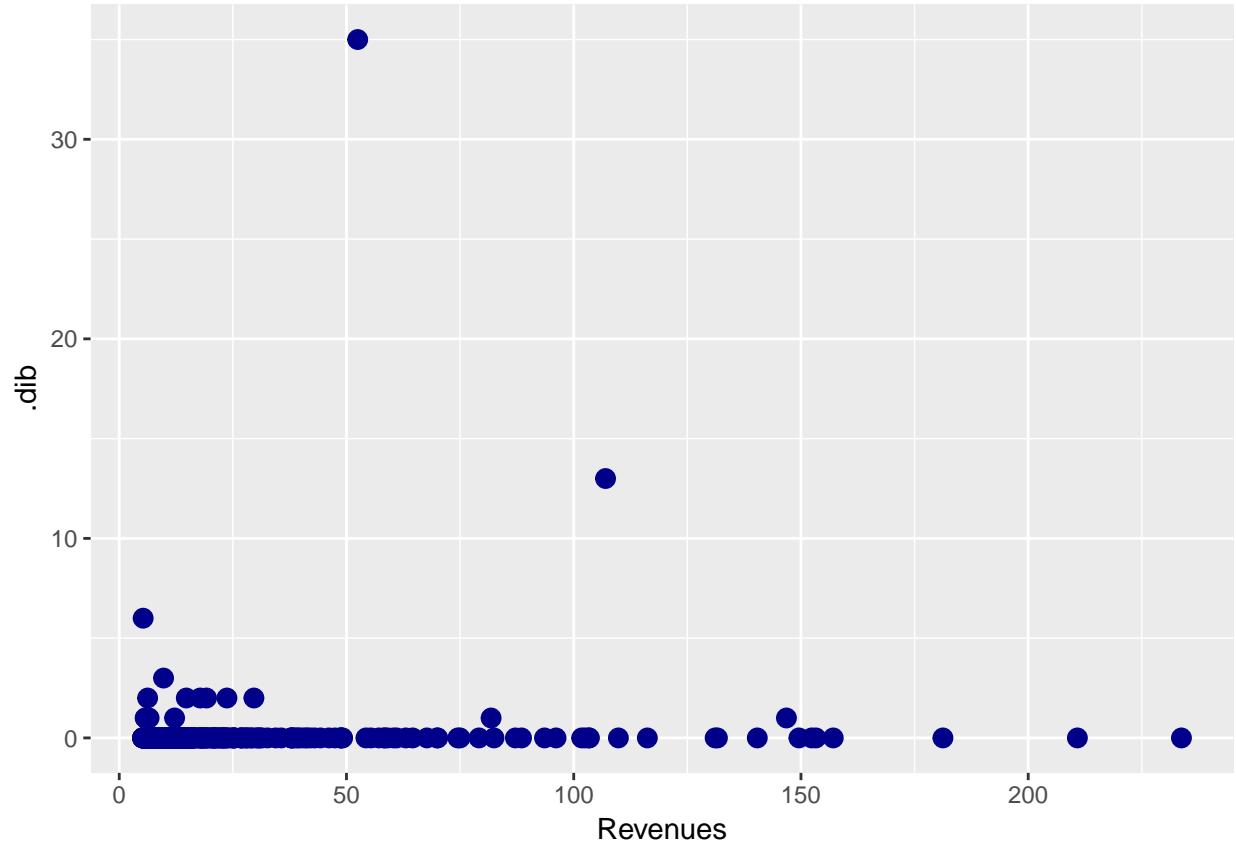


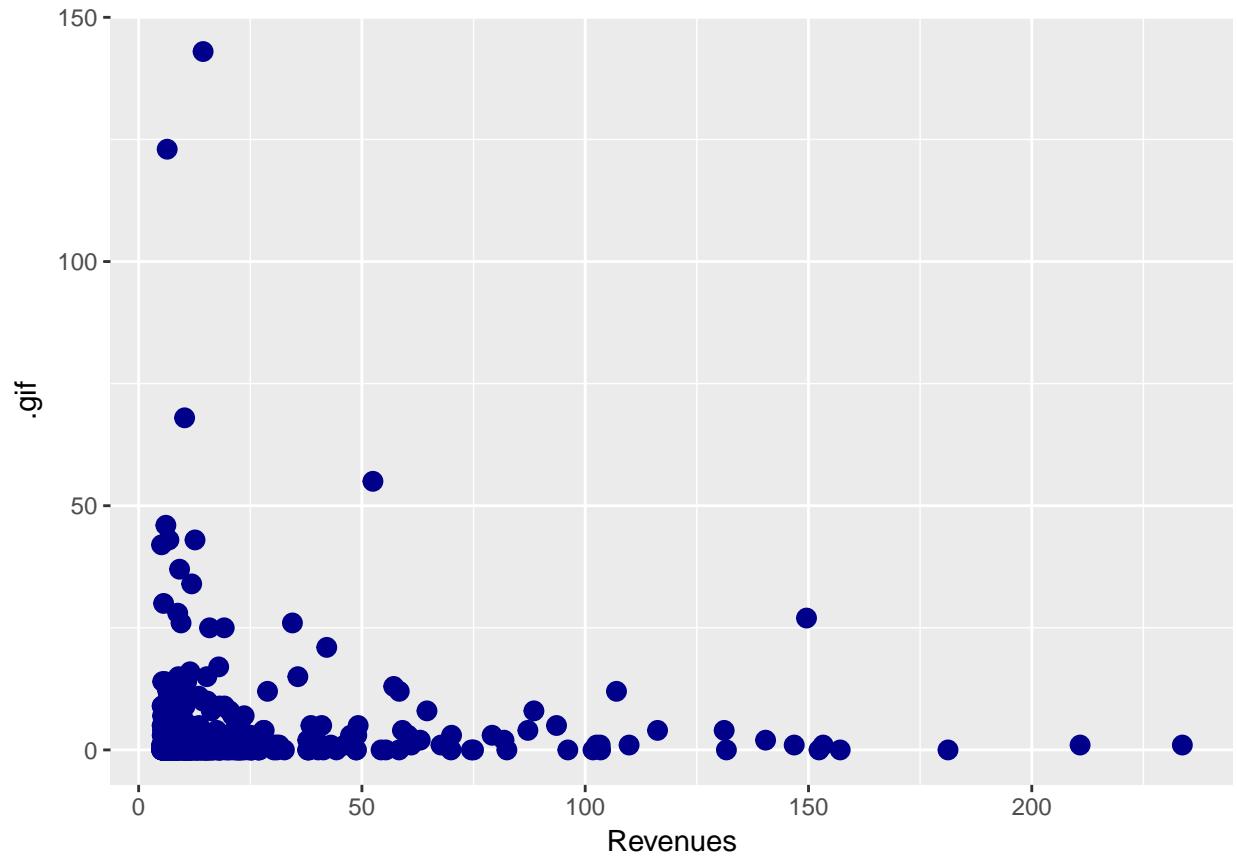




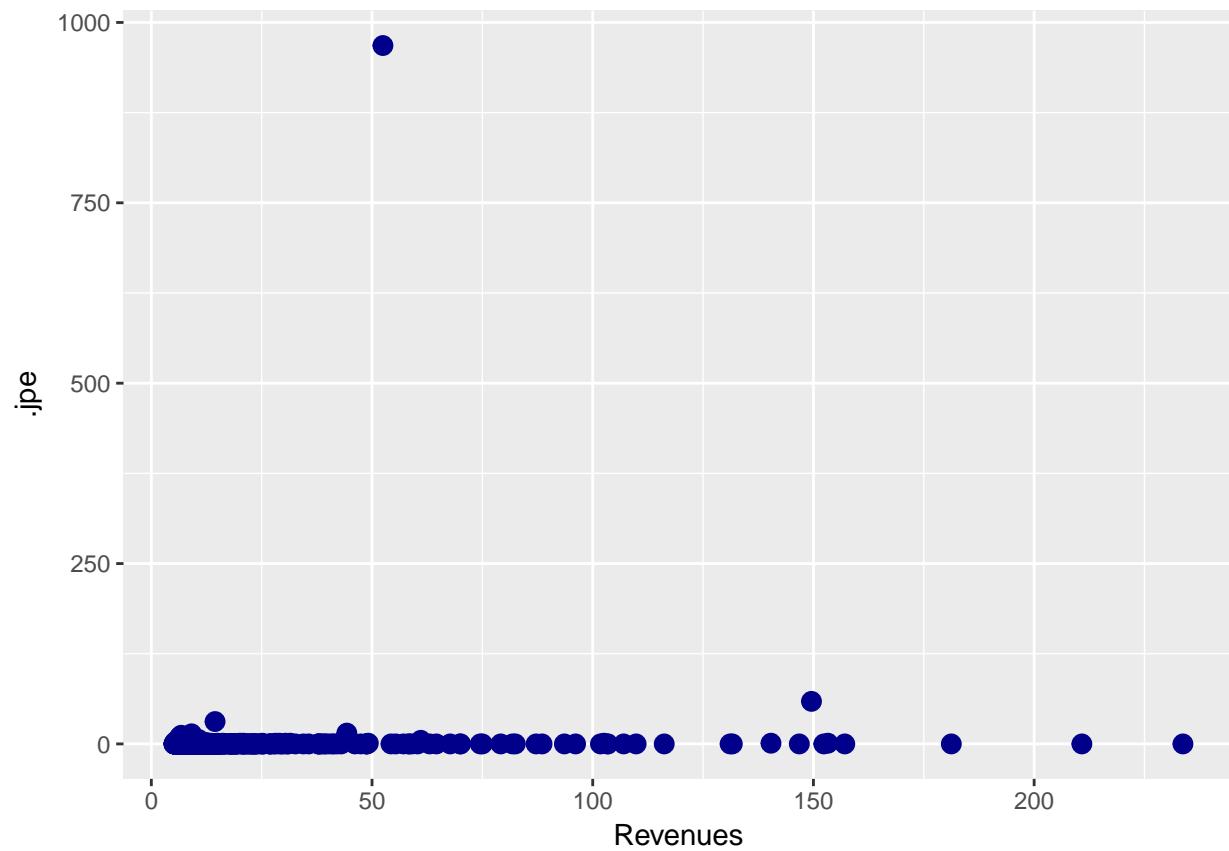
```
#It is obvious that the most common images type are .jpg, gif and .png  
#We will check now the types in relationship with the revenues  
ggplot(total_500_final, aes(Revenues, .bmp)) + geom_point(size=3, colour = "dark blue")
```

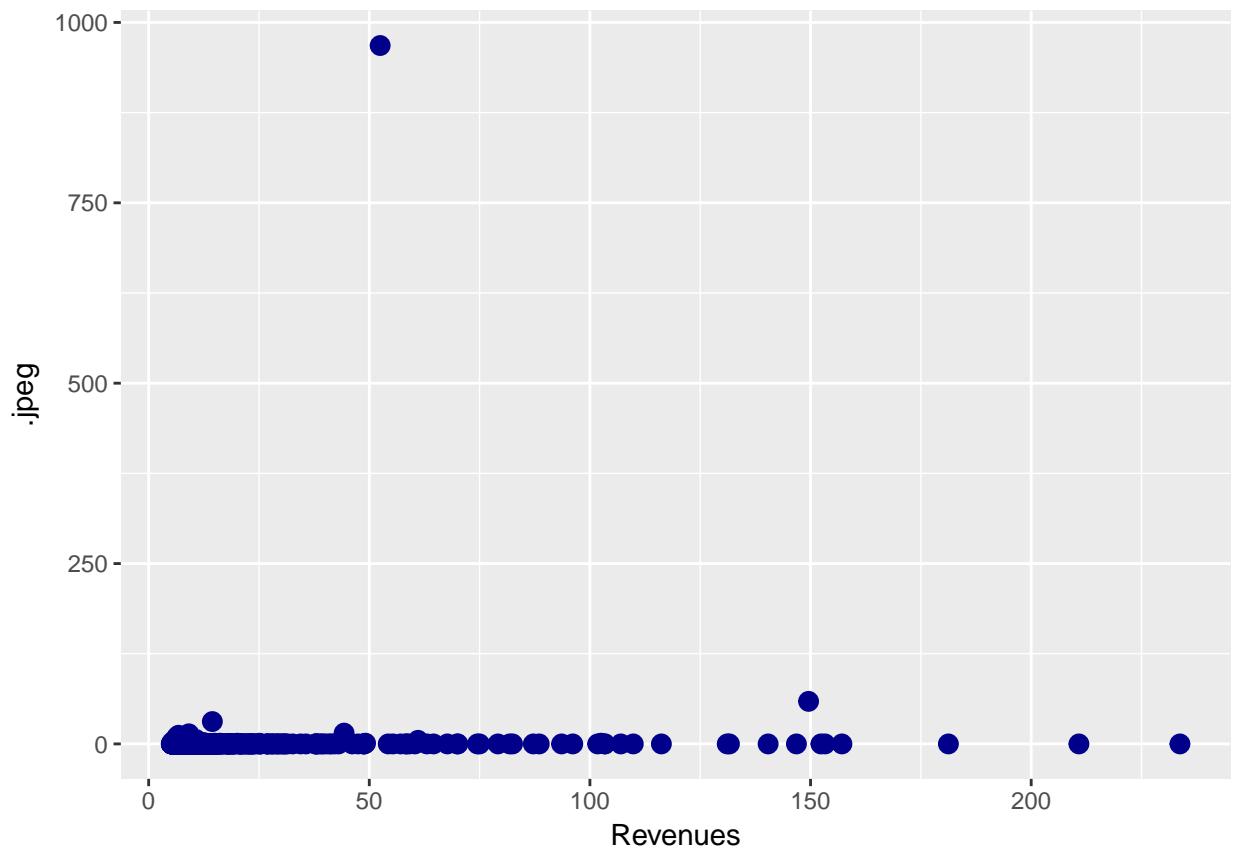


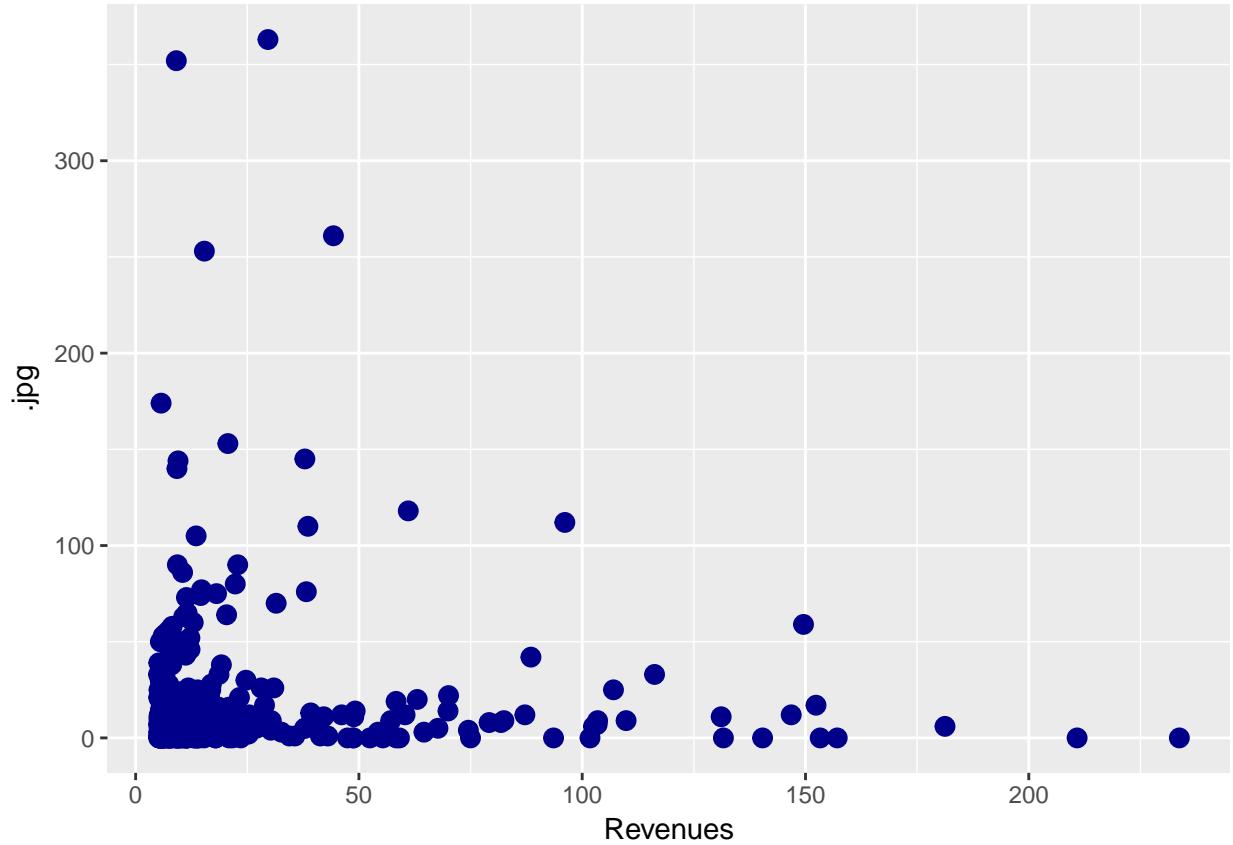


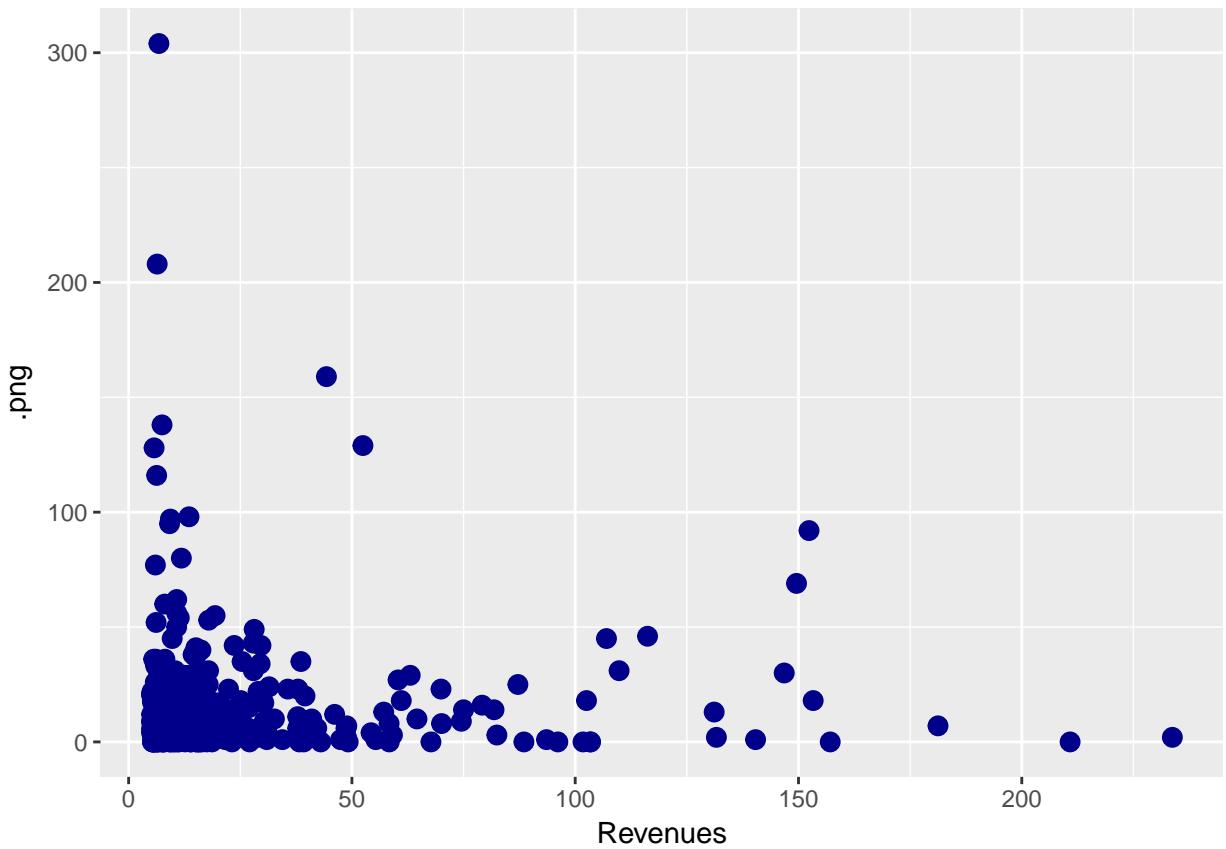


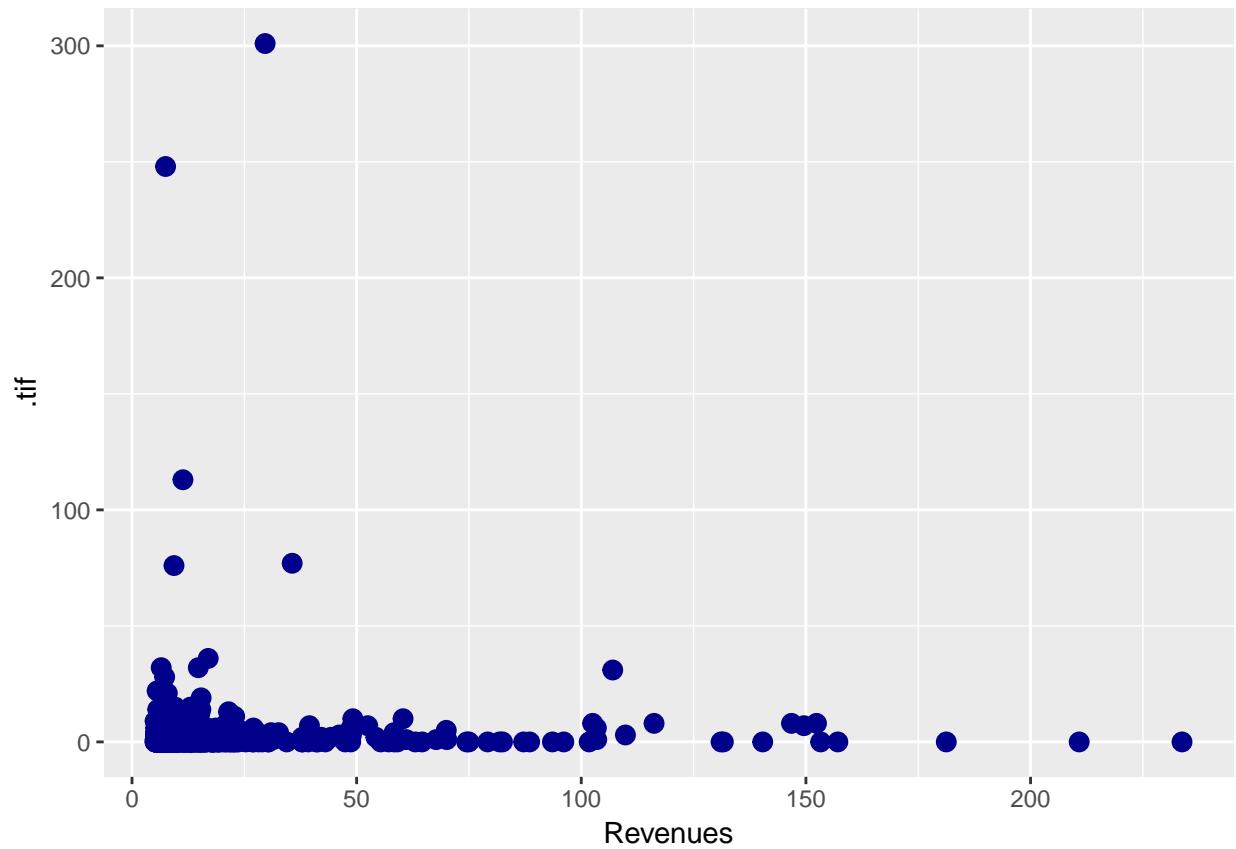
```
ggplot(total_500_final, aes(Revenues, .jpe)) + geom_point(size=3, colour = "dark blue")
```

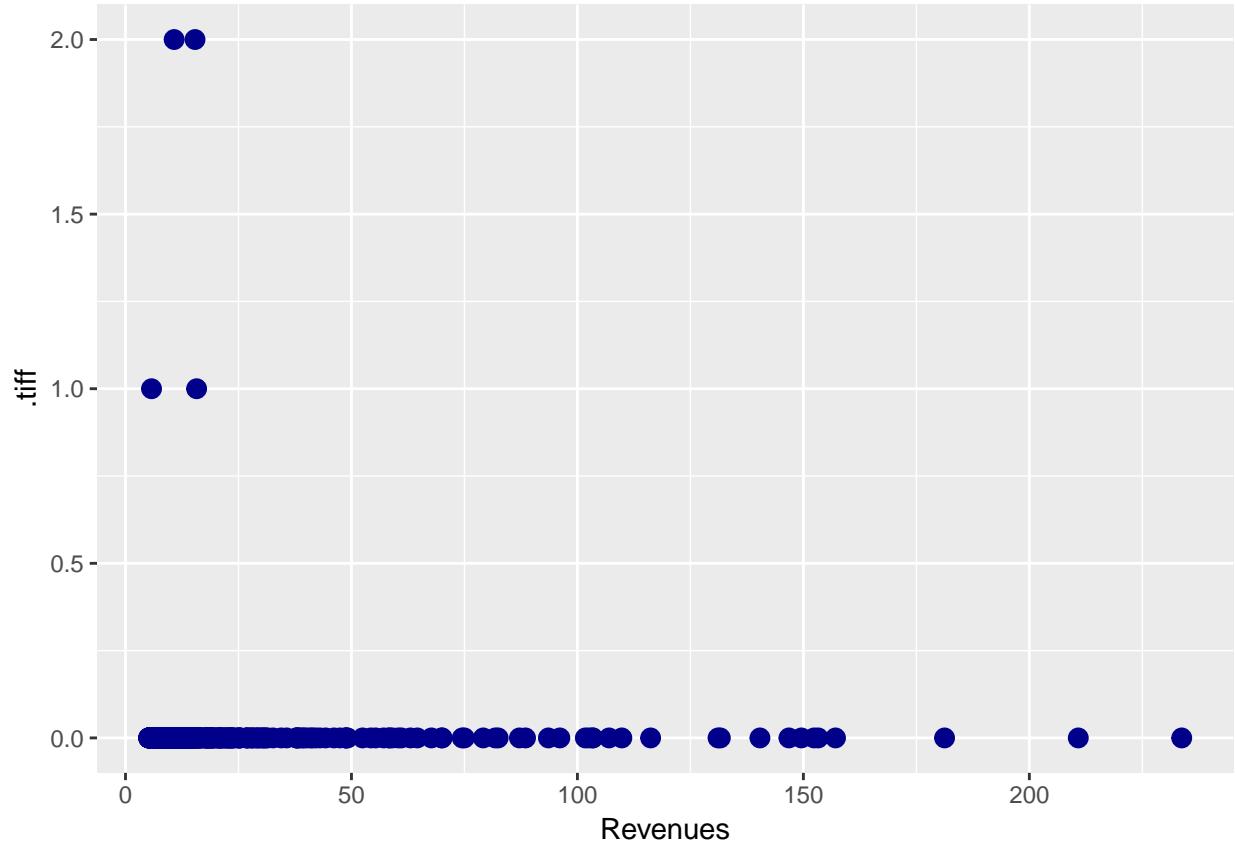












```
#And we can also see for correlations
total_500_im<- total_500_final[,c(4,717:726)]
library(corrplot)
library(caret)
tl <- cor(total_500_im)
tl

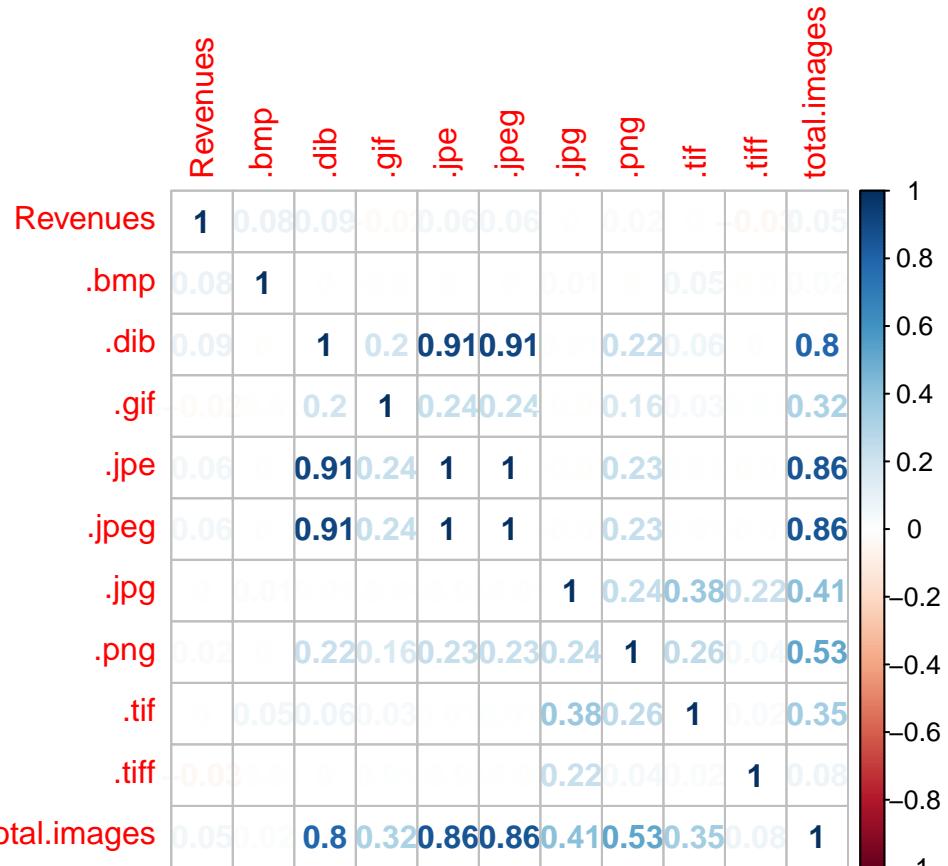
##          Revenues      .bmp      .dib      .gif
## Revenues 1.000000000 0.083489281 0.0877047067 -0.020930575
## .bmp      0.083489281 1.000000000 -0.0013011275 -0.005780172
## .dib      0.087704707 -0.001301127 1.0000000000 0.196433219
## .gif     -0.020930575 -0.005780172 0.1964332192 1.0000000000
## .jpe      0.059660288 -0.003534780 0.9108504455 0.235836802
## .jpeg     0.059427022 -0.003482535 0.9108963639 0.236016392
## .jpg      0.004515870  0.013995542 0.0081915788 -0.006939606
## .png      0.022053238 -0.001301001 0.2204575537 0.164538688
## .tif      -0.002466003  0.050383371 0.0628334455 0.031167772
## .tiff    -0.030165817 -0.005659409 -0.0008571371 0.007252077
## total.images 0.051605339  0.018395189 0.7961253175 0.319131935
##          .jpe      .jpeg      .jpg      .png
## Revenues 0.059660288 0.059427022 0.004515870 0.022053238
## .bmp     -0.003534780 -0.003482535 0.013995542 -0.001301001
## .dib      0.910850445 0.910896364 0.008191579 0.220457554
## .gif      0.235836802 0.236016392 -0.006939606 0.164538688
## .jpe      1.000000000 0.999991326 -0.008220505 0.231367993
## .jpeg     0.999991326 1.000000000 -0.008242305 0.231422560
## .jpg     -0.008220505 -0.008242305 1.000000000 0.244033499
```

```

## .png          0.231367993  0.231422560  0.244033499  1.000000000
## .tif          0.007431086  0.007587510  0.375321187  0.259001712
## .tiff         -0.005630343 -0.005548159  0.224429140  0.040635288
## total.images  0.855175392  0.855225420  0.413319367  0.529706228
##               .tif          .tiff      total.images
## Revenues     -0.002466003 -0.0301658169  0.05160534
## .bmp          0.050383371 -0.0056594093  0.01839519
## .dib          0.062833445 -0.0008571371  0.79612532
## .gif          0.031167772  0.0072520772  0.31913194
## .jpe          0.007431086 -0.0056303426  0.85517539
## .jpeg         0.007587510 -0.0055481589  0.85522542
## .jpg          0.375321187  0.2244291400  0.41331937
## .png          0.259001712  0.0406352880  0.52970623
## .tif          1.0000000000  0.0222123897  0.34887215
## .tiff         0.022212390  1.0000000000  0.07827891
## total.images  0.348872154  0.0782789113  1.000000000

```

```
corrplot(cor(total_500_im),method="number")
```



```
#We will see now the frequency of image sizes that is being used
```

```

k = c()
#Check for sizes that are half and half divided in existing and not
for(i in 24:716){
  image_size<- round(table(total_500_final[,i]))
  if ((image_size[[1]]==408)==TRUE){
    k <- union(k, c(i))
  }
}
```

```

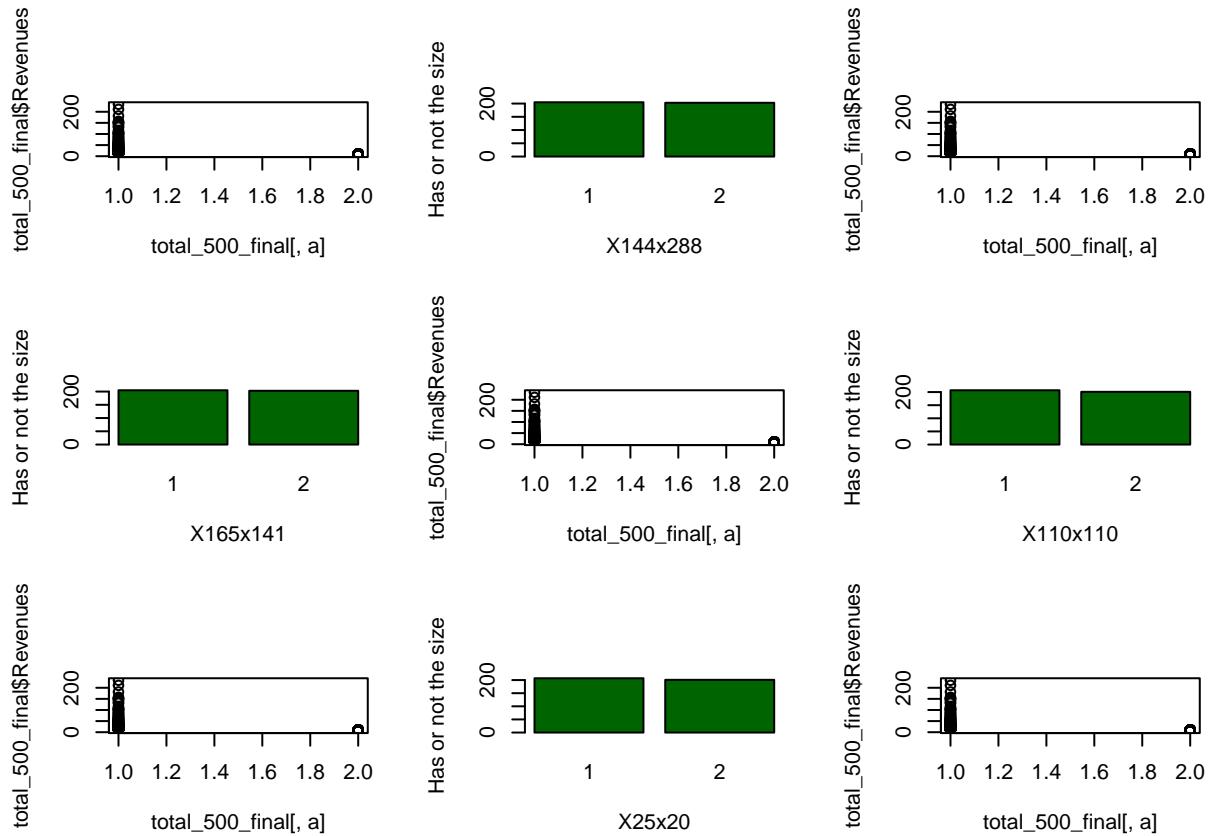
#####
#Number 24 is all own price so we want use it
names(total_500_final)[24]

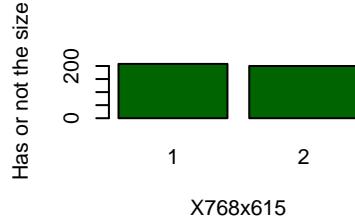
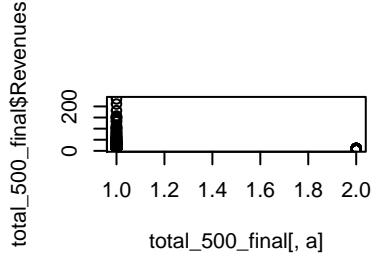
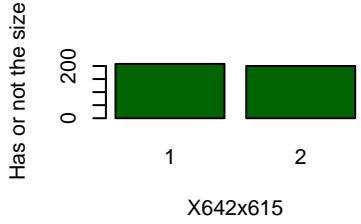
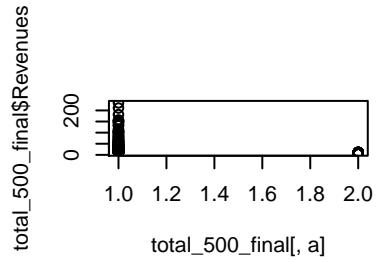
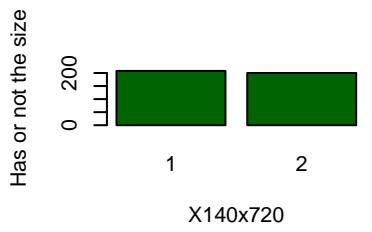
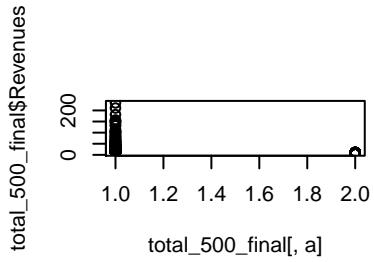
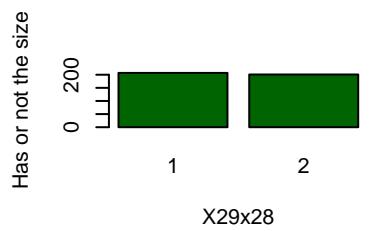
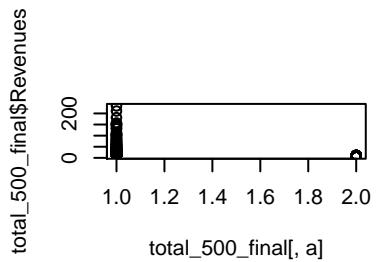
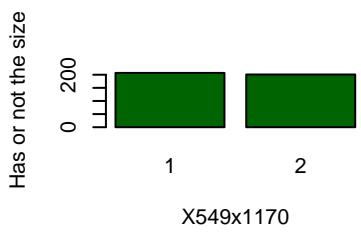
## [1] "X144x144"

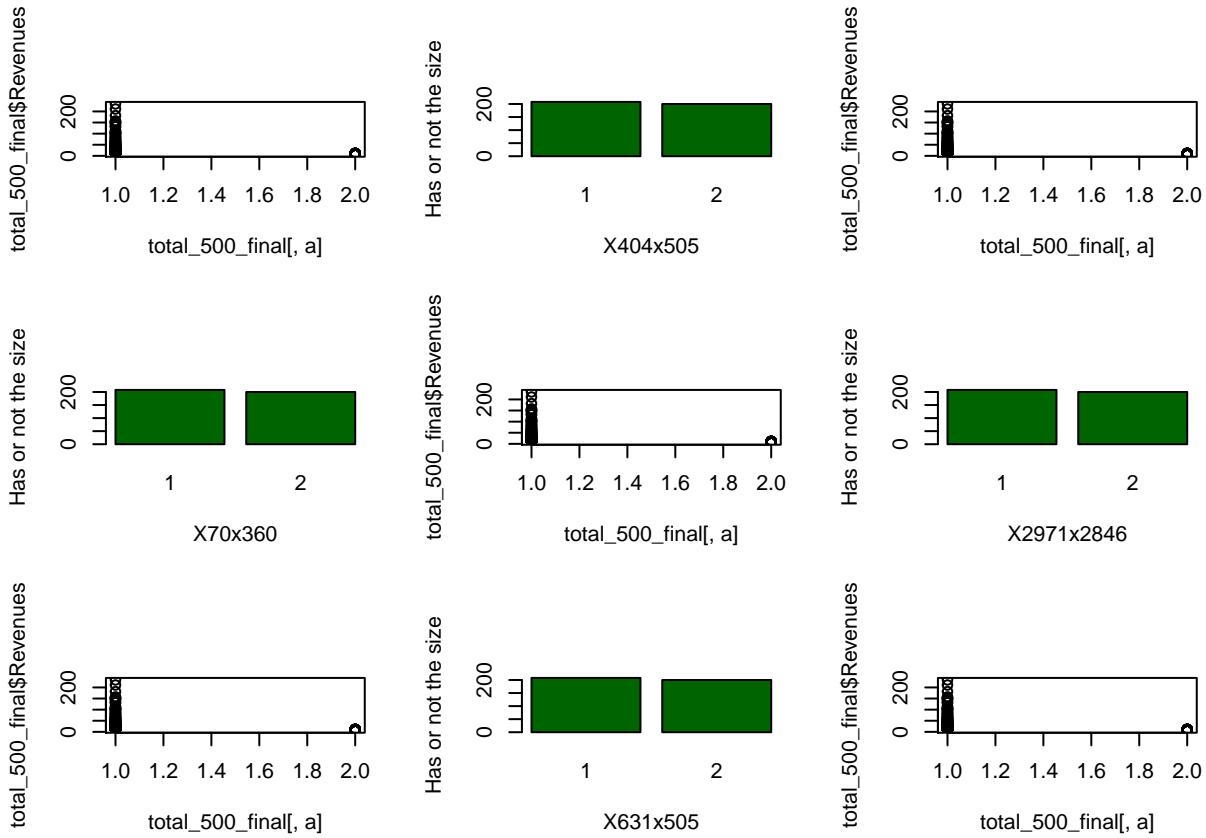
total_500_final$X144x144 <- NULL

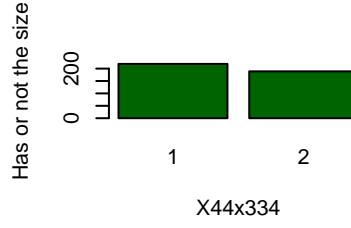
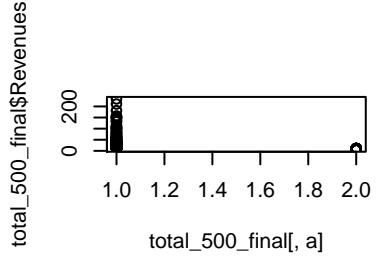
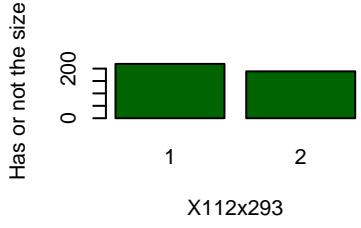
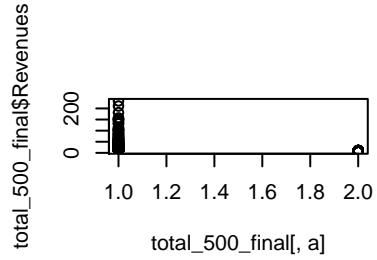
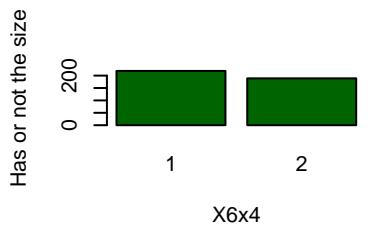
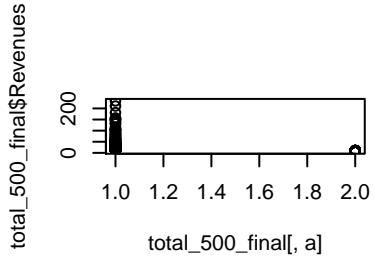
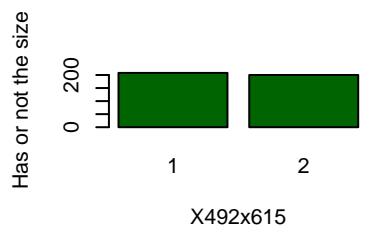
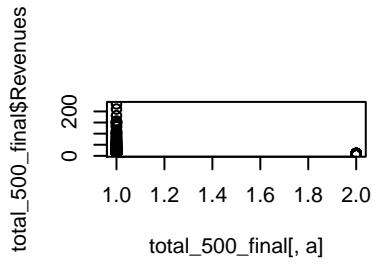
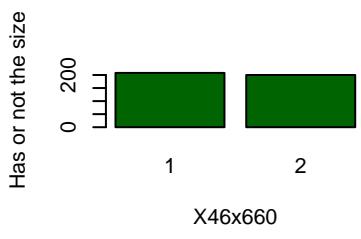
false_not_existing = c()
#Check for sizes that are less than half divided in existing and not
for(i in 24:715){
  image_size<- round(table(total_500_final[,i]))
  if ((image_size[2]<204)==TRUE){
    false_not_existing <- union(false_not_existing, c(i))
  }
#####
#Now we will take the sizes that exist in less than half the instances and check graphically the deviation
par(mfrow=c(3,3))
for(i in 1:416){
  a = false_not_existing[i]
  plot(total_500_final[,a],total_500_final$Revenues)
  image_size<- round(table(total_500_final[,a]))
  barplot(image_size,xlab=names(total_500_final)[a],ylab = "Has or not the size", col = "dark green")
}

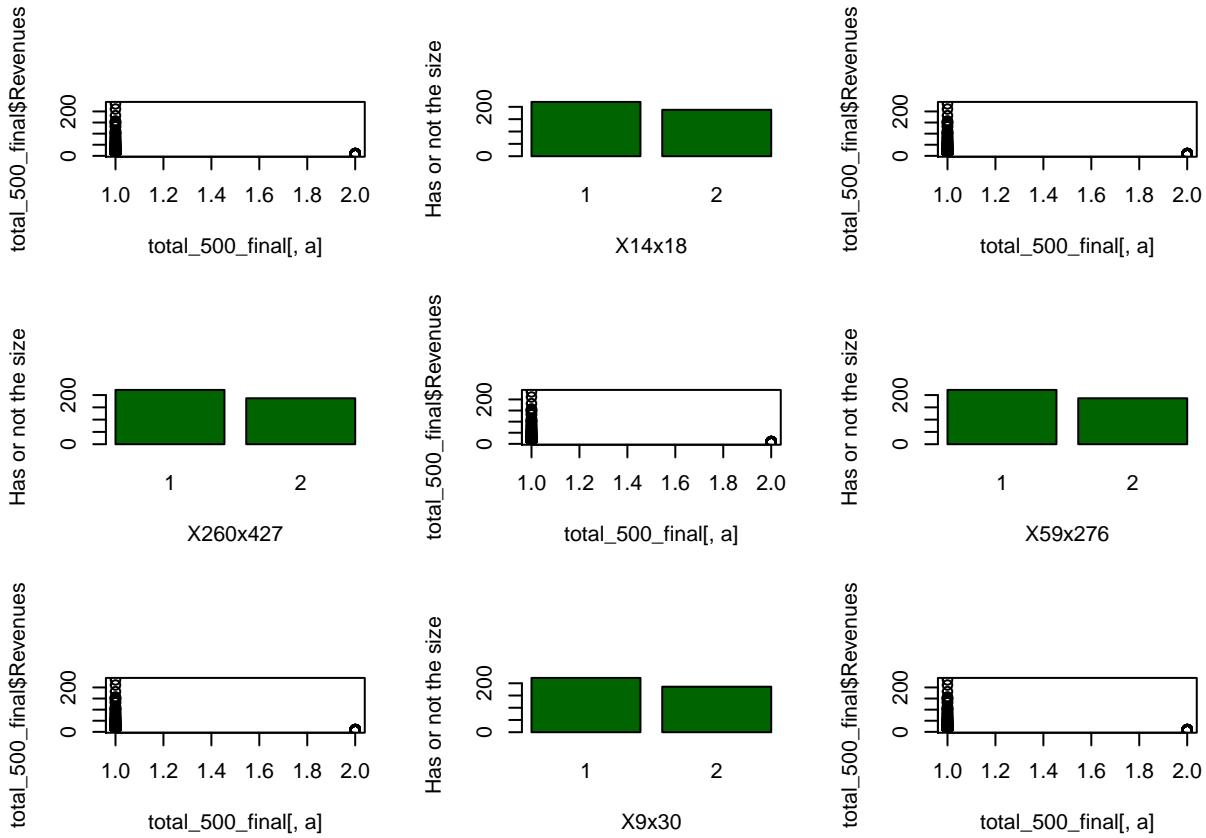
```

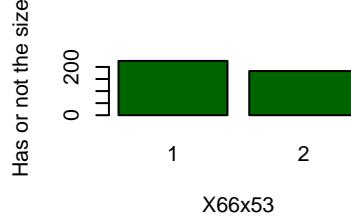
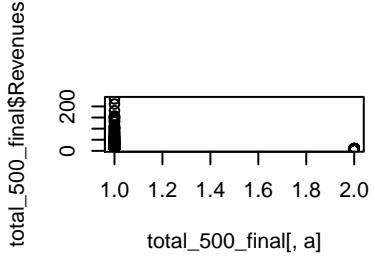
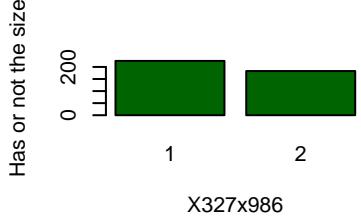
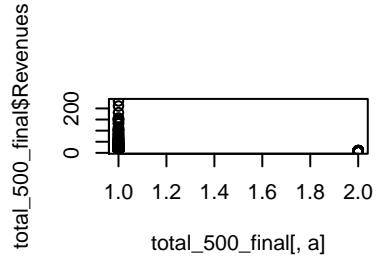
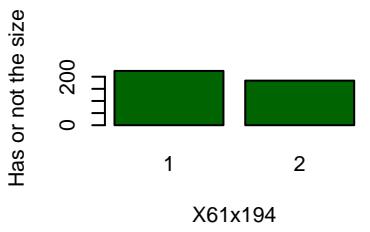
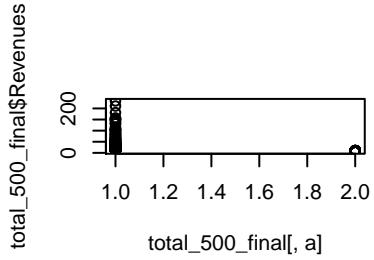
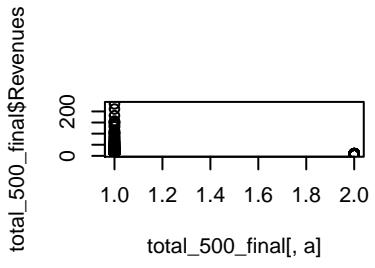
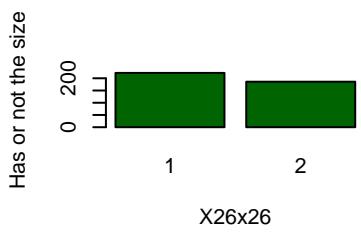


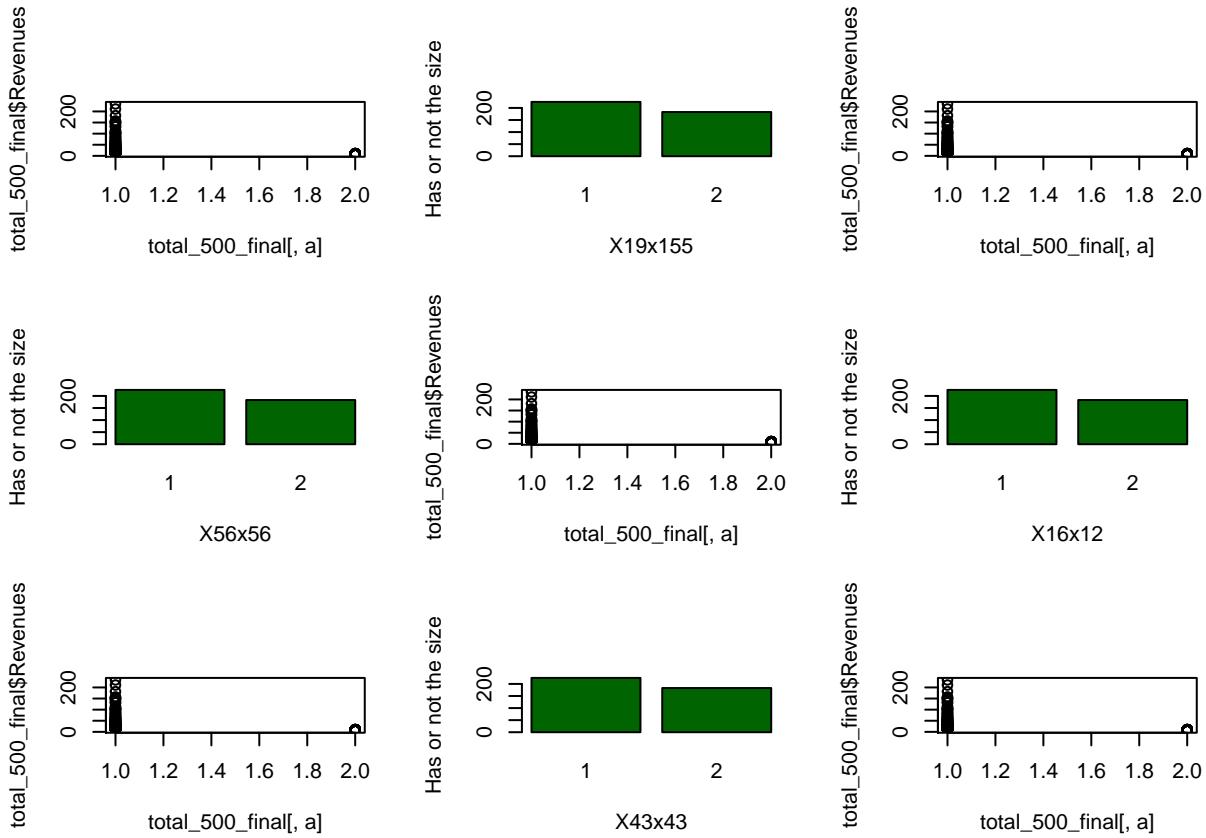


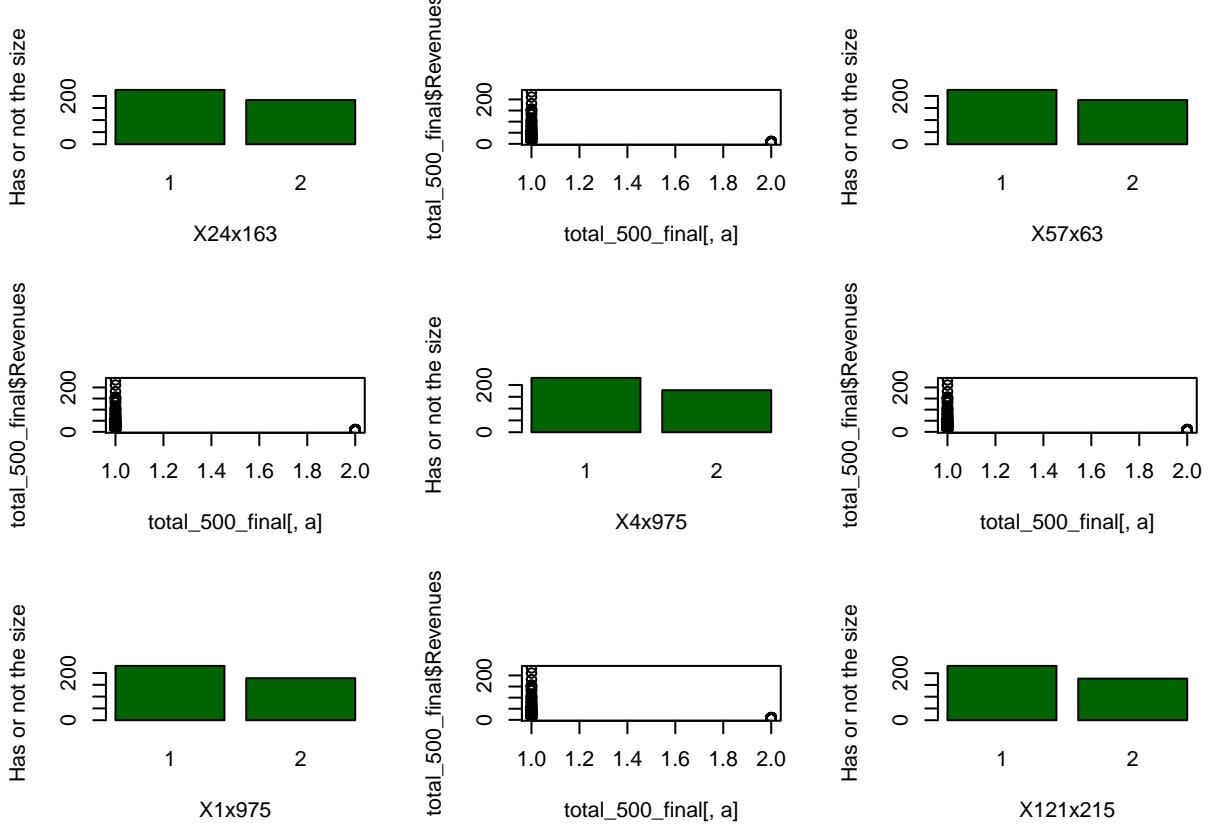


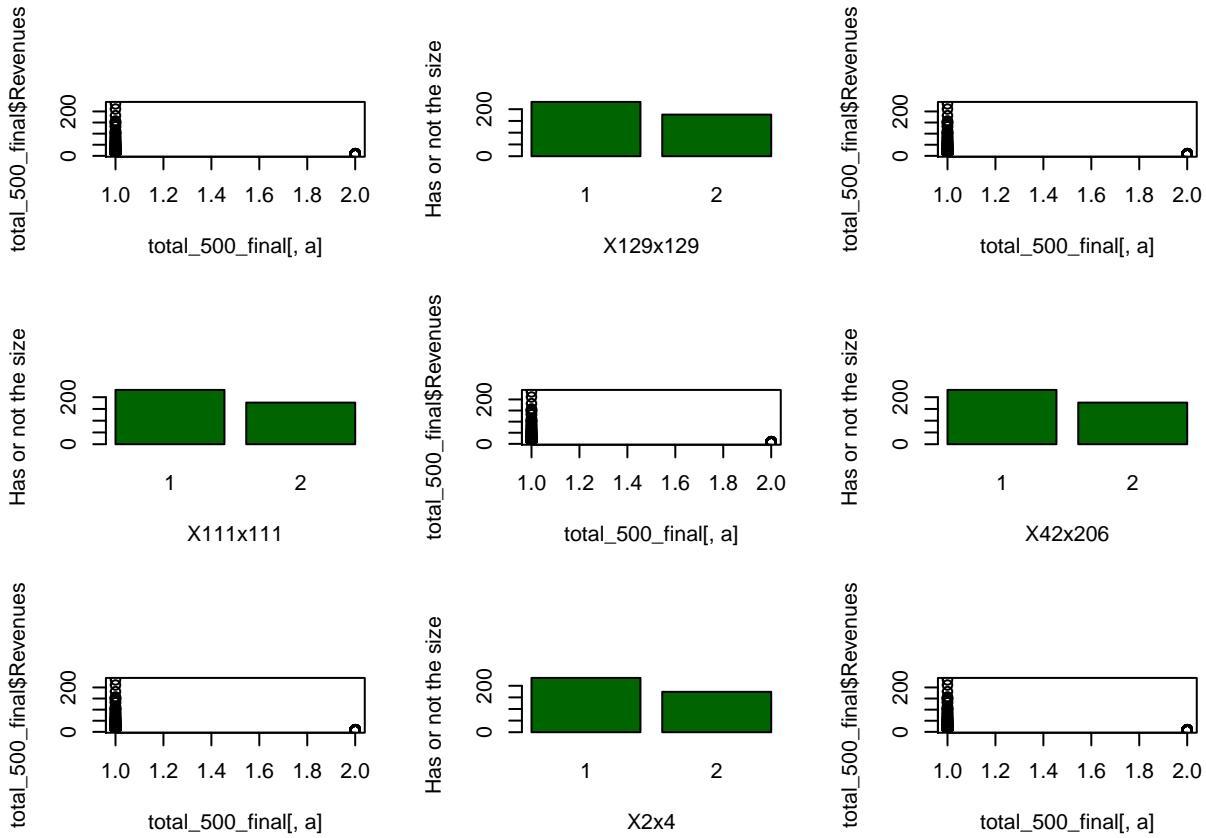


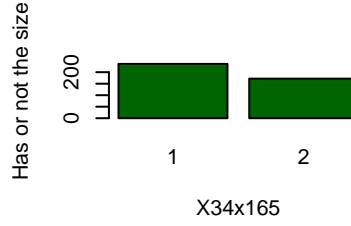
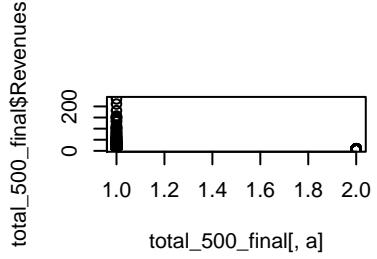
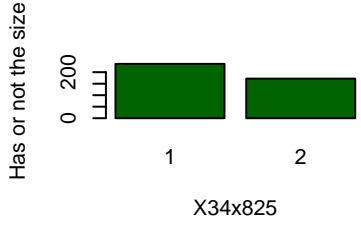
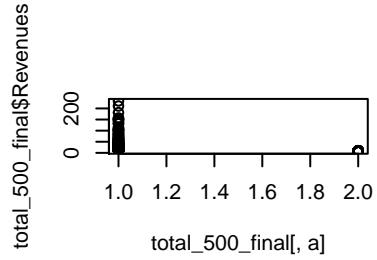
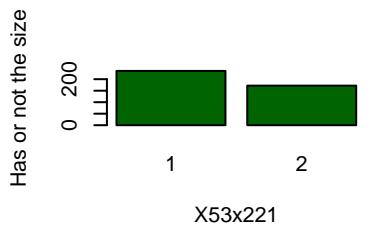
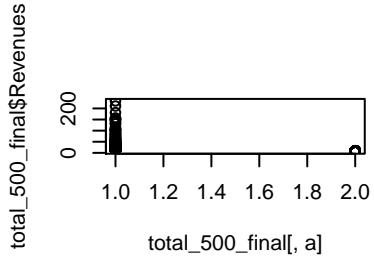
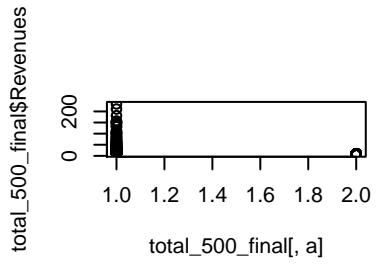
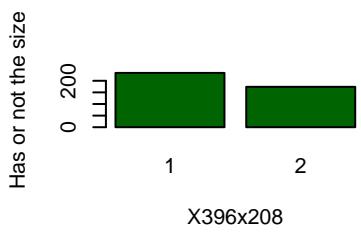


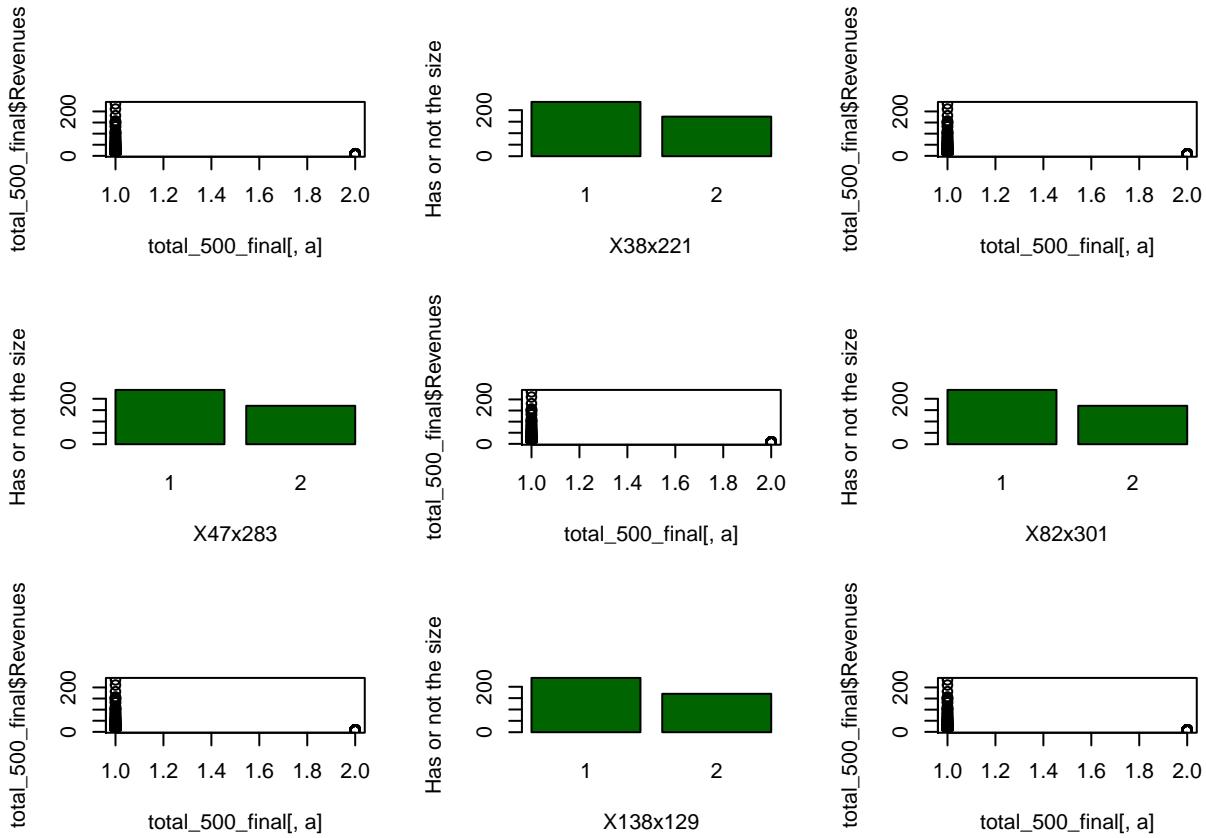


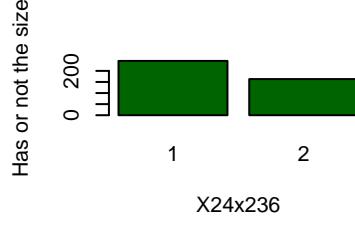
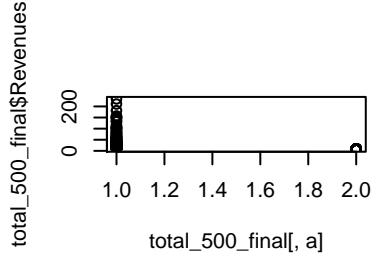
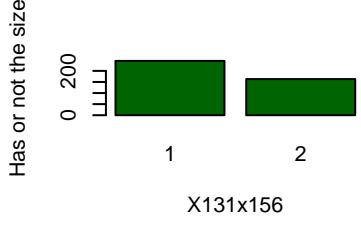
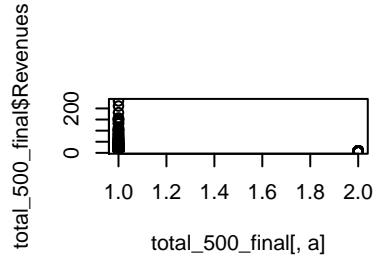
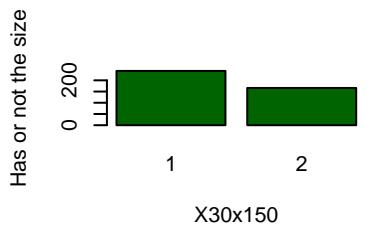
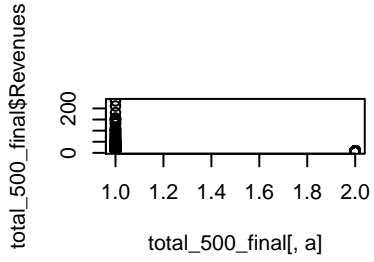
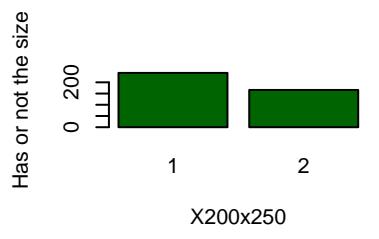
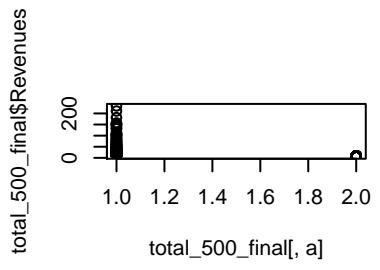
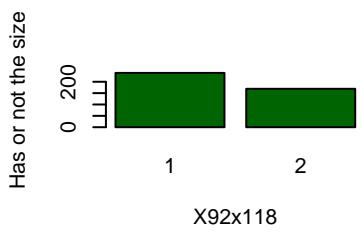


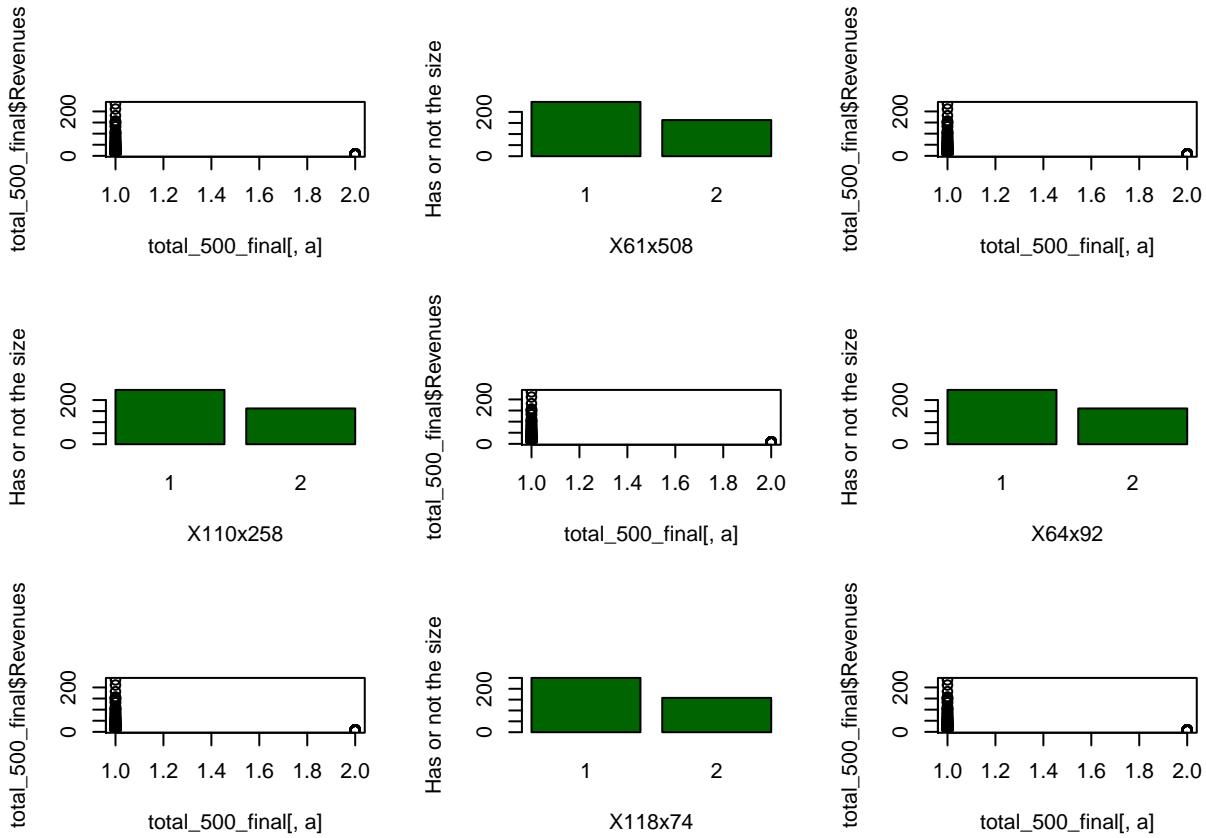


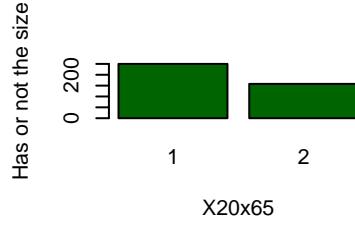
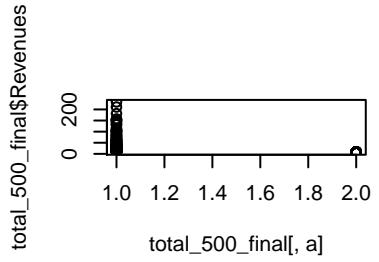
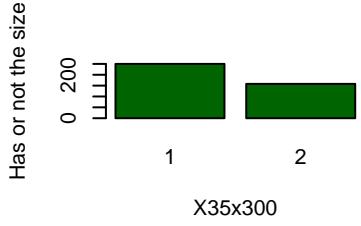
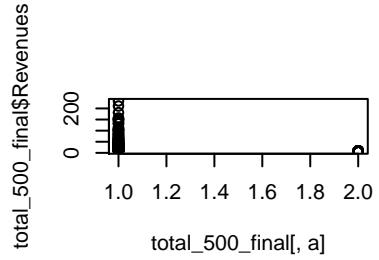
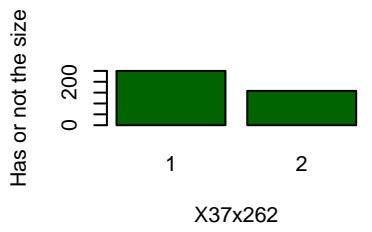
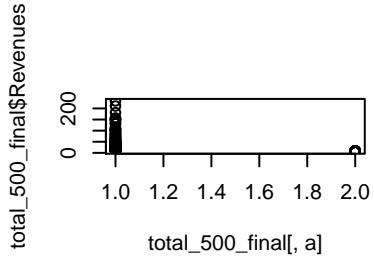
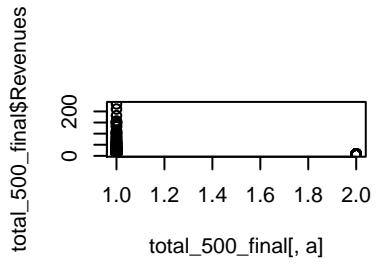
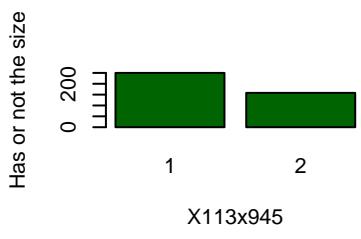


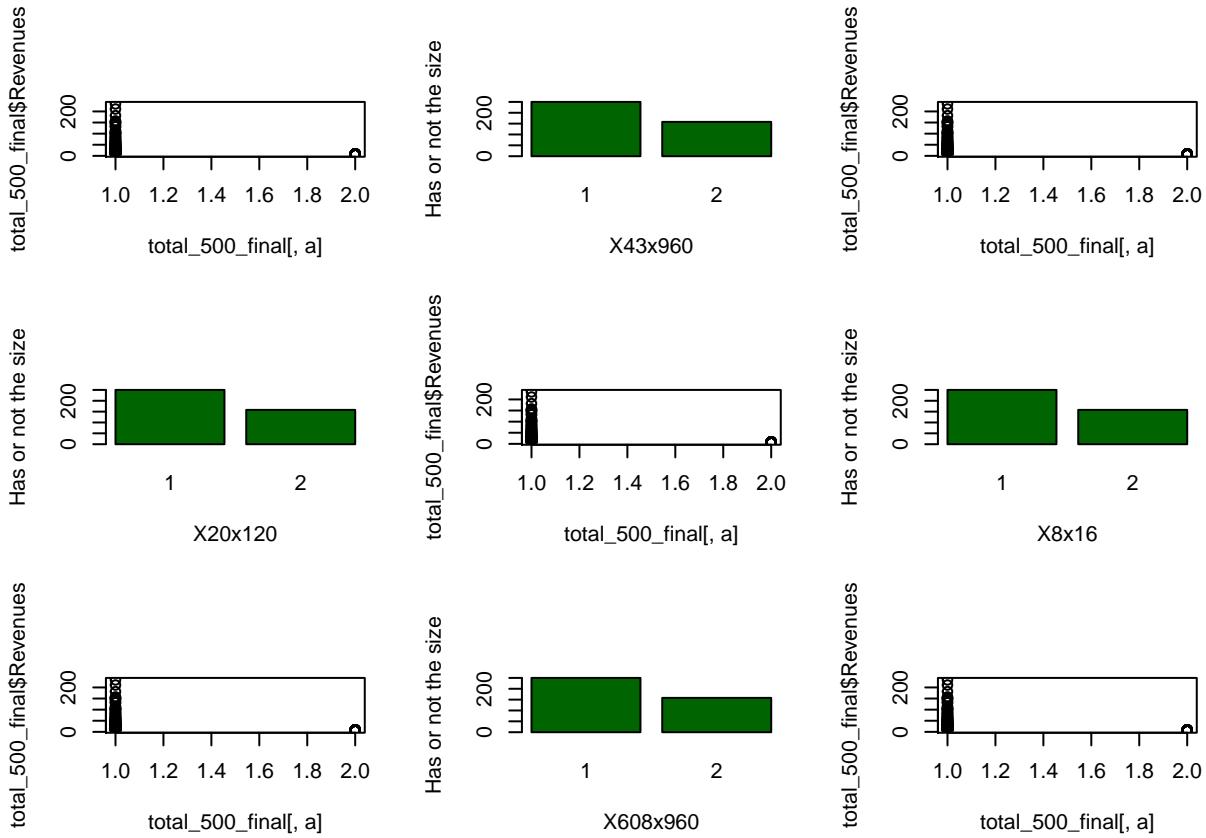


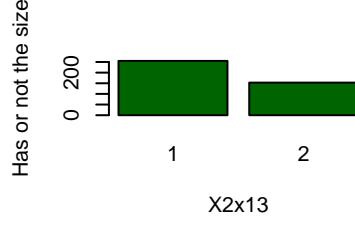
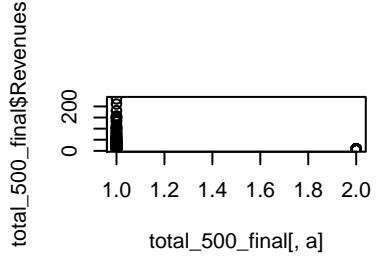
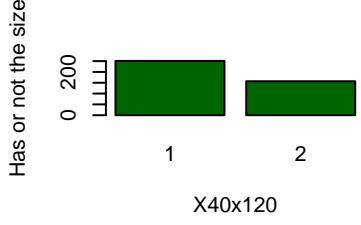
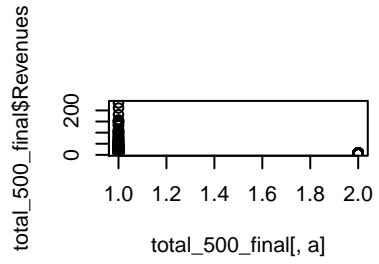
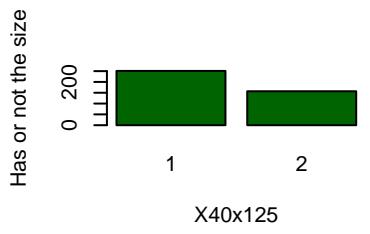
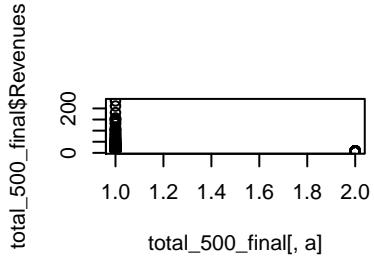
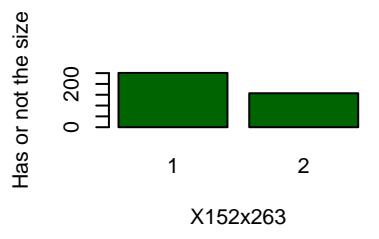
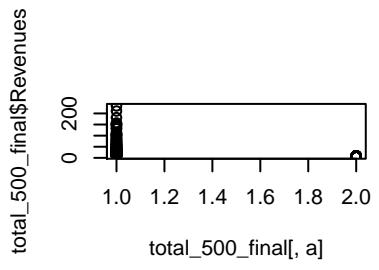
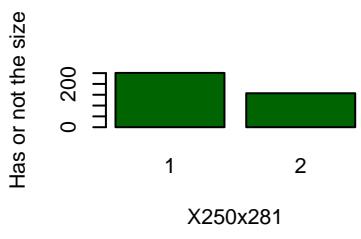


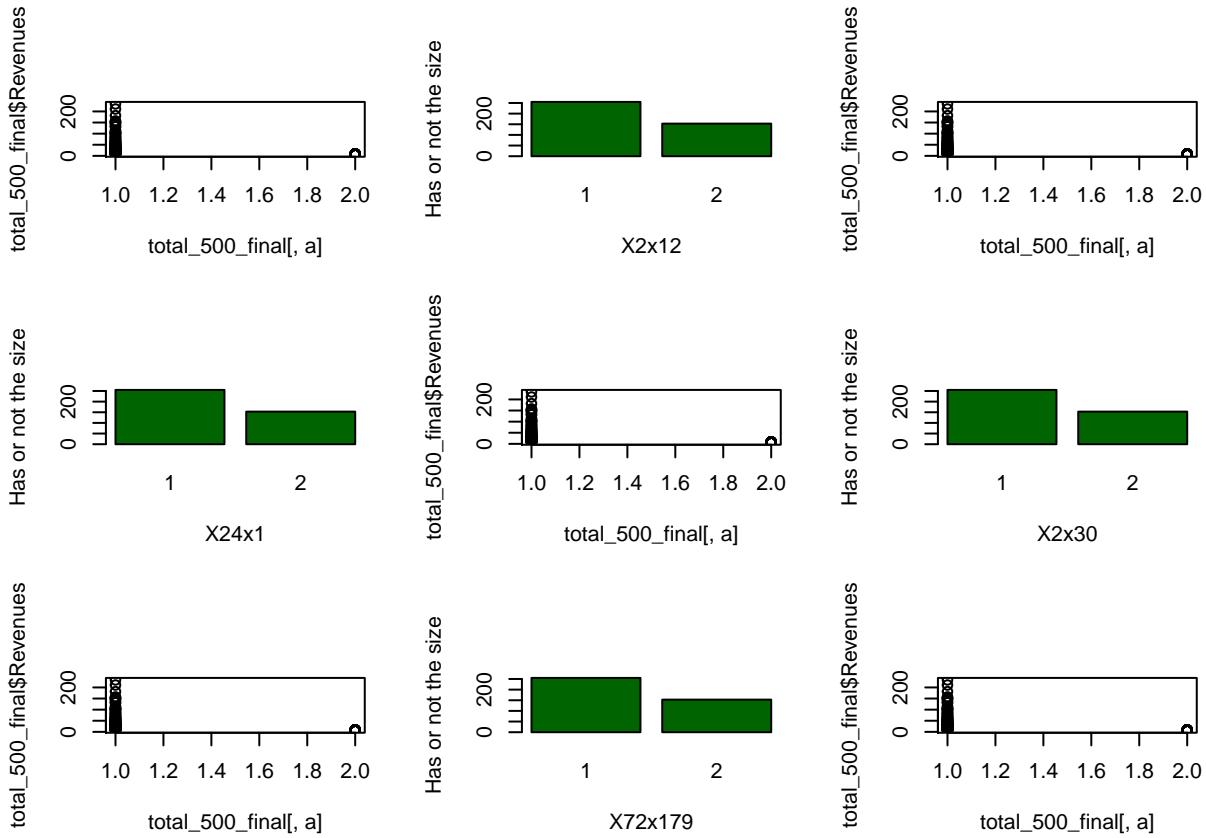


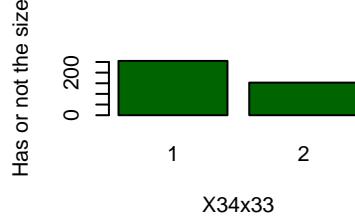
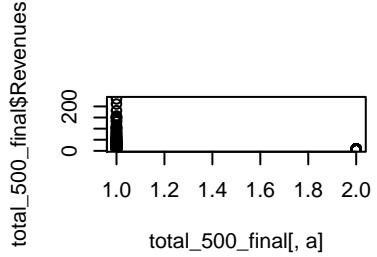
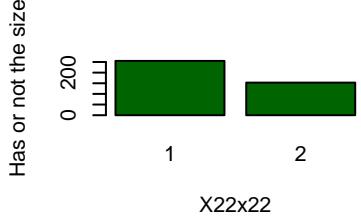
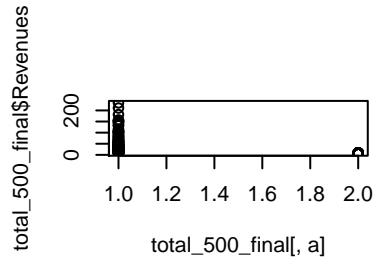
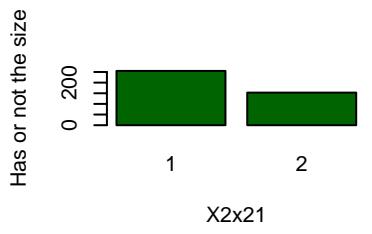
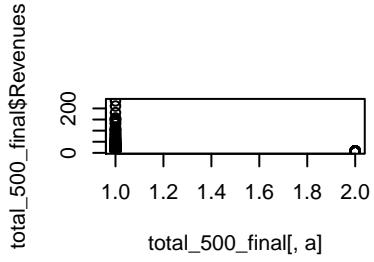
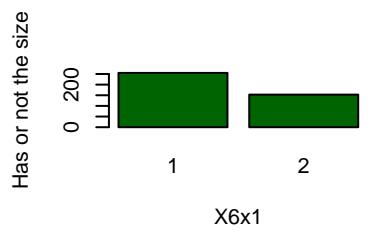
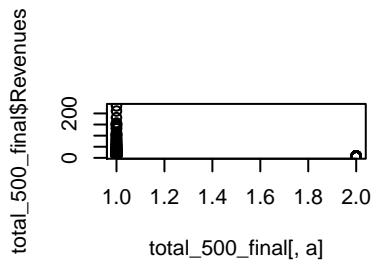
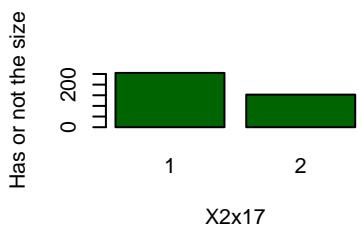


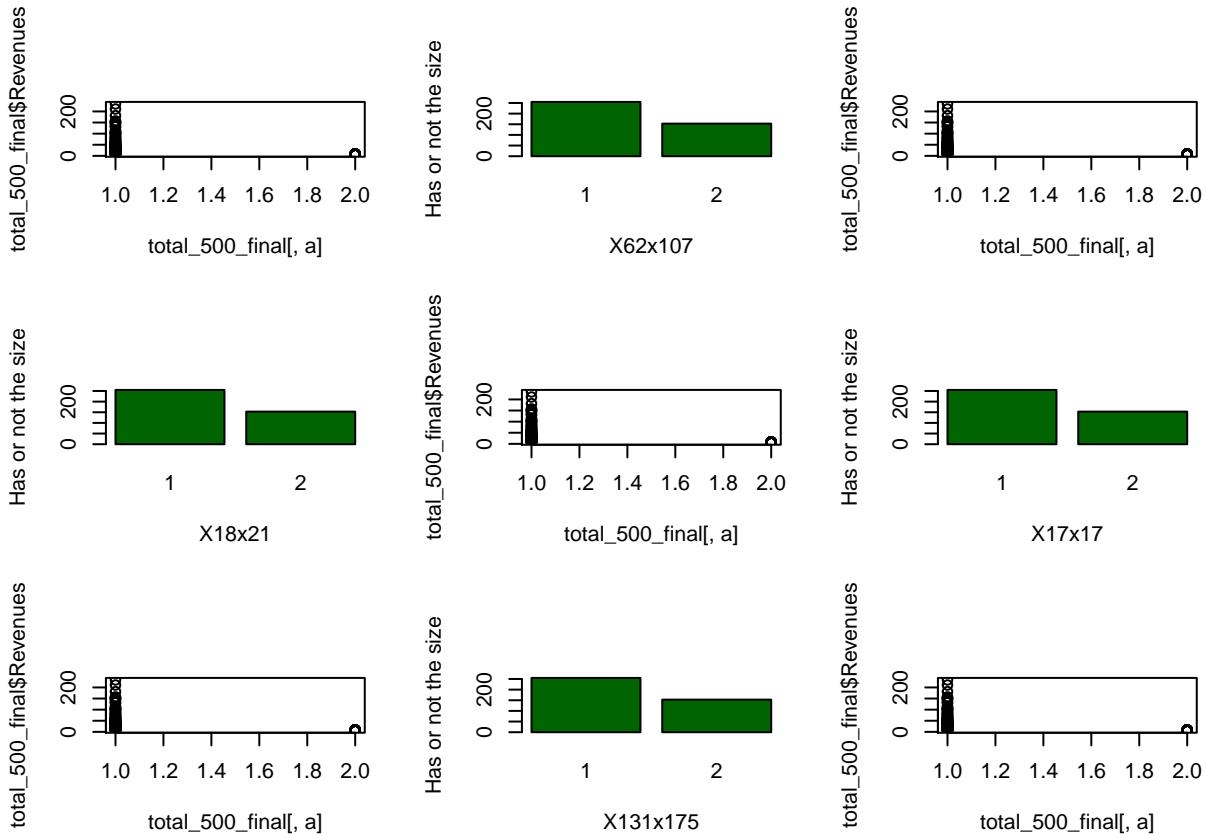


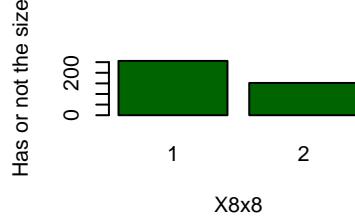
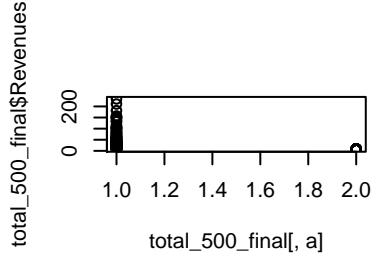
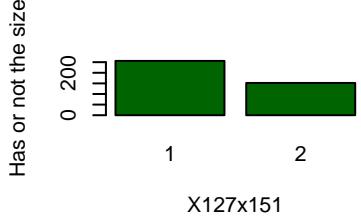
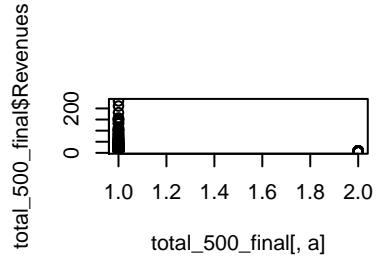
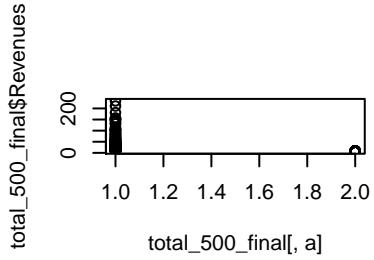
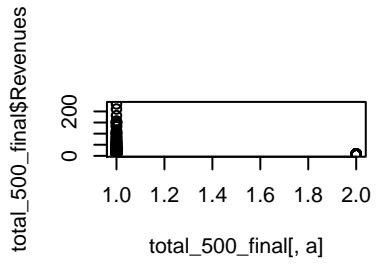
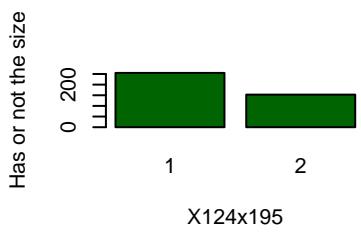


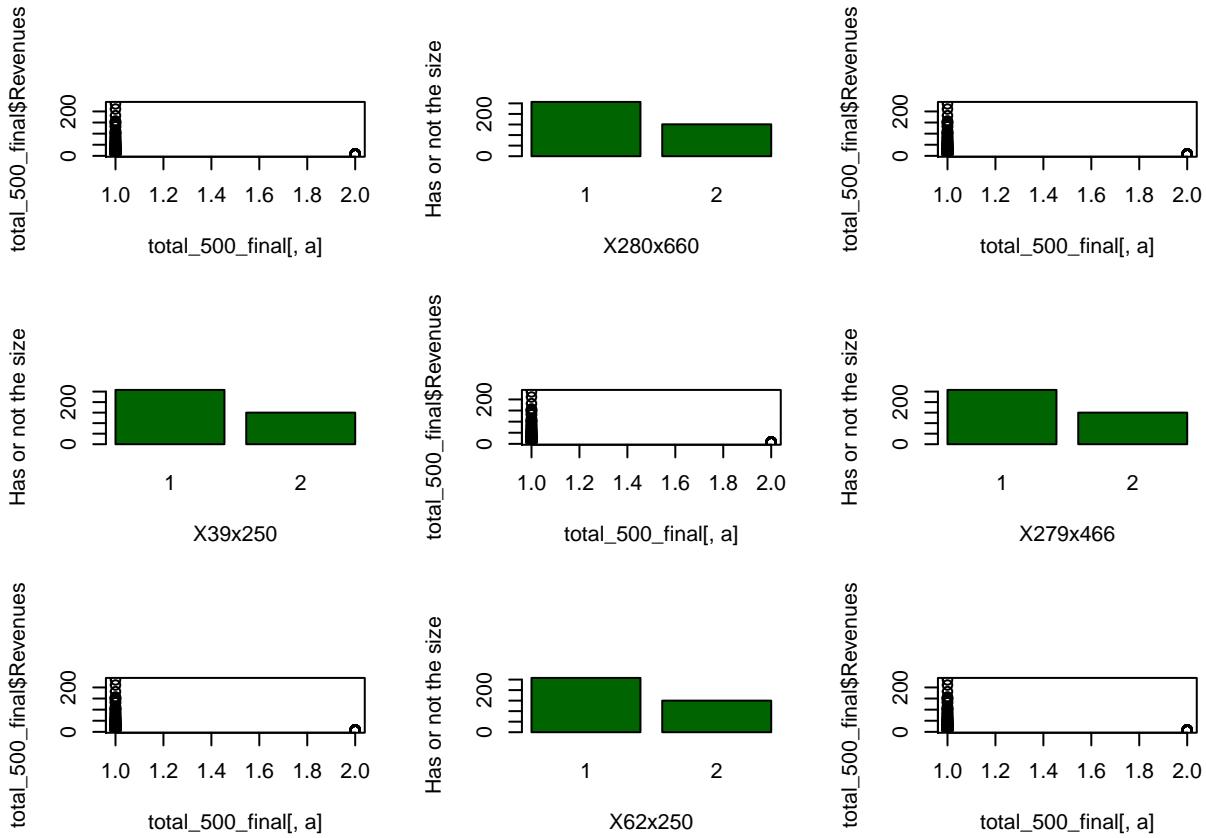


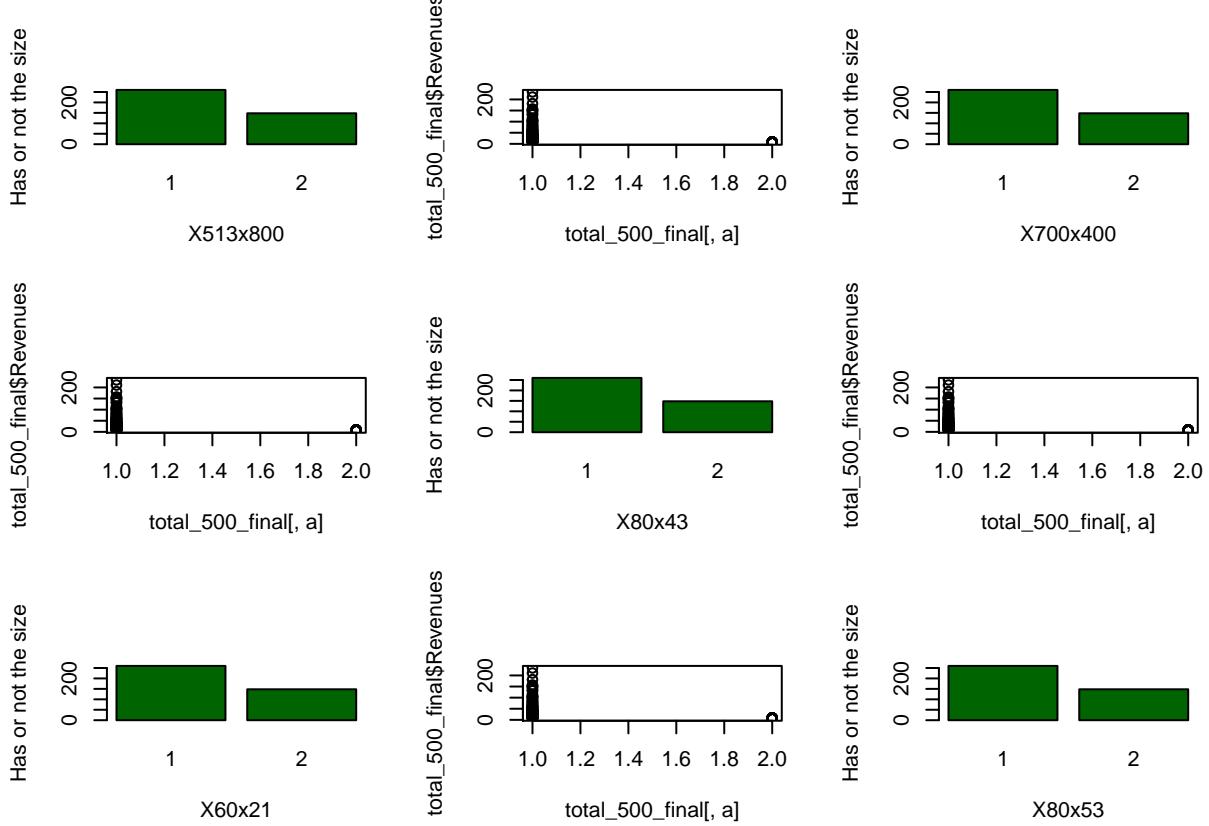


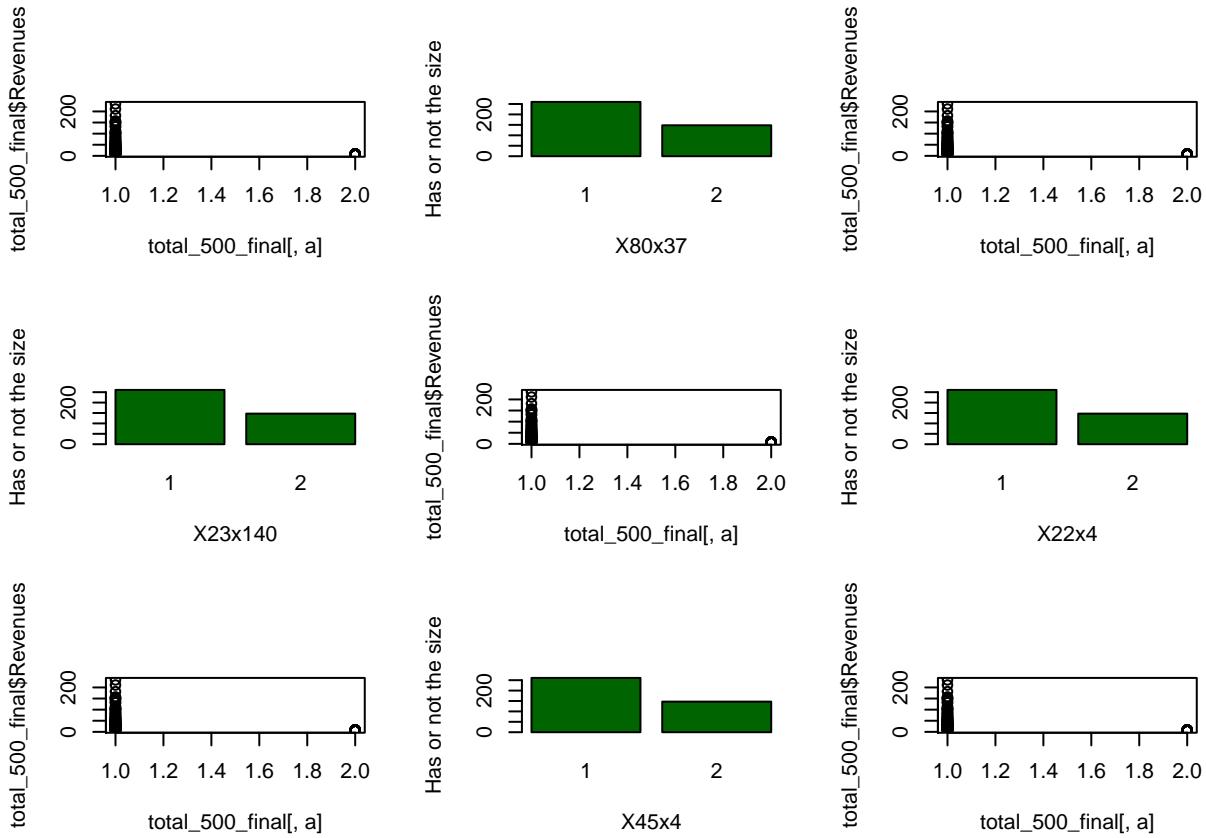


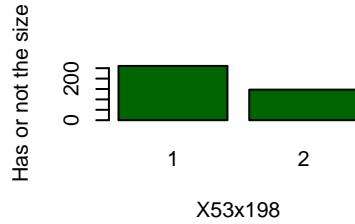
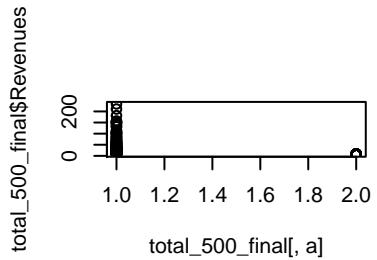
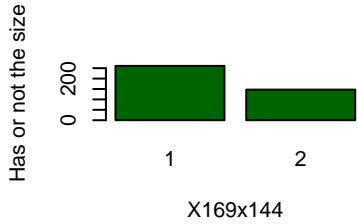
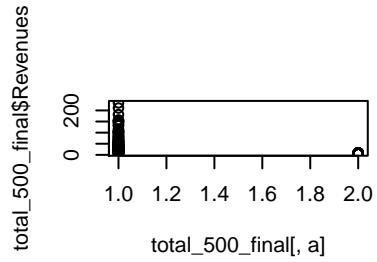
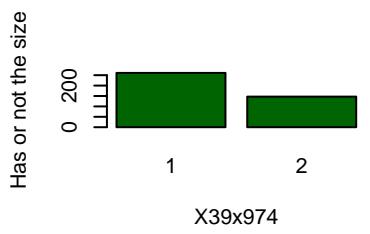
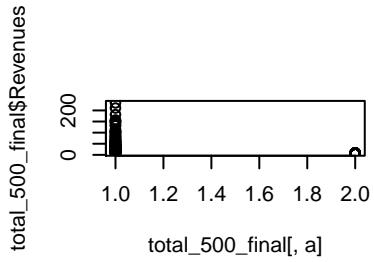
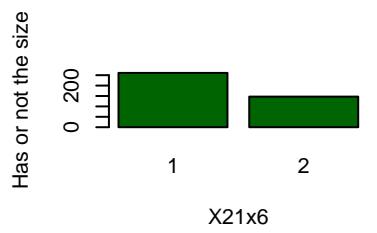
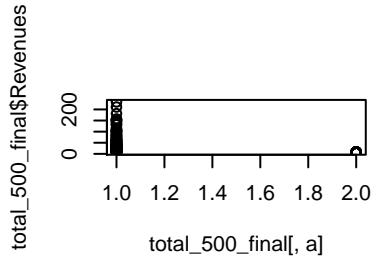
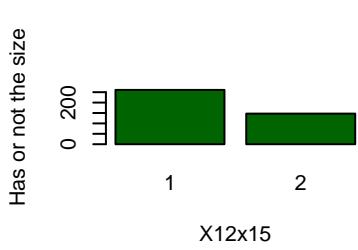


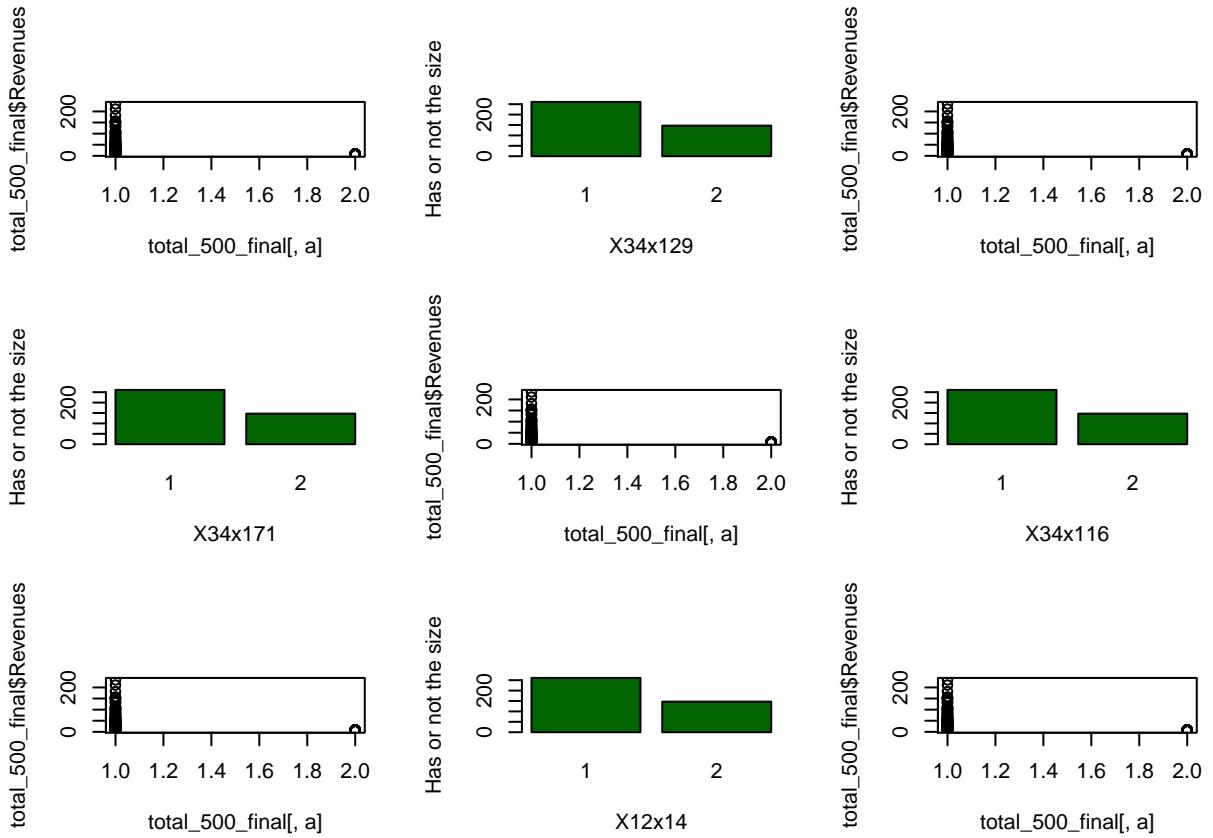


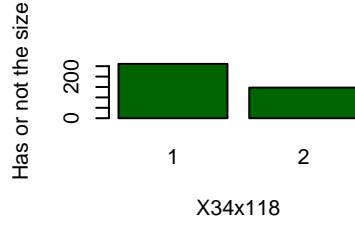
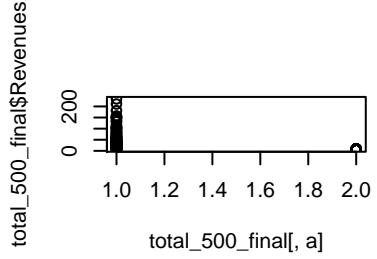
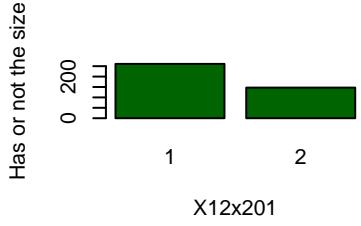
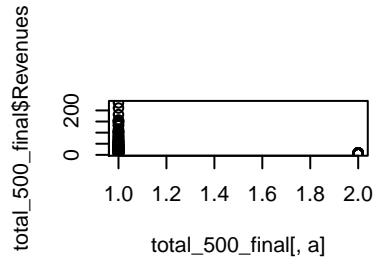
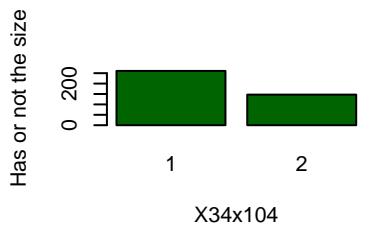
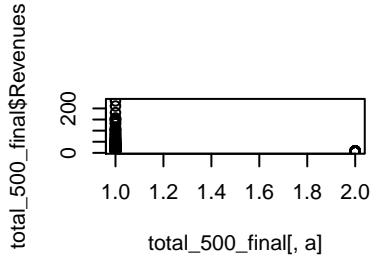
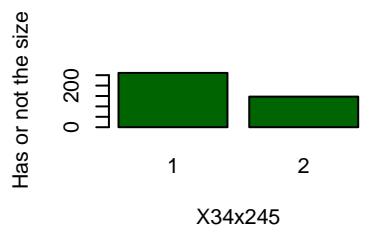
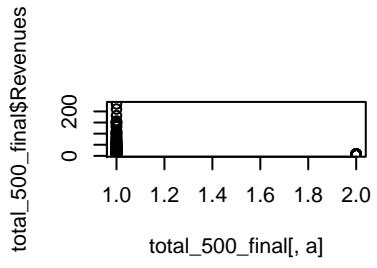
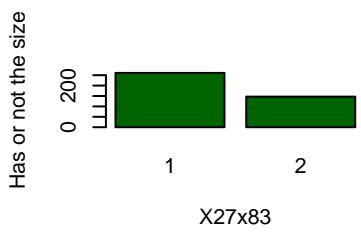


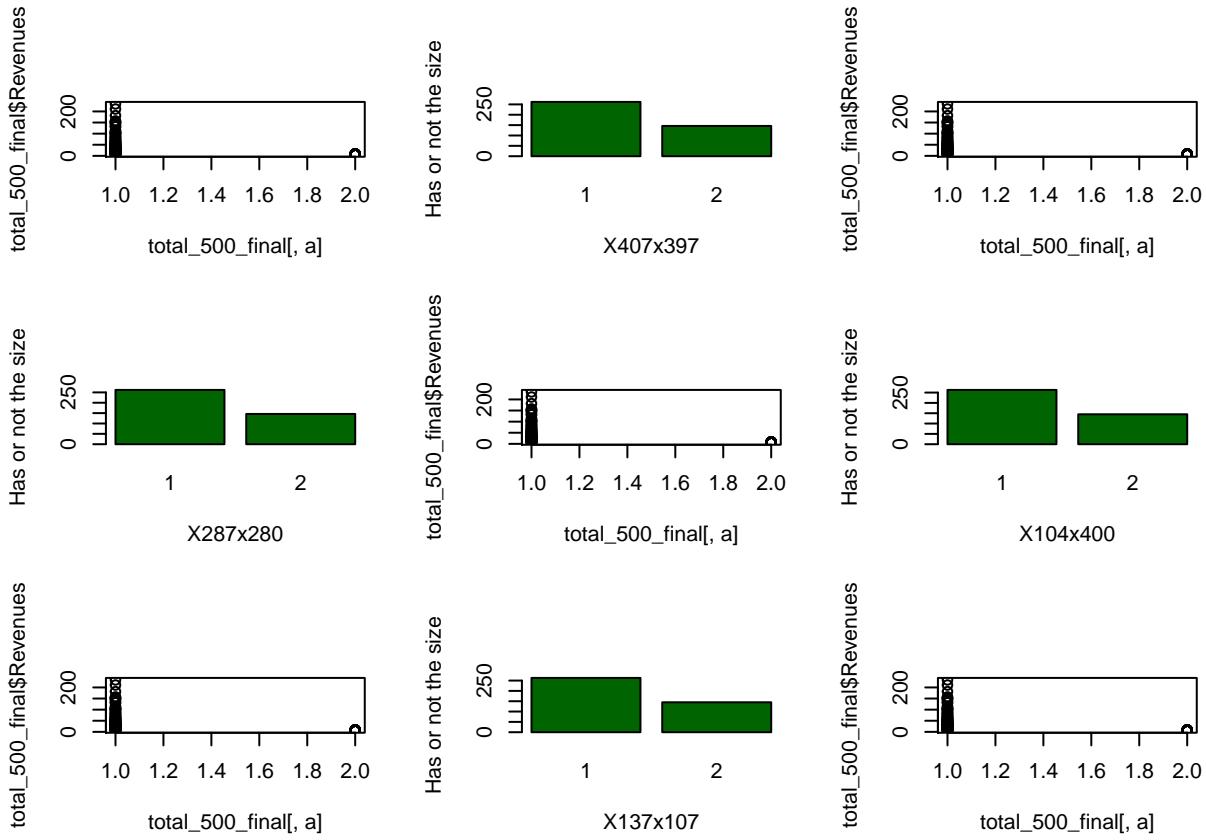


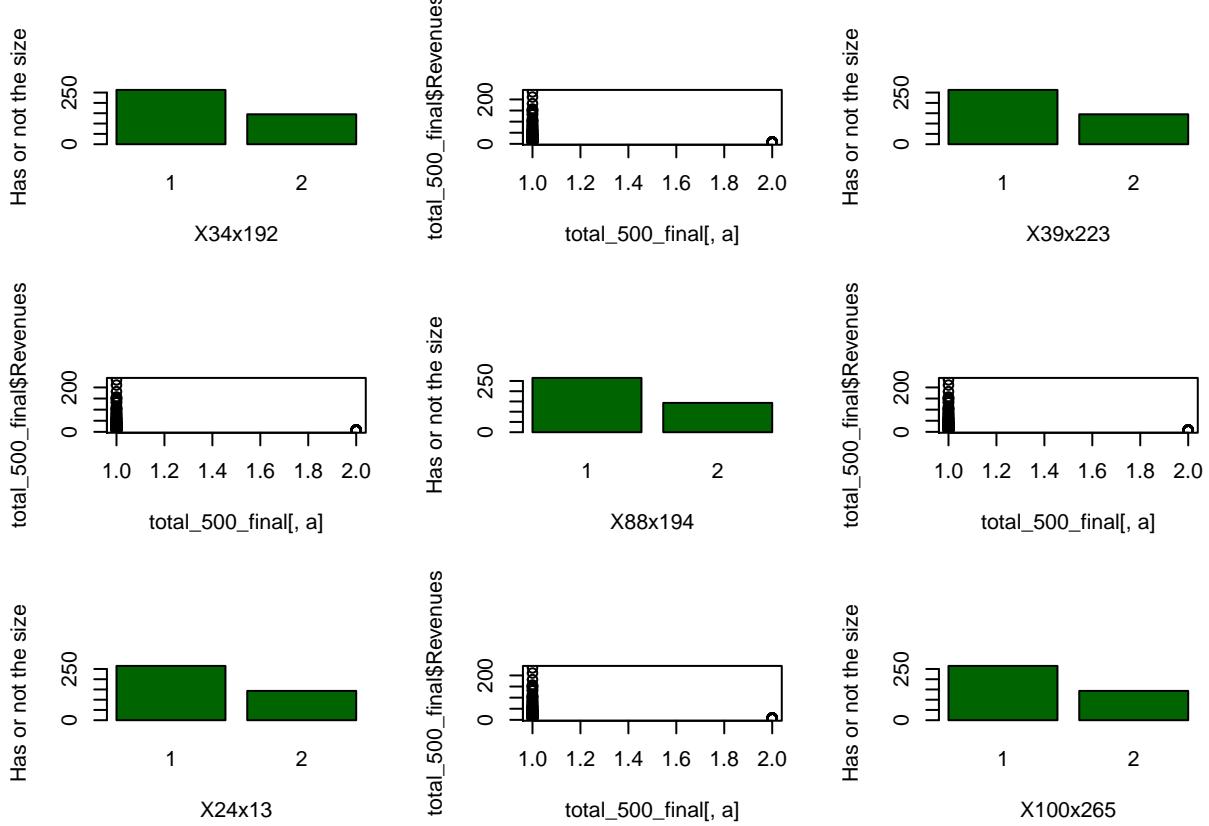


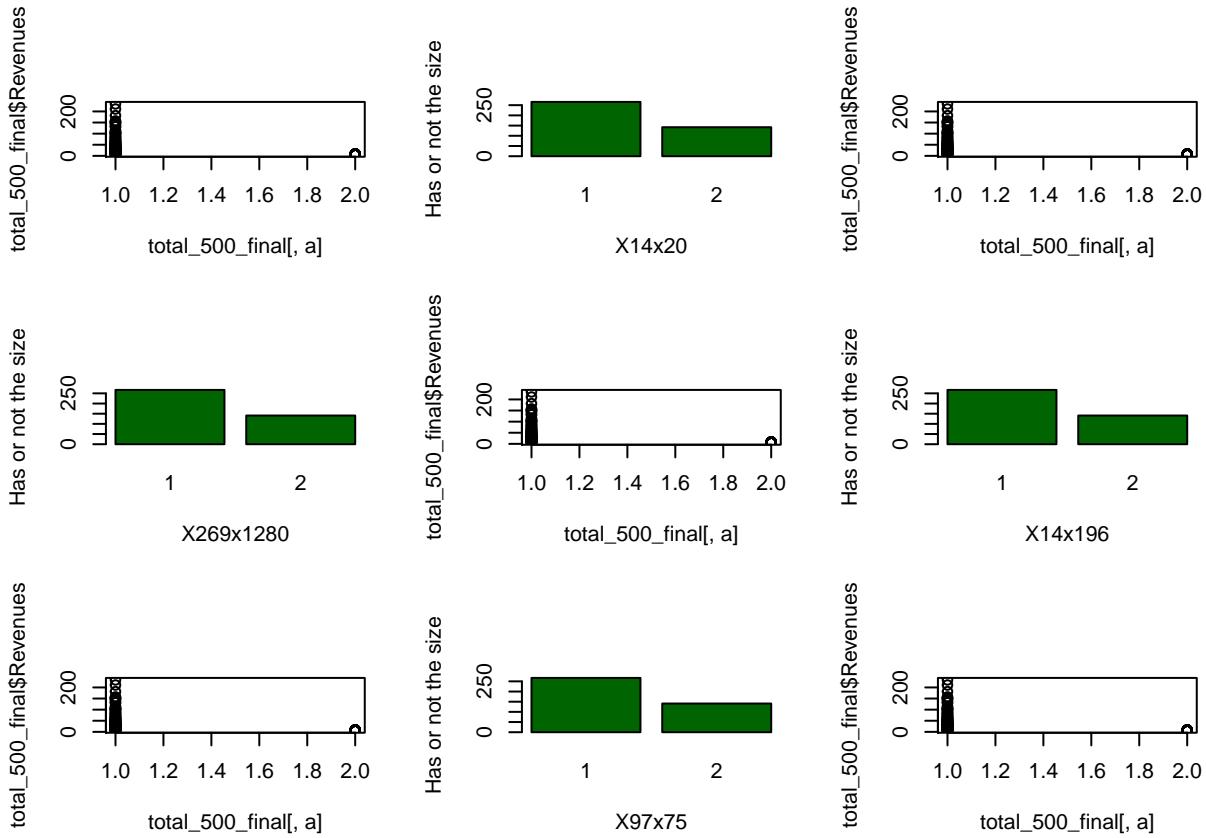


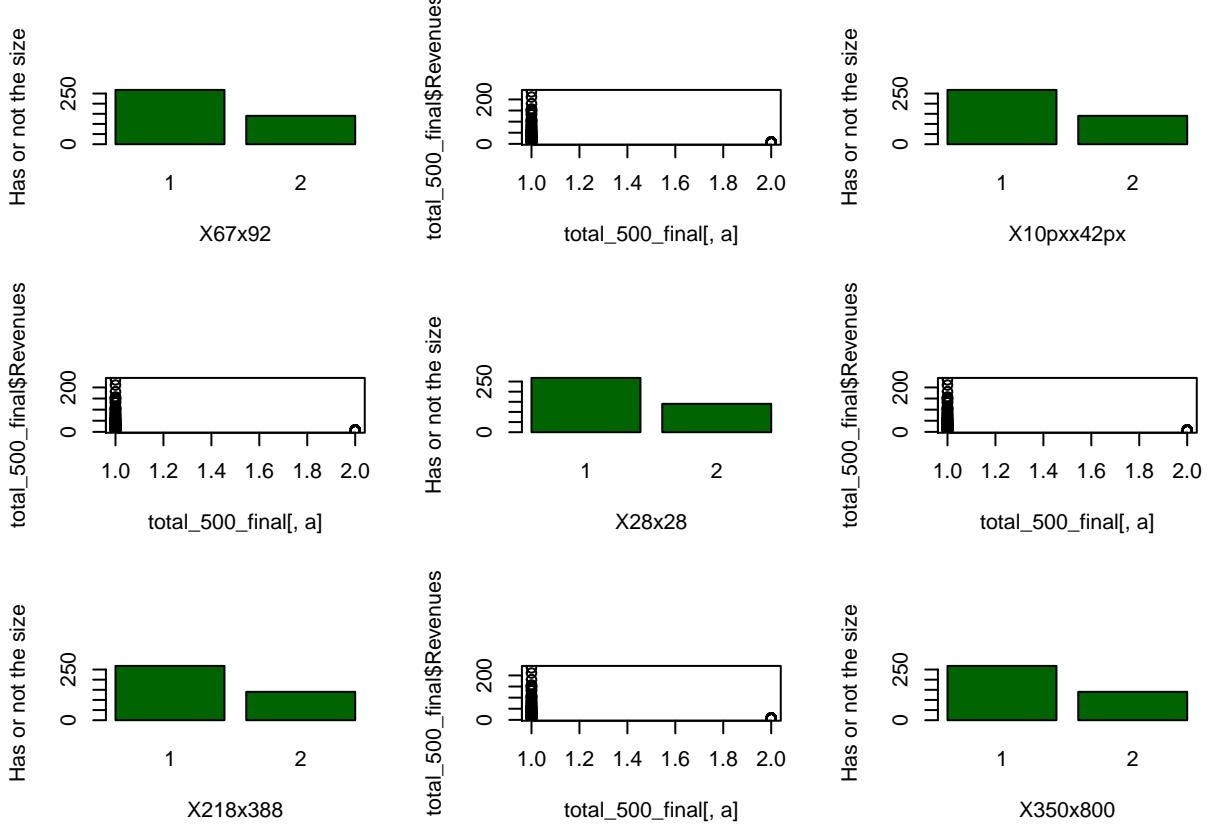


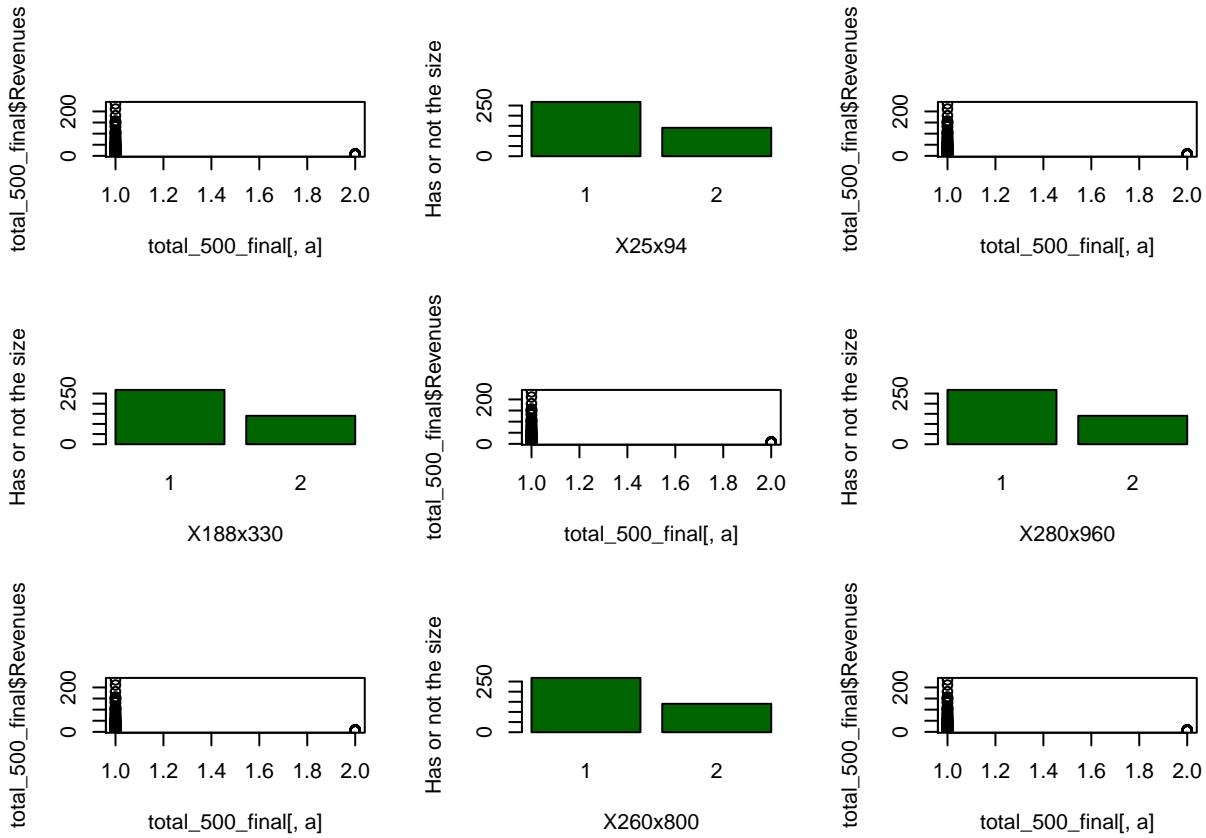


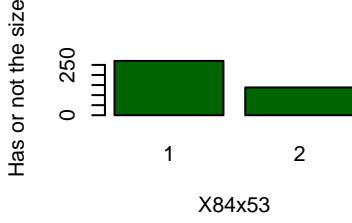
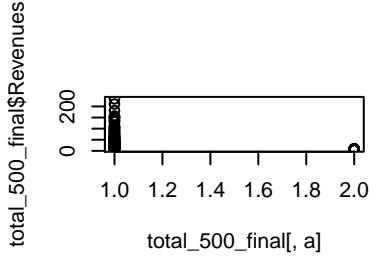
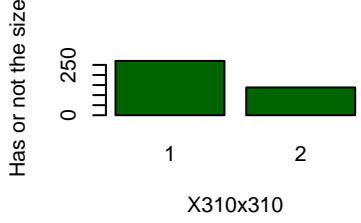
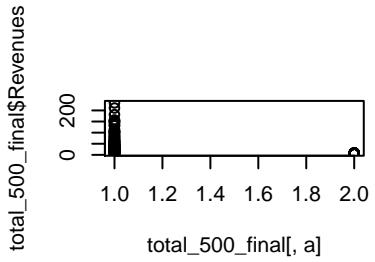
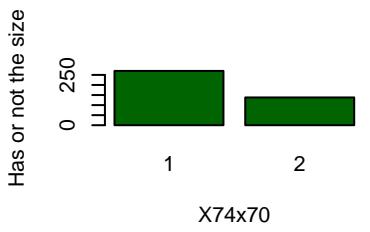
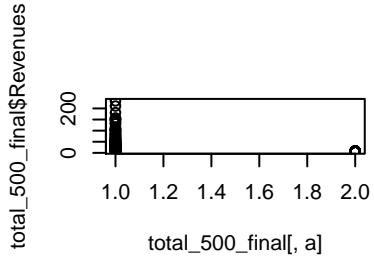
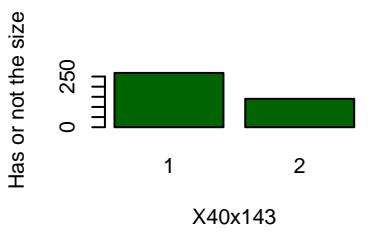
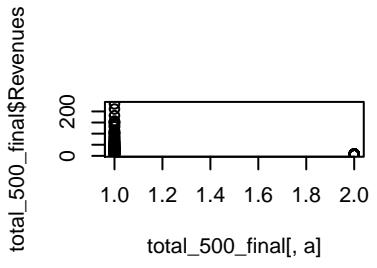
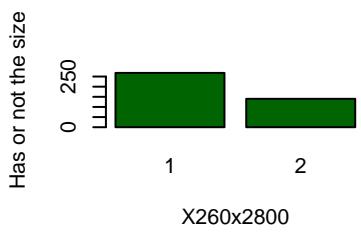


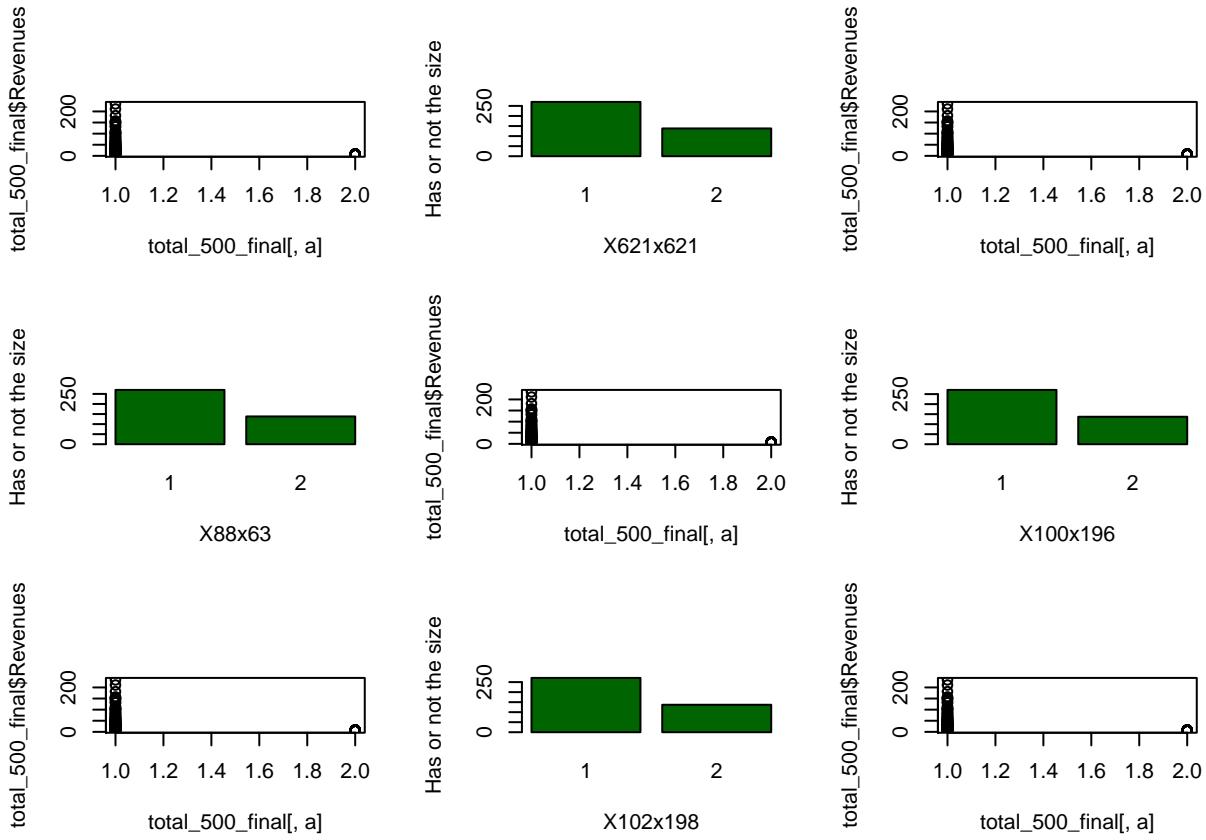


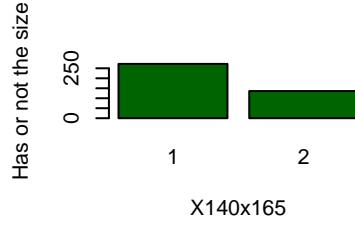
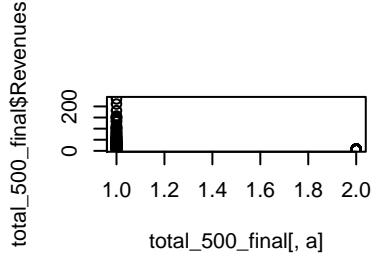
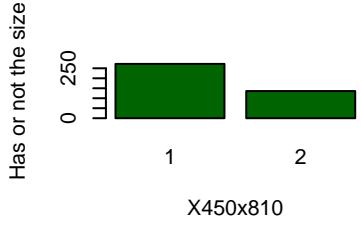
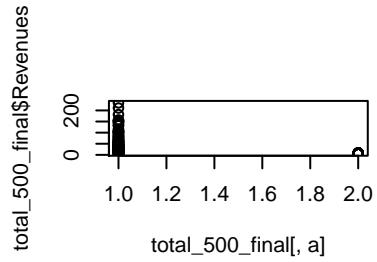
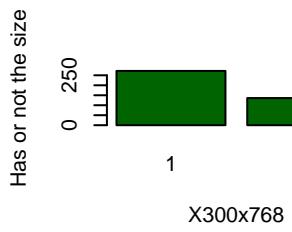
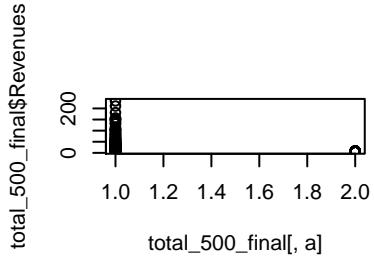
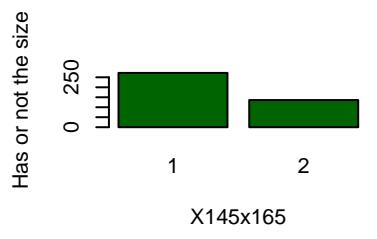
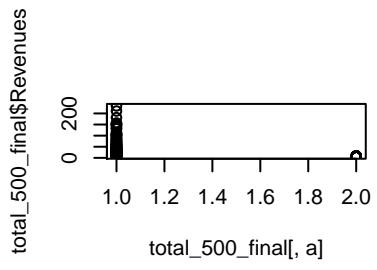
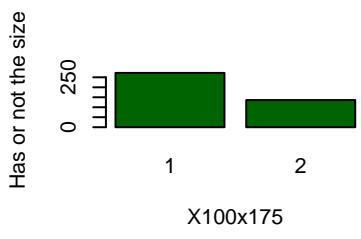


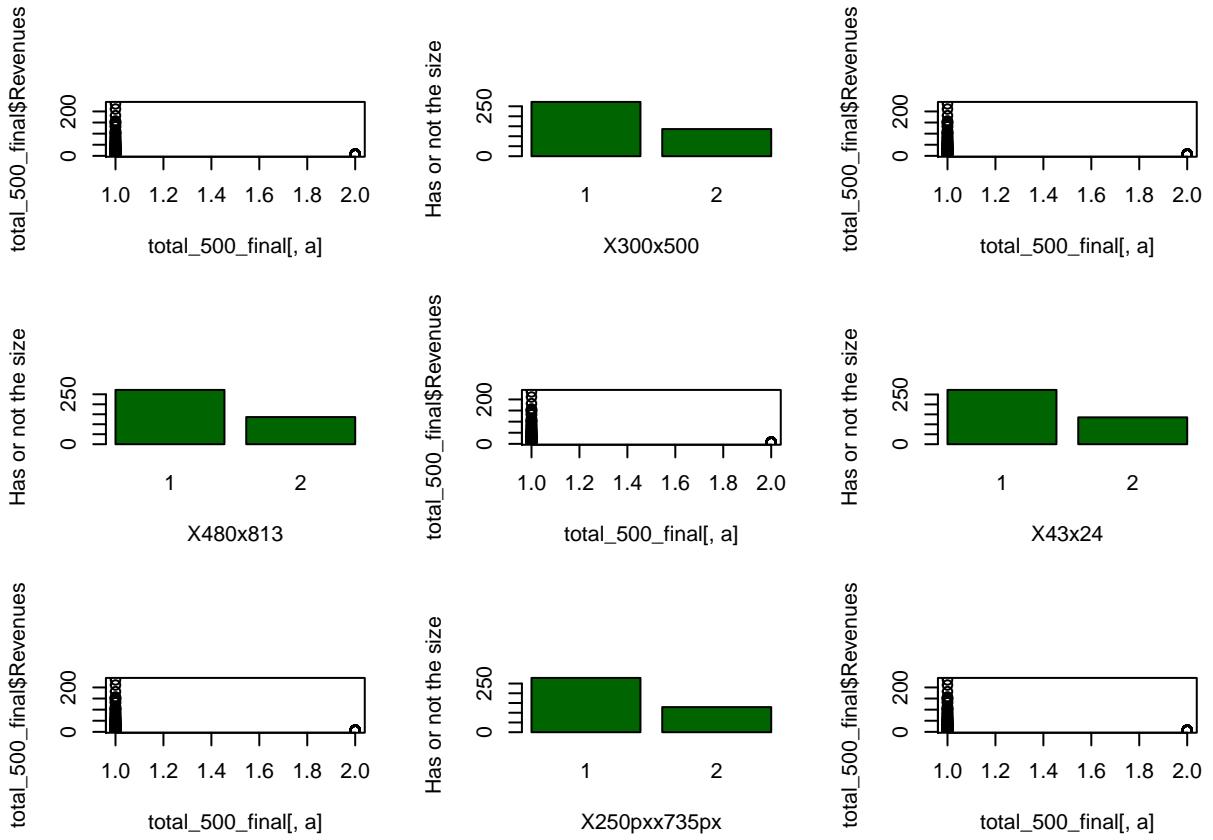


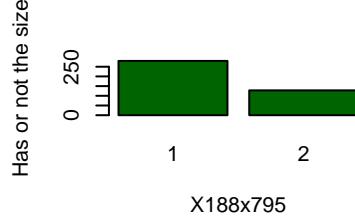
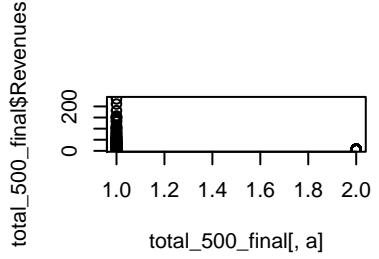
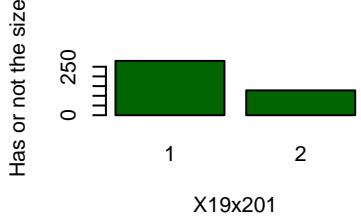
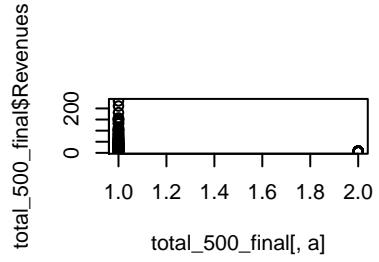
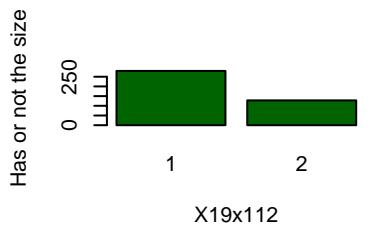
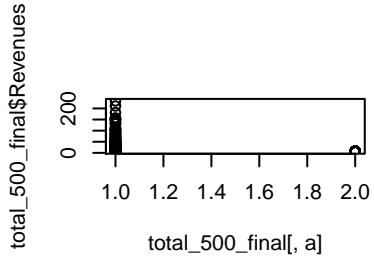
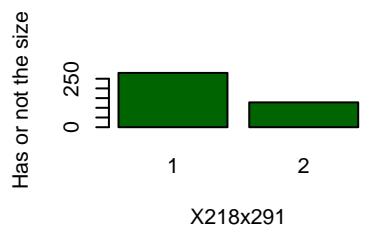
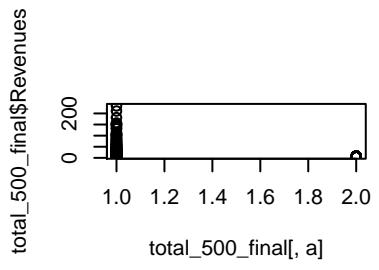
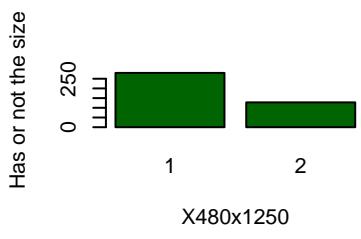


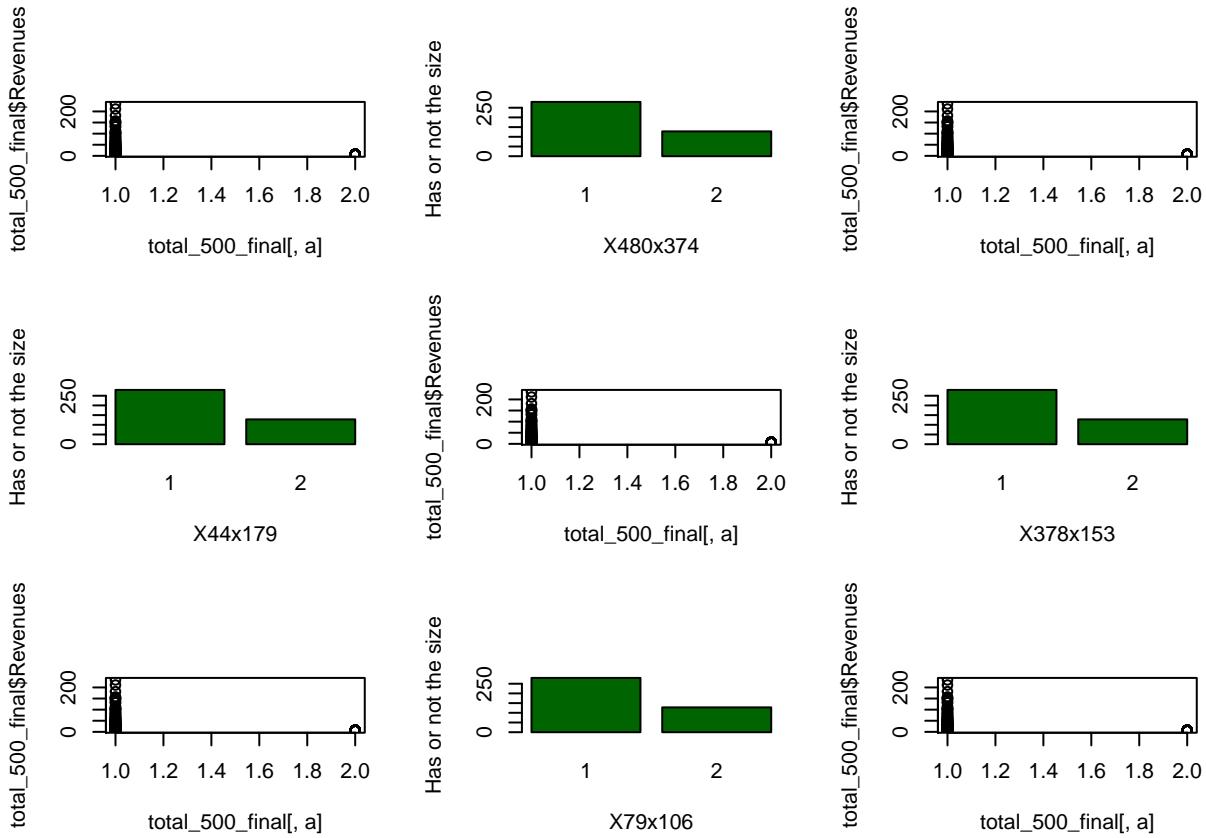


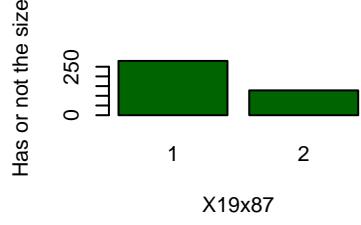
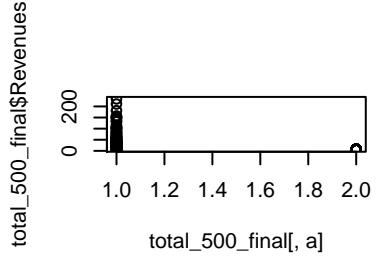
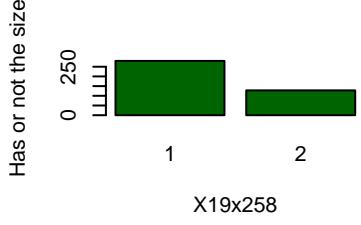
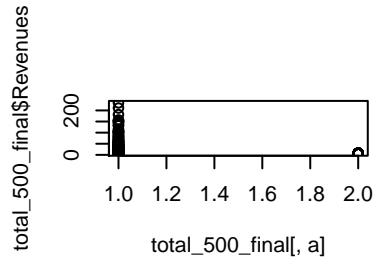
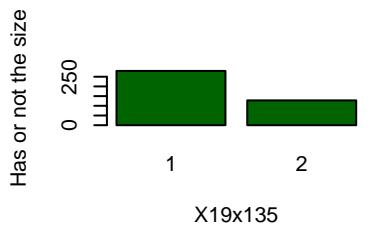
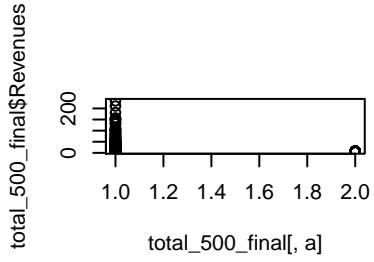
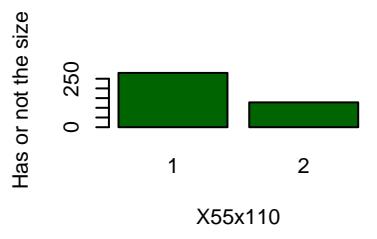
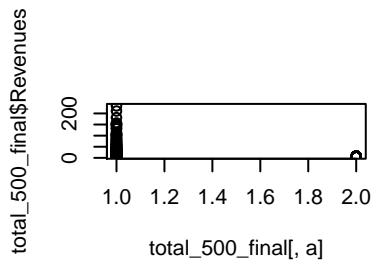
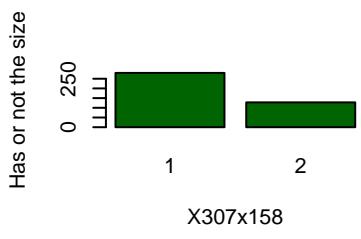


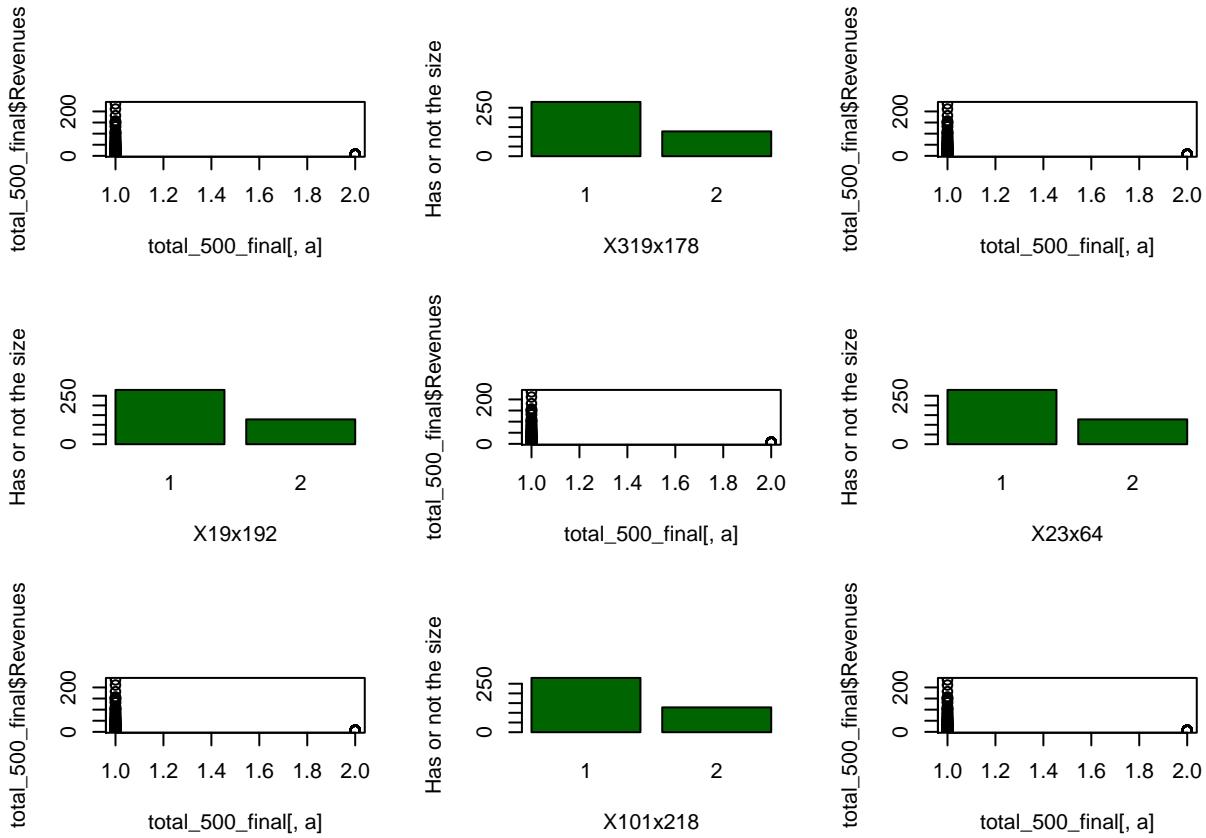


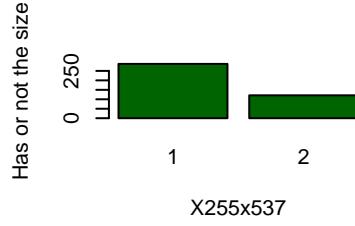
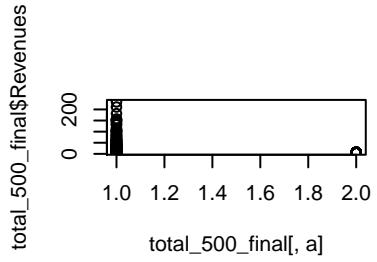
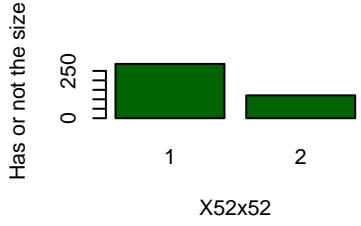
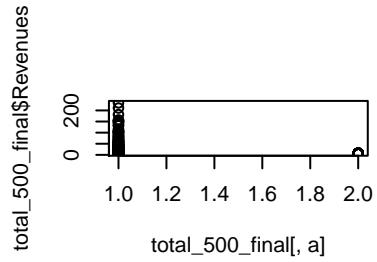
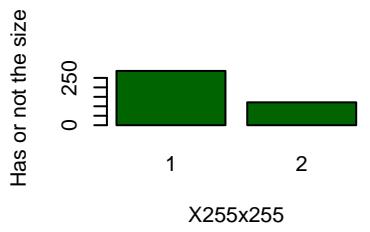
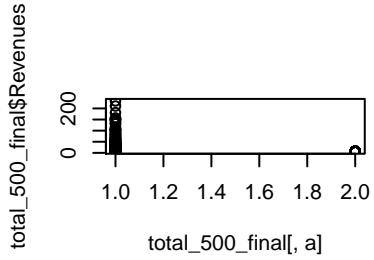
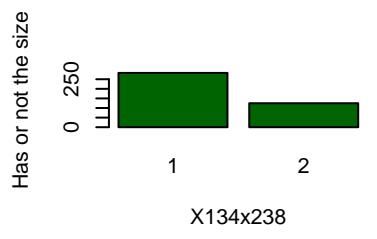
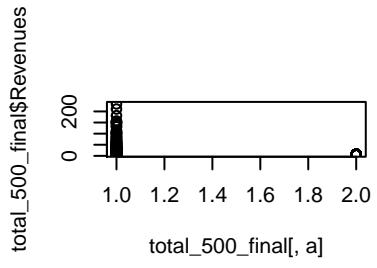
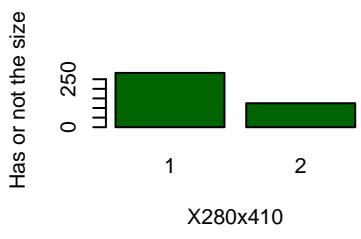


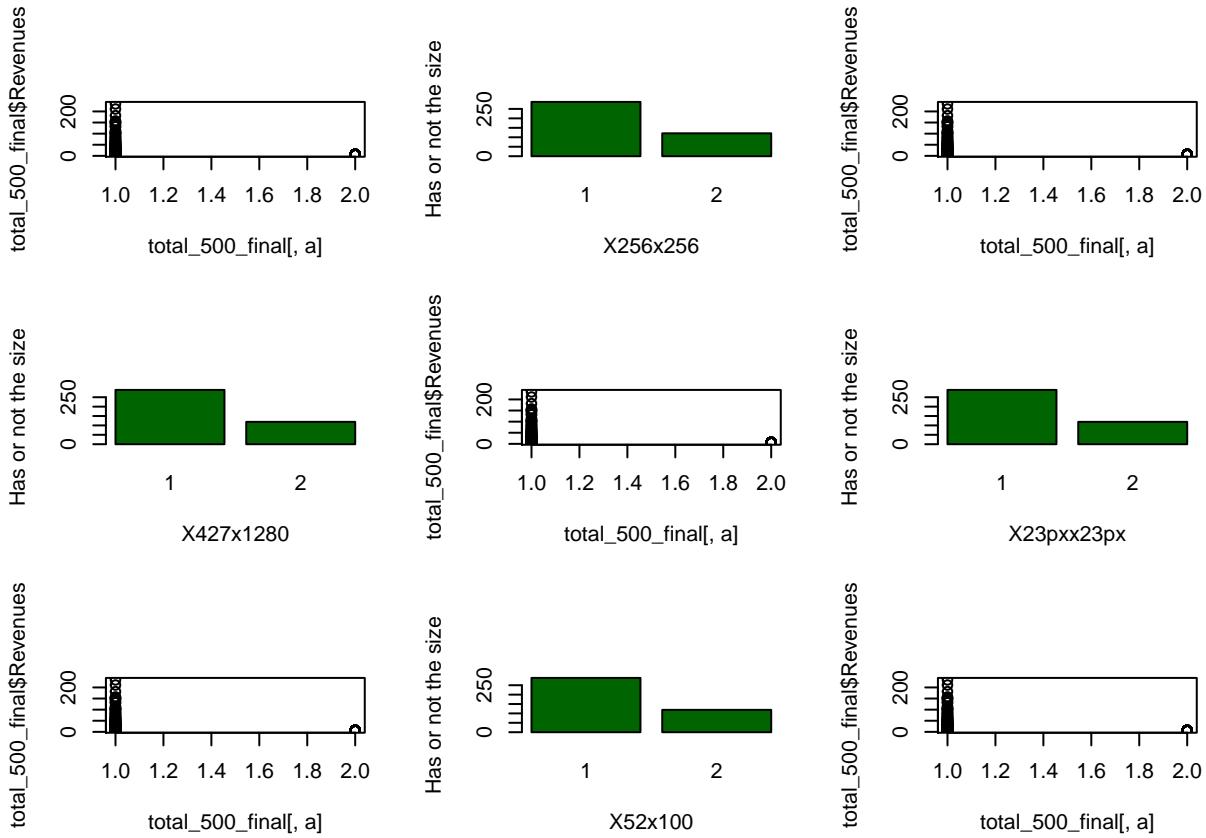


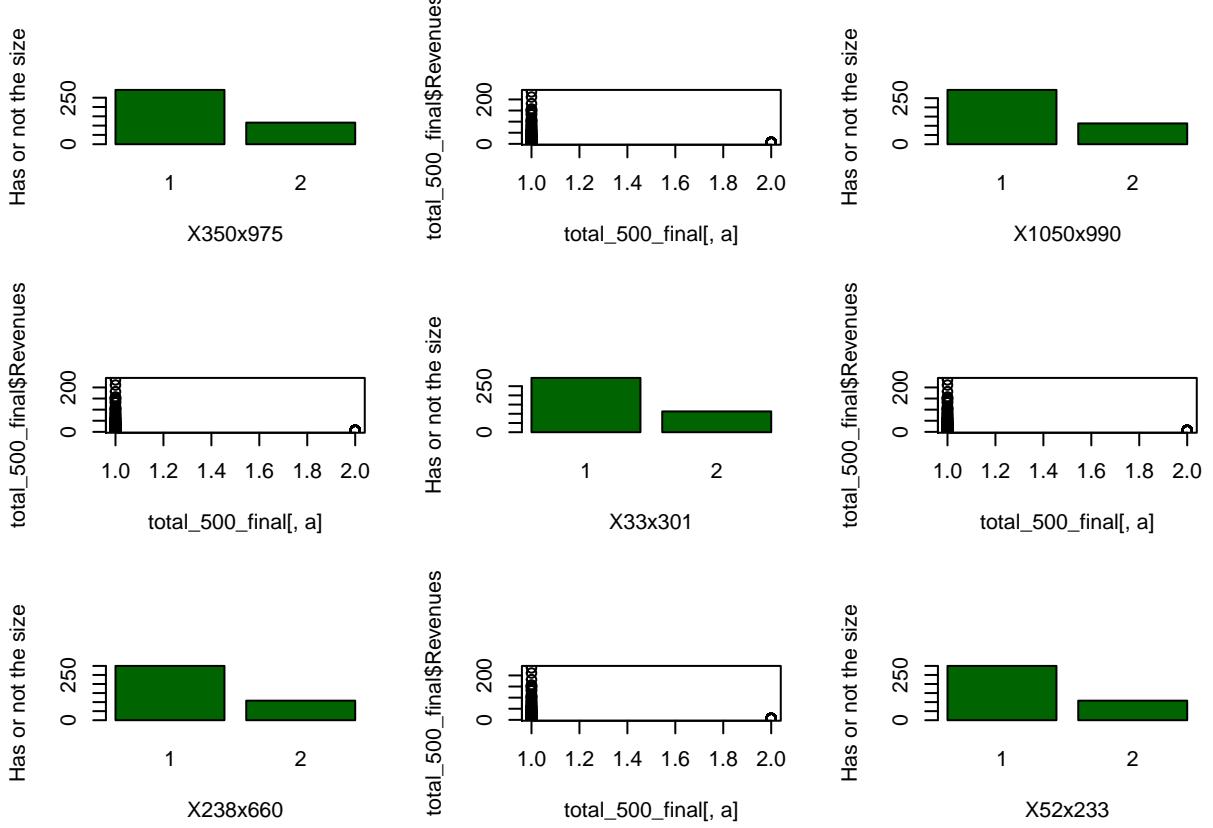


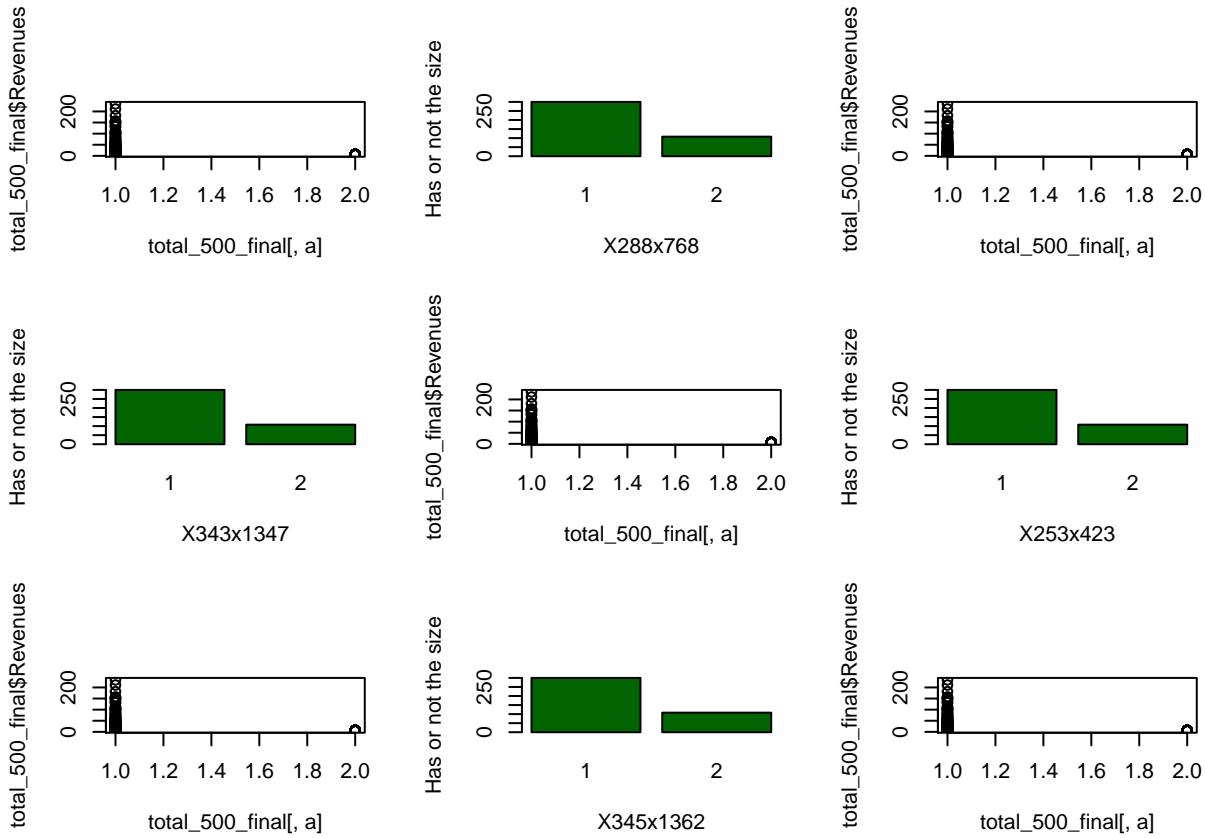


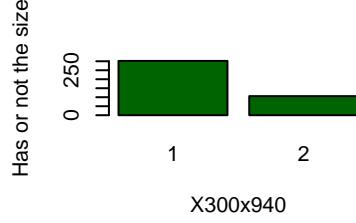
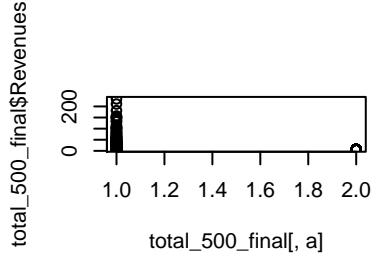
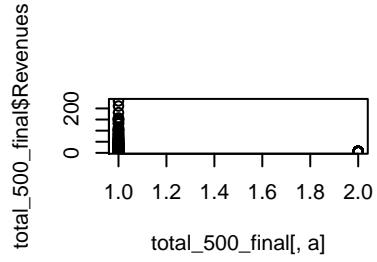
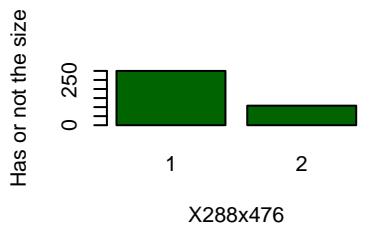
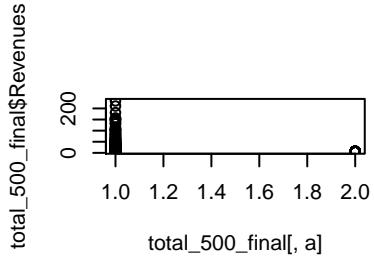
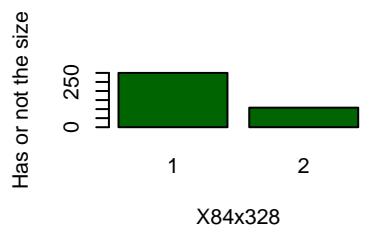
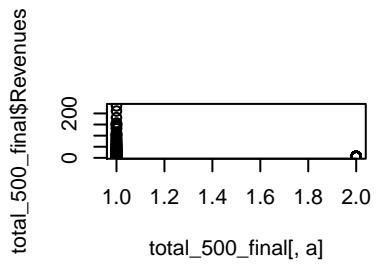
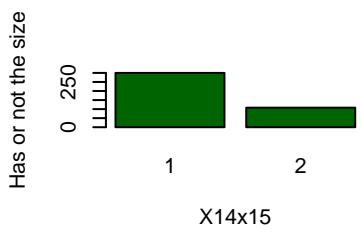


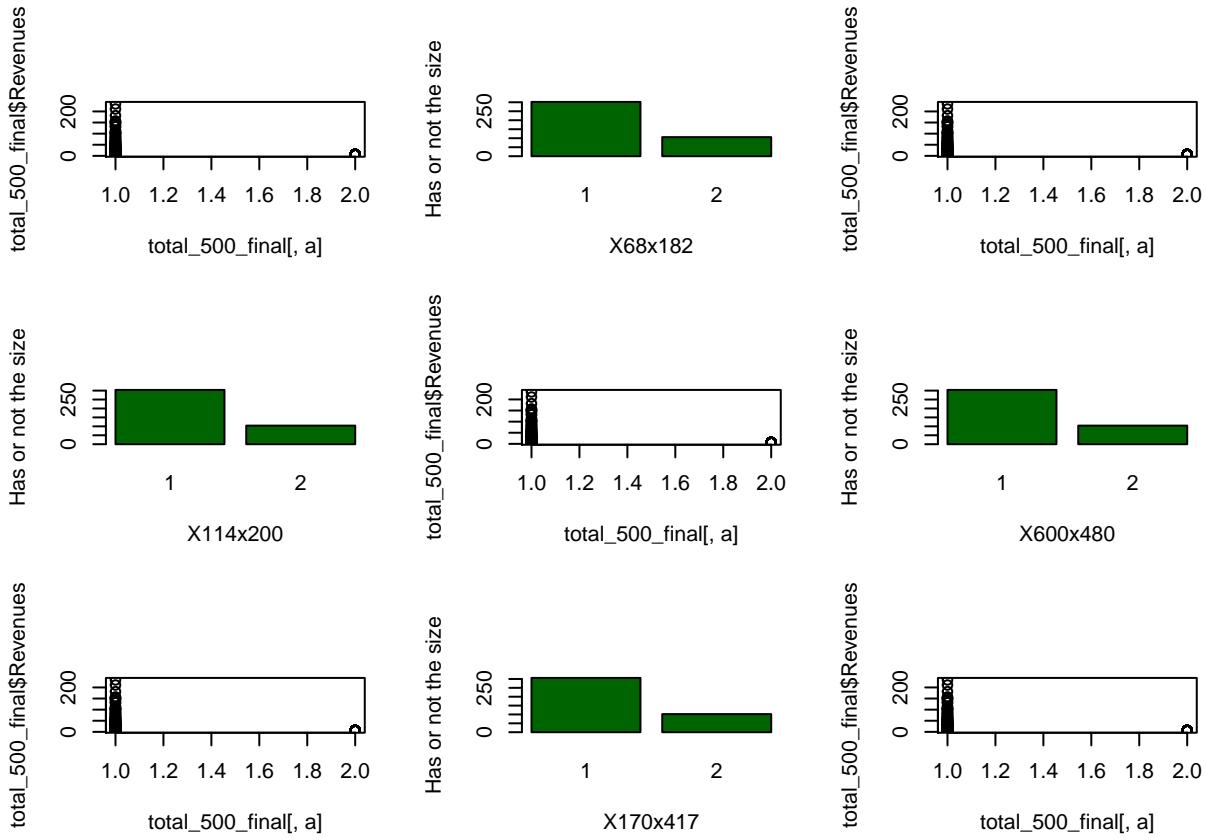


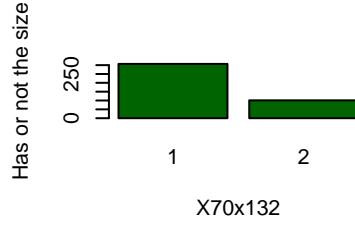
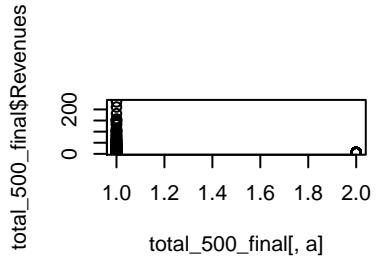
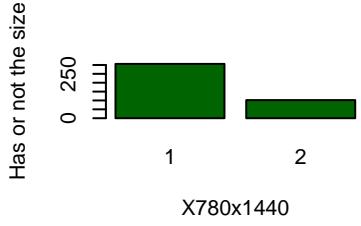
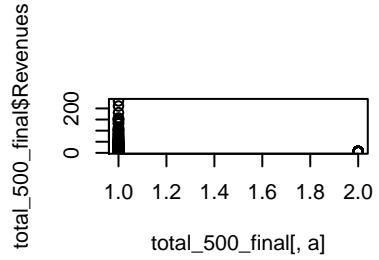
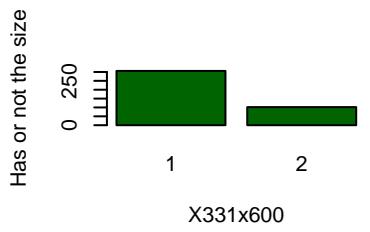
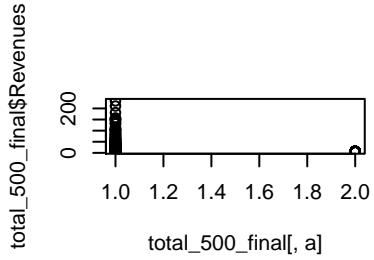
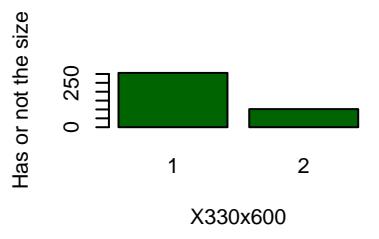
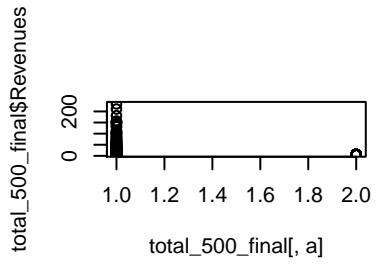
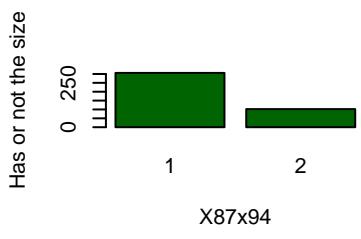


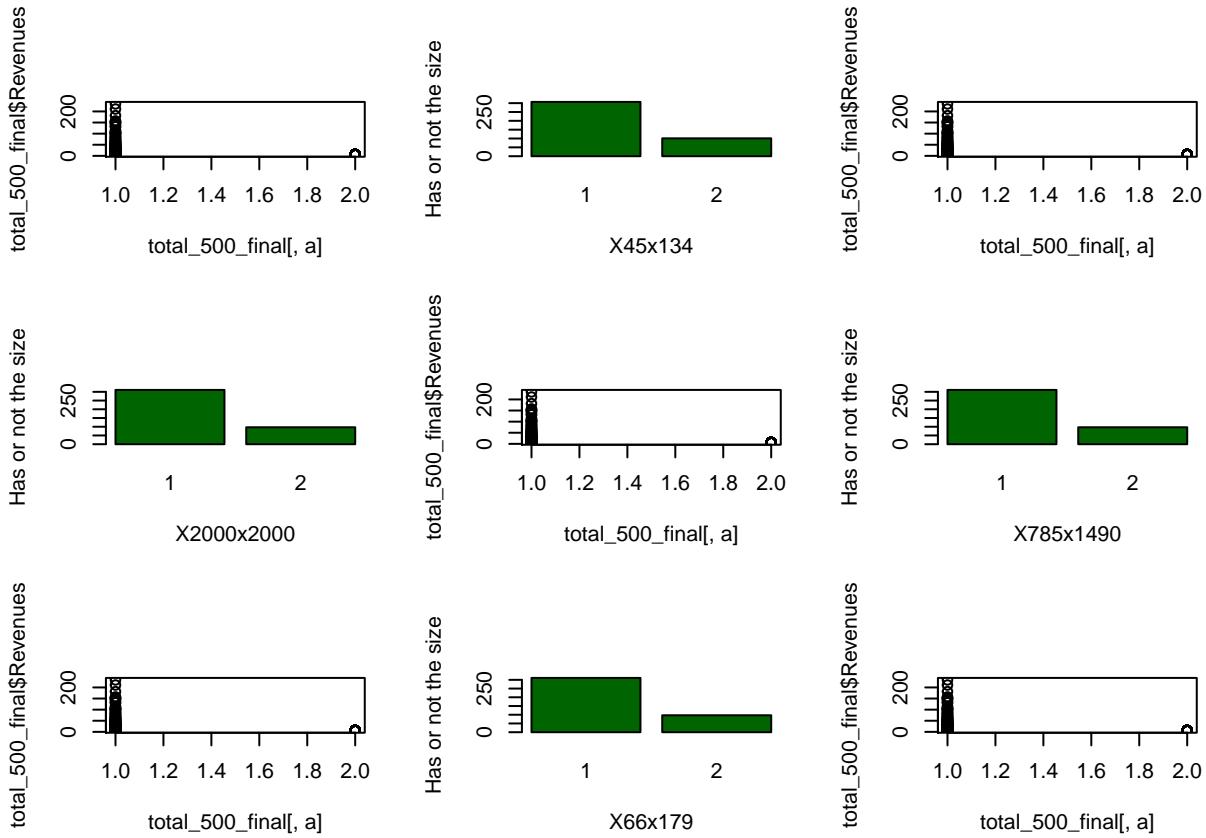


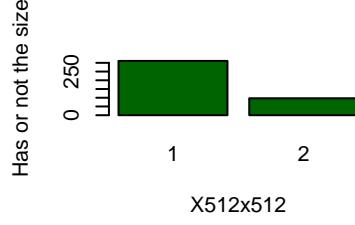
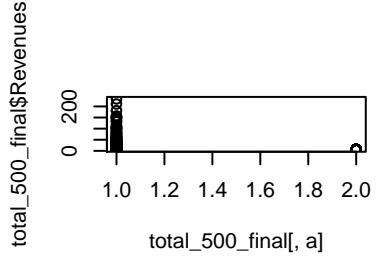
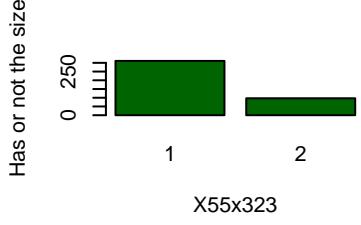
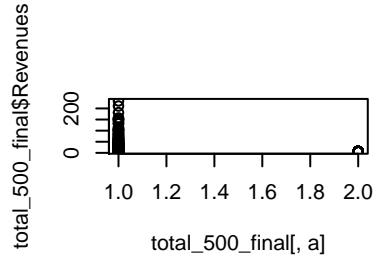
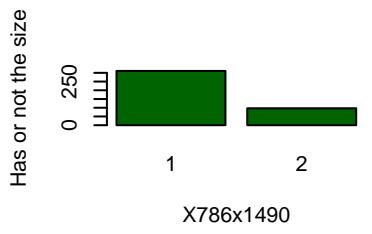
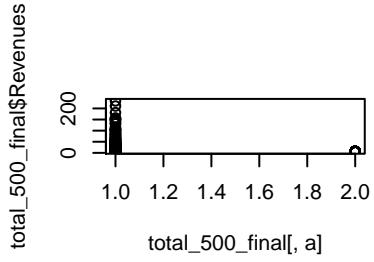
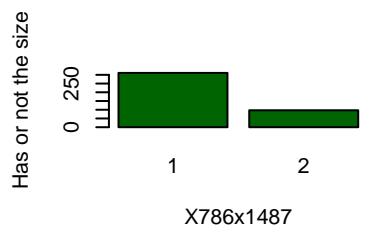
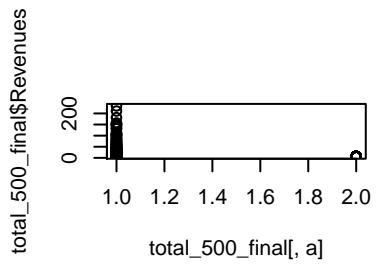
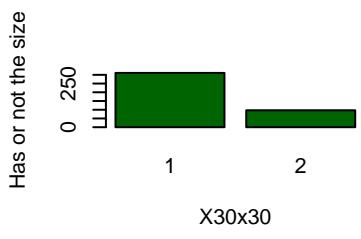


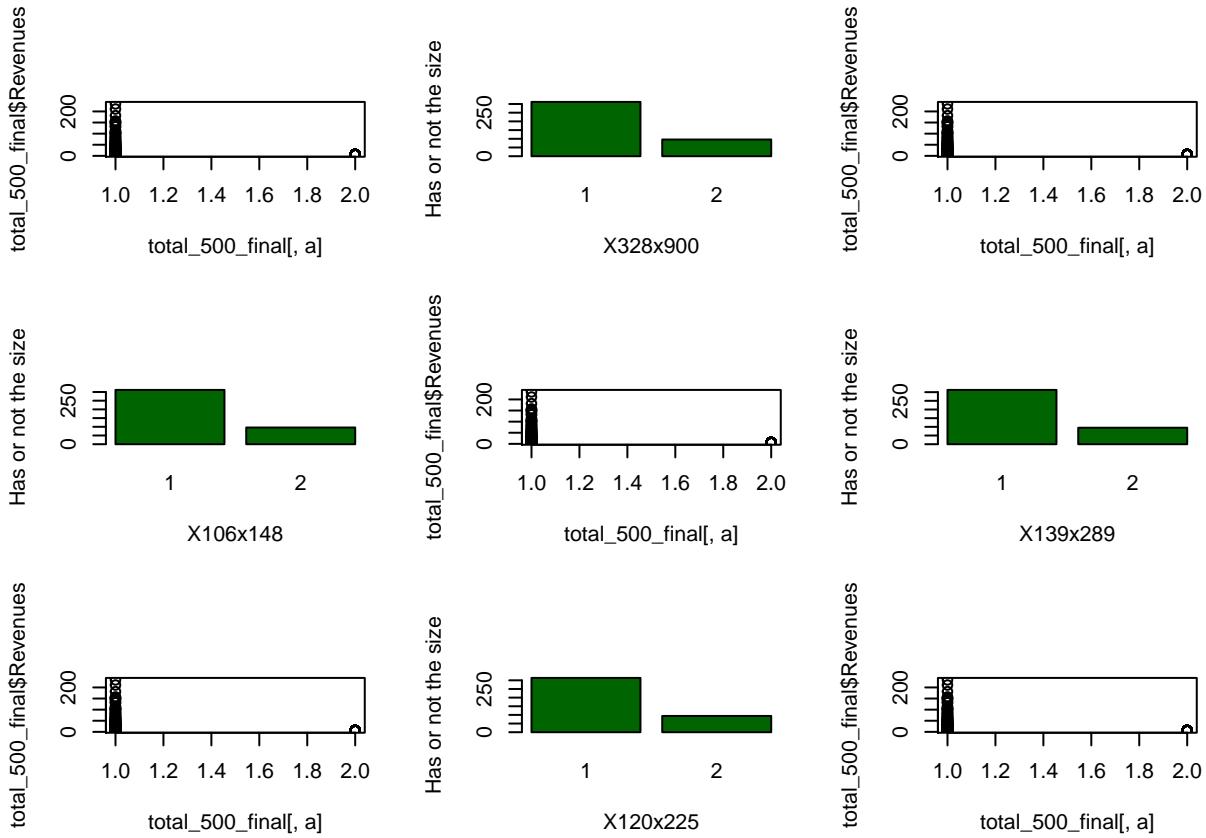


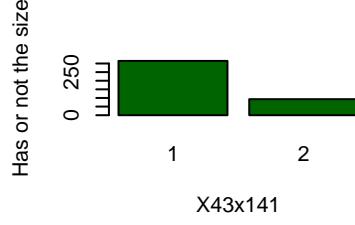
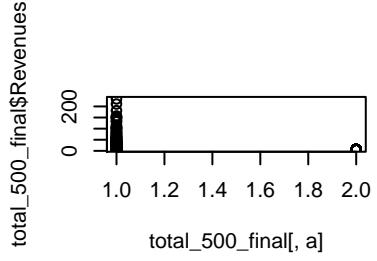
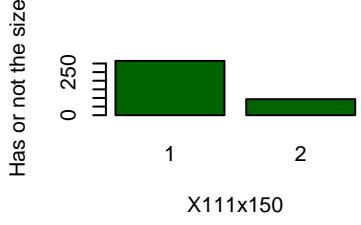
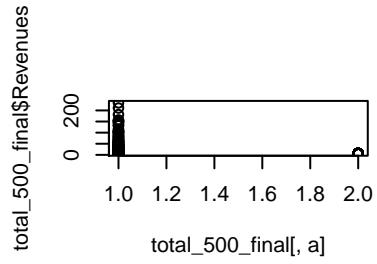
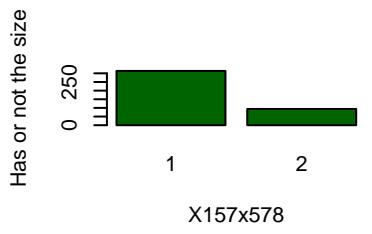
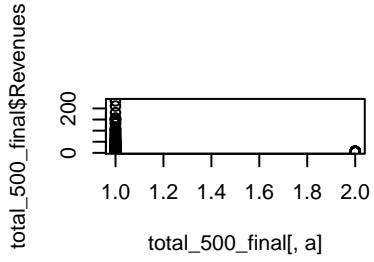
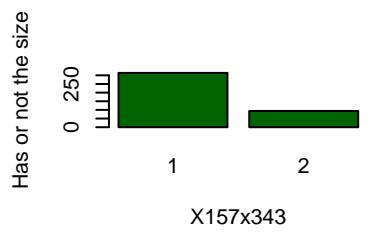
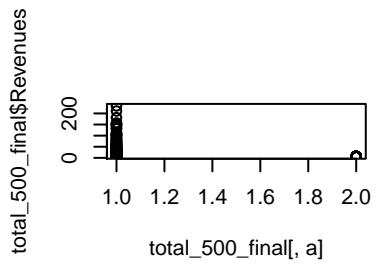
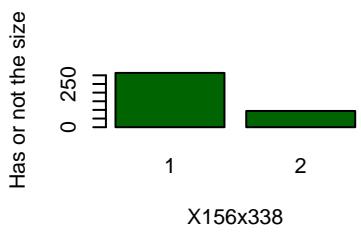


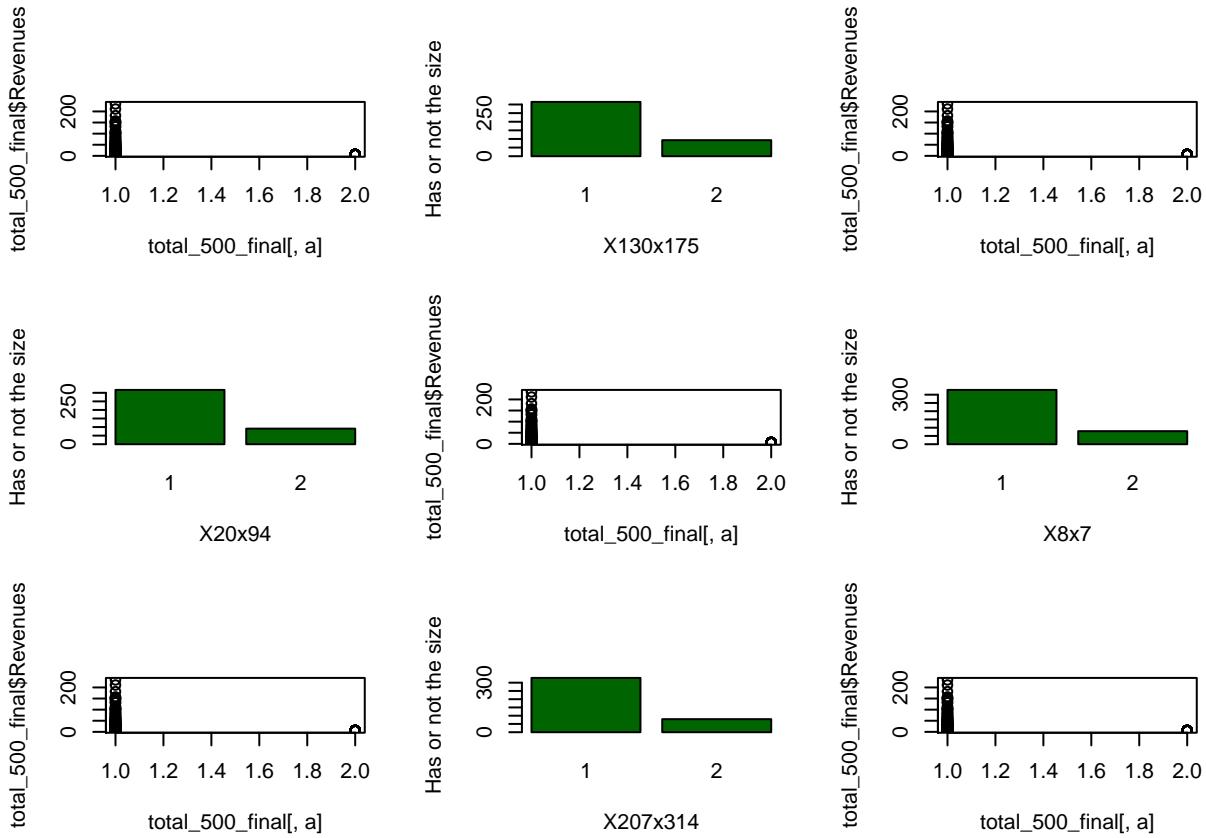


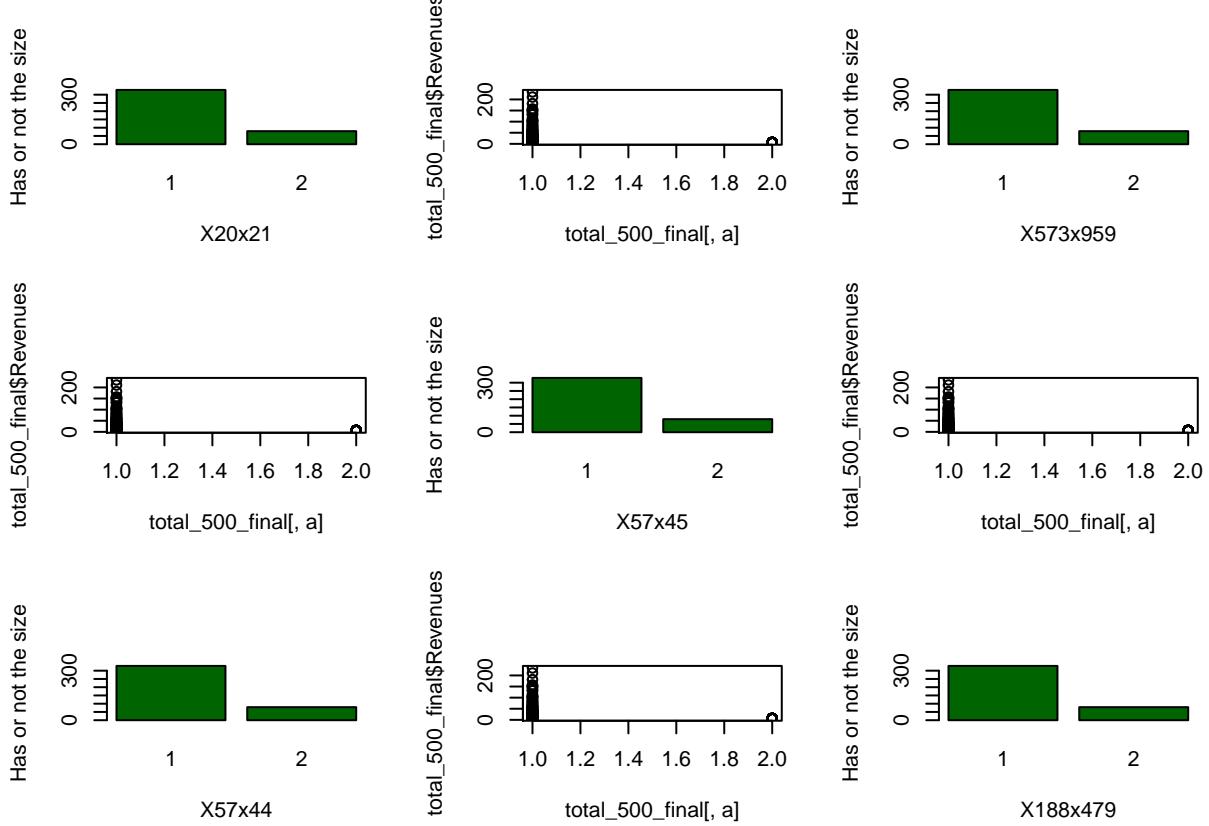


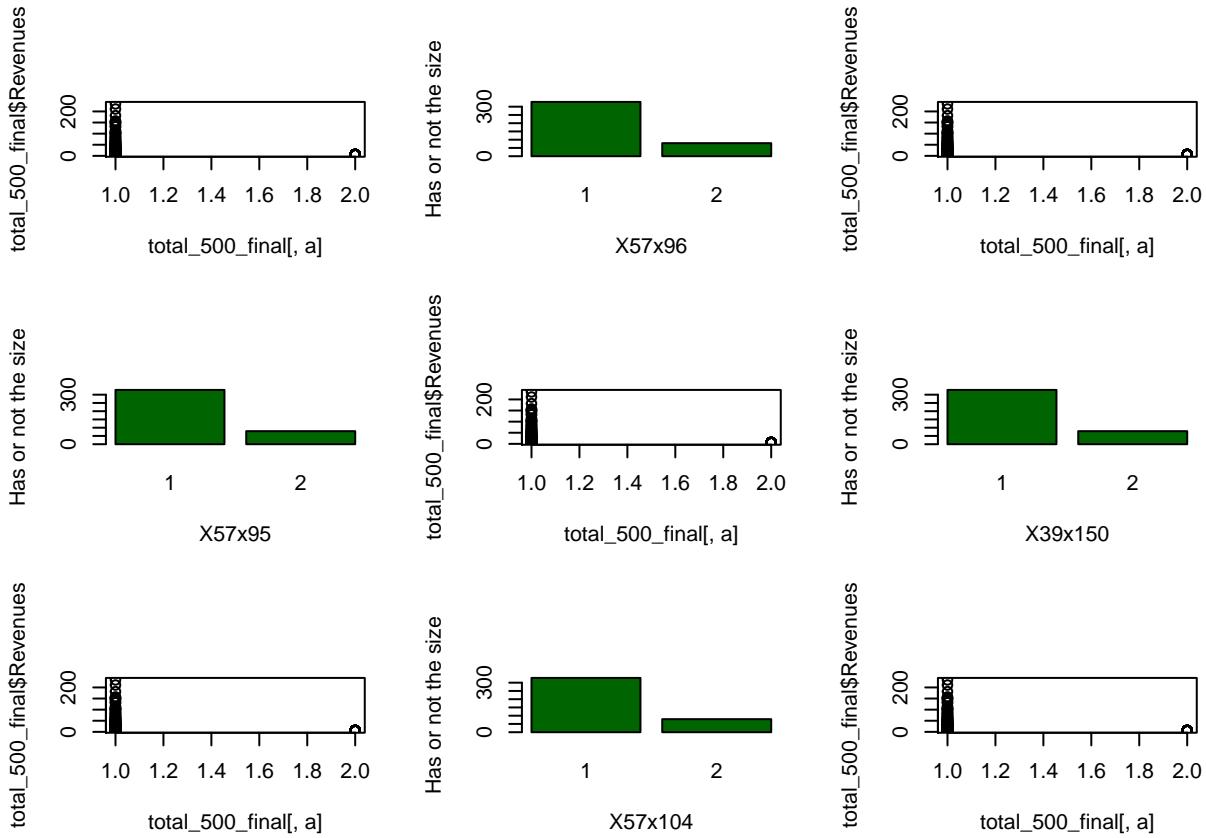


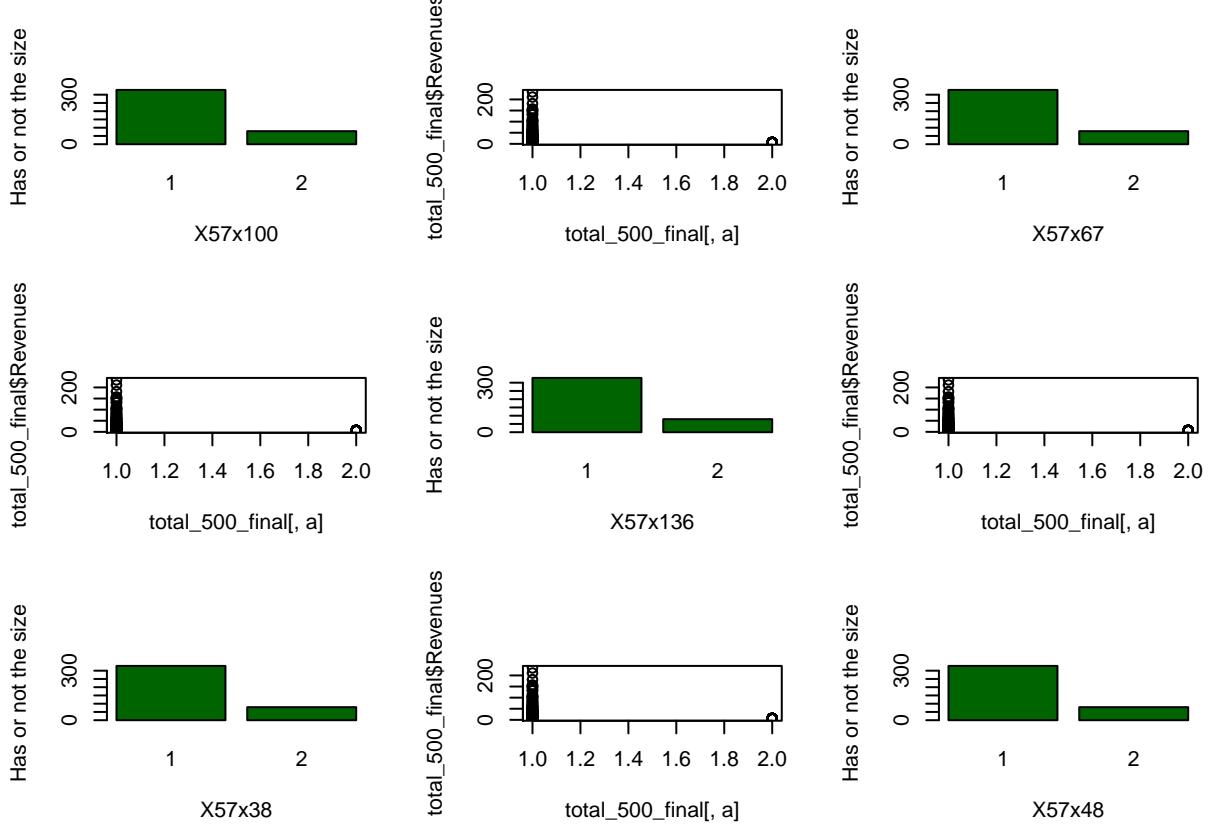


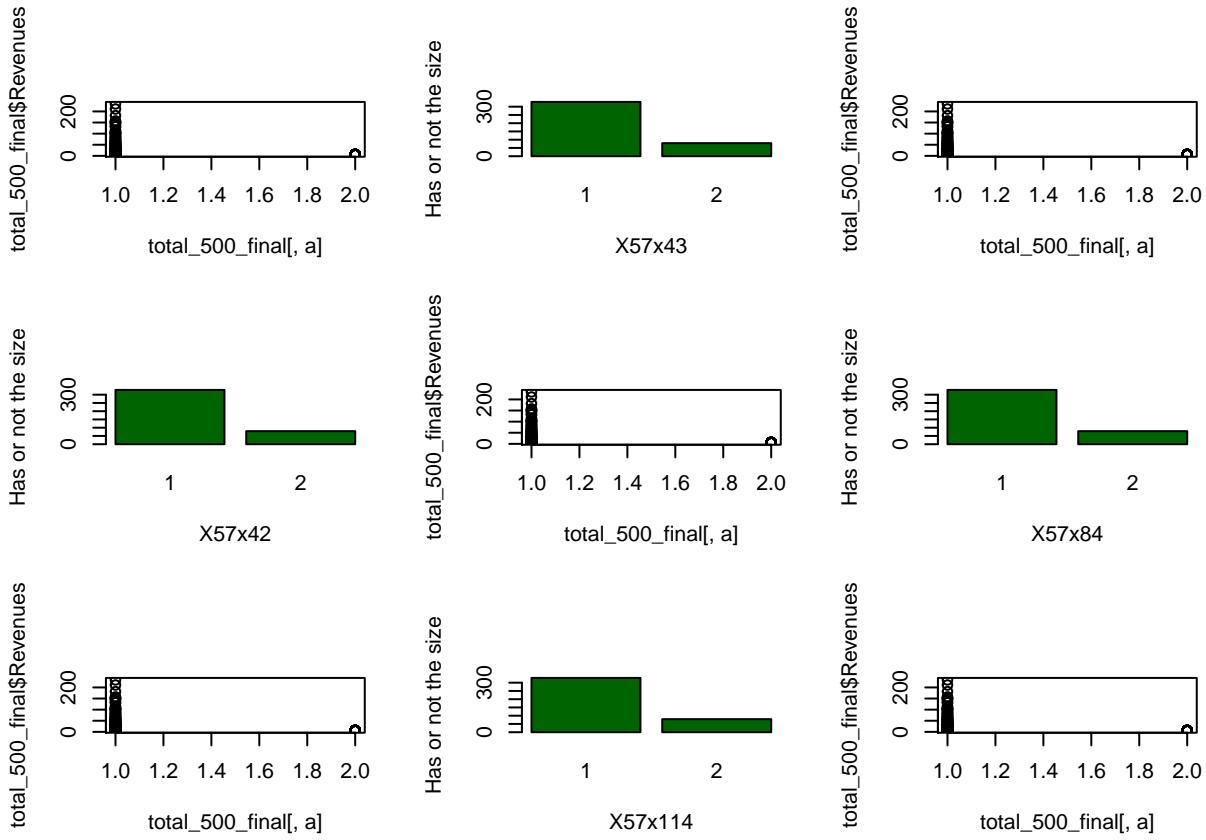


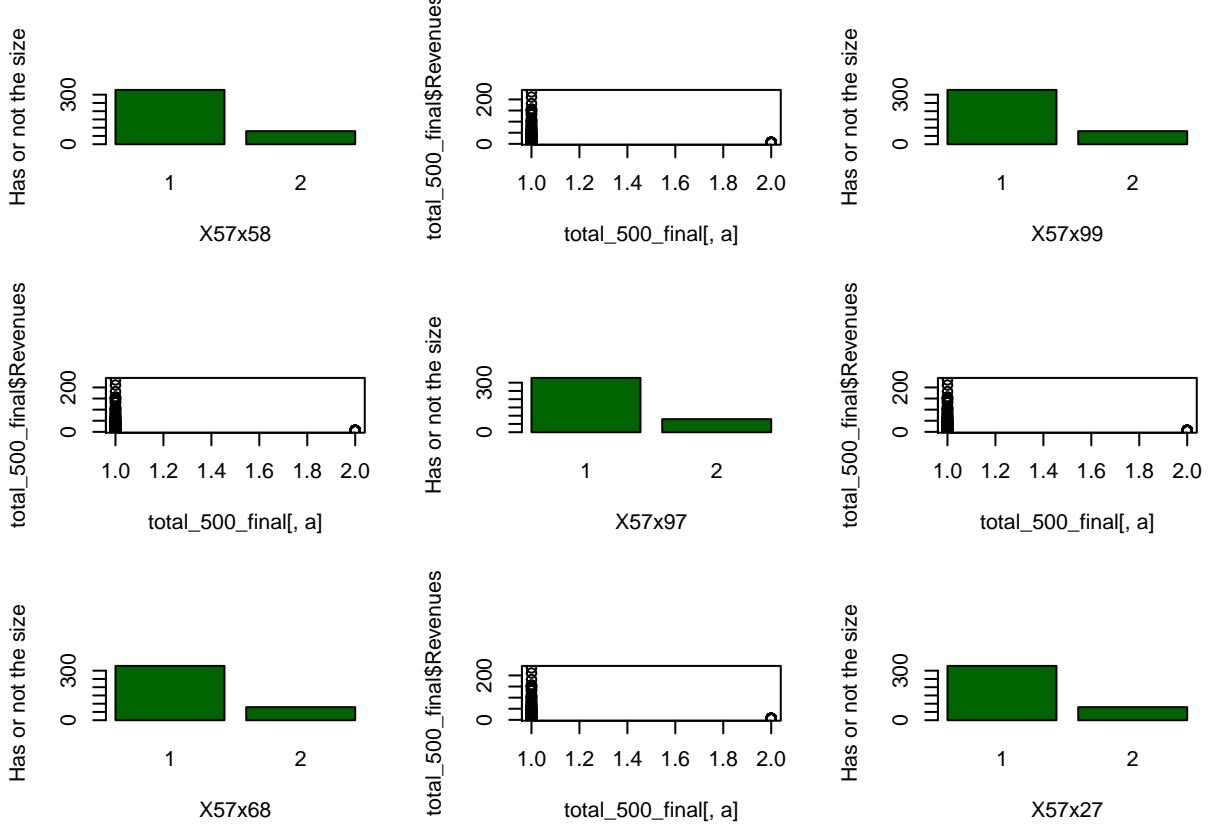


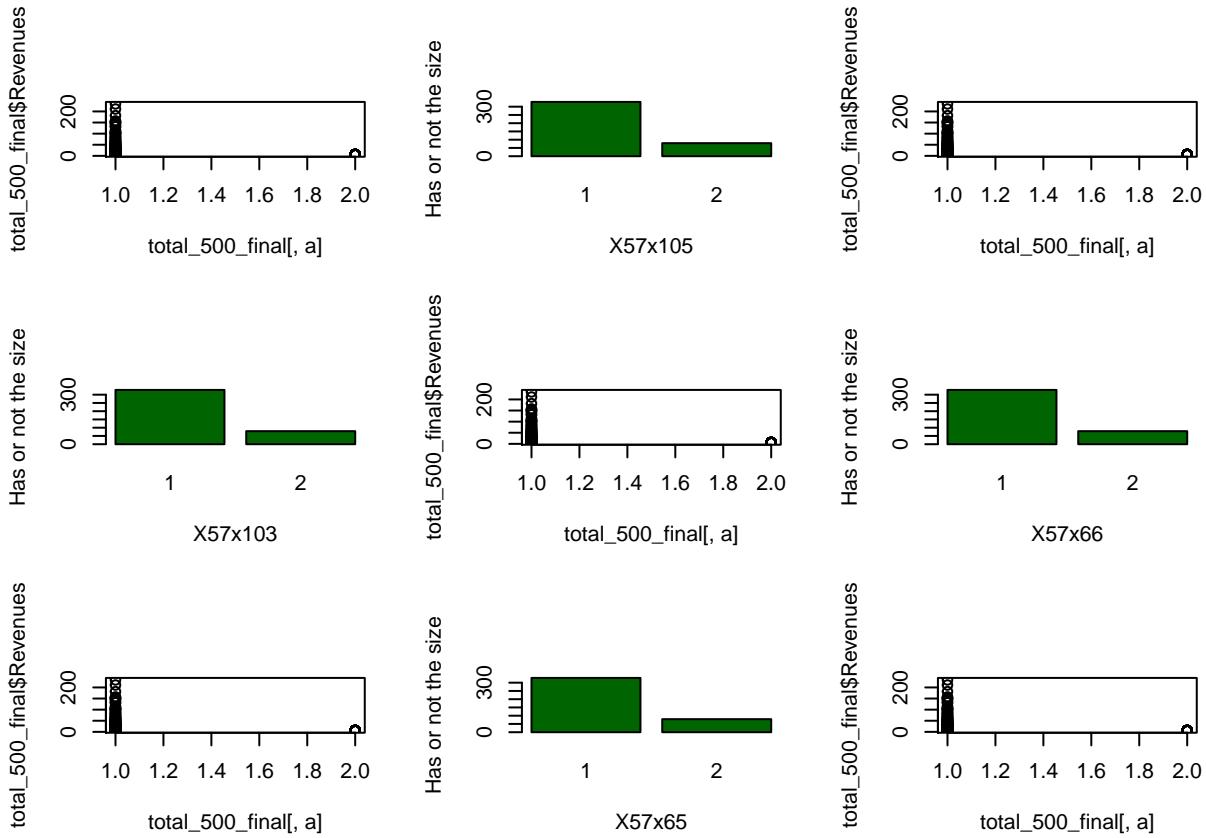


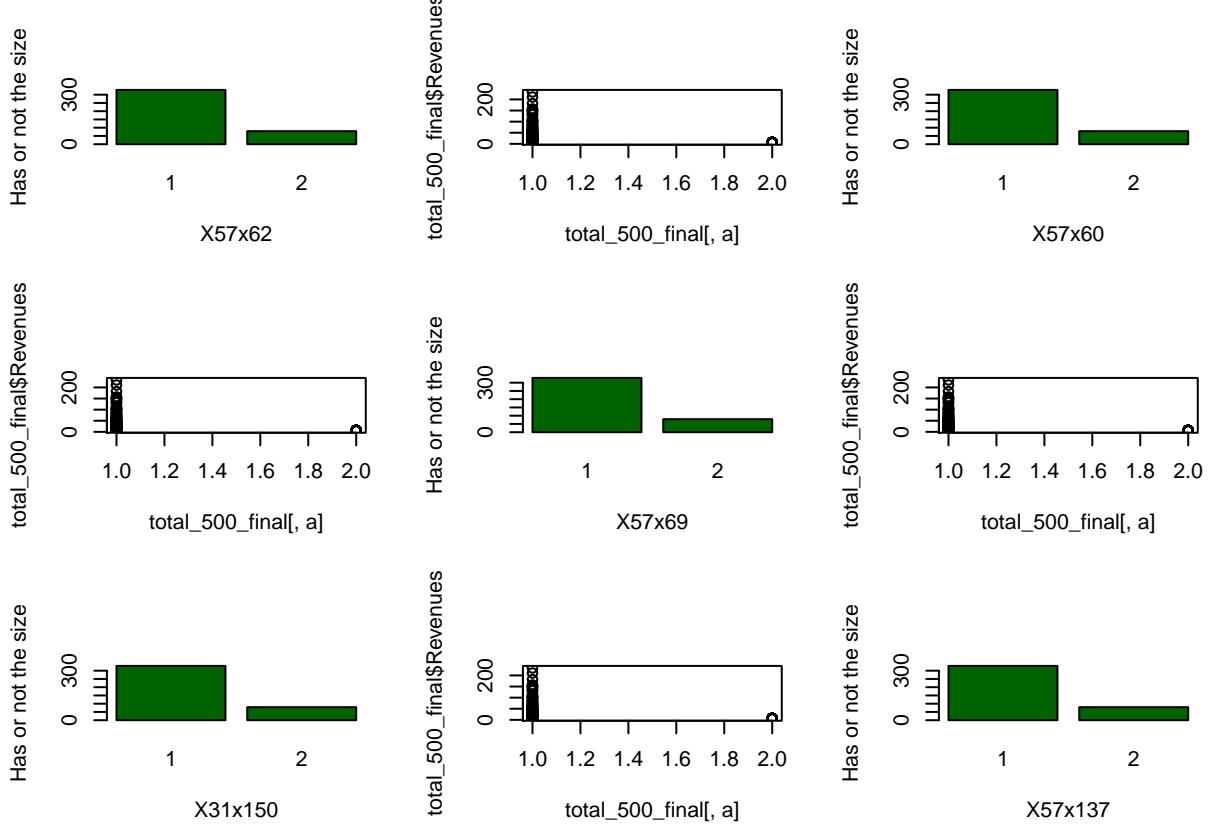


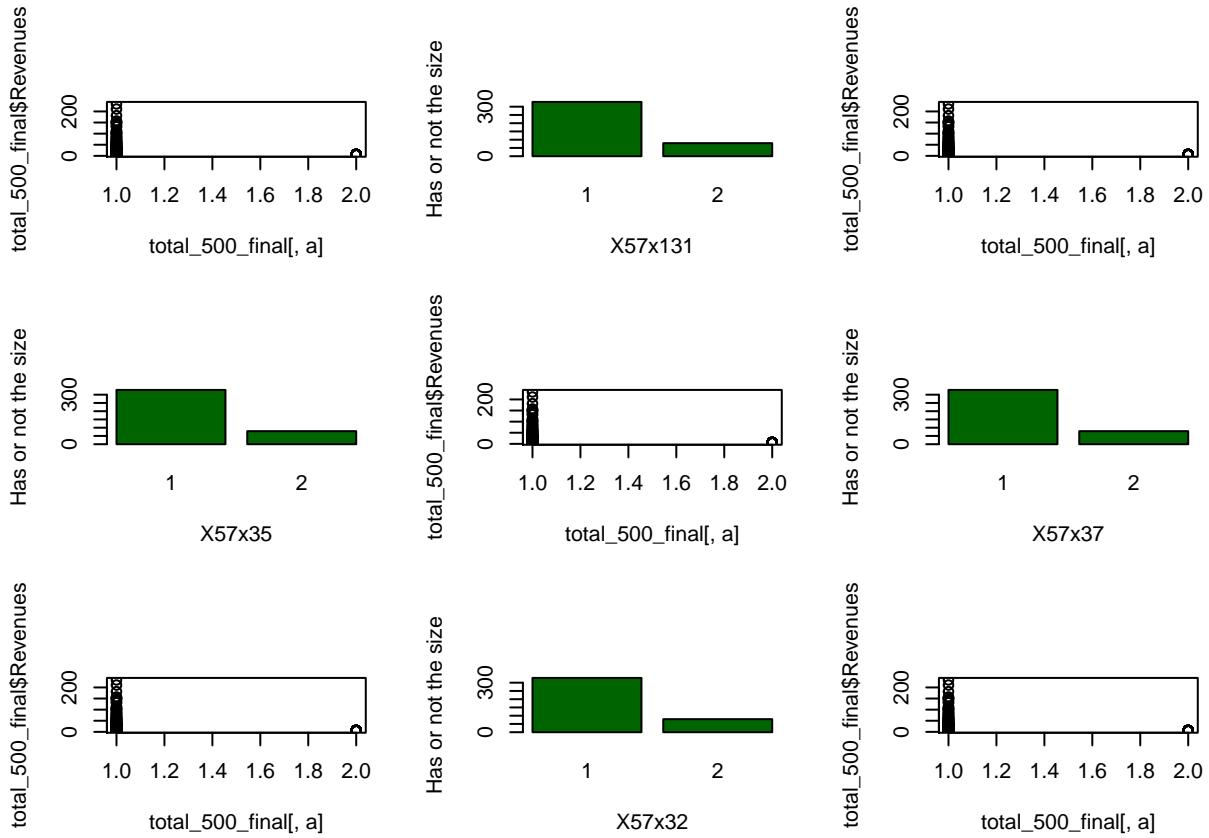


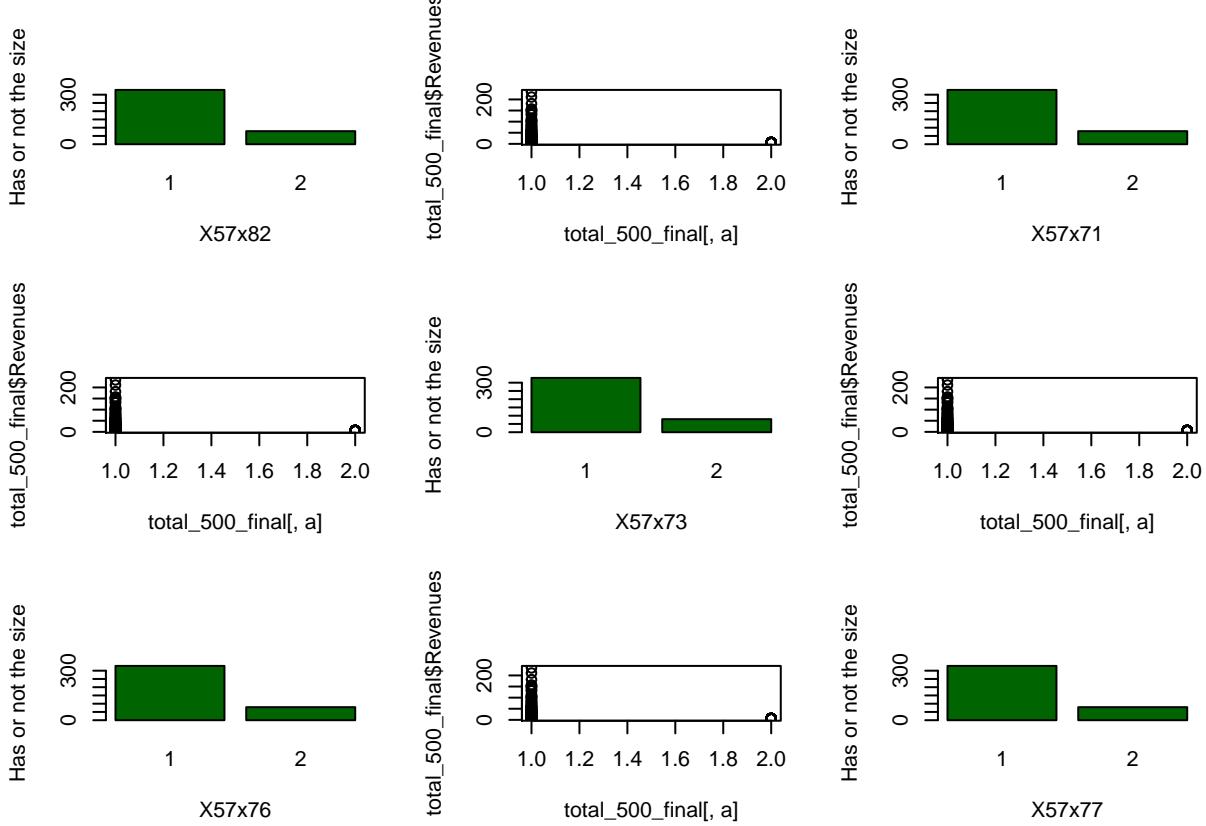


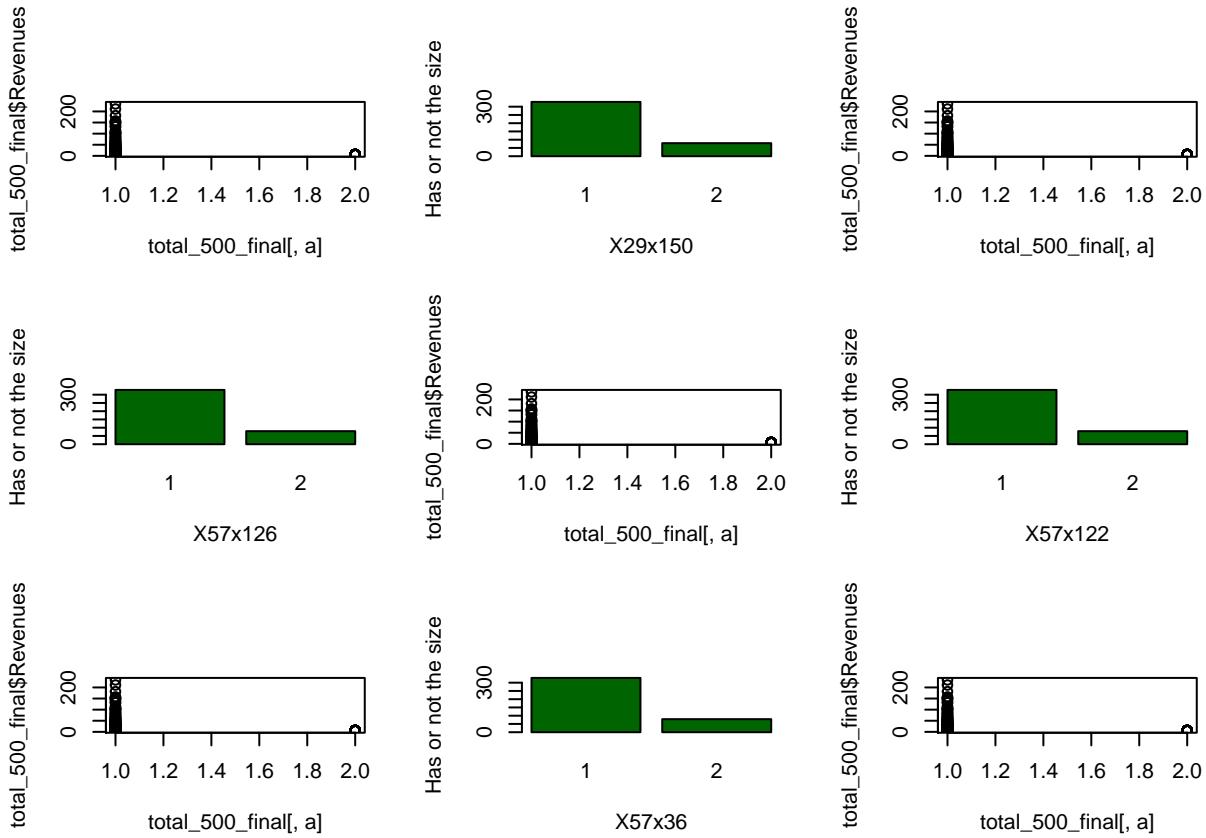


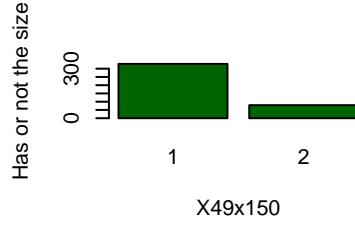
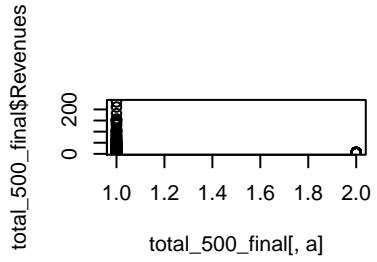
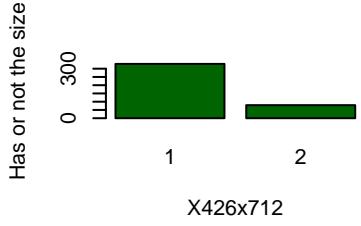
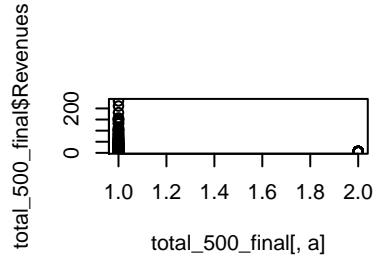
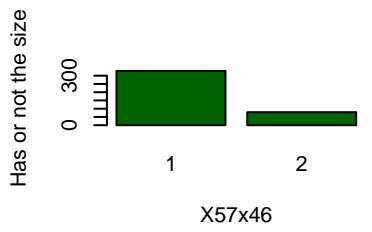
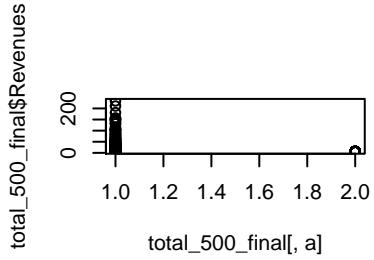
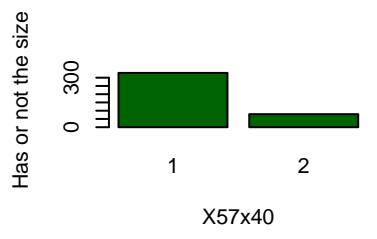
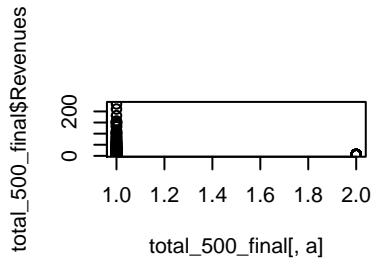
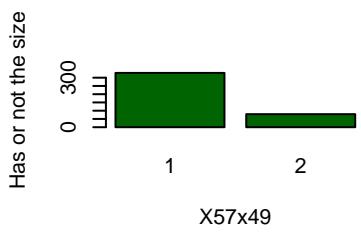


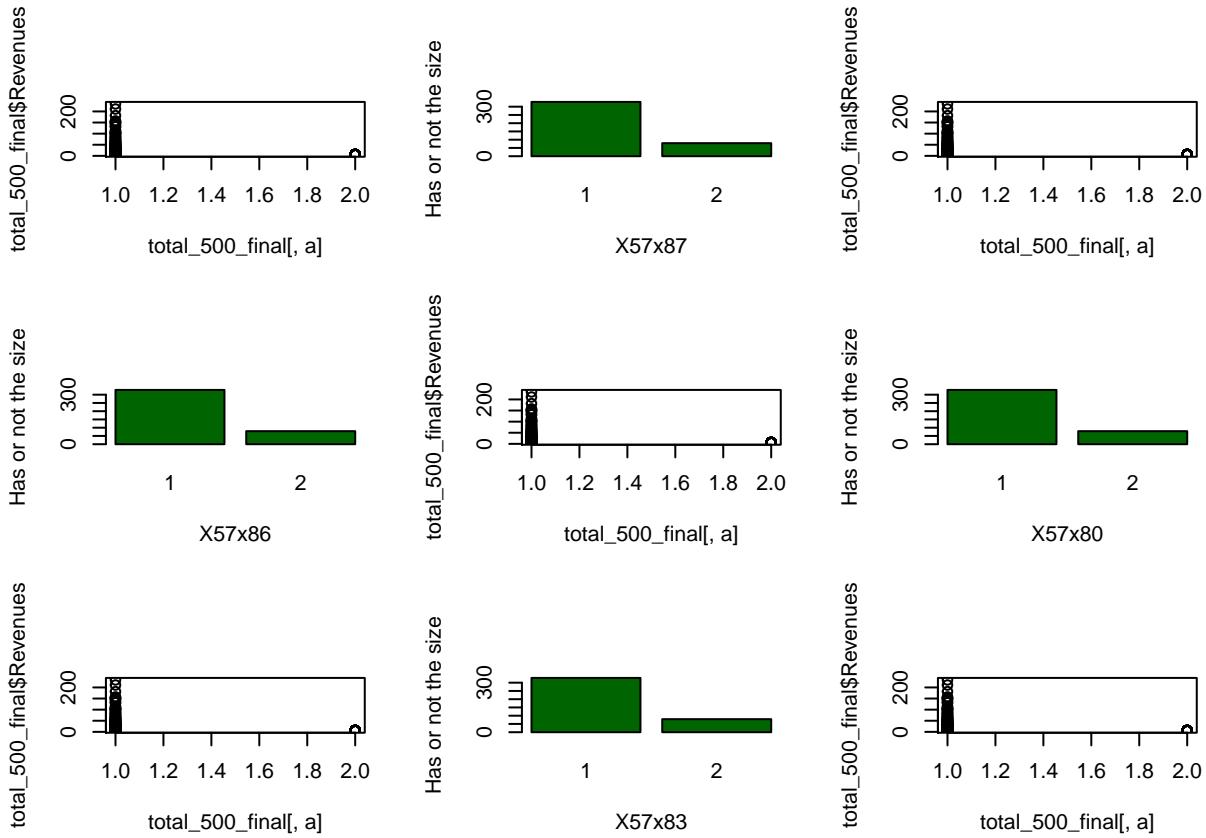


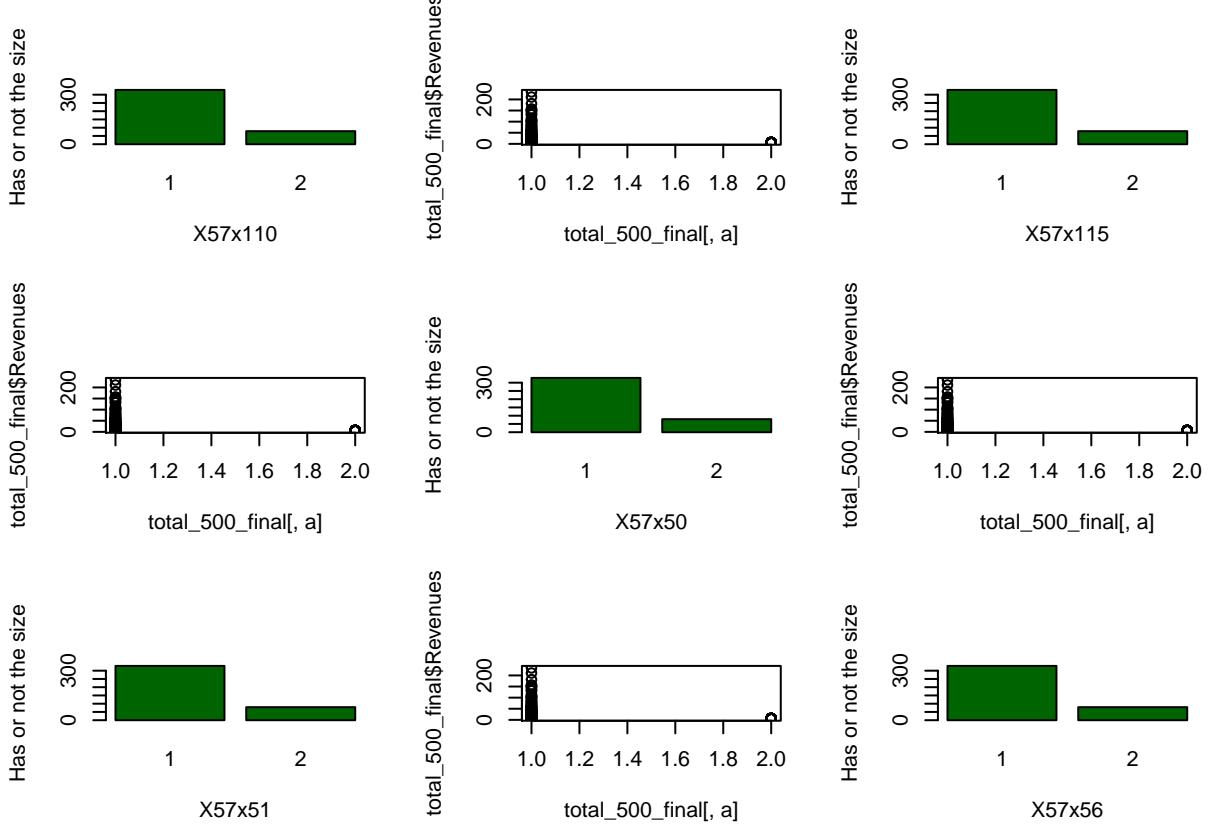


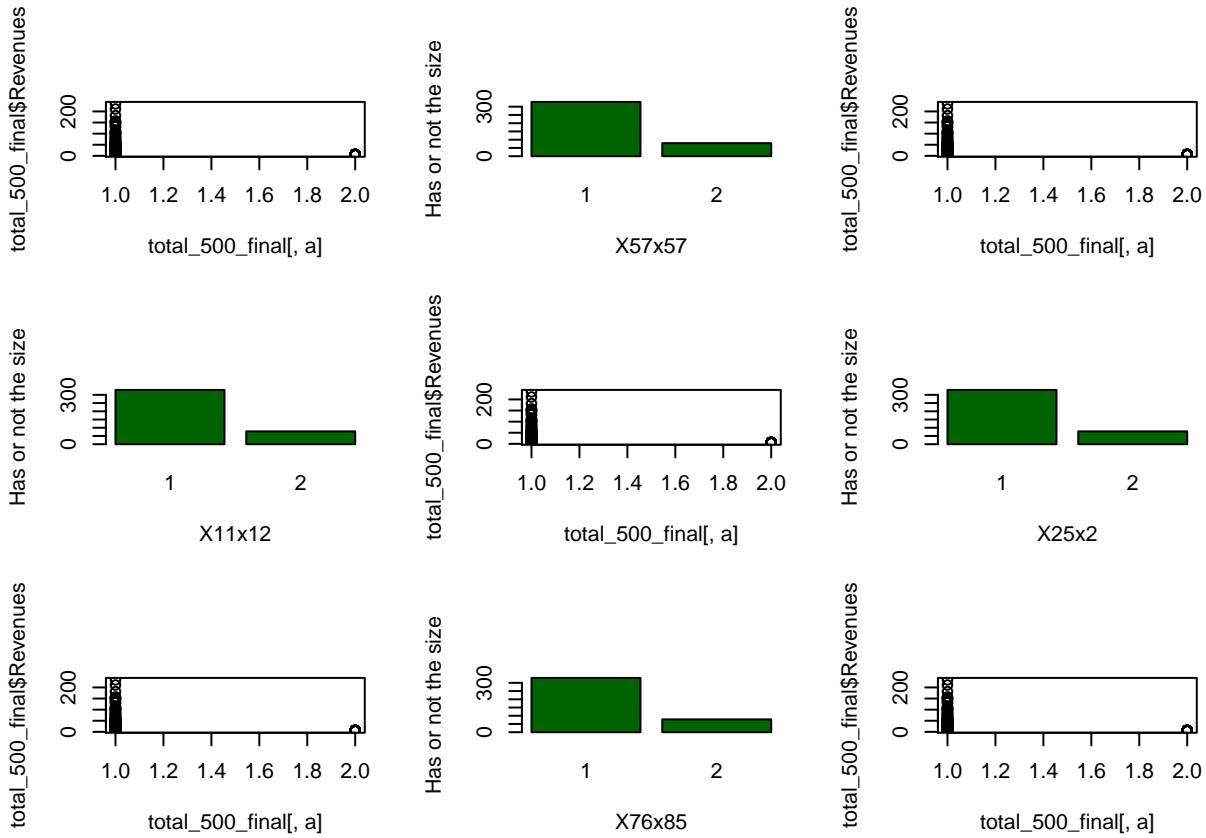


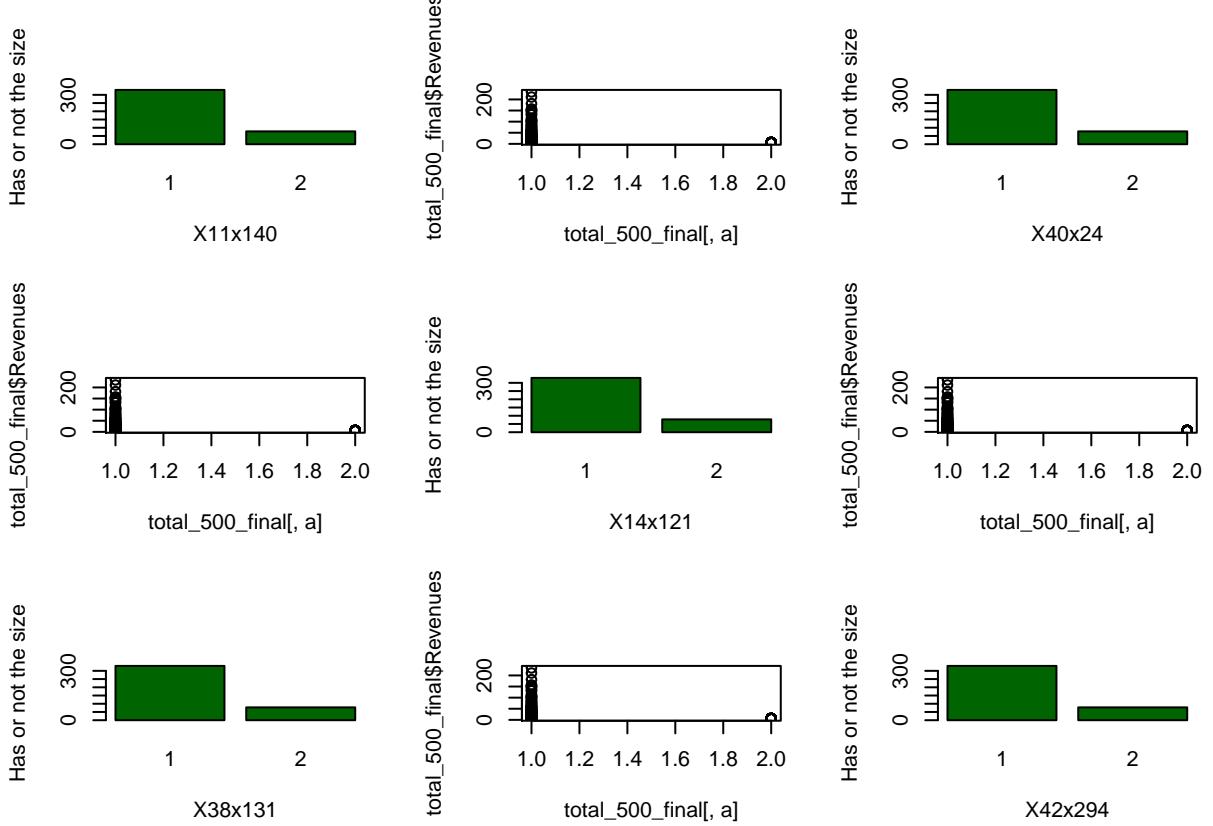


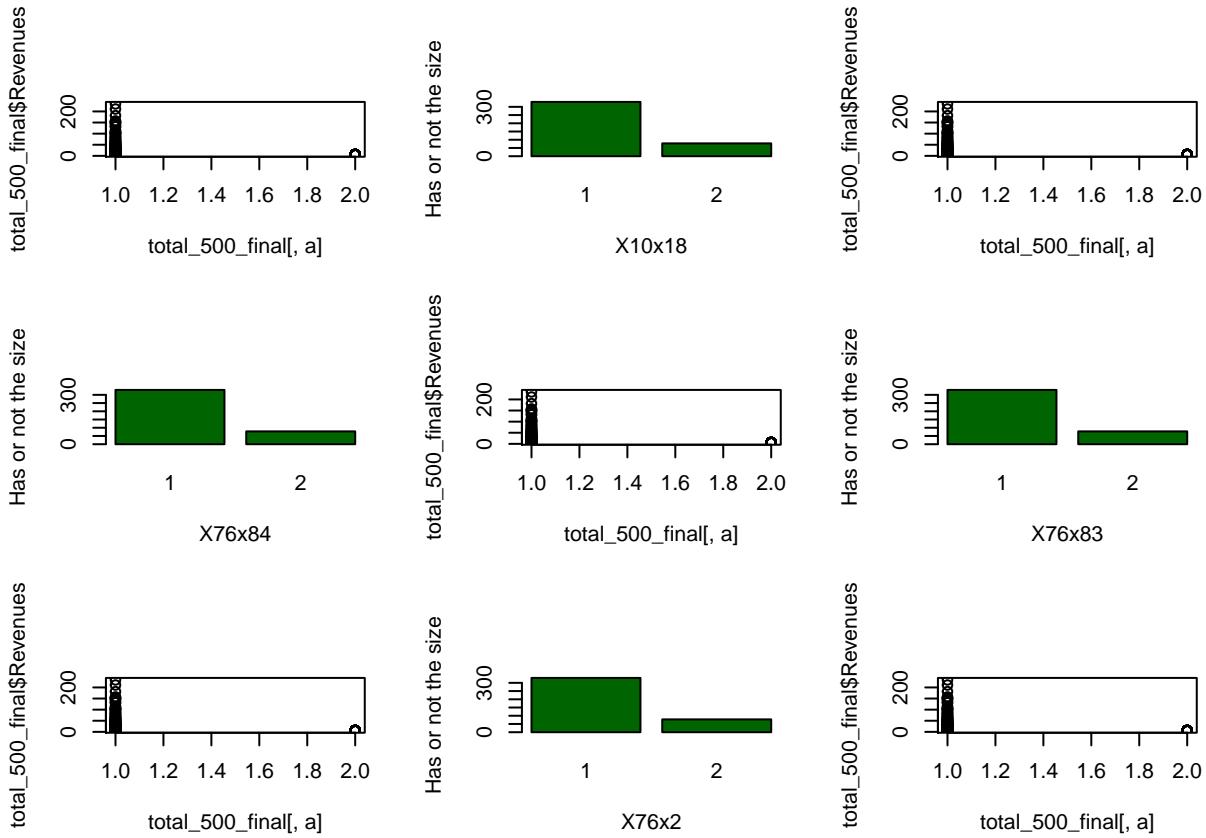


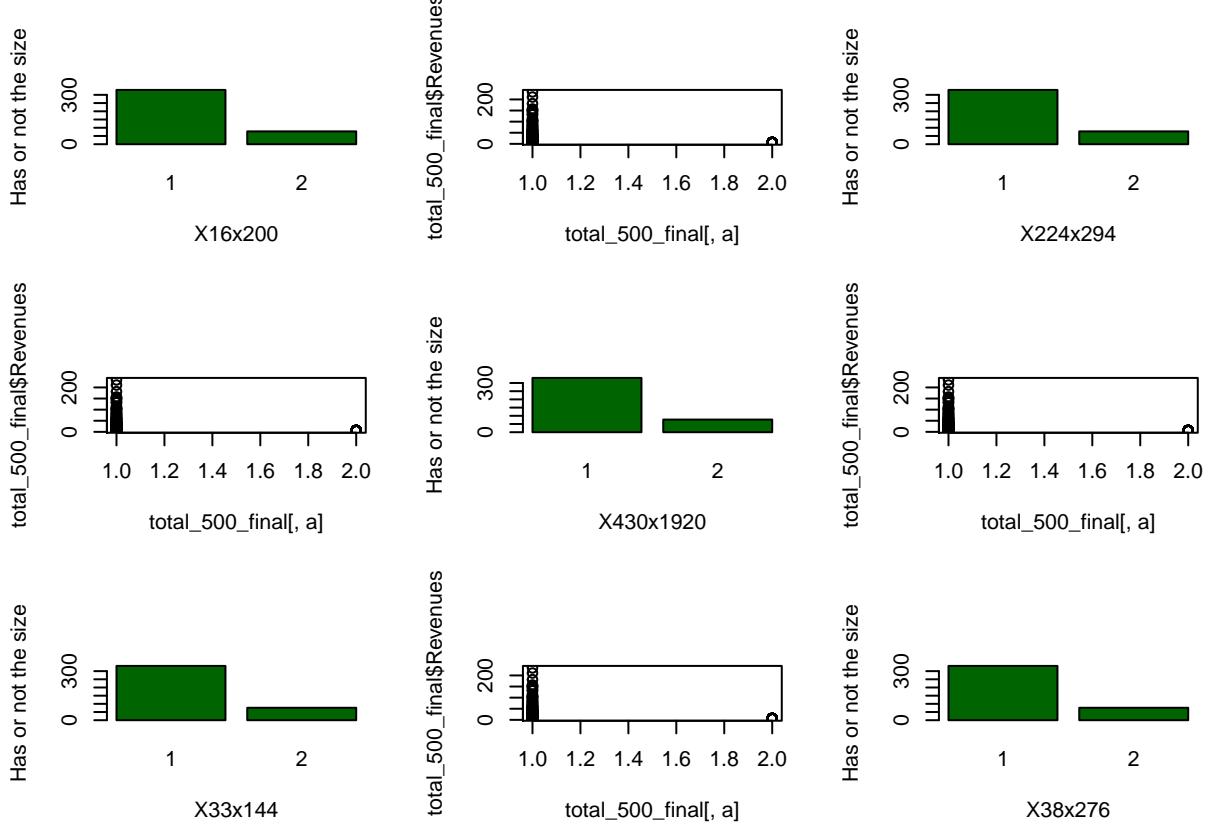


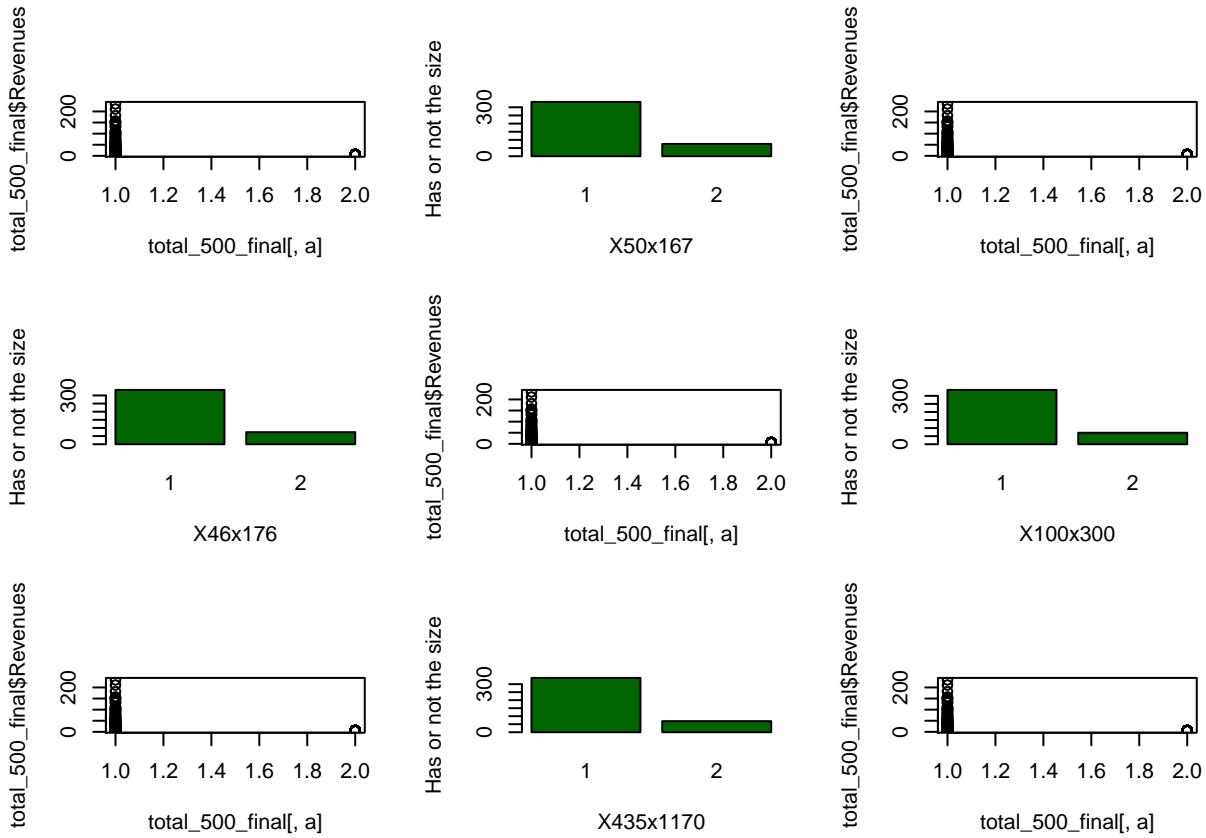


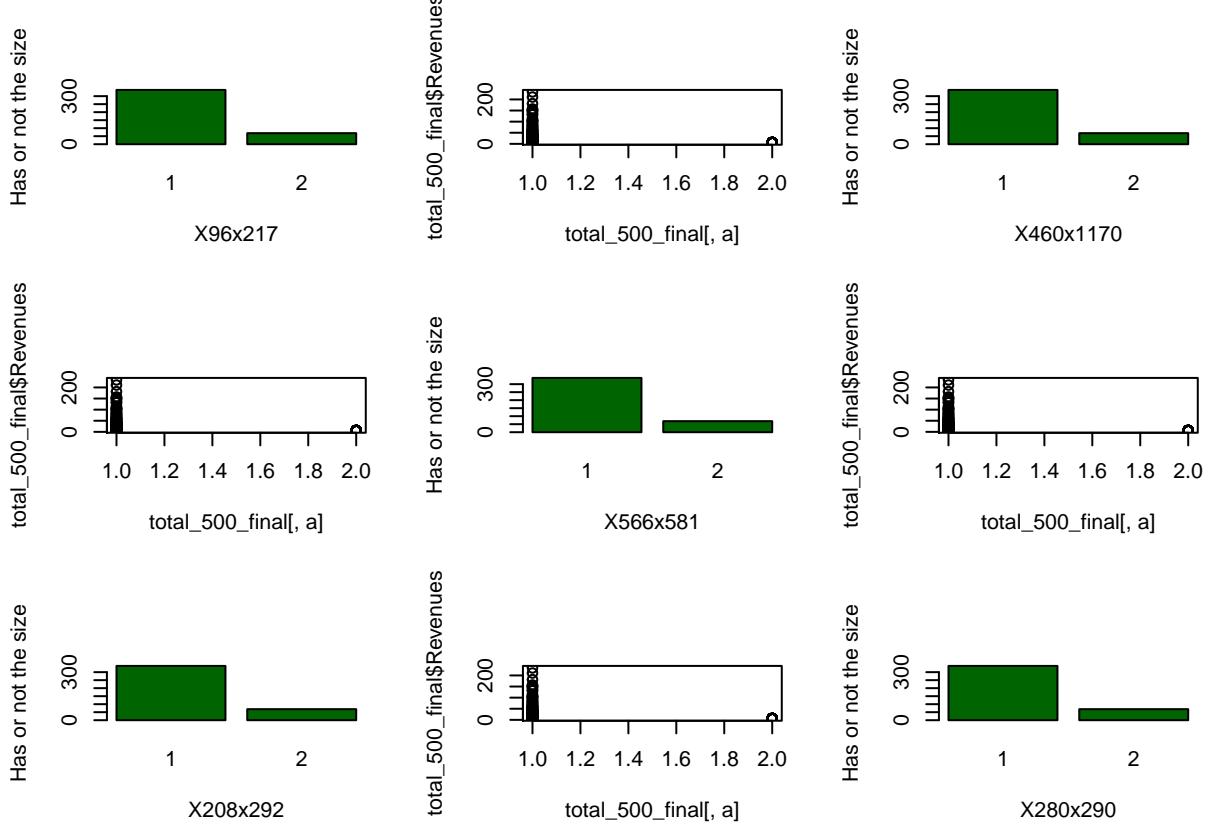


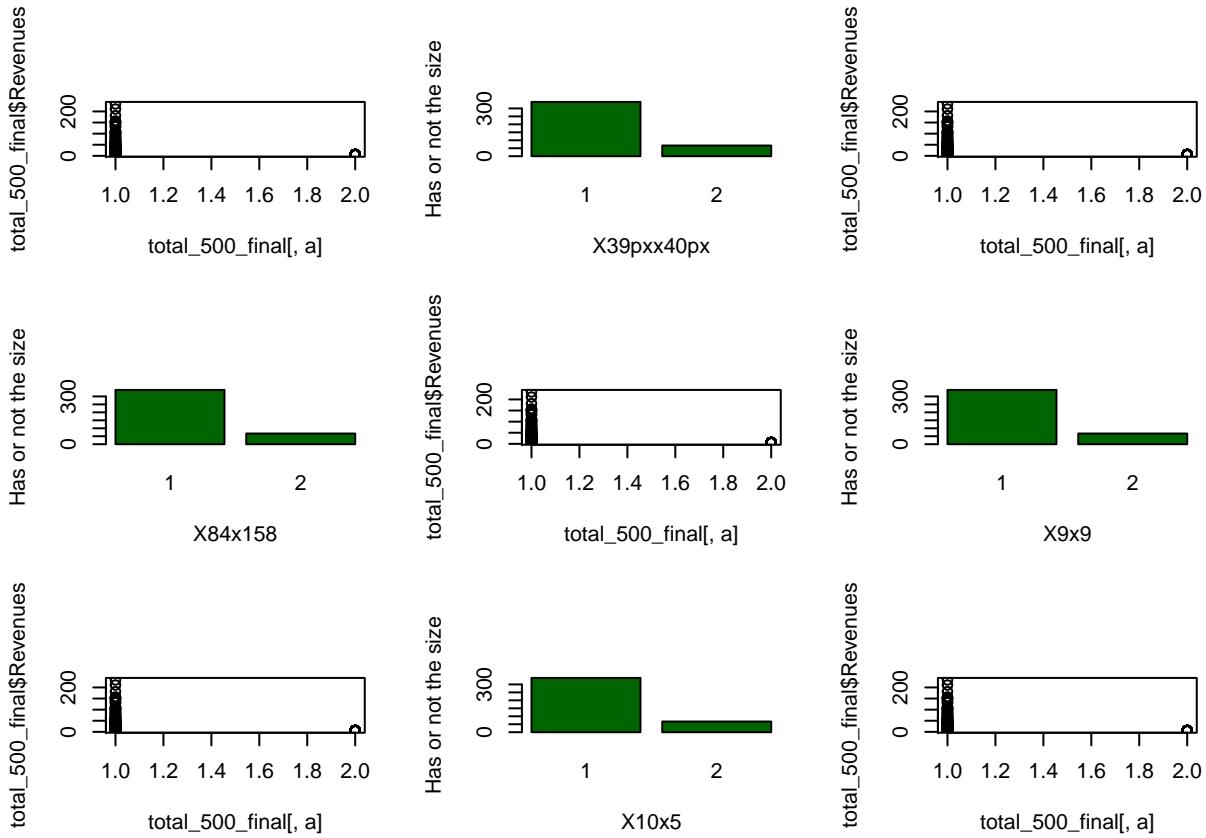


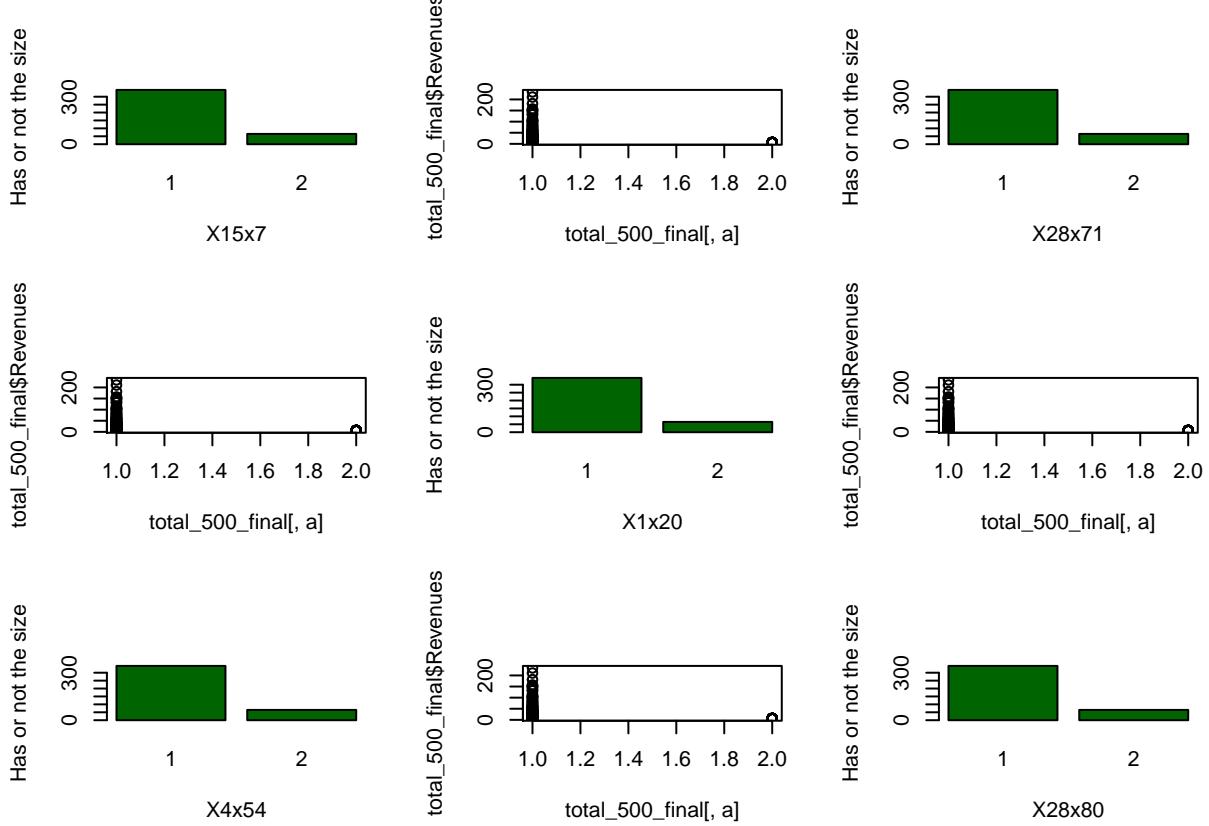


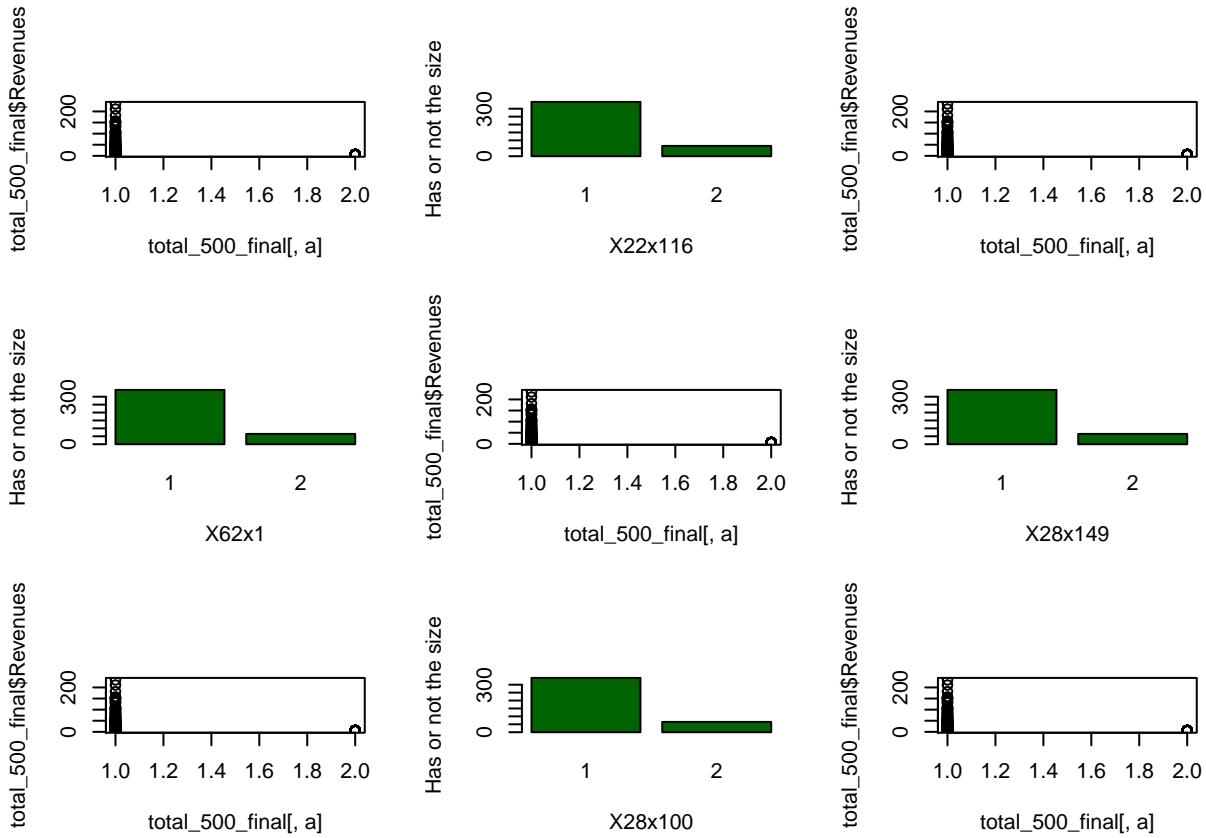


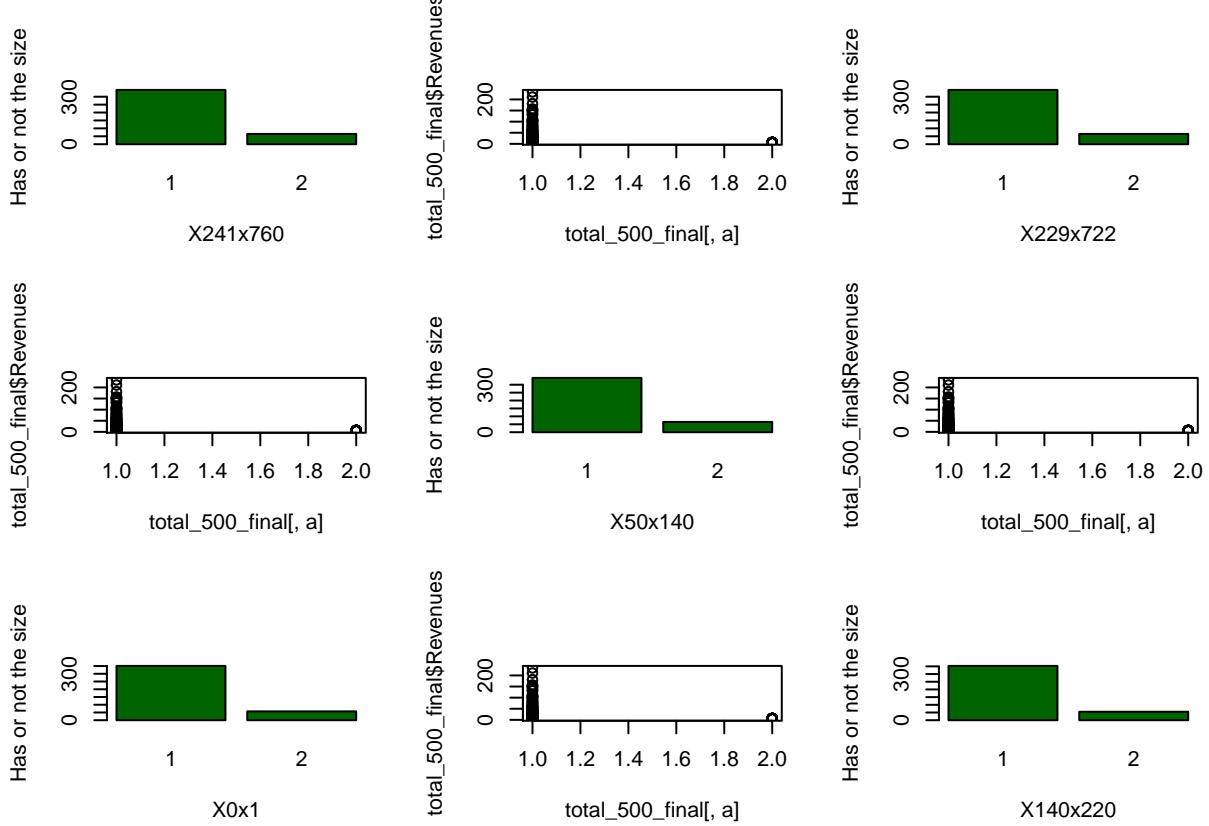


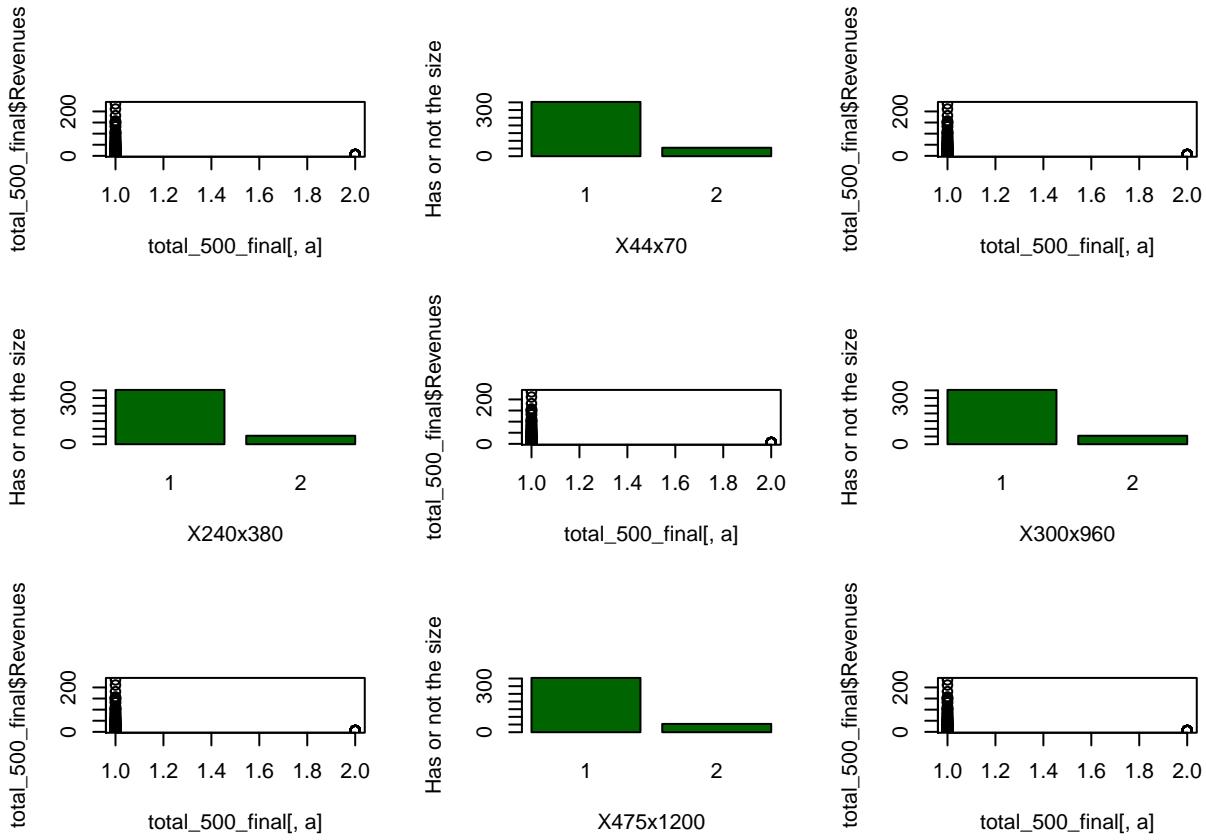


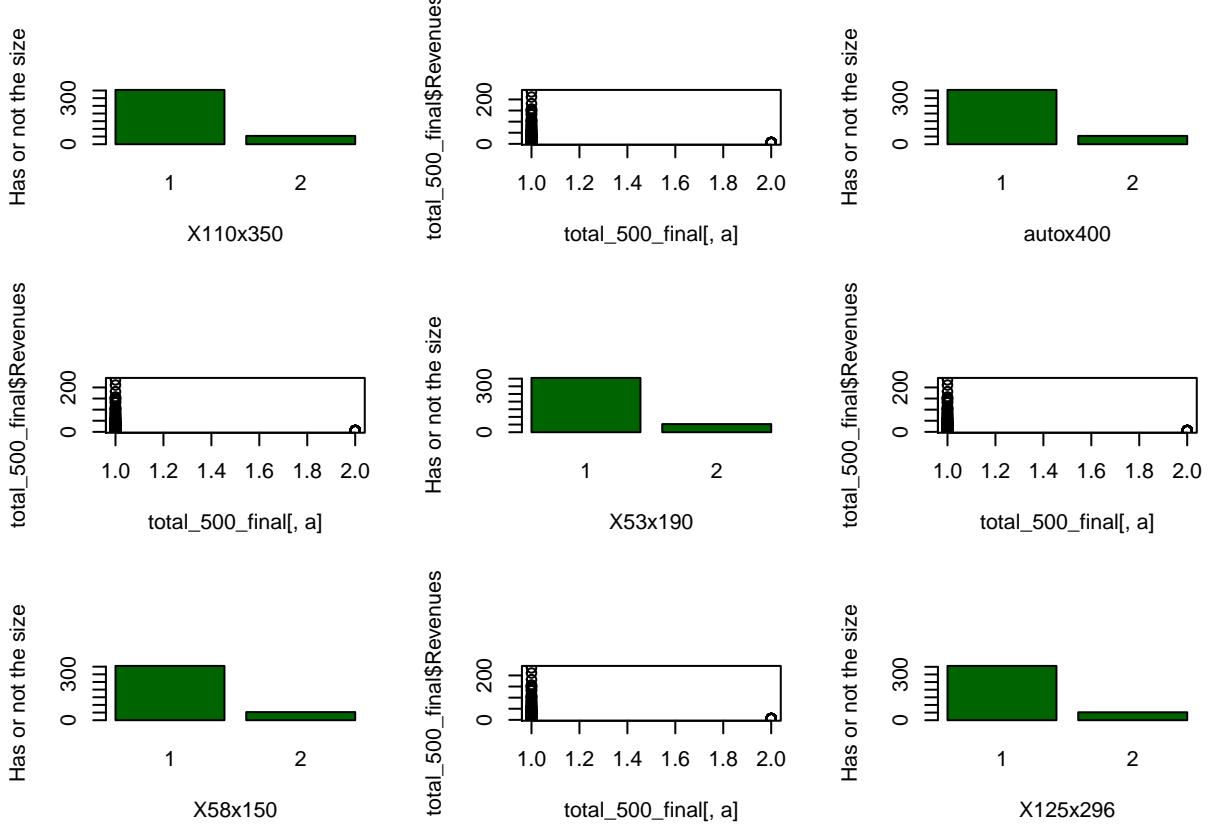


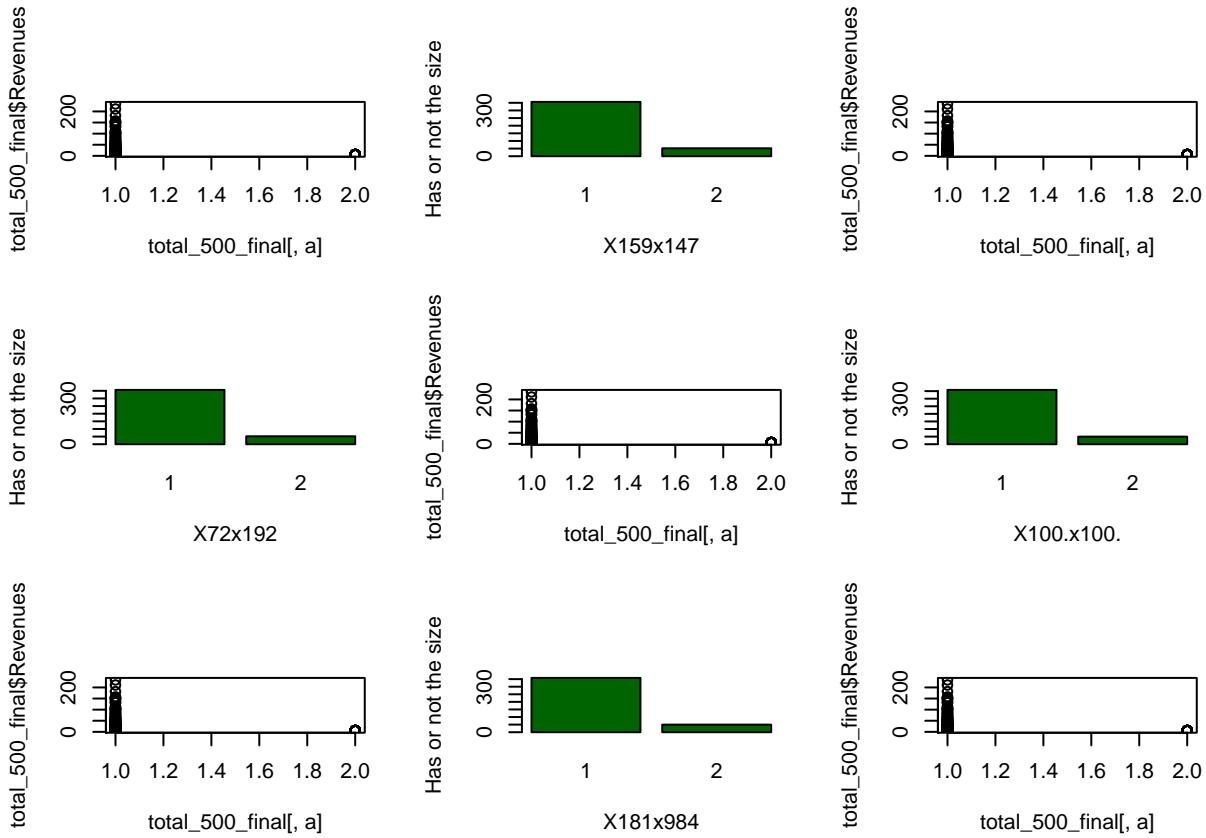


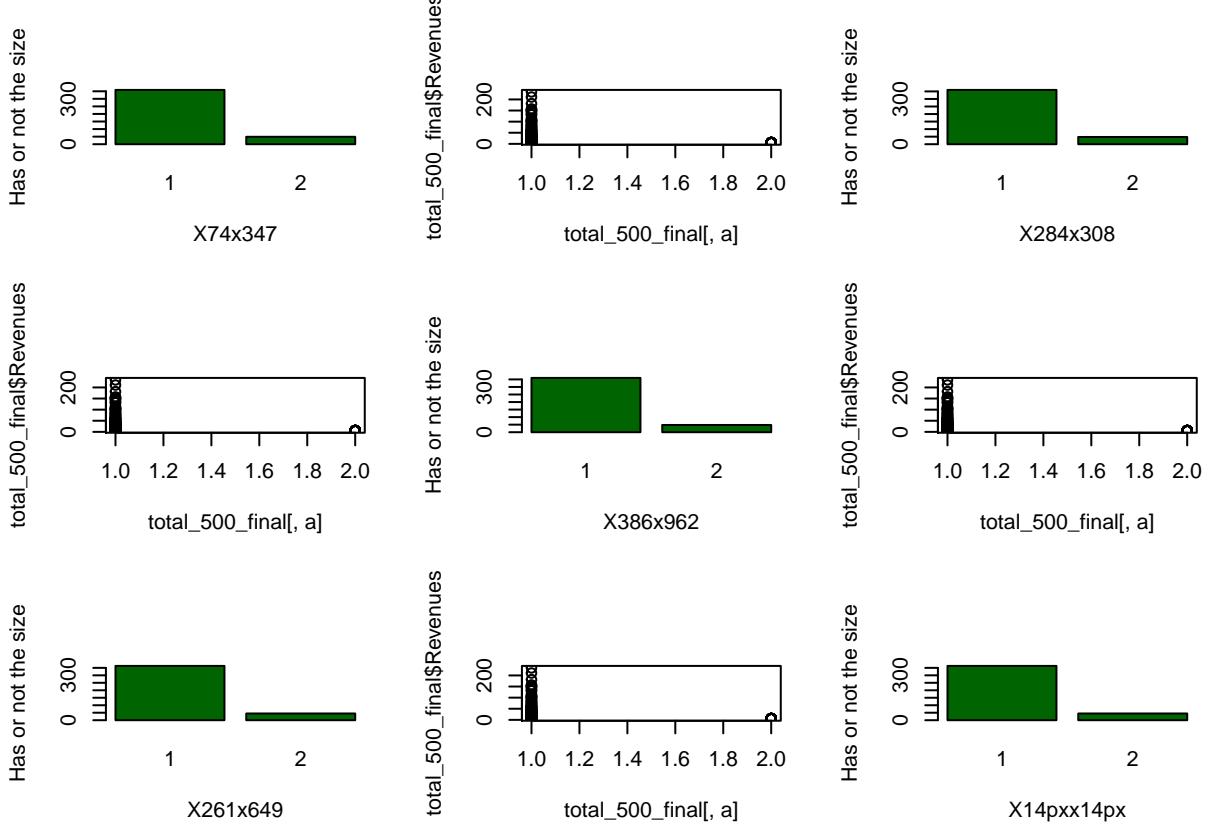


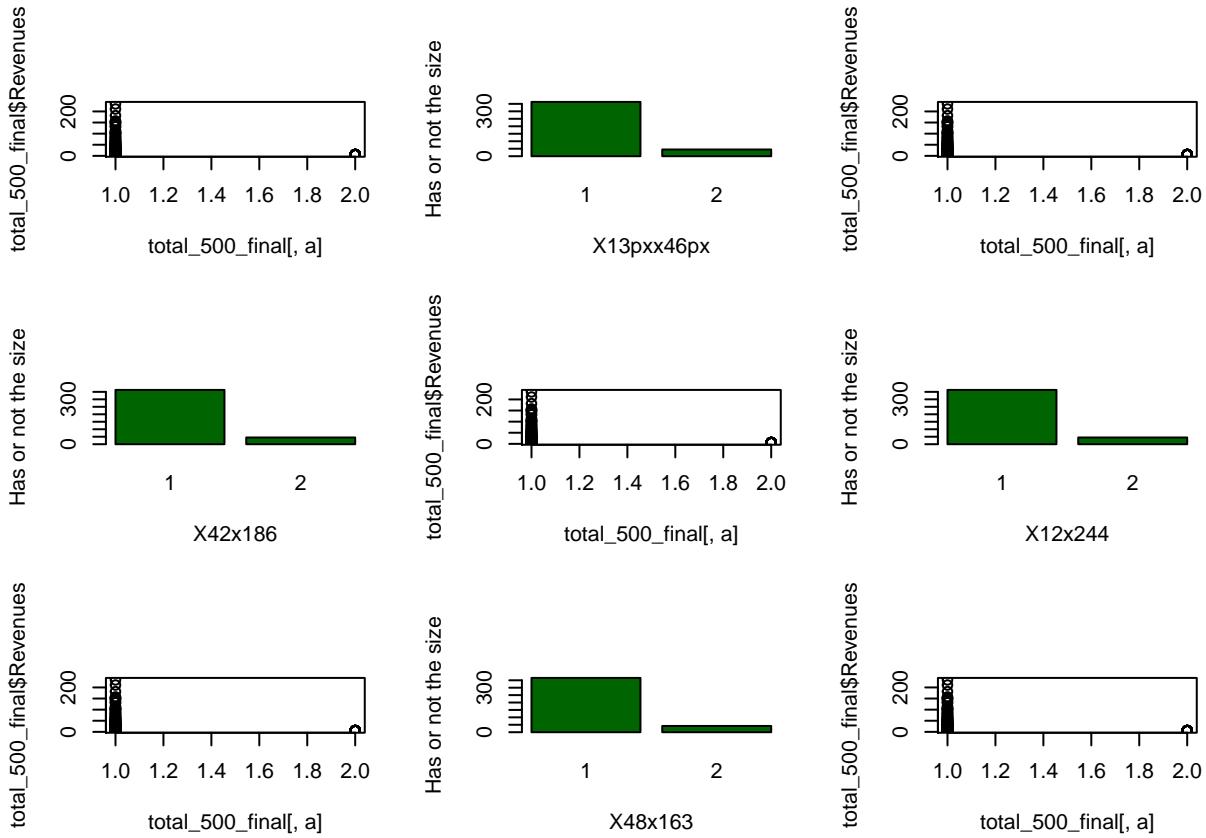


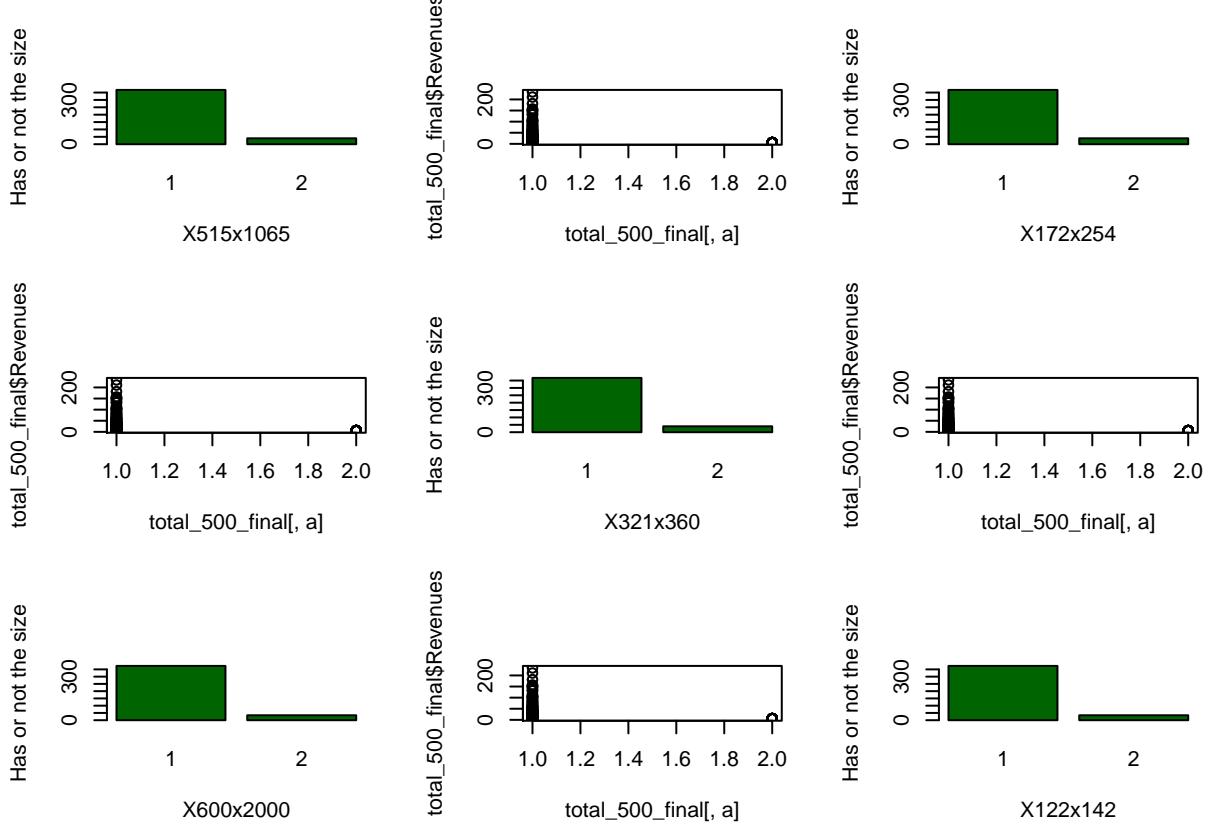


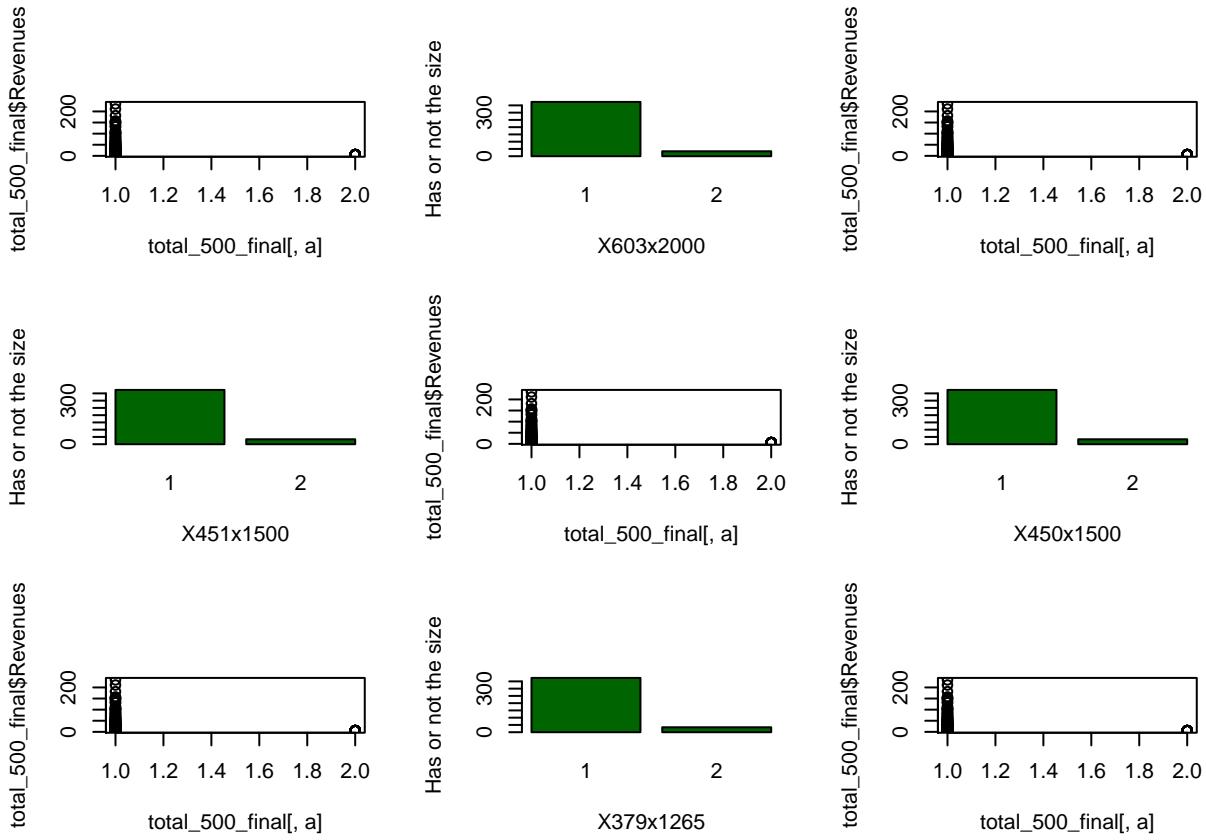


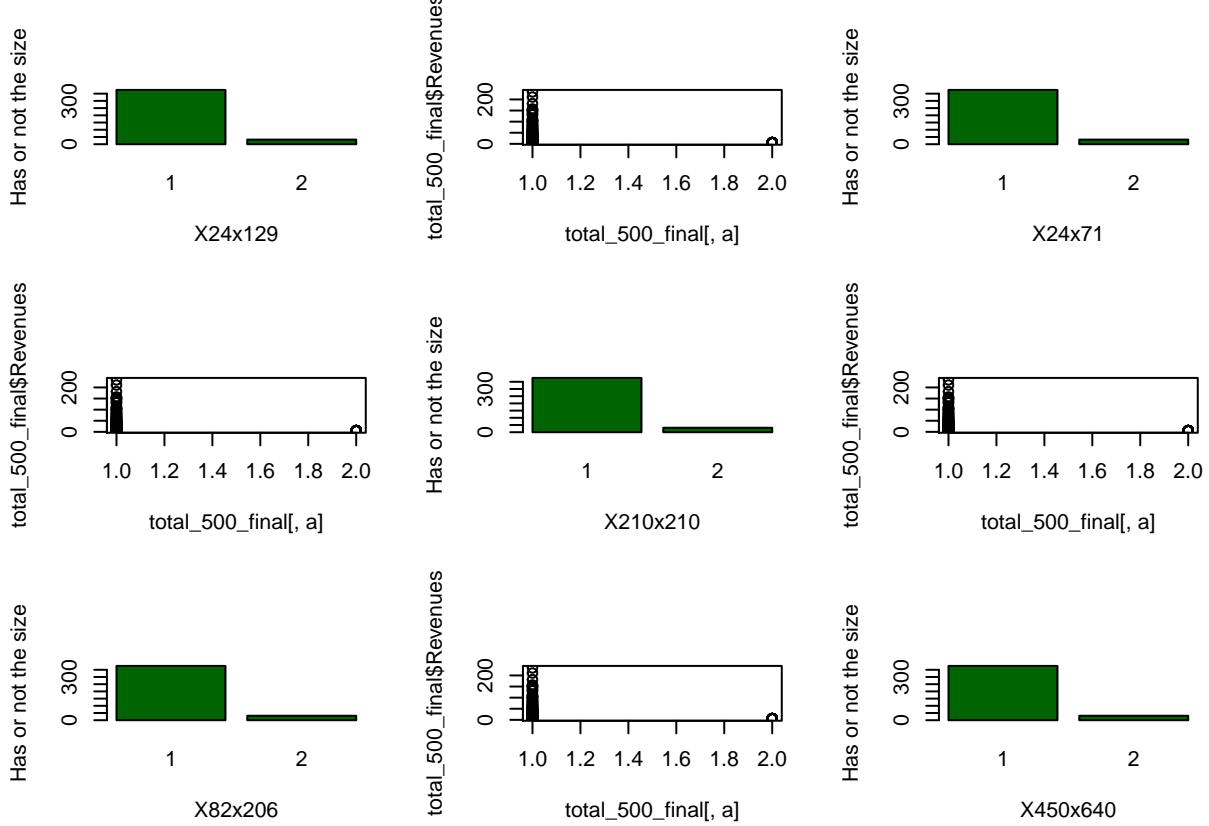


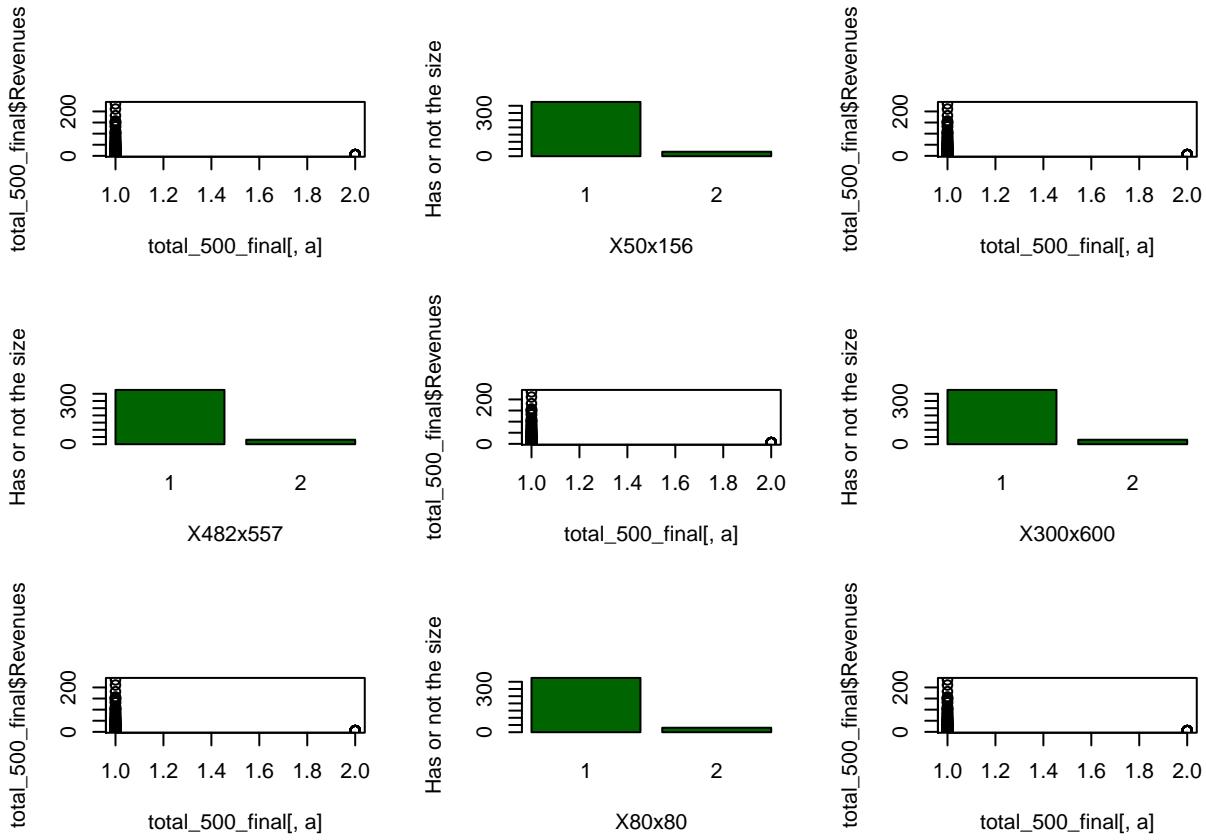


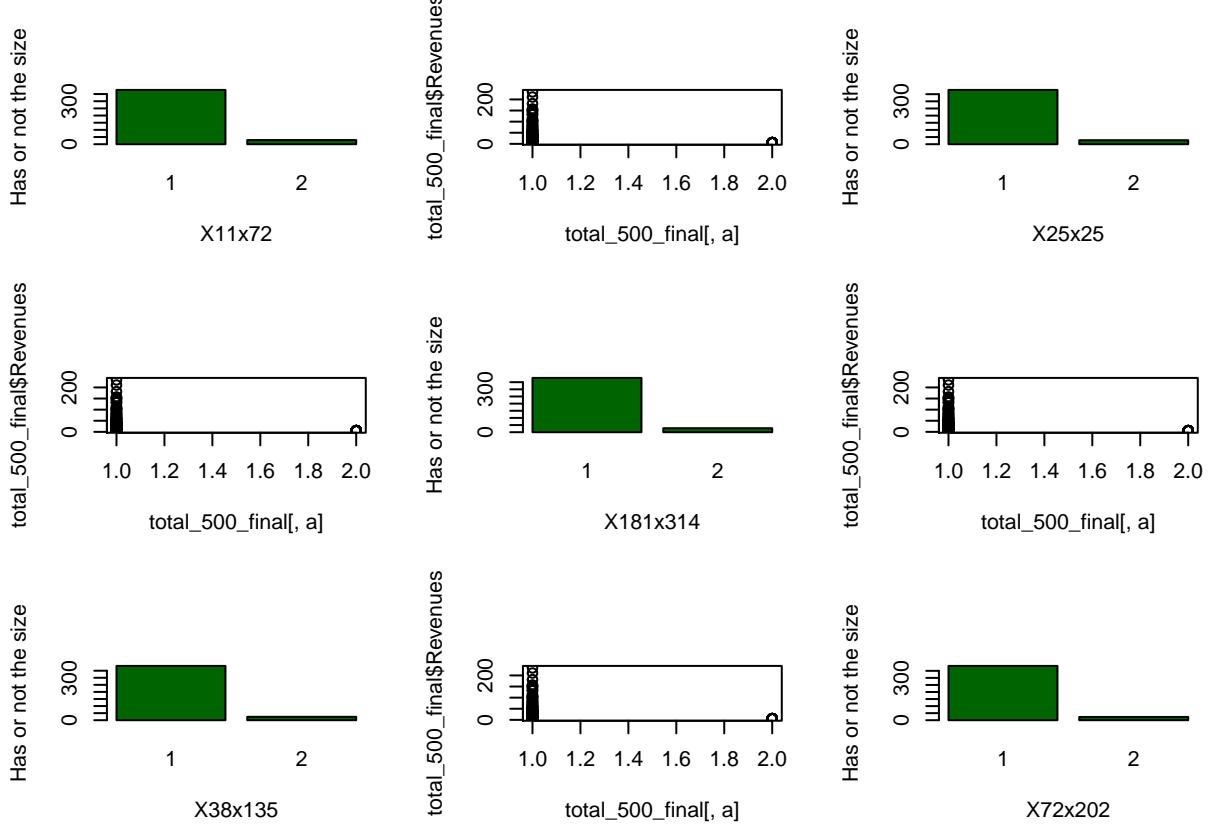


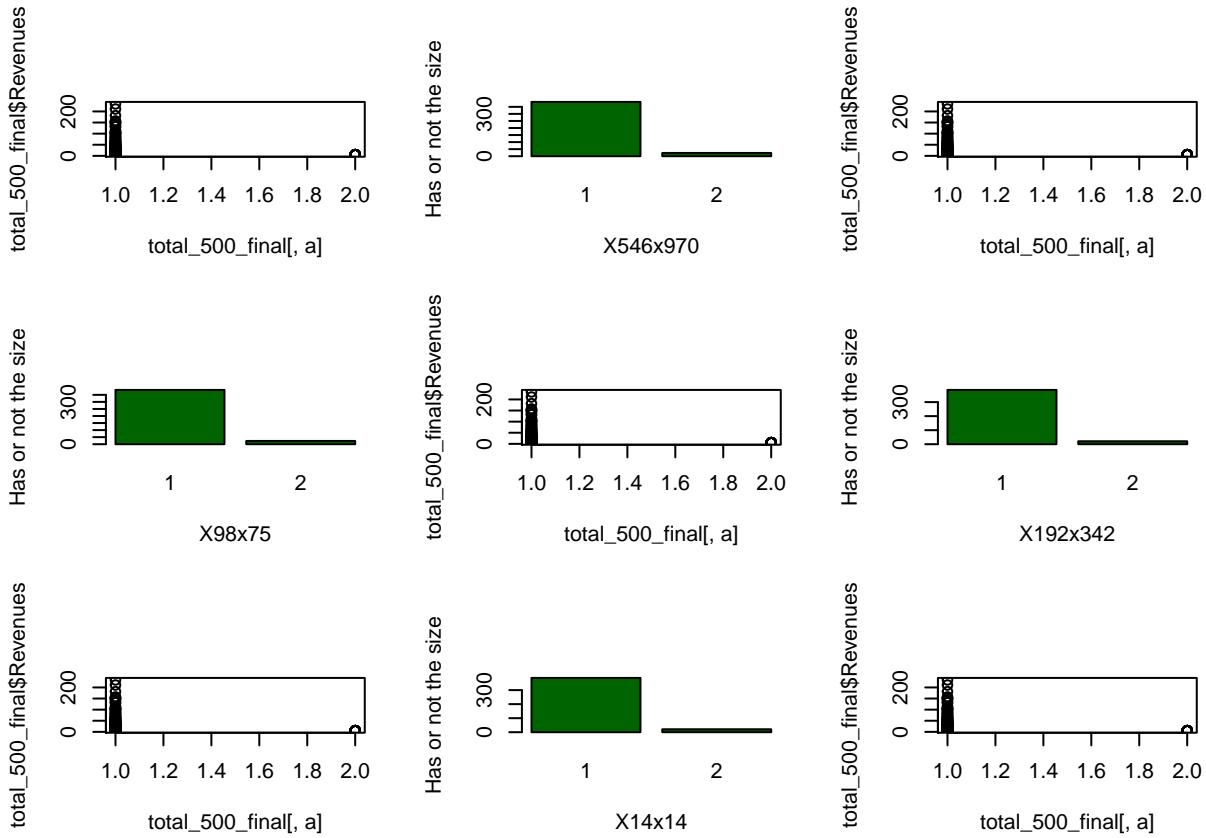


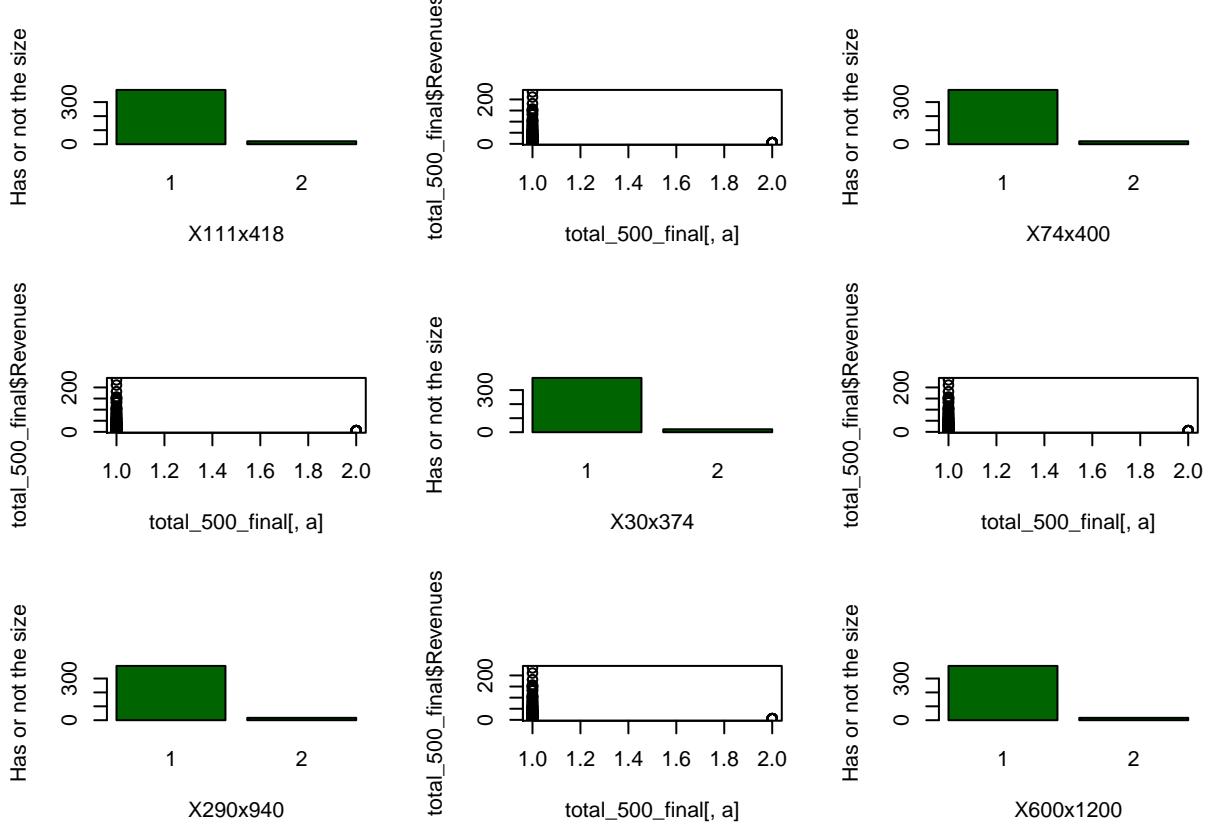


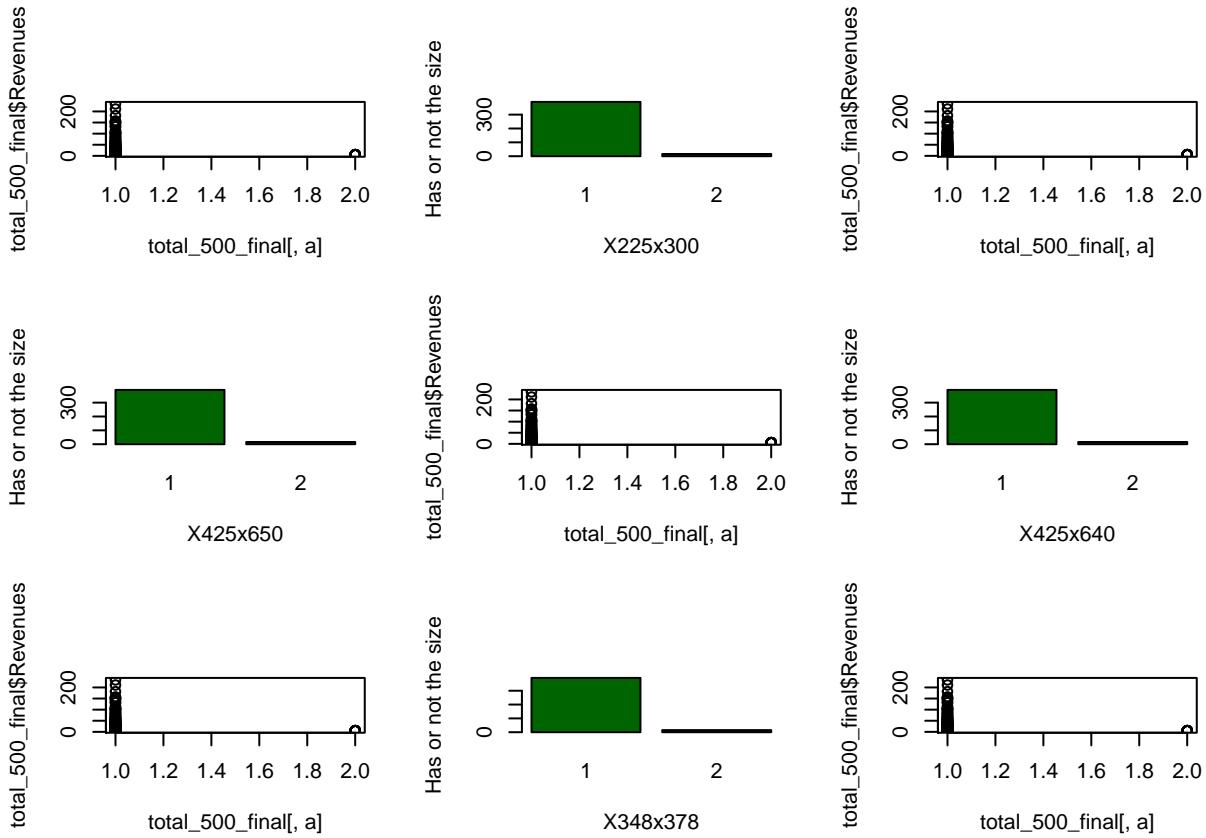


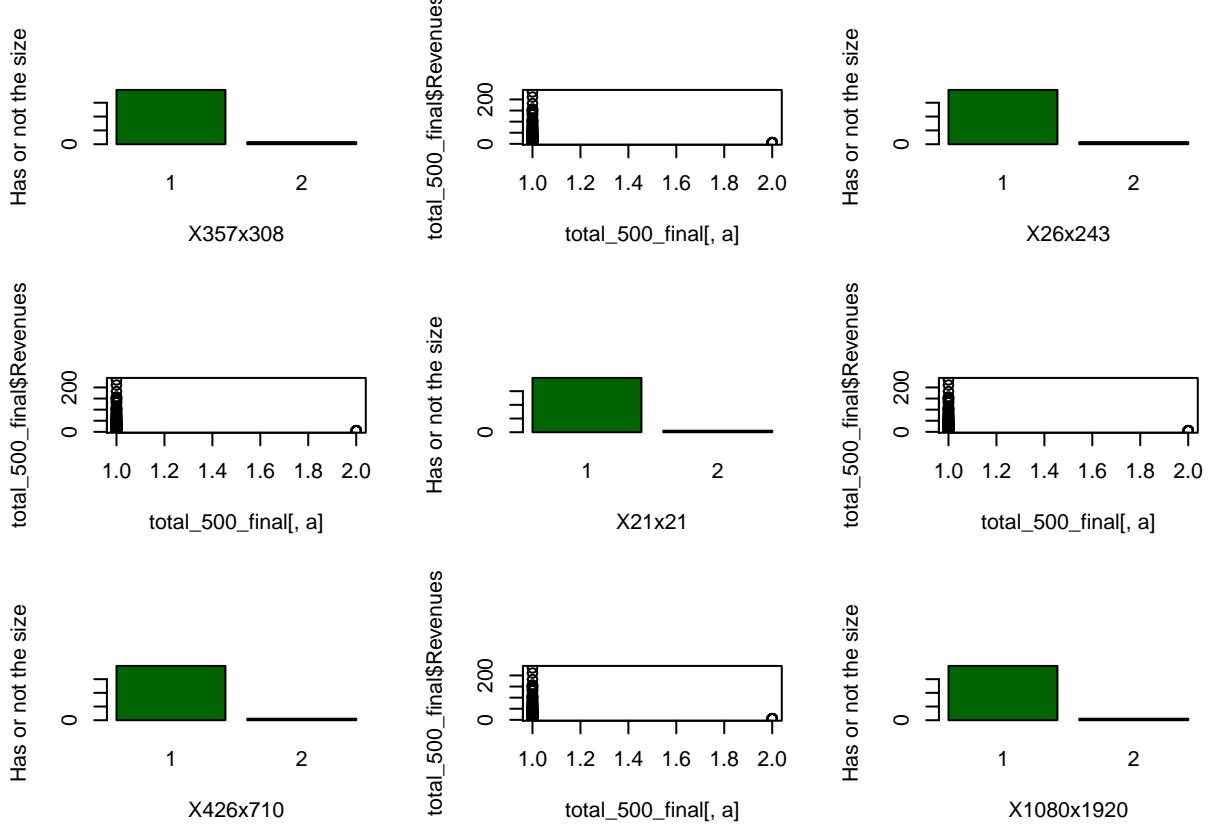


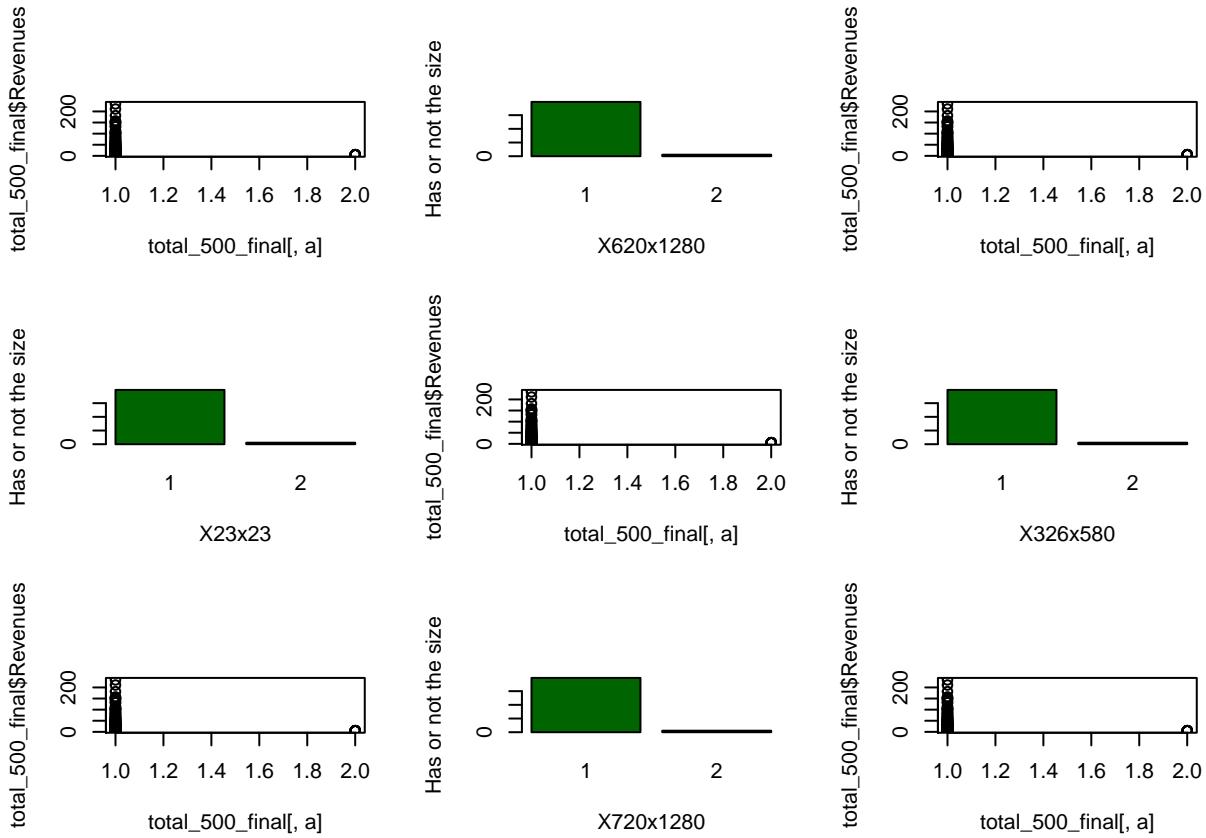


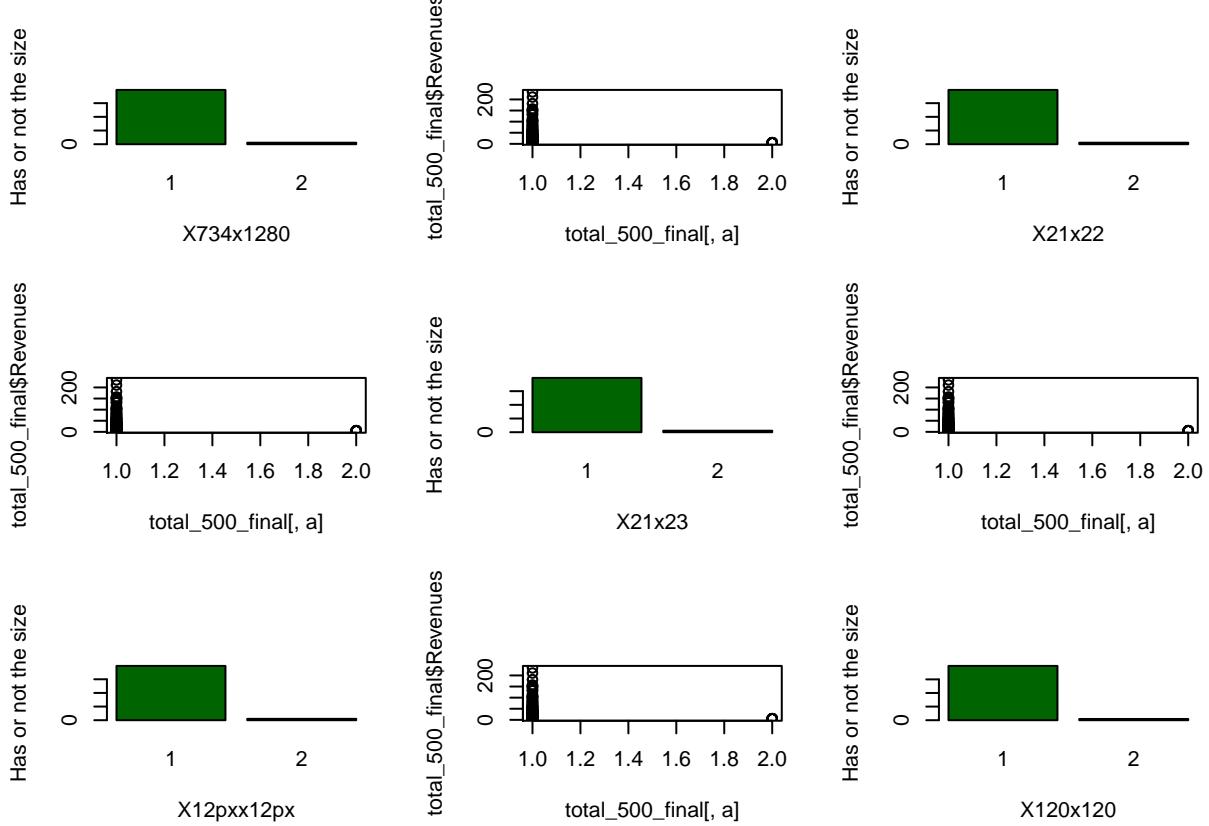


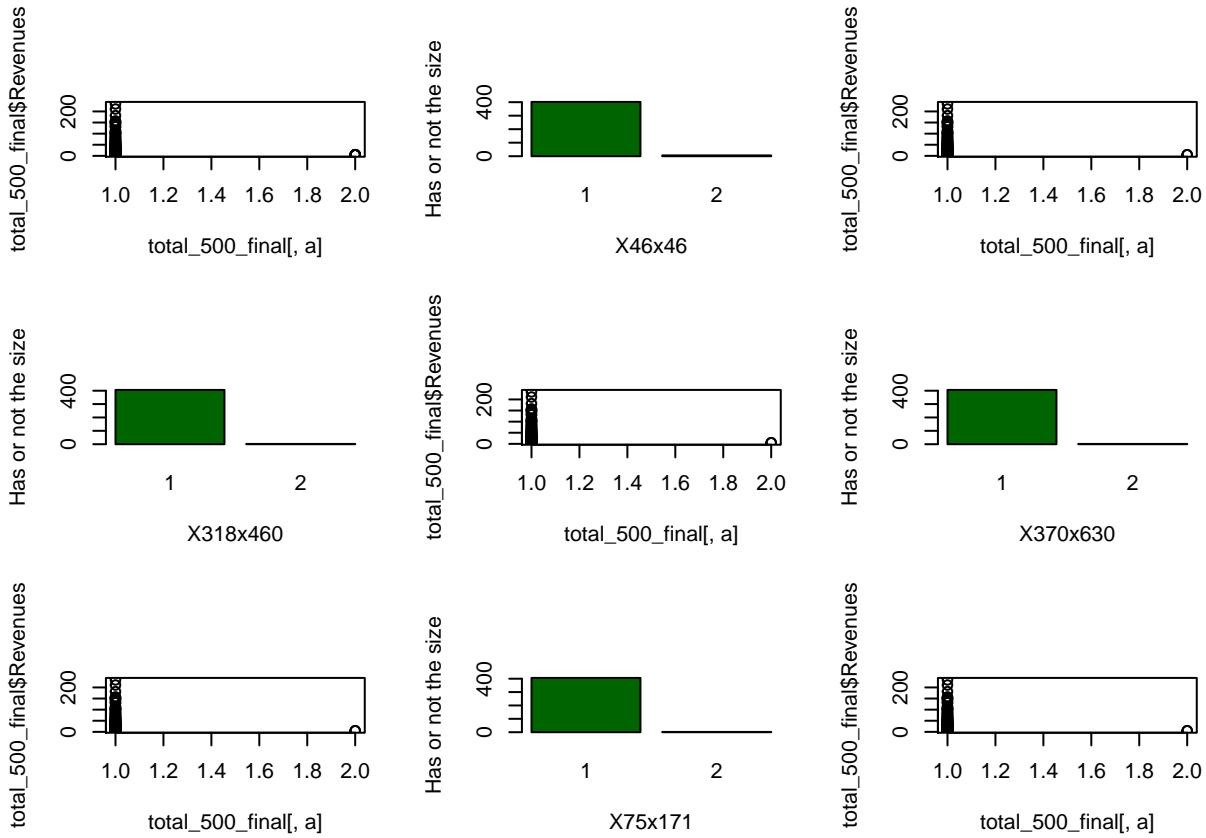


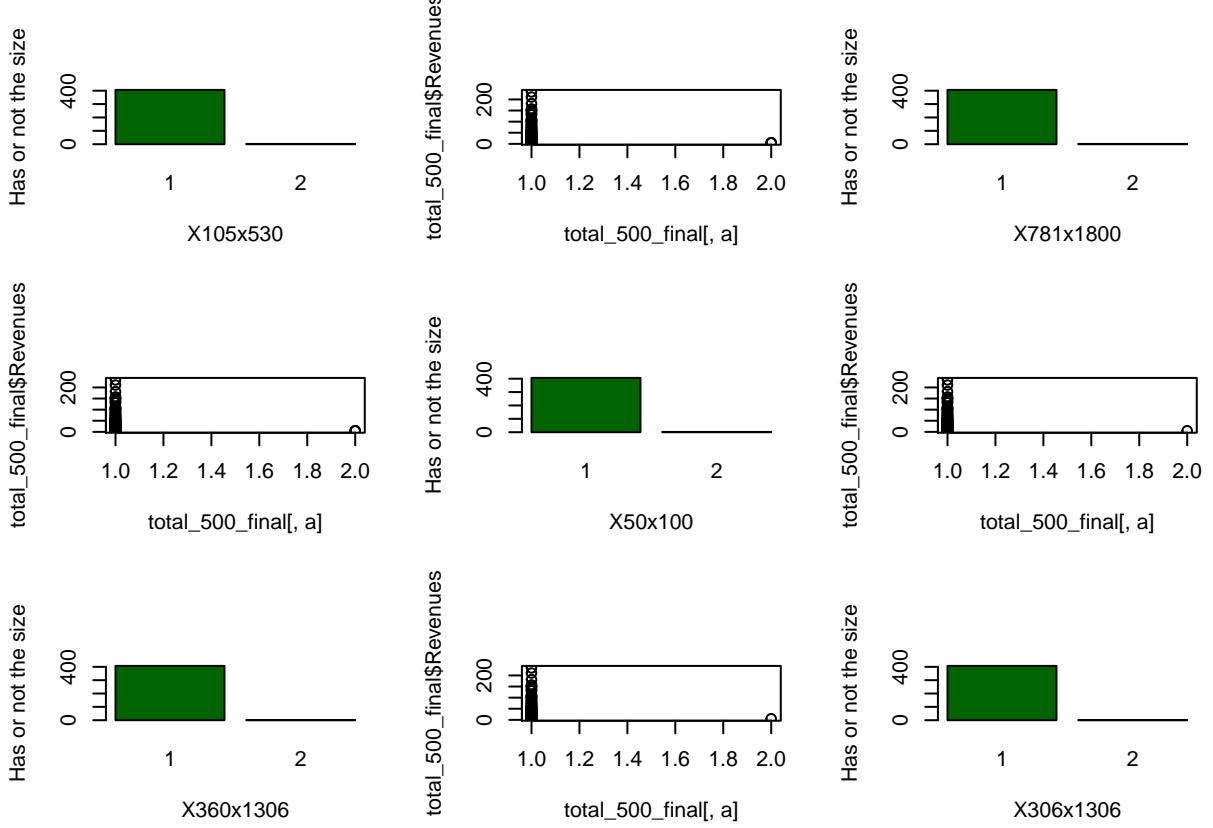


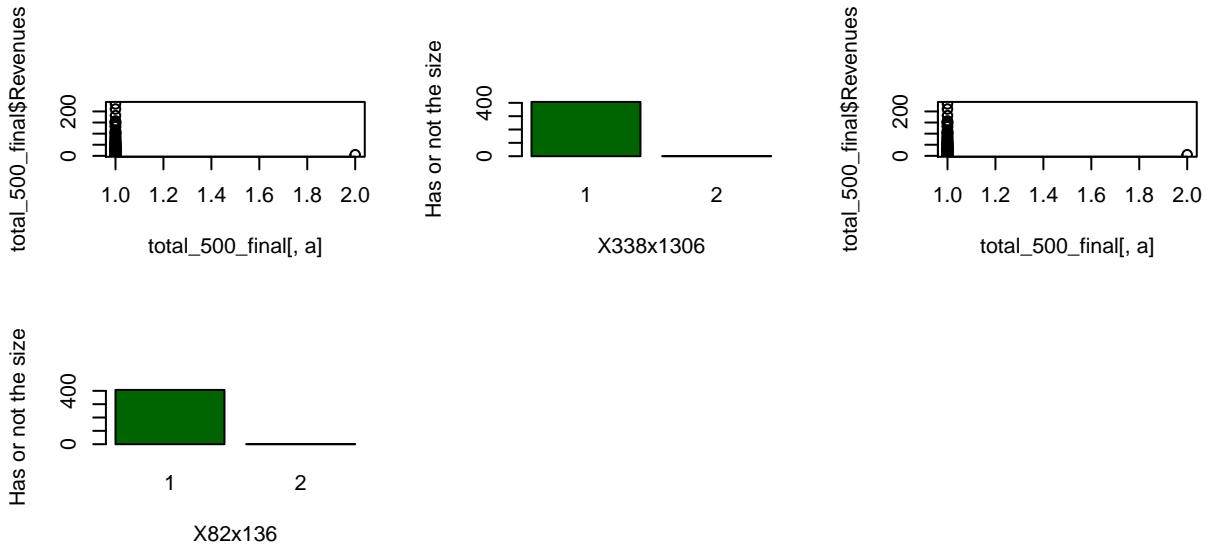










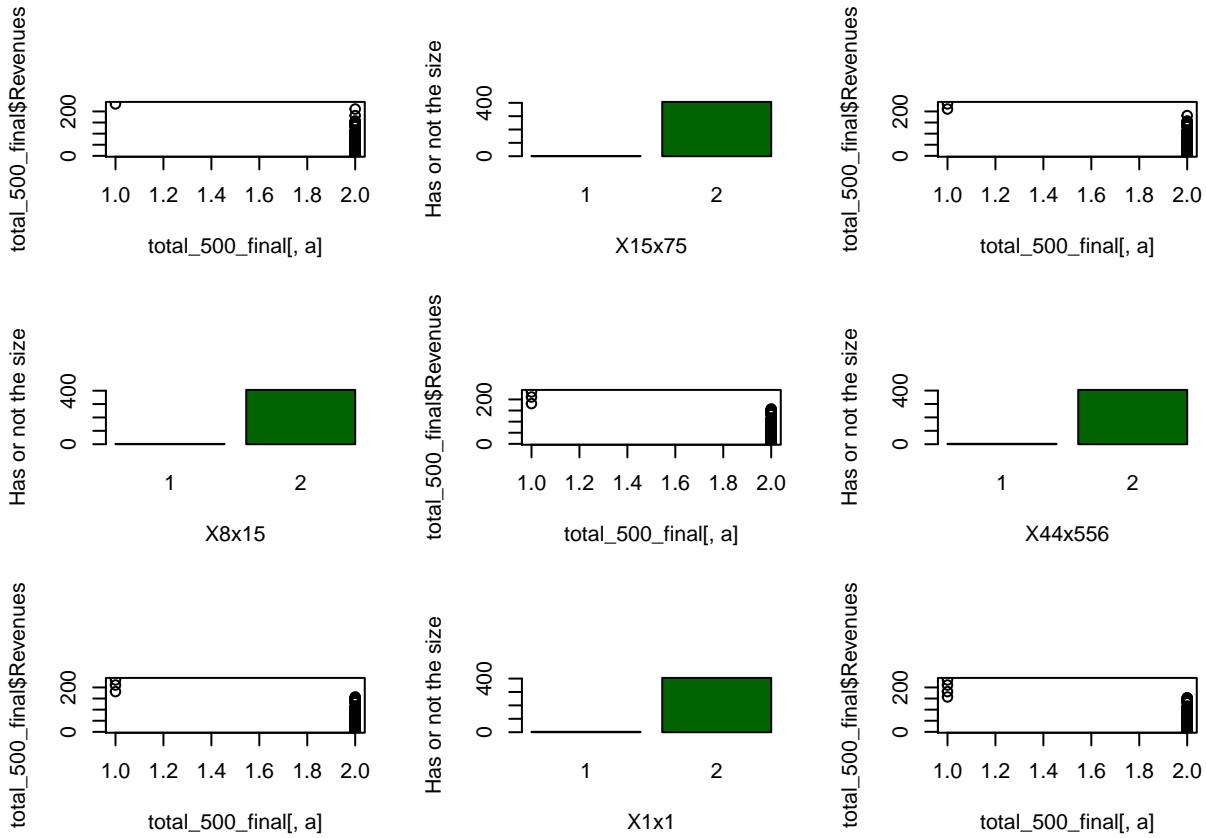


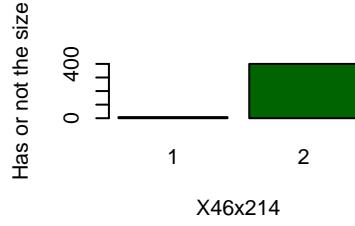
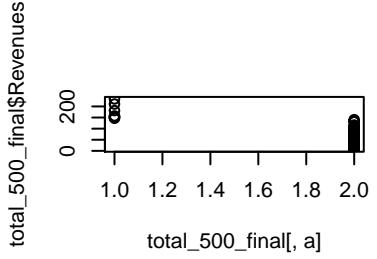
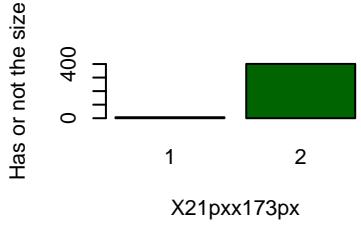
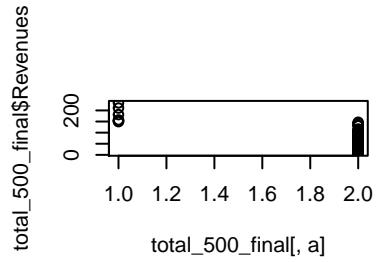
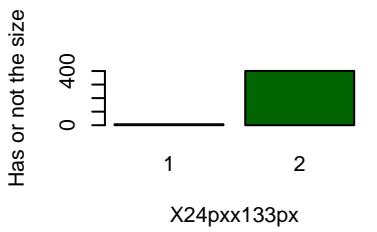
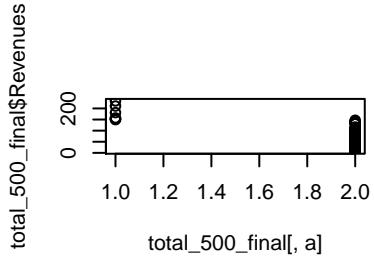
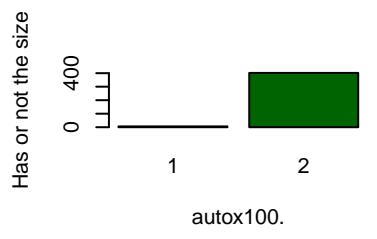
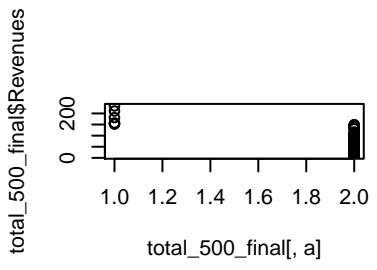
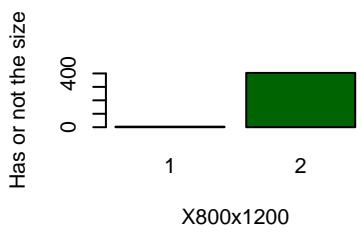
```

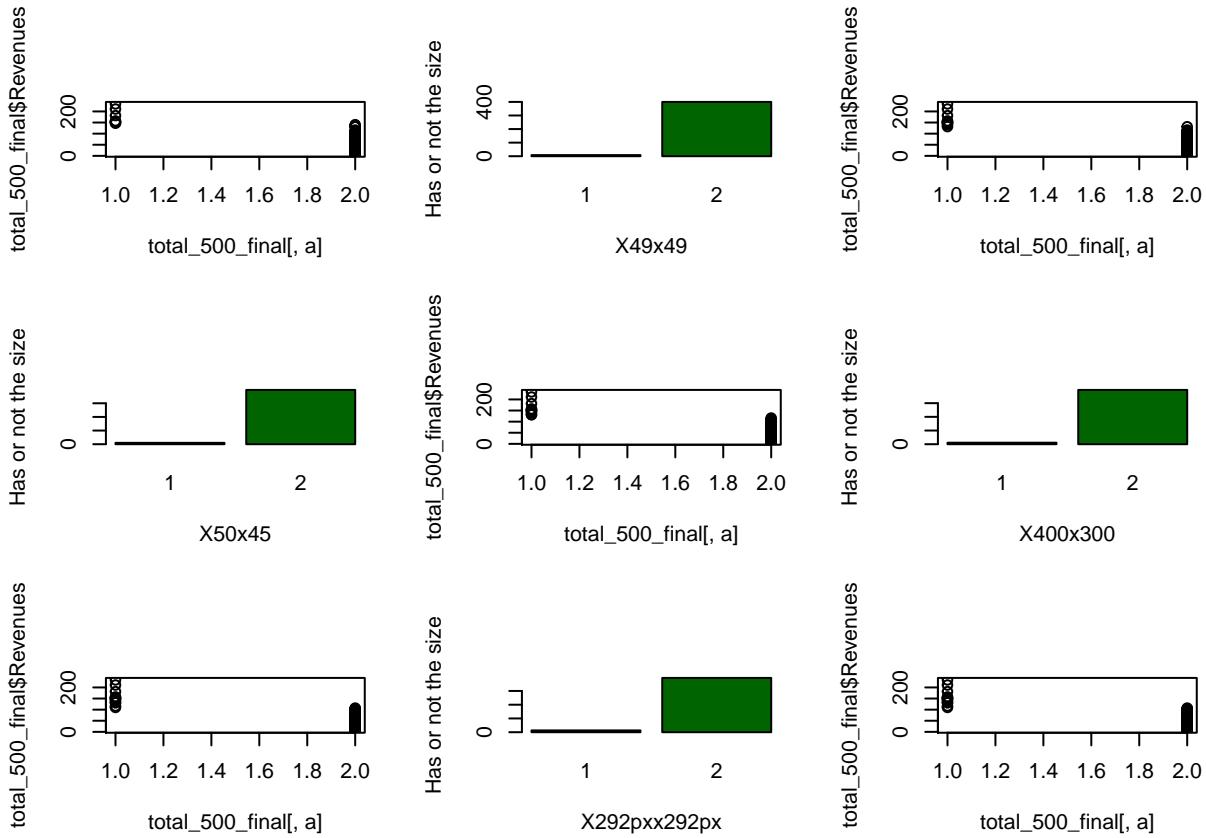
true_existing = c()
#Check for sizes that are more than half divided in existing and not
for(i in 24:715){
  image_size<- round(table(total_500_final[,i]))
  if ((image_size[[2]]>204)==TRUE){
    true_existing <- union(true_existing, c(i))
  }
}

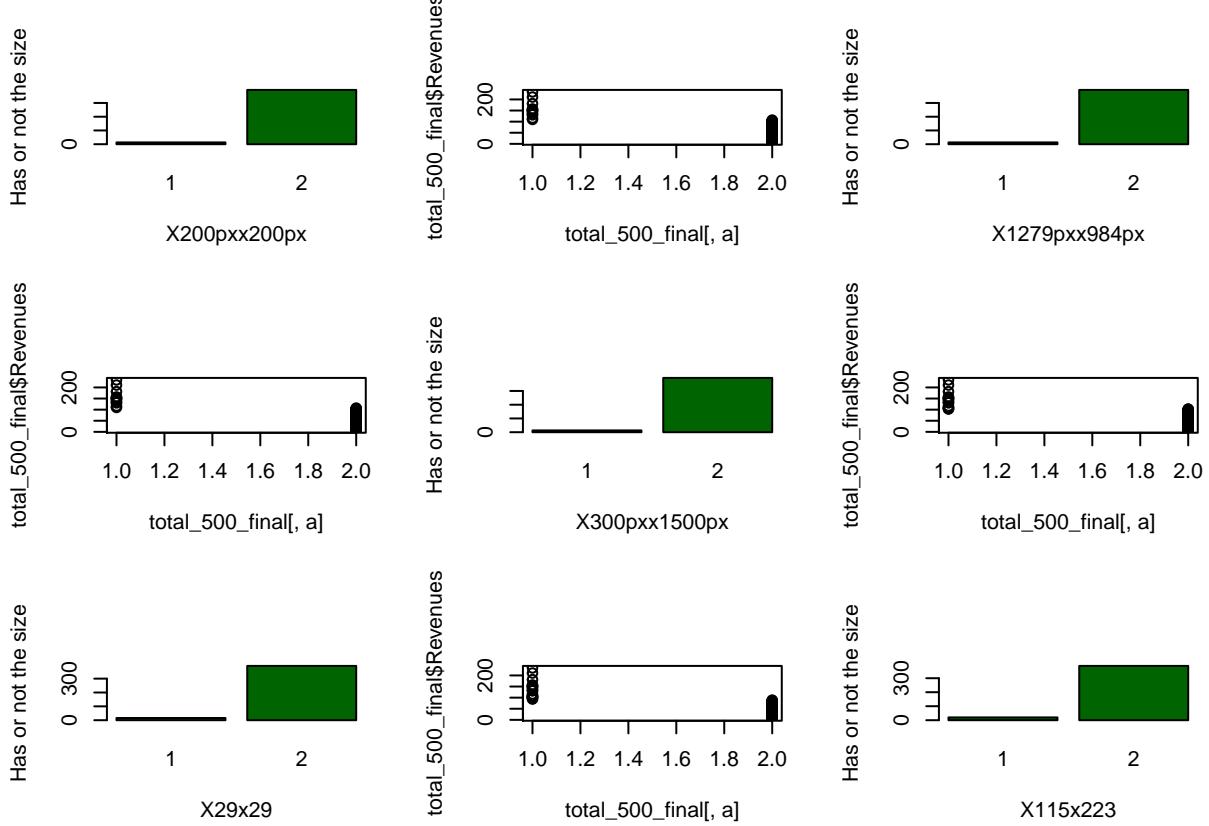
#Now we will take the sizes that exist in more than half the instances and check graphically the deviation
par(mfrow=c(3,3))
for(i in 1:276){
  a = true_existing[i]
  image_size<- round(table(total_500_final[,a]))
  plot(total_500_final[,a],total_500_final$Revenues)
  barplot(image_size,xlab=names(total_500_final)[a],ylab = "Has or not the size", col = "dark green")
}

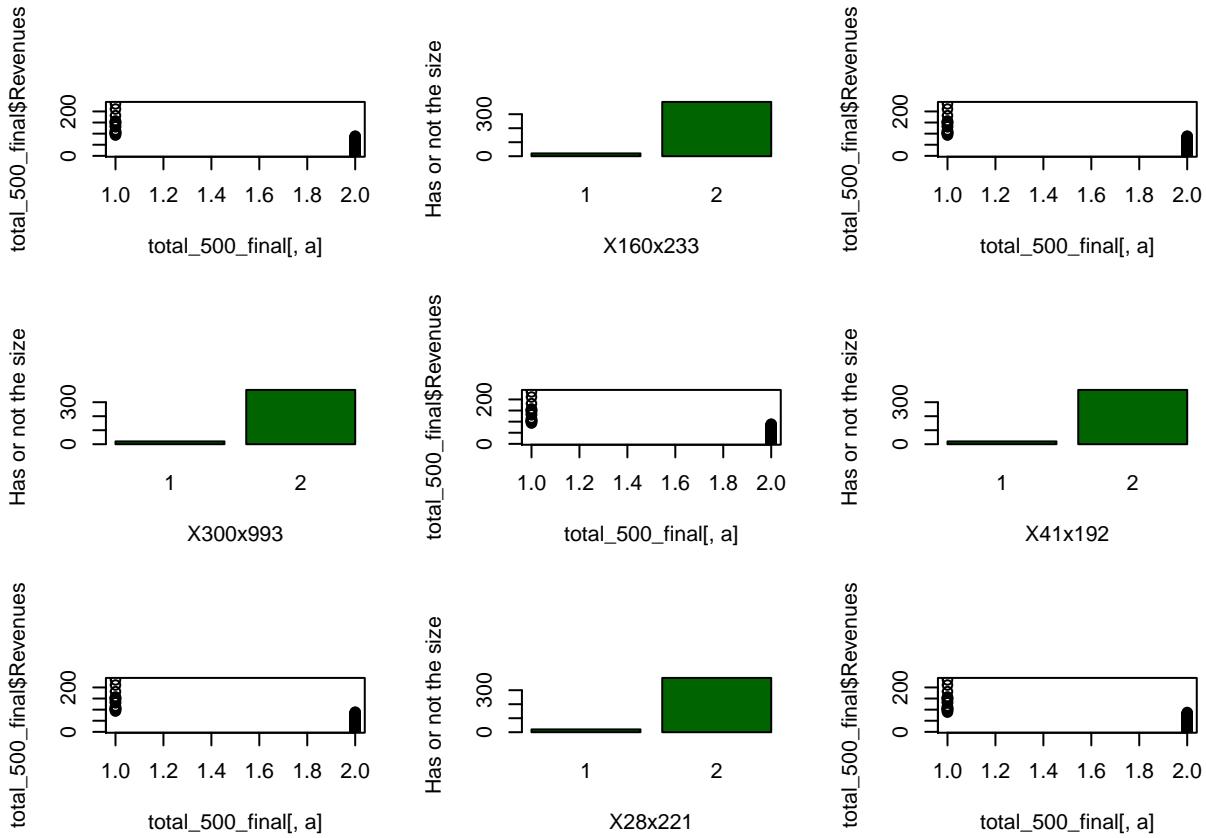
```

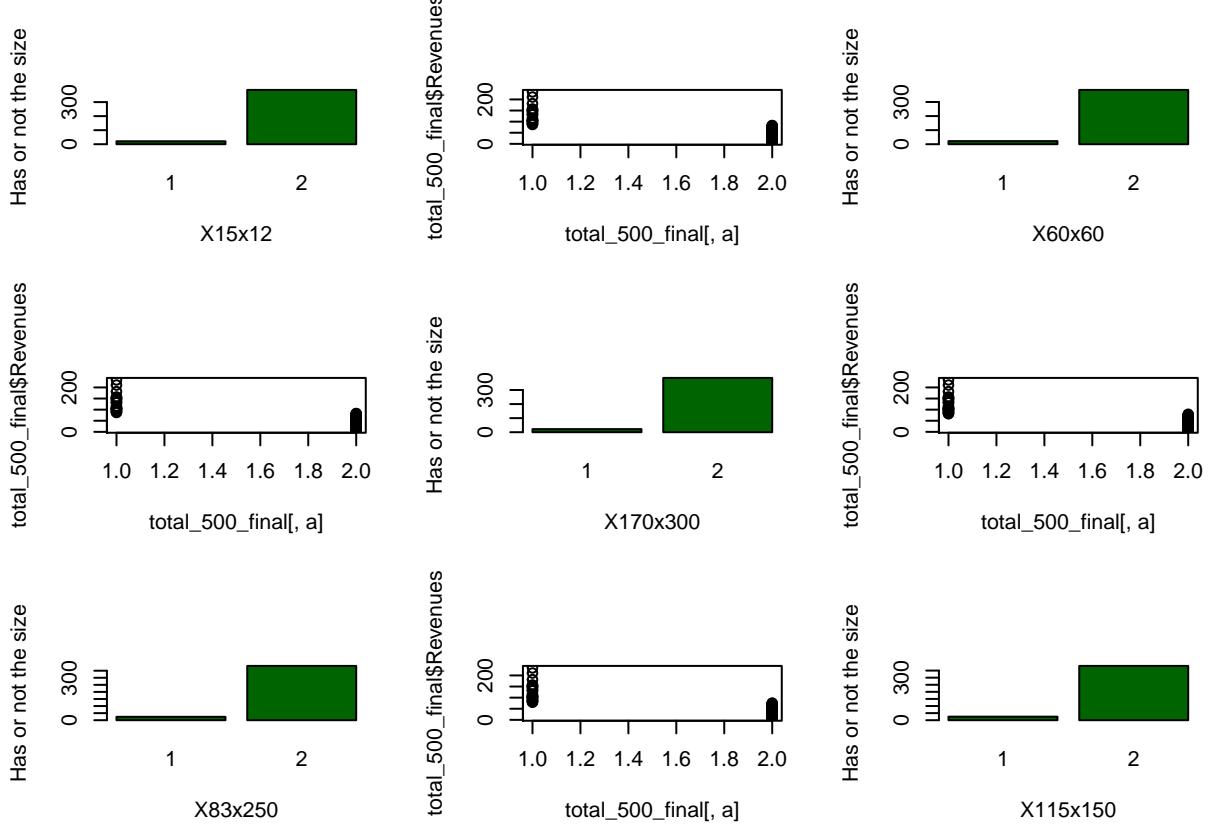


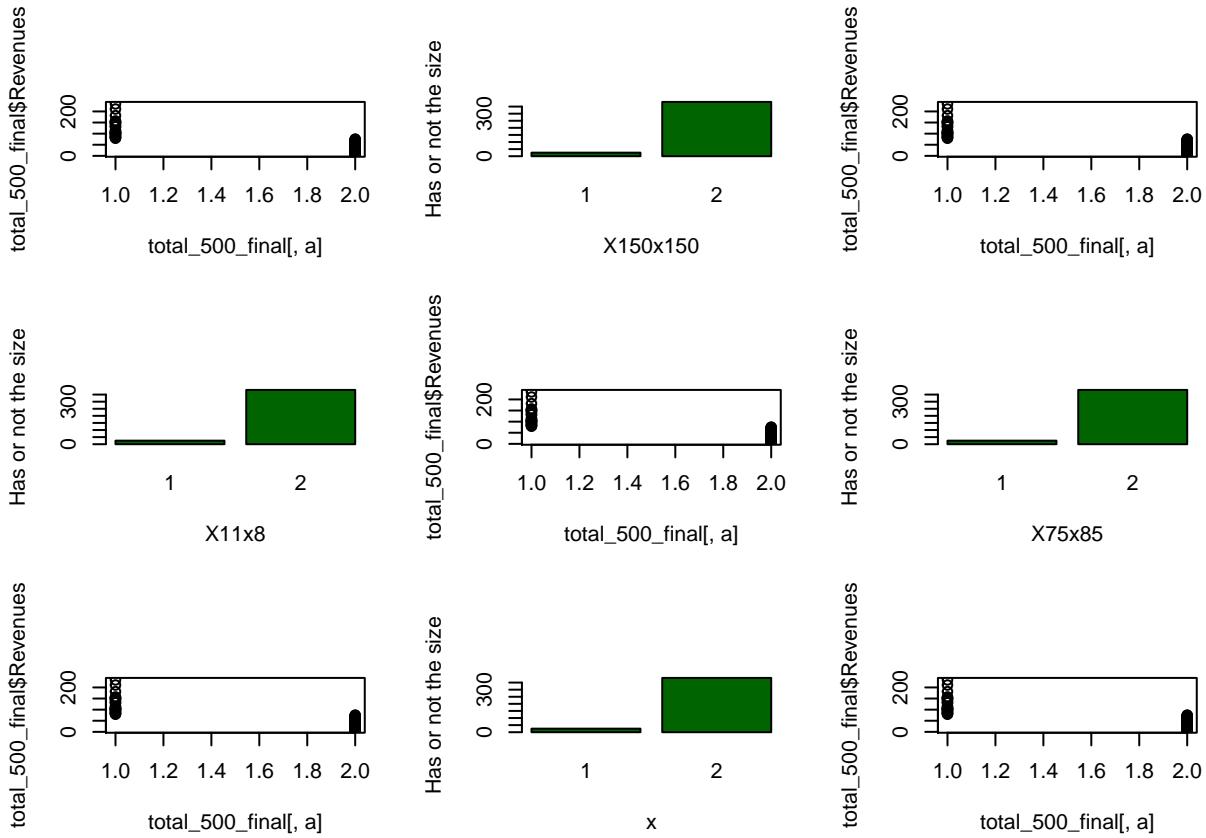


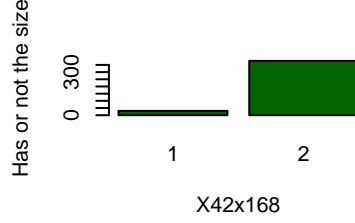
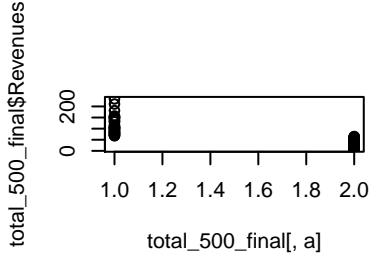
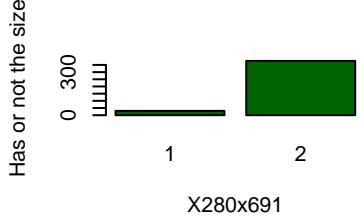
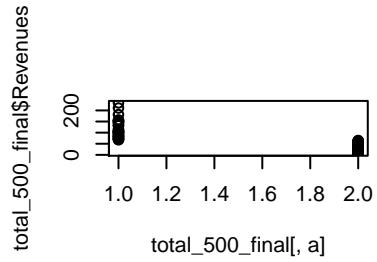
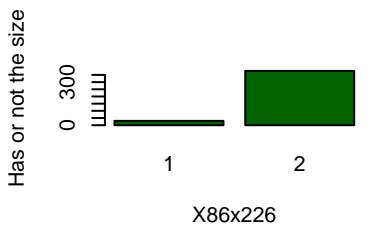
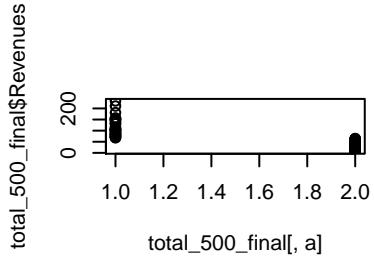
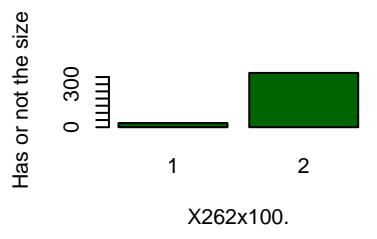
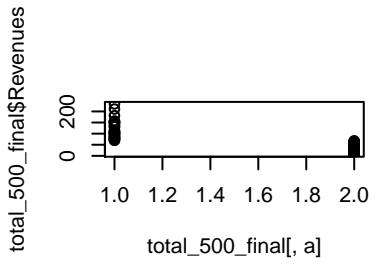
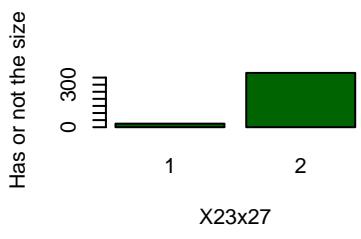


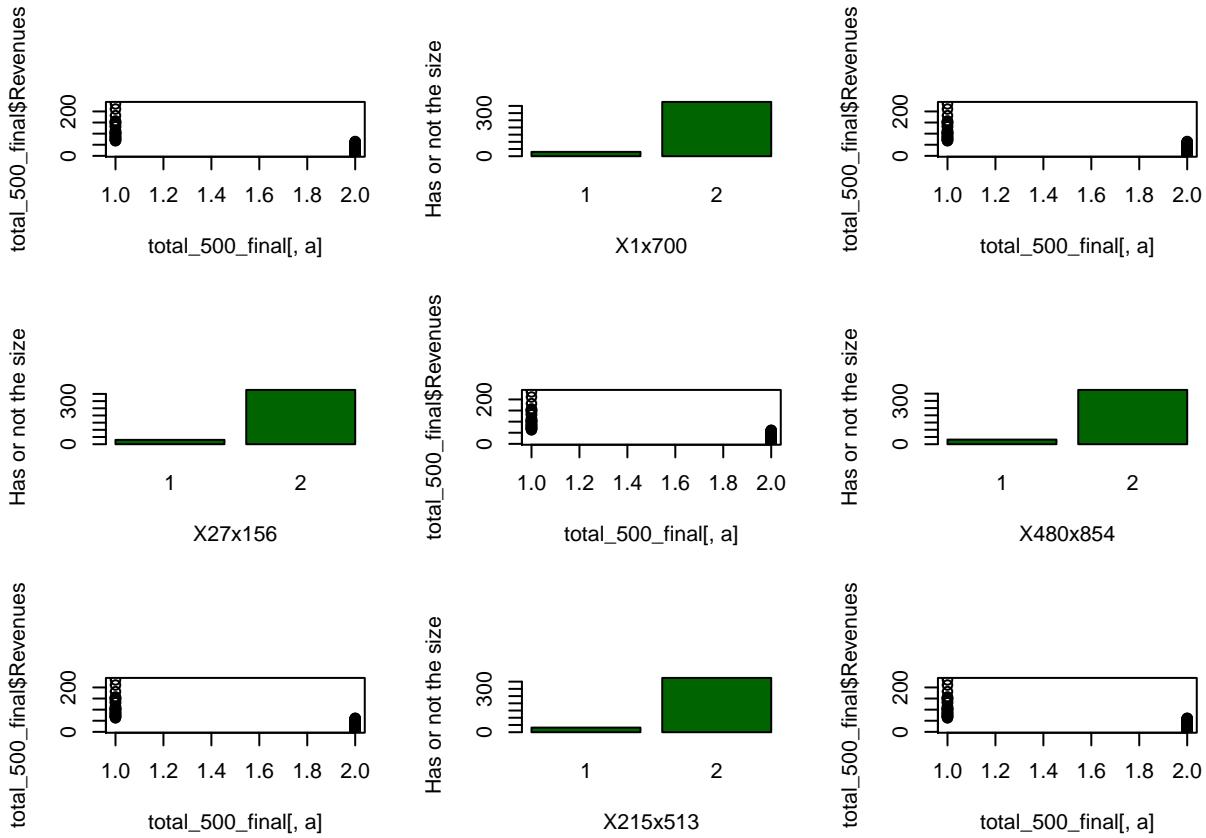


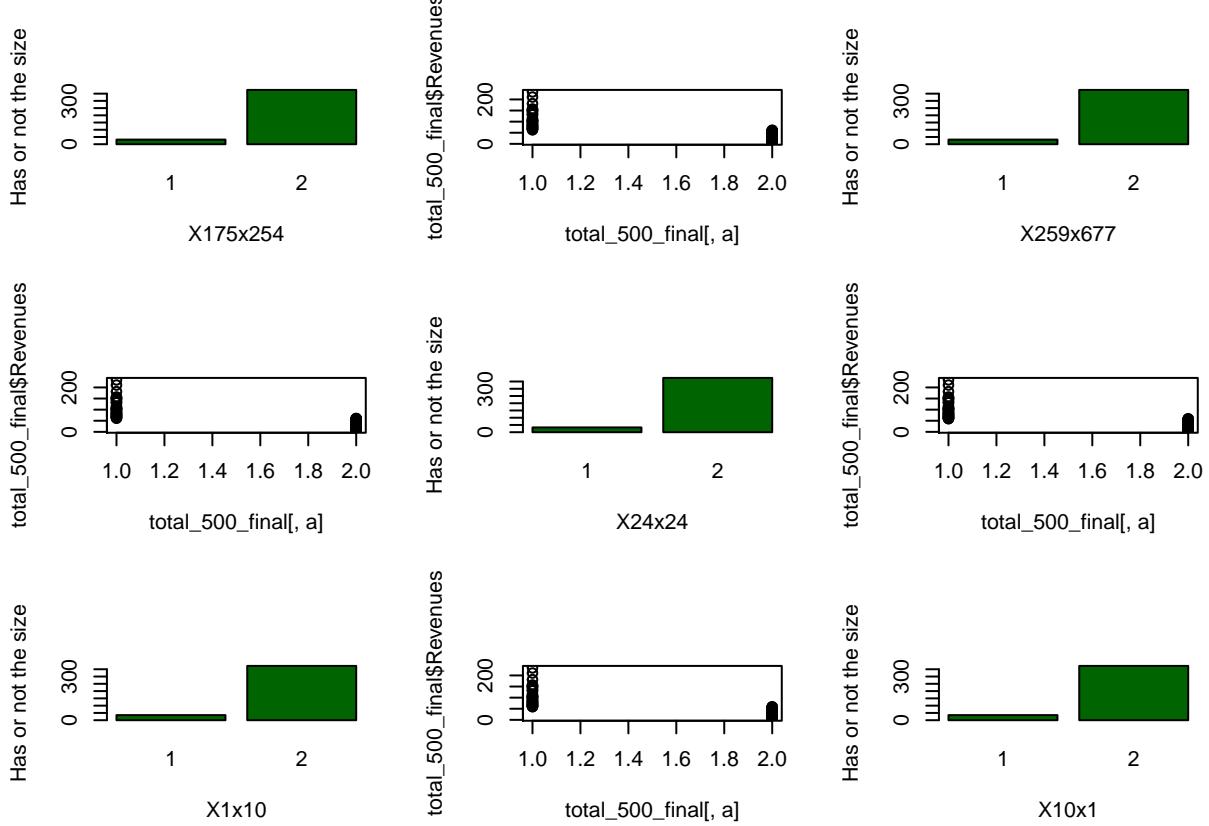


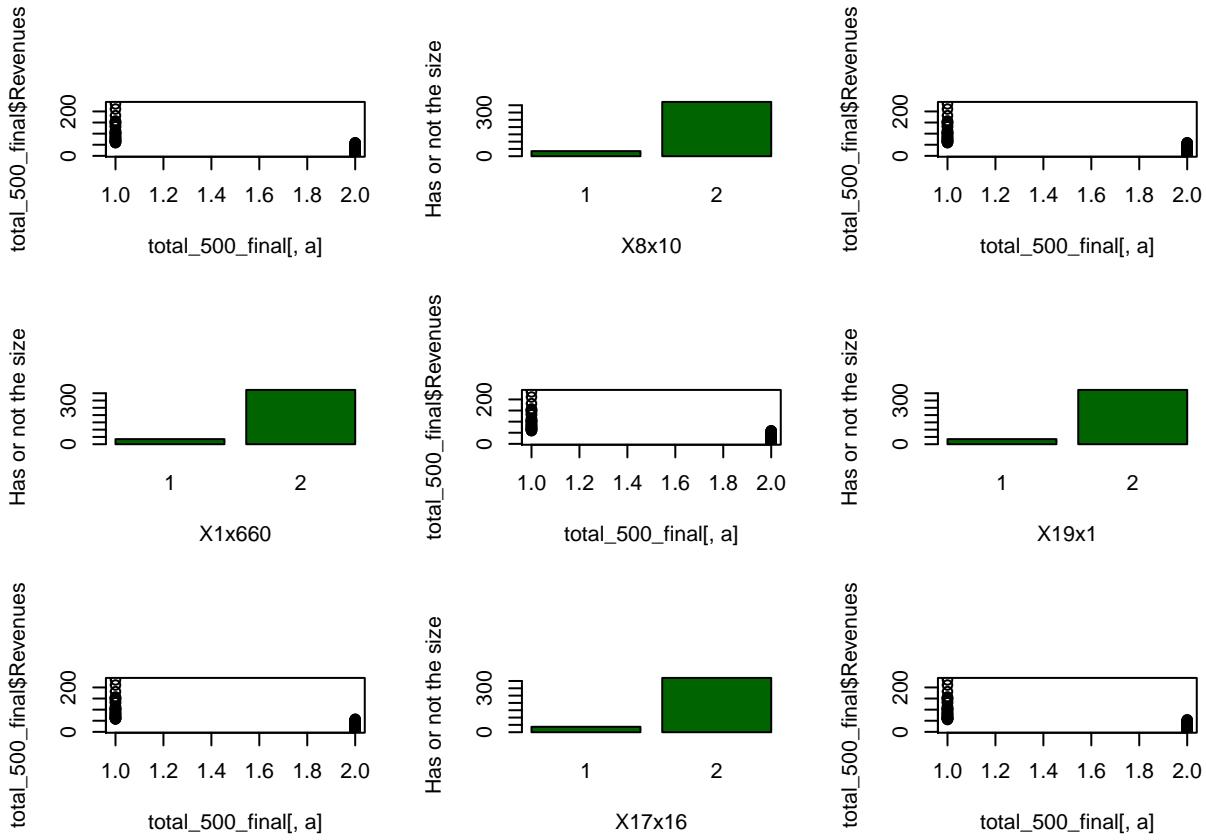


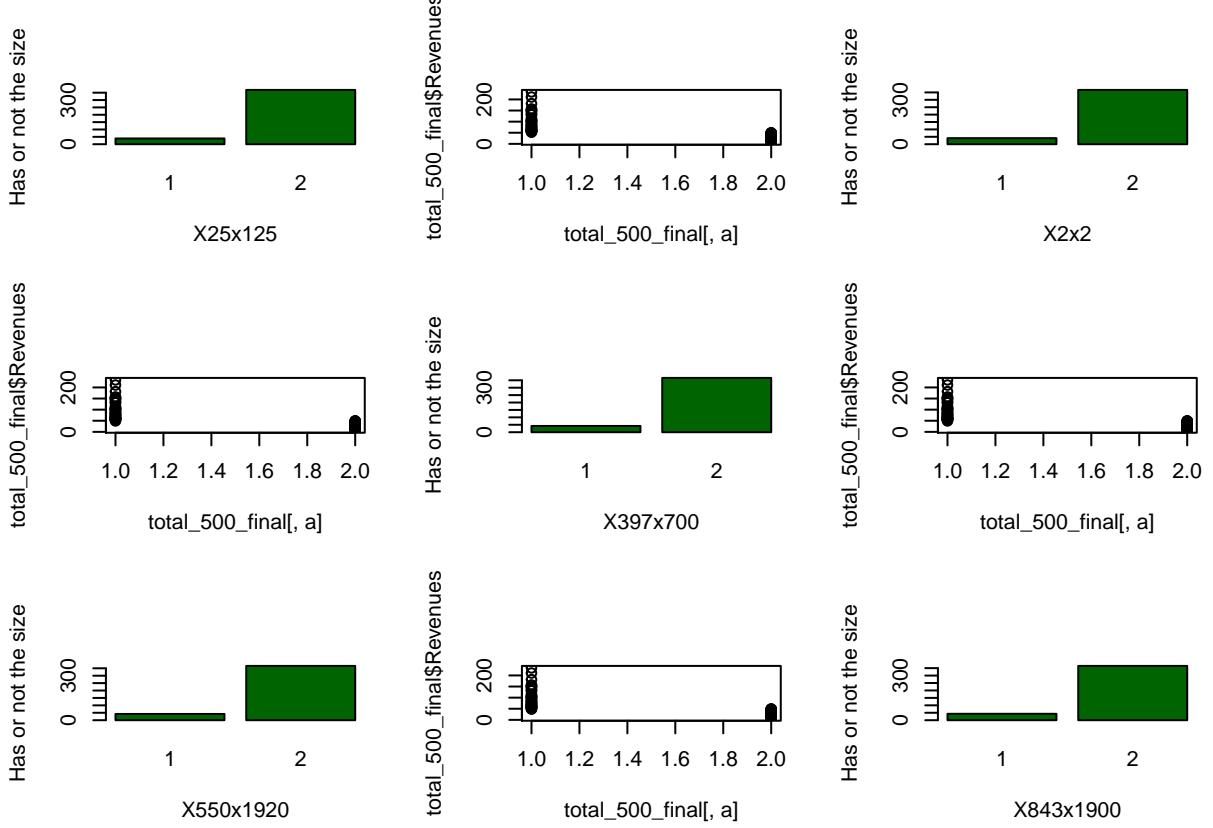


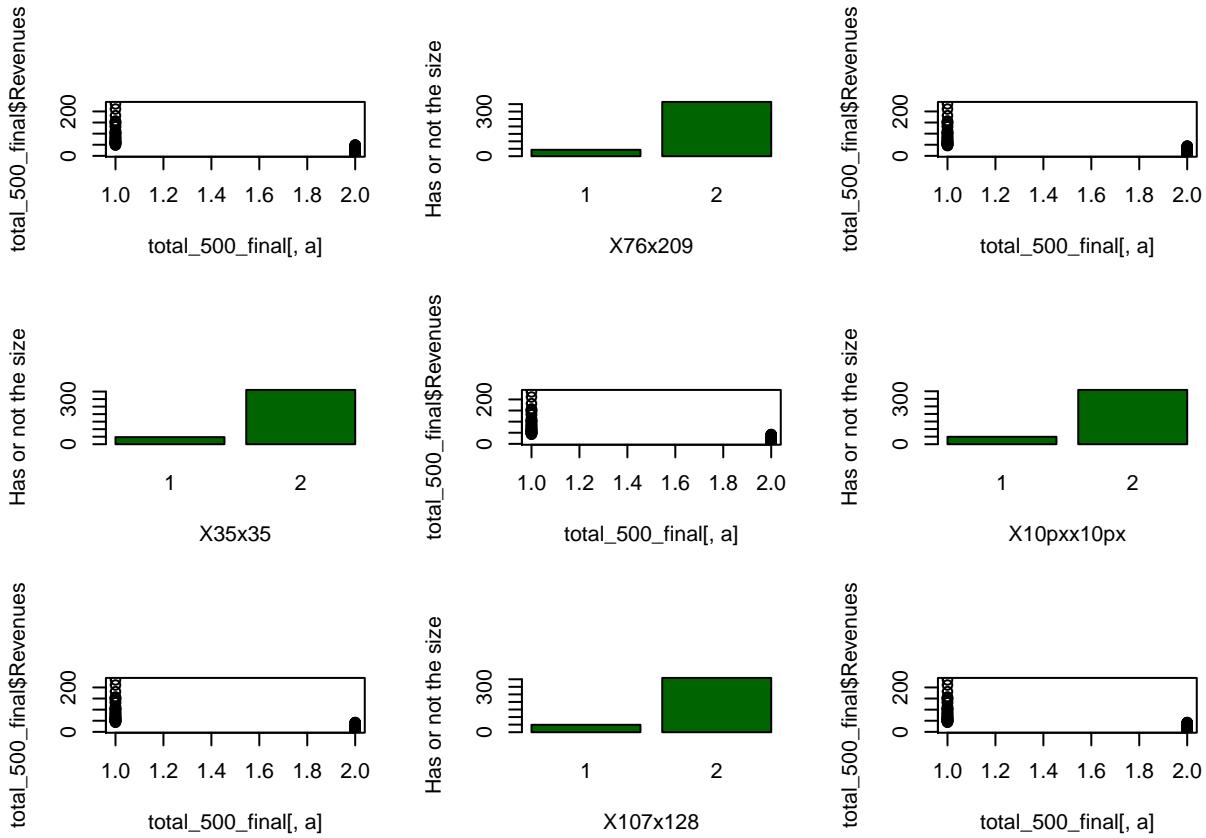


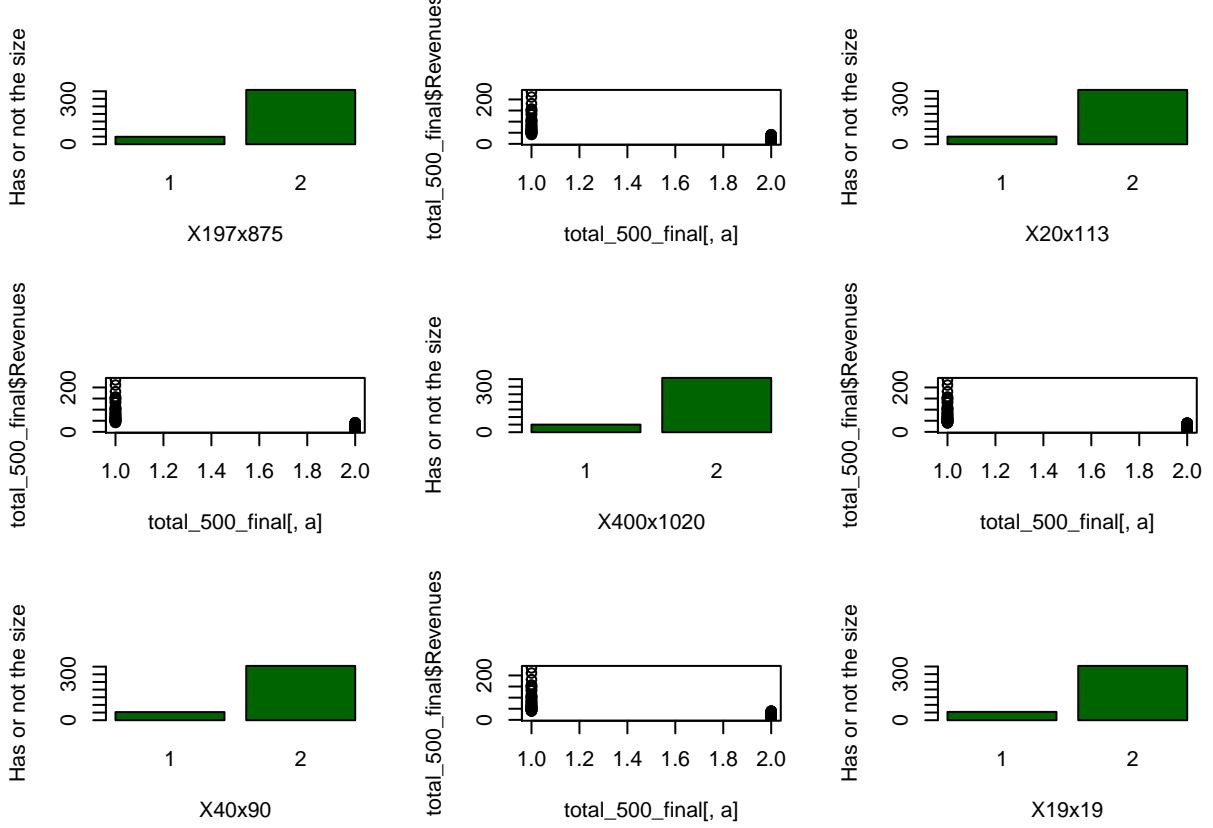


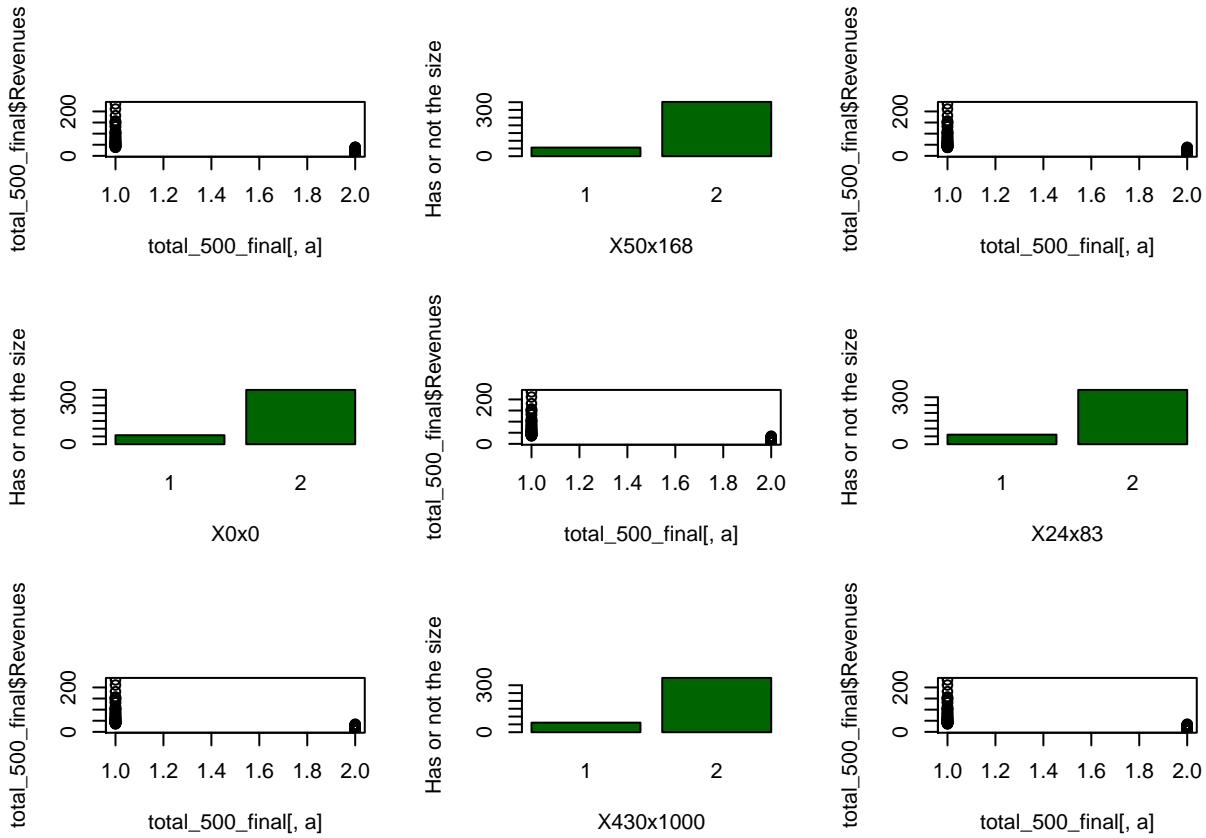


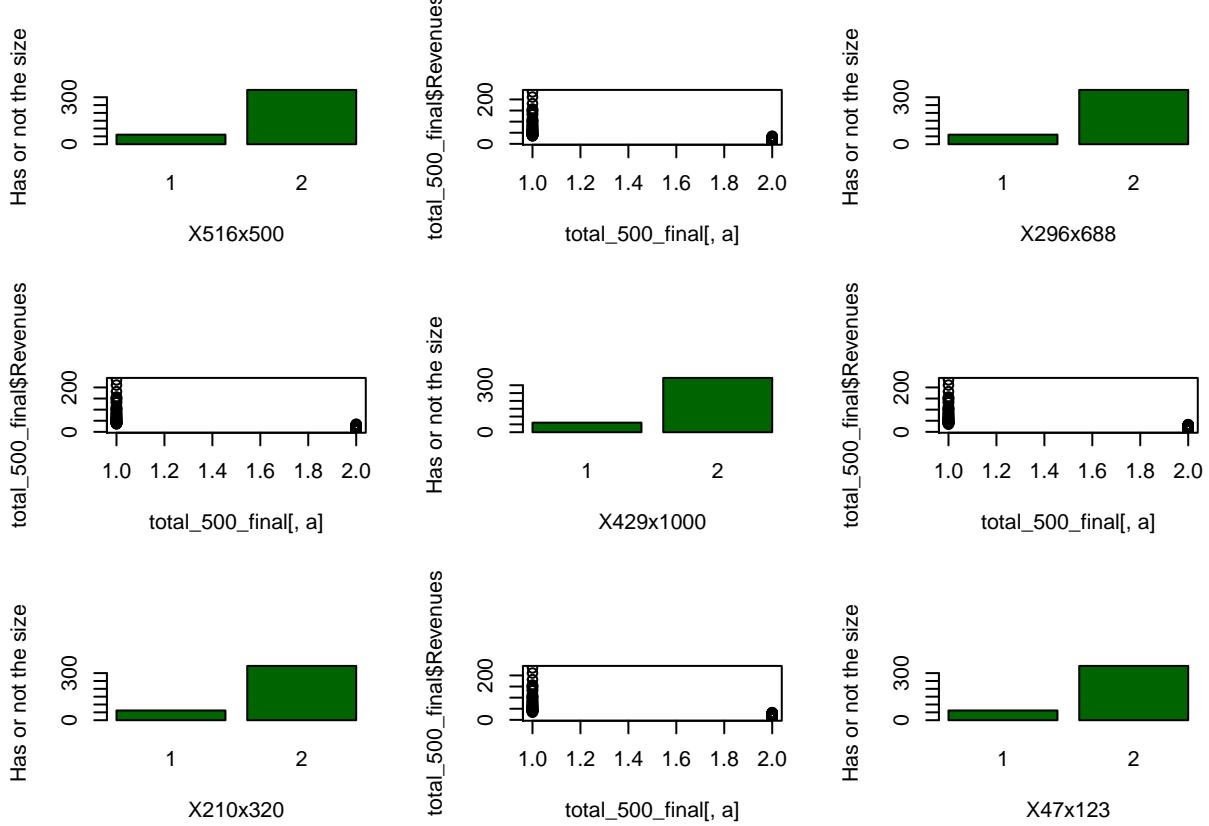


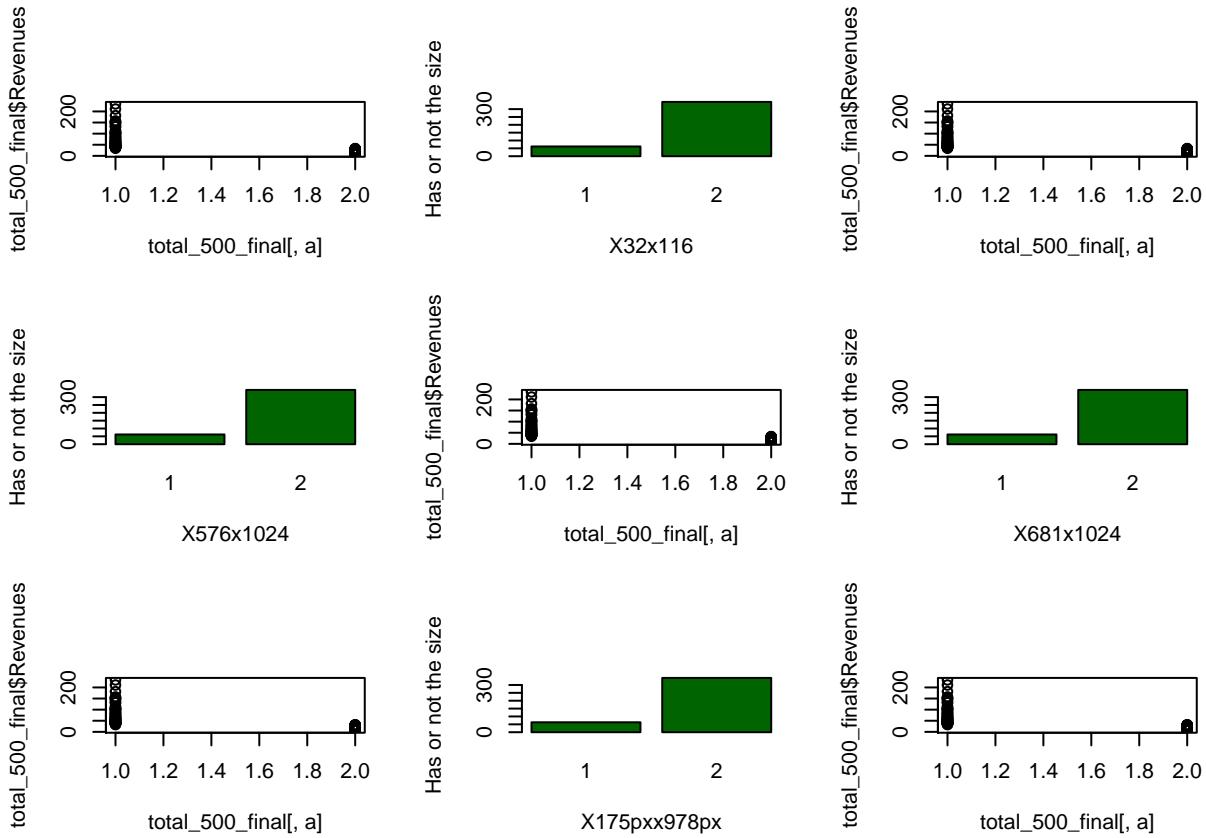


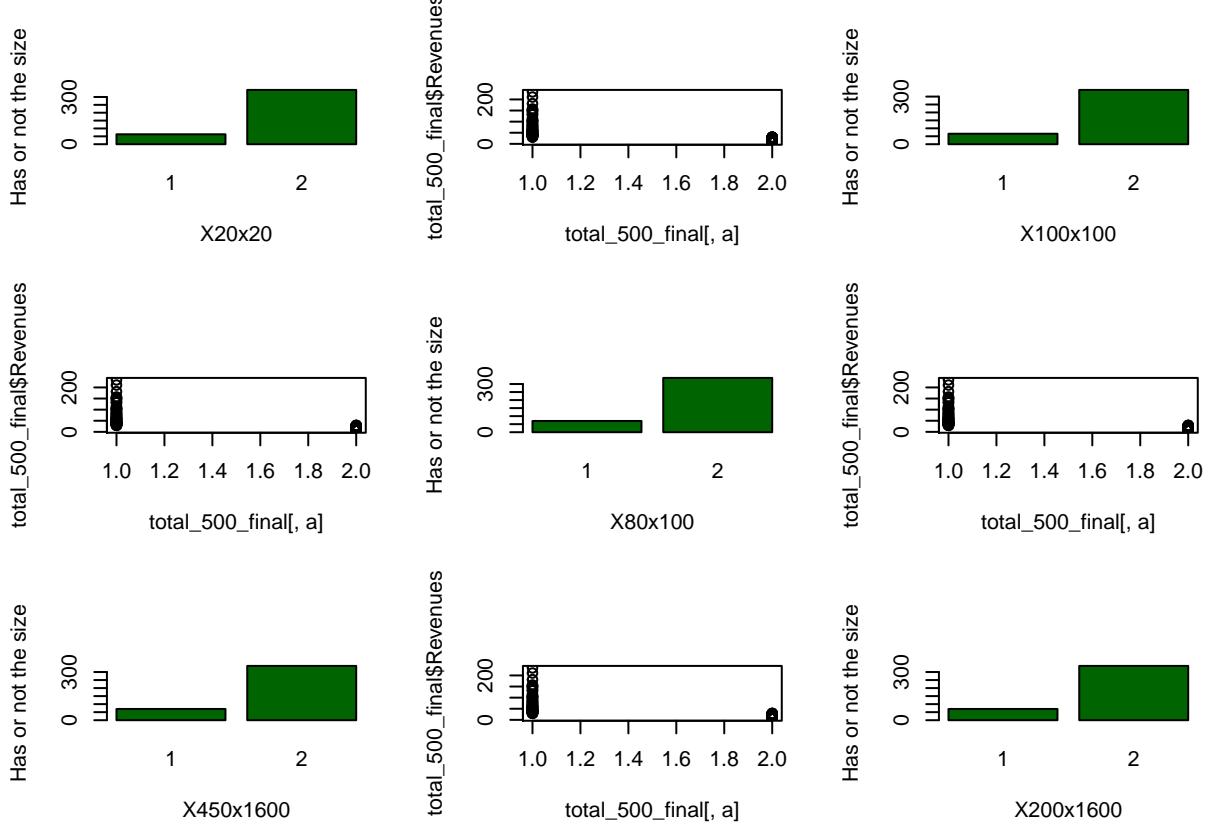


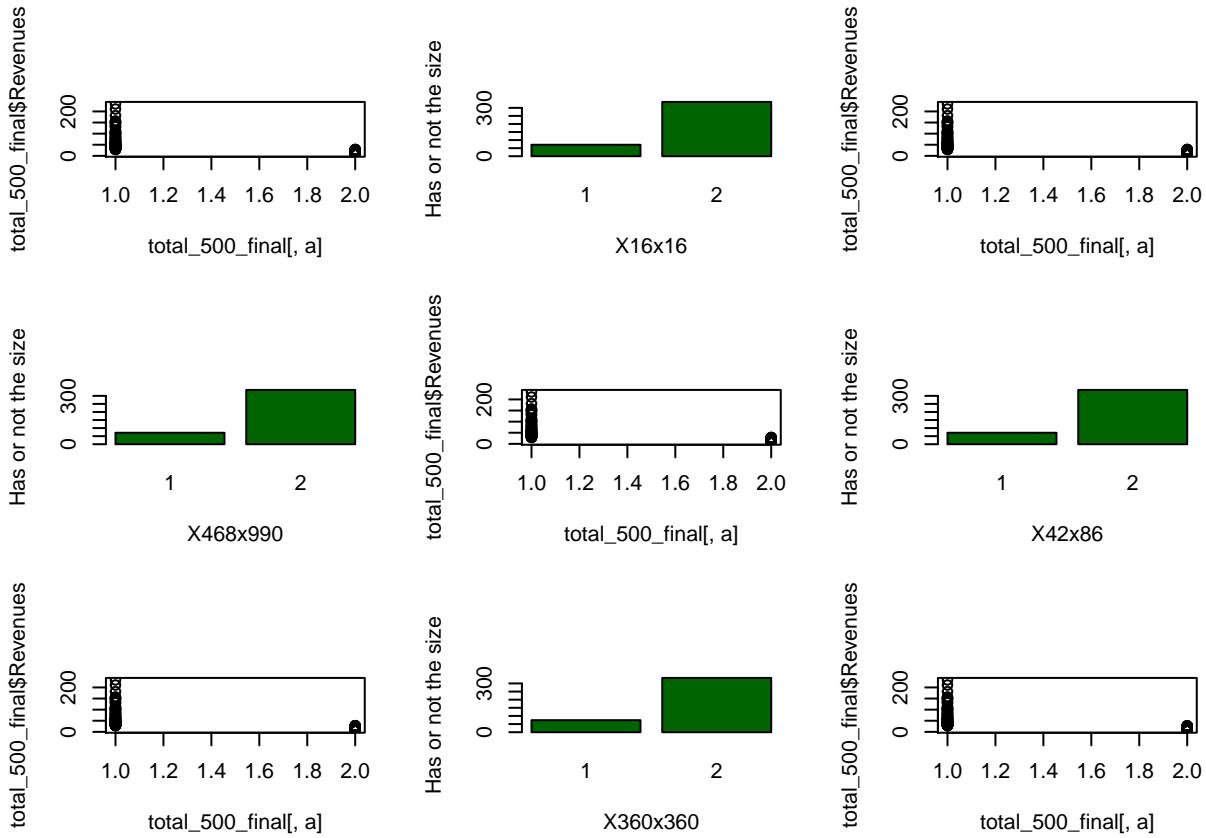


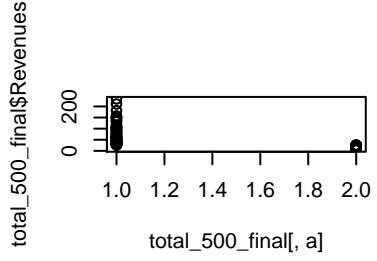
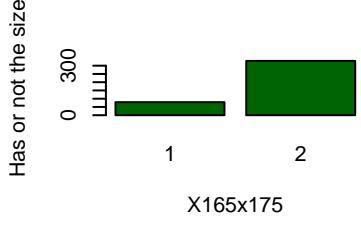
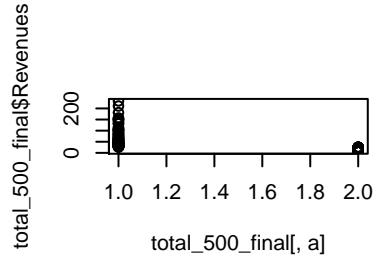
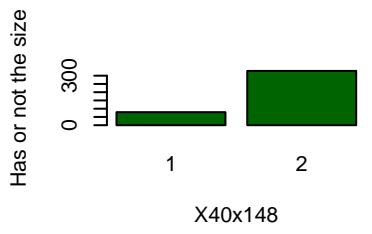
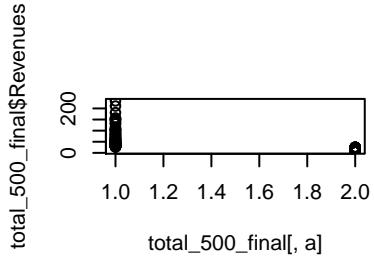
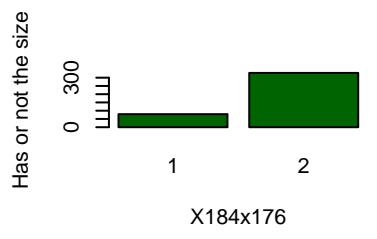
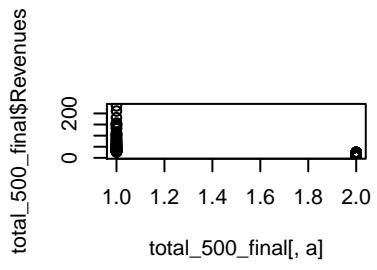
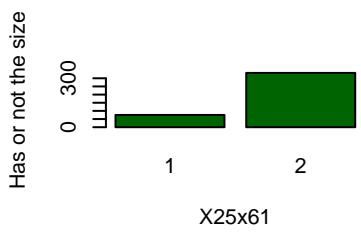


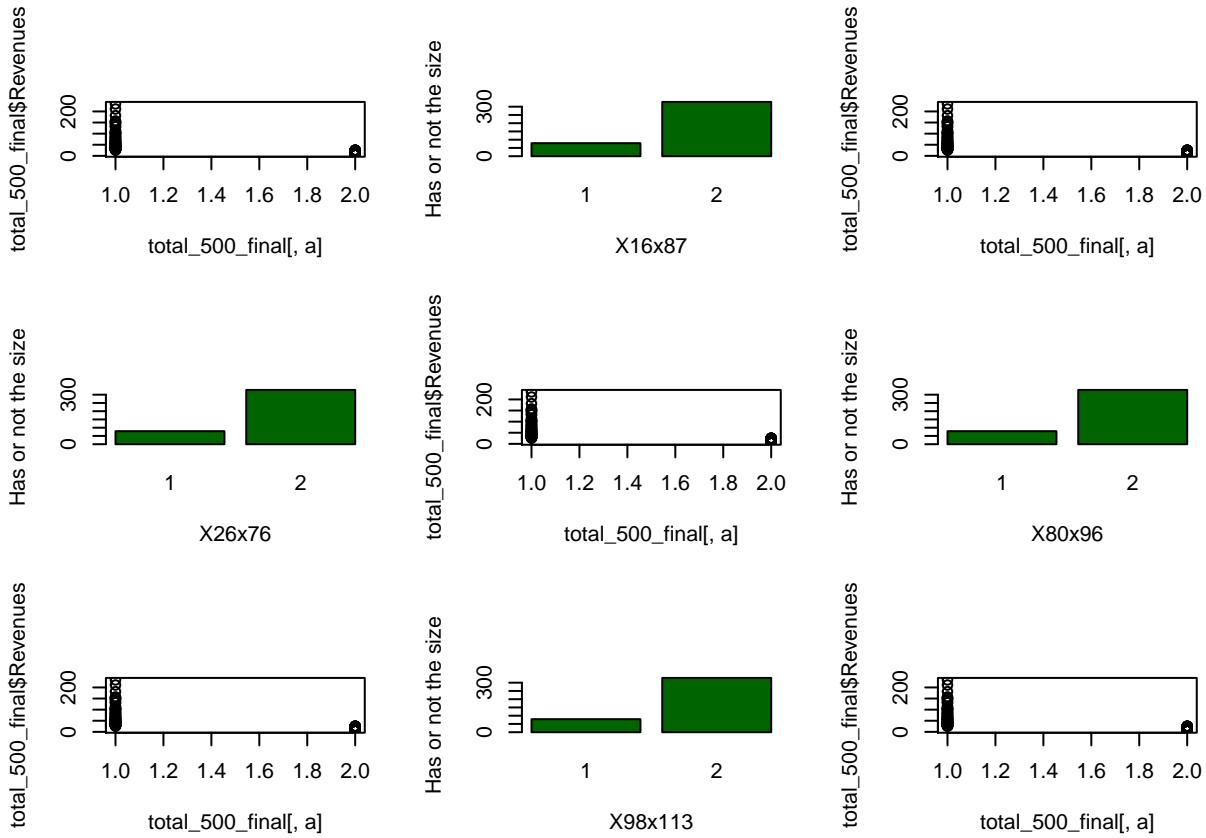


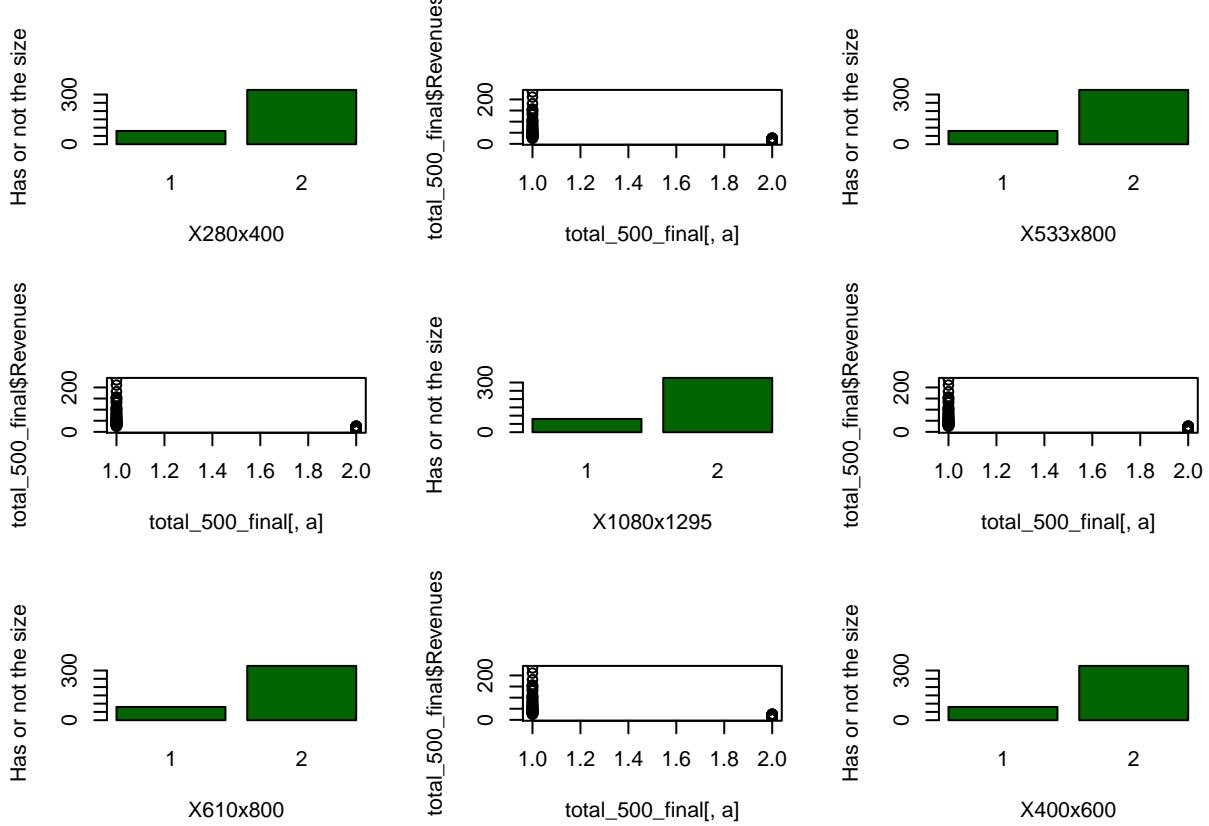


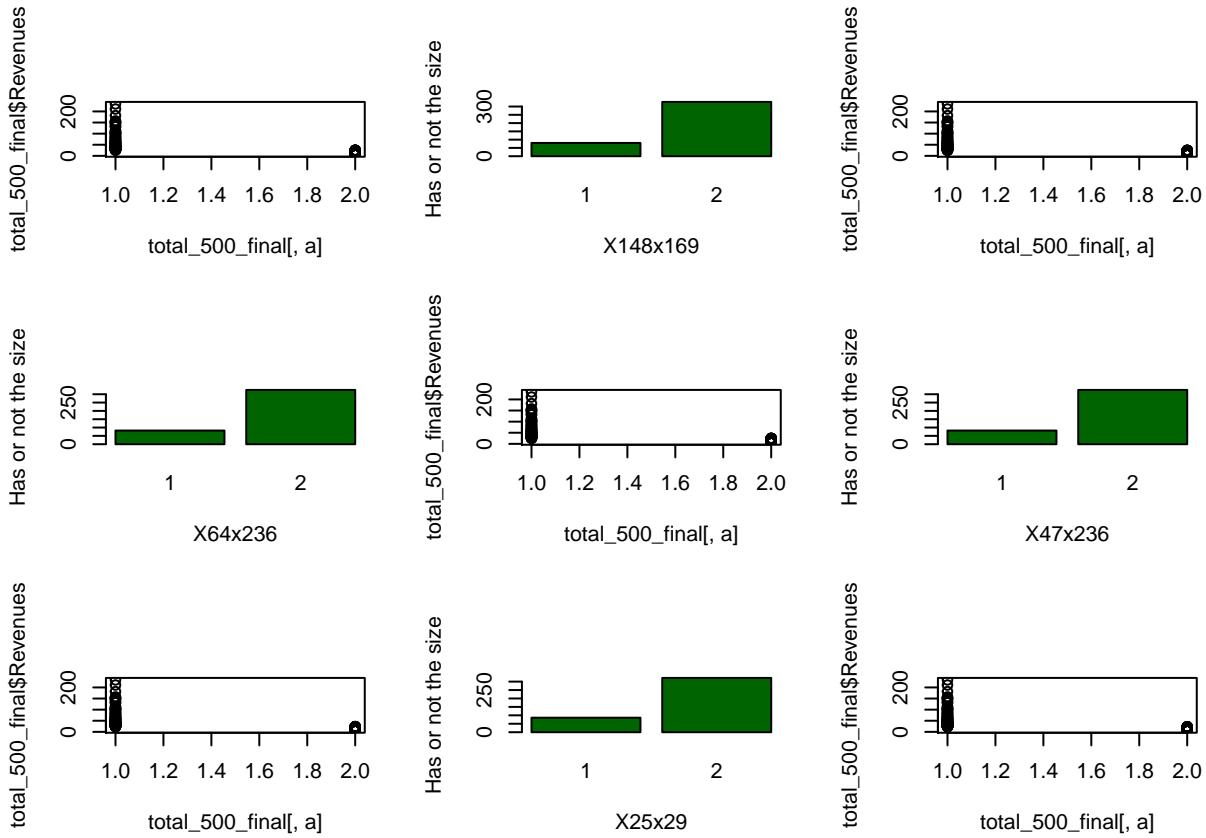


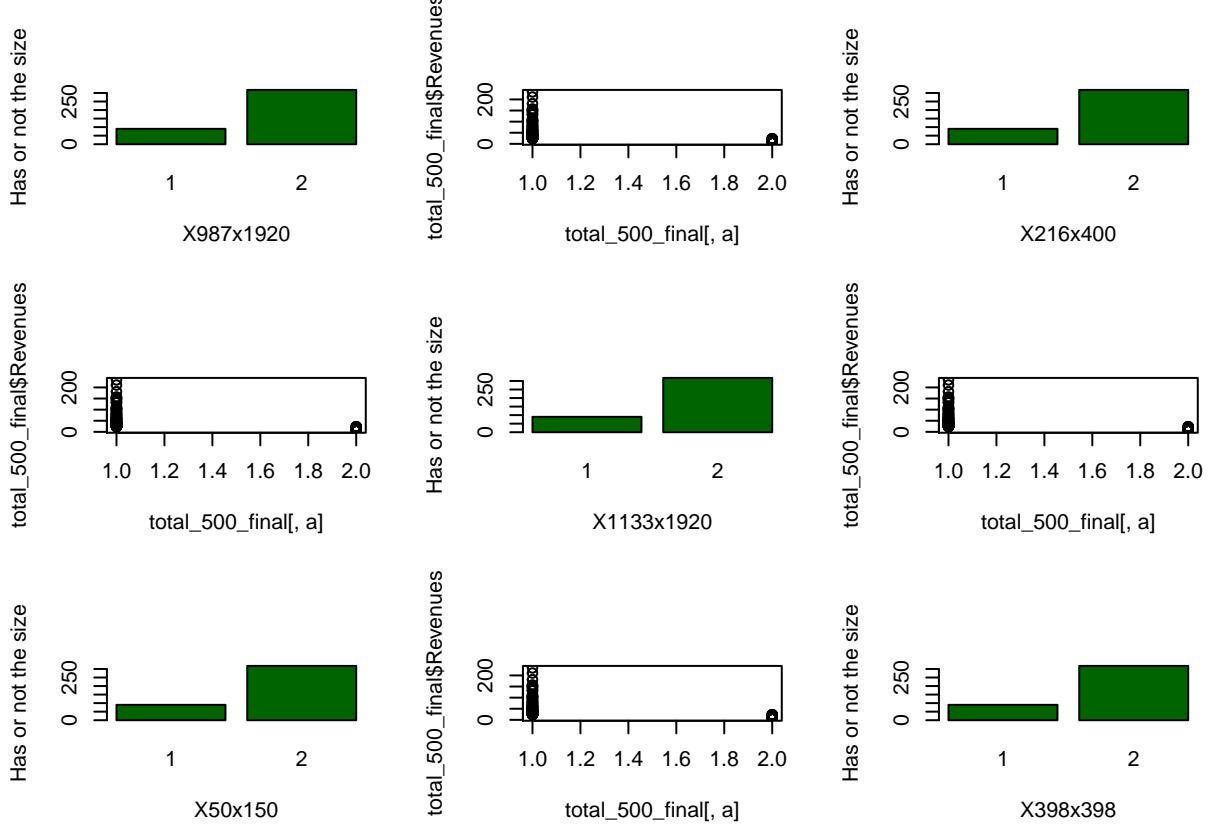


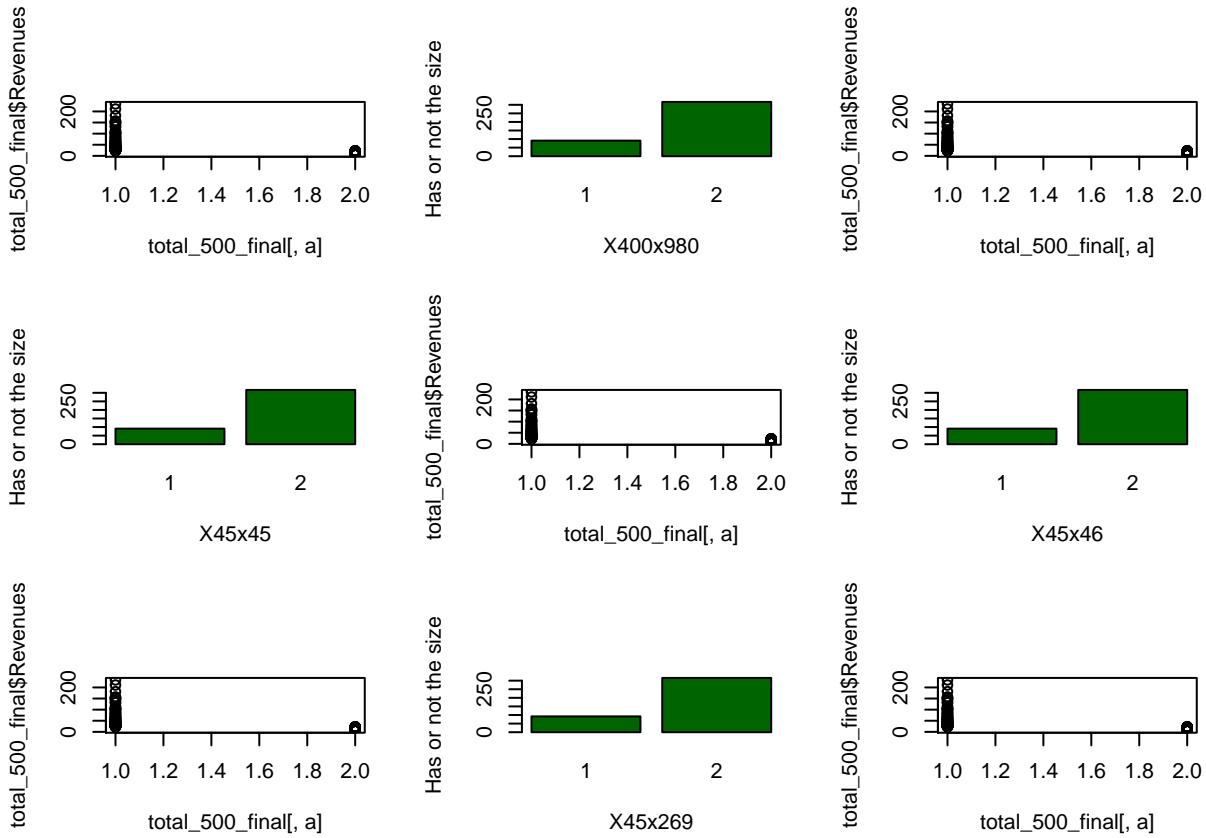


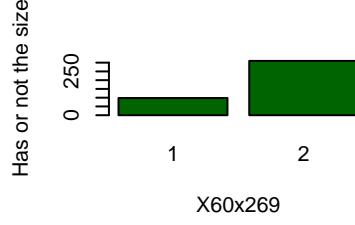
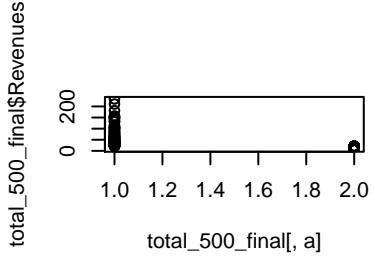
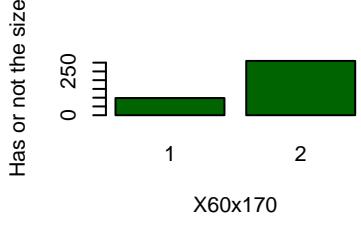
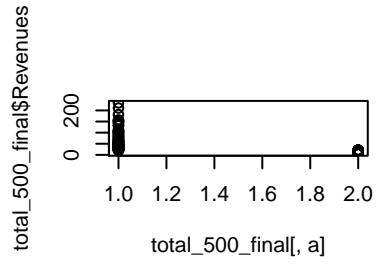
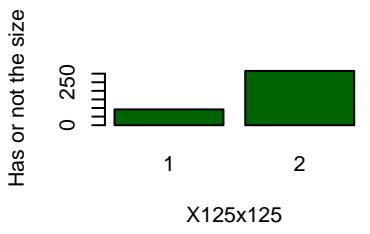
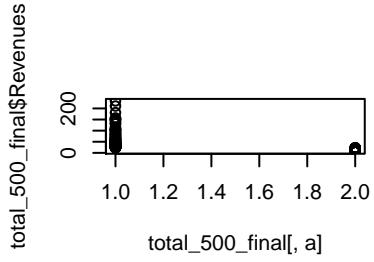
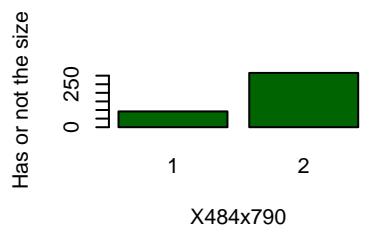
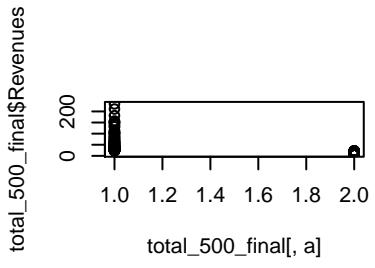
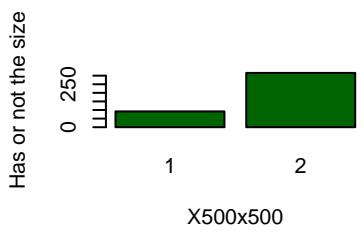


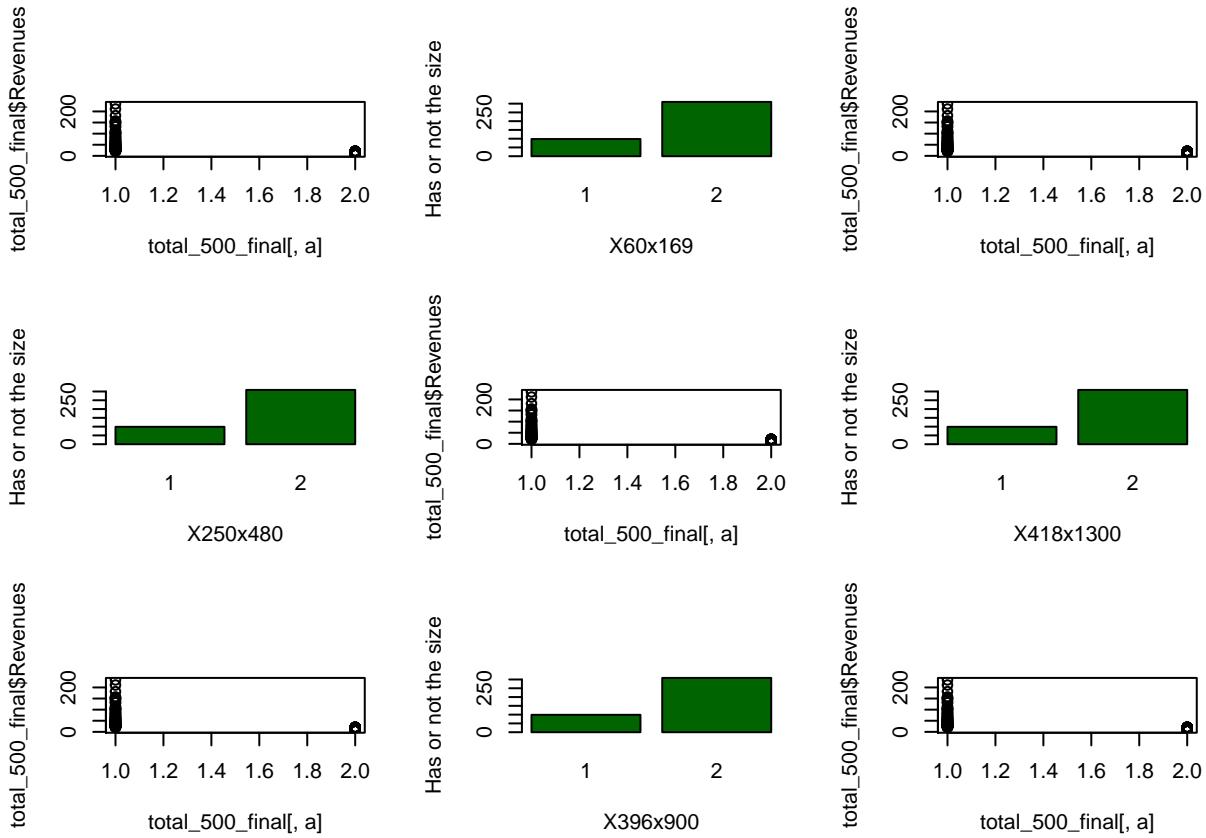


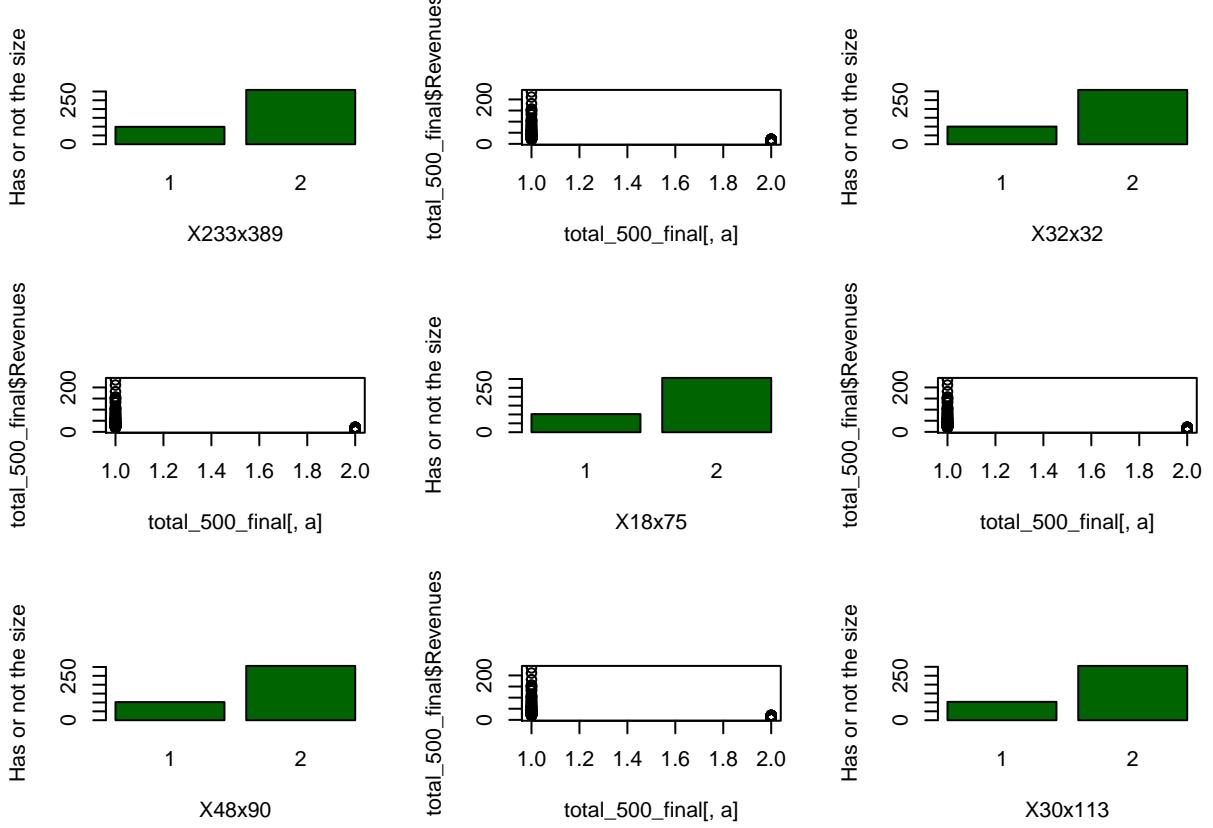


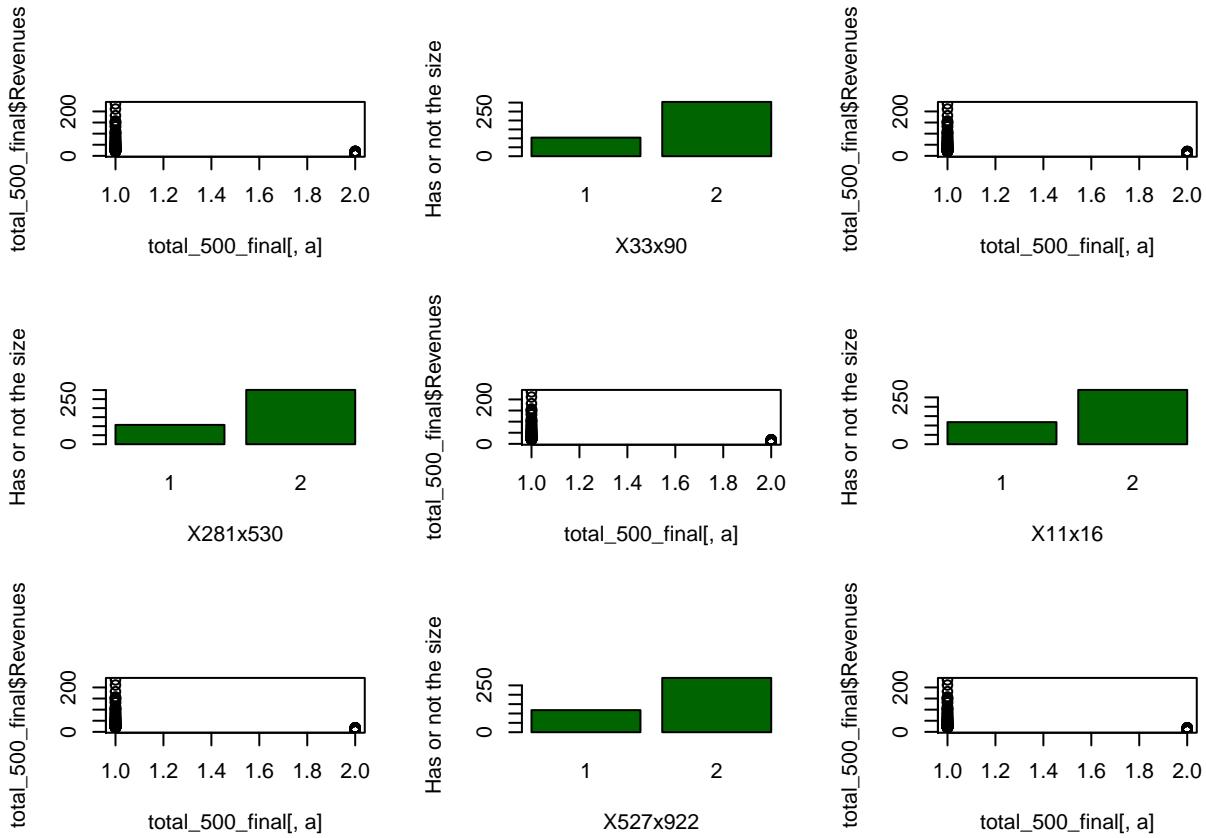


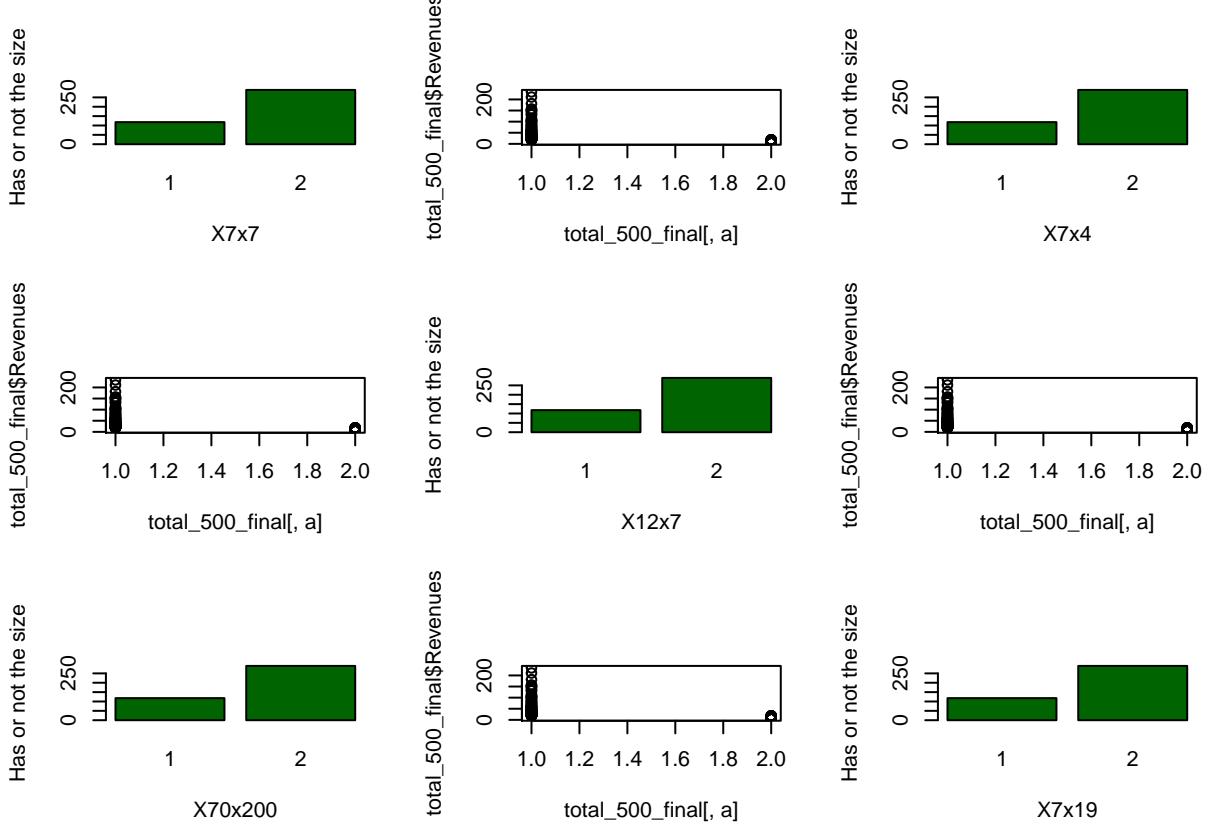


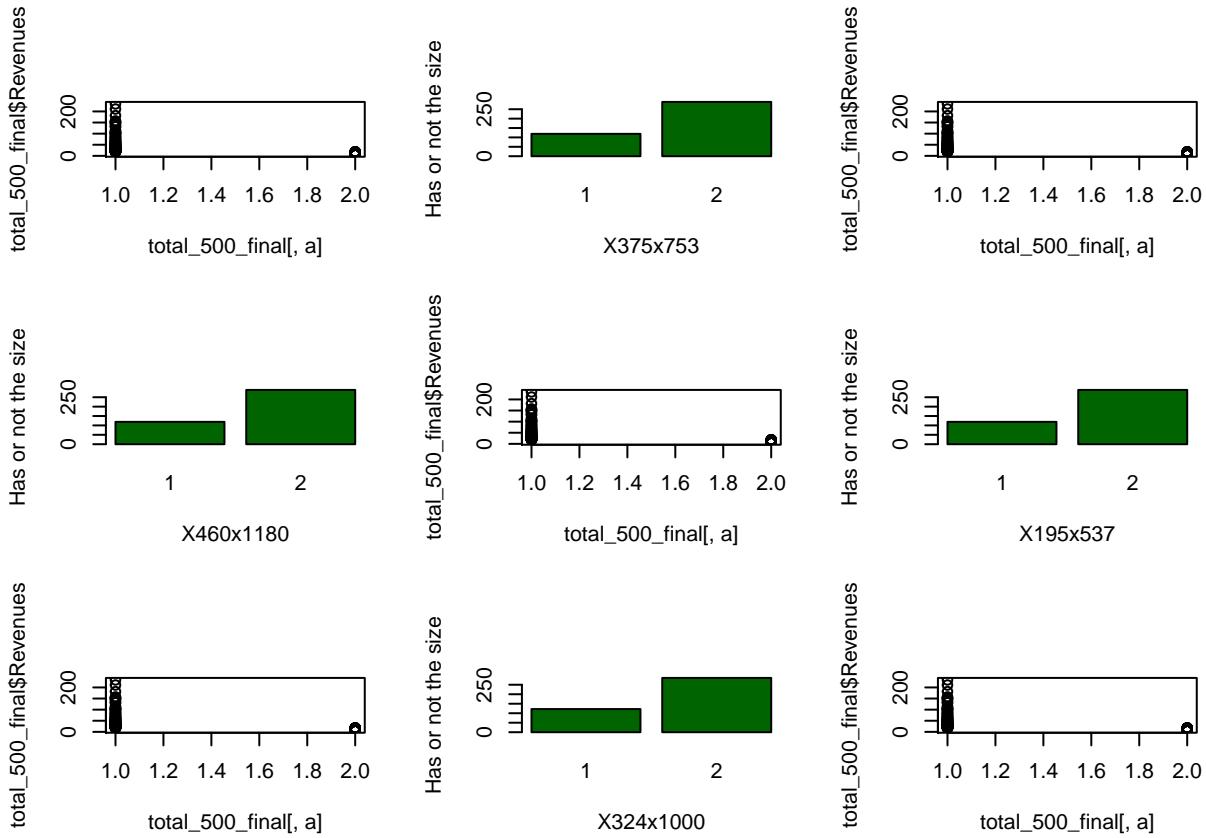


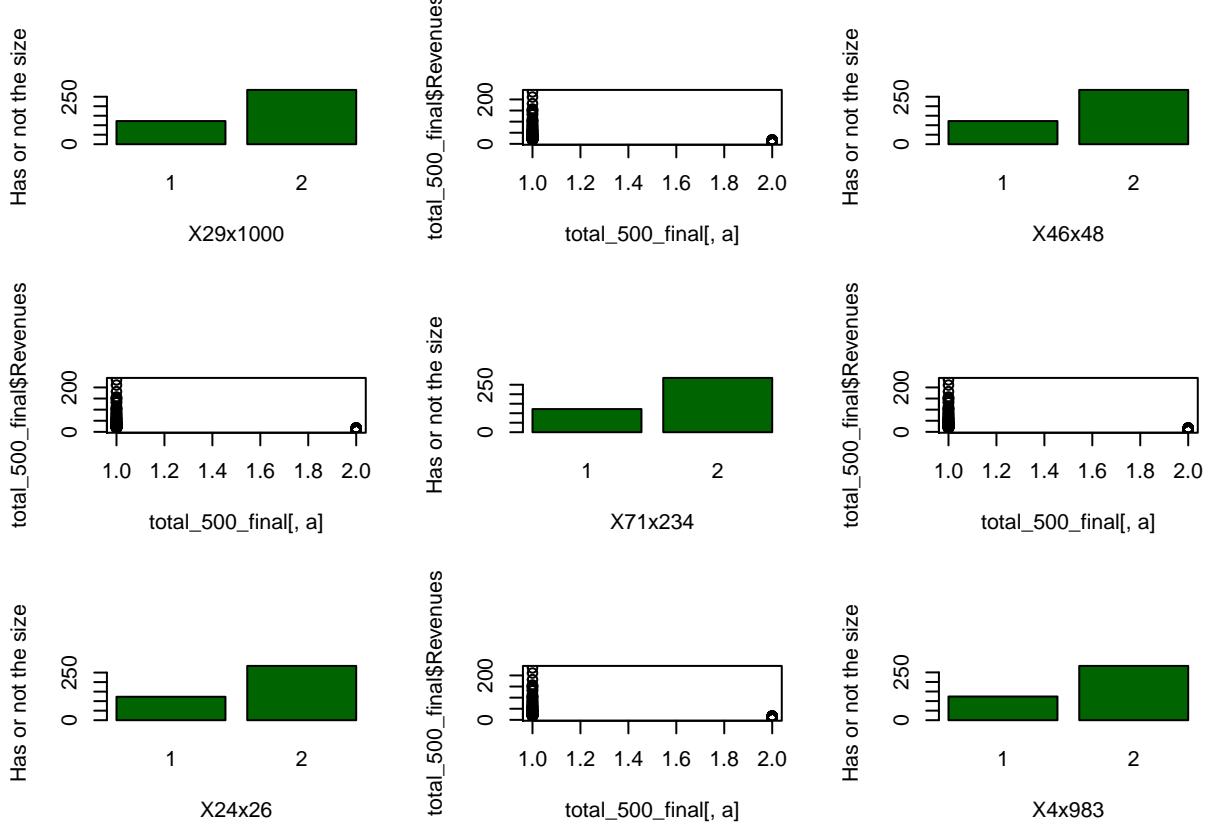


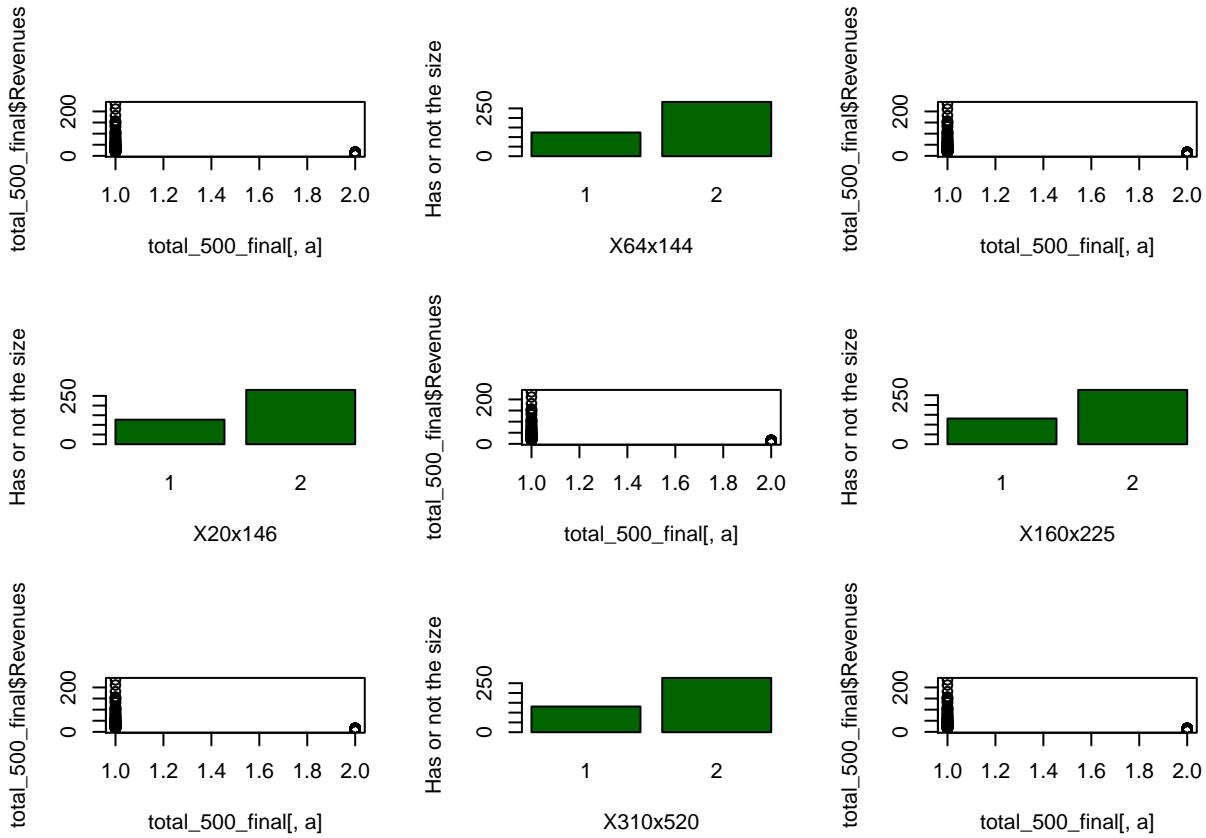


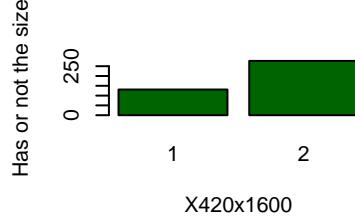
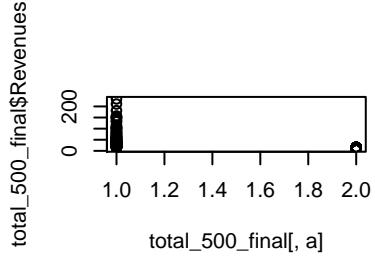
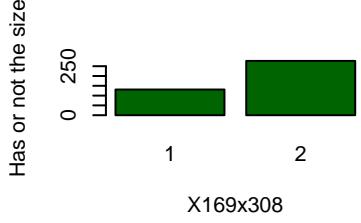
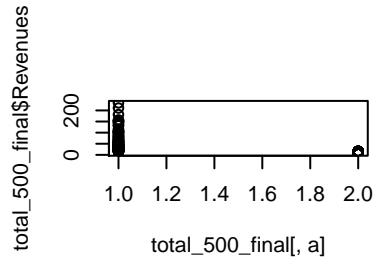
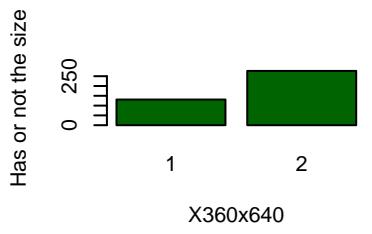
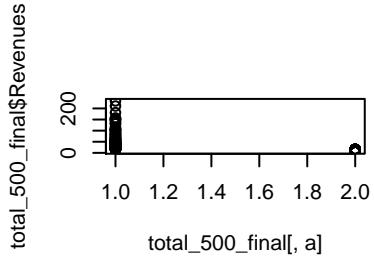
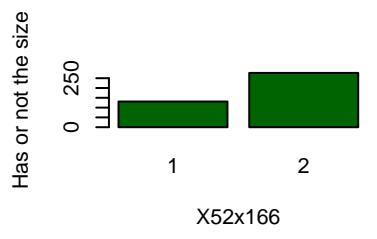
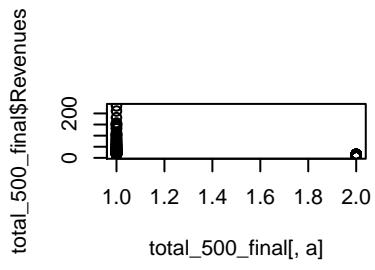
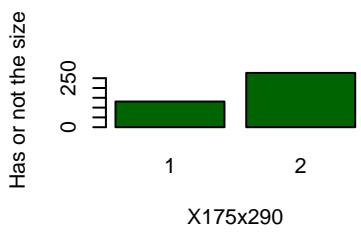


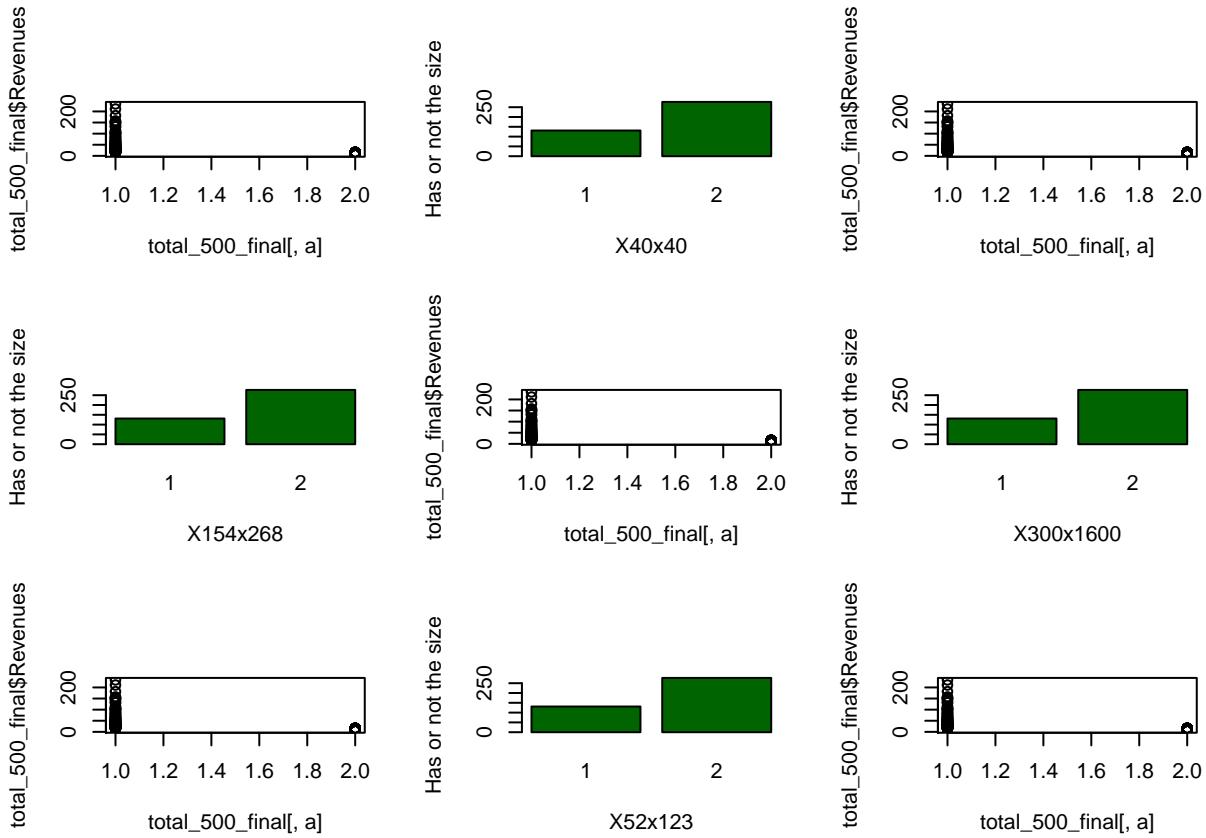


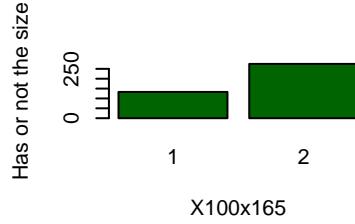
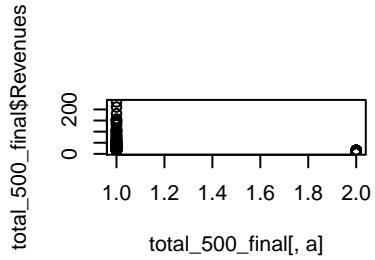
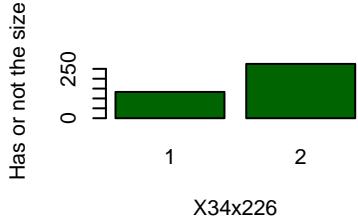
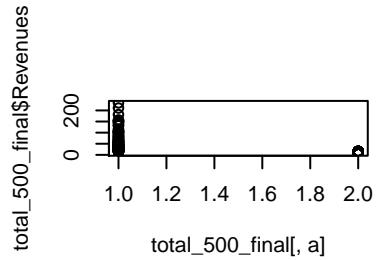
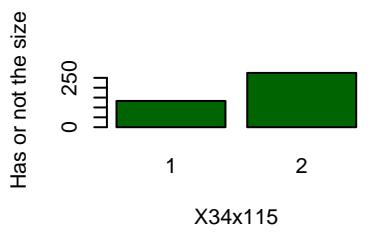
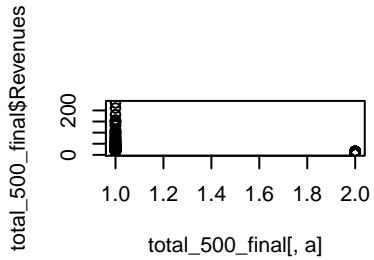
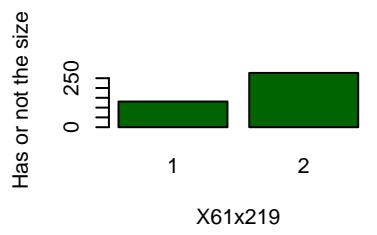
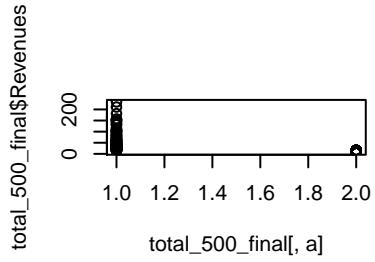
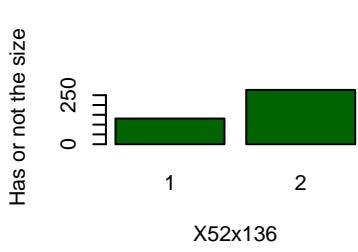


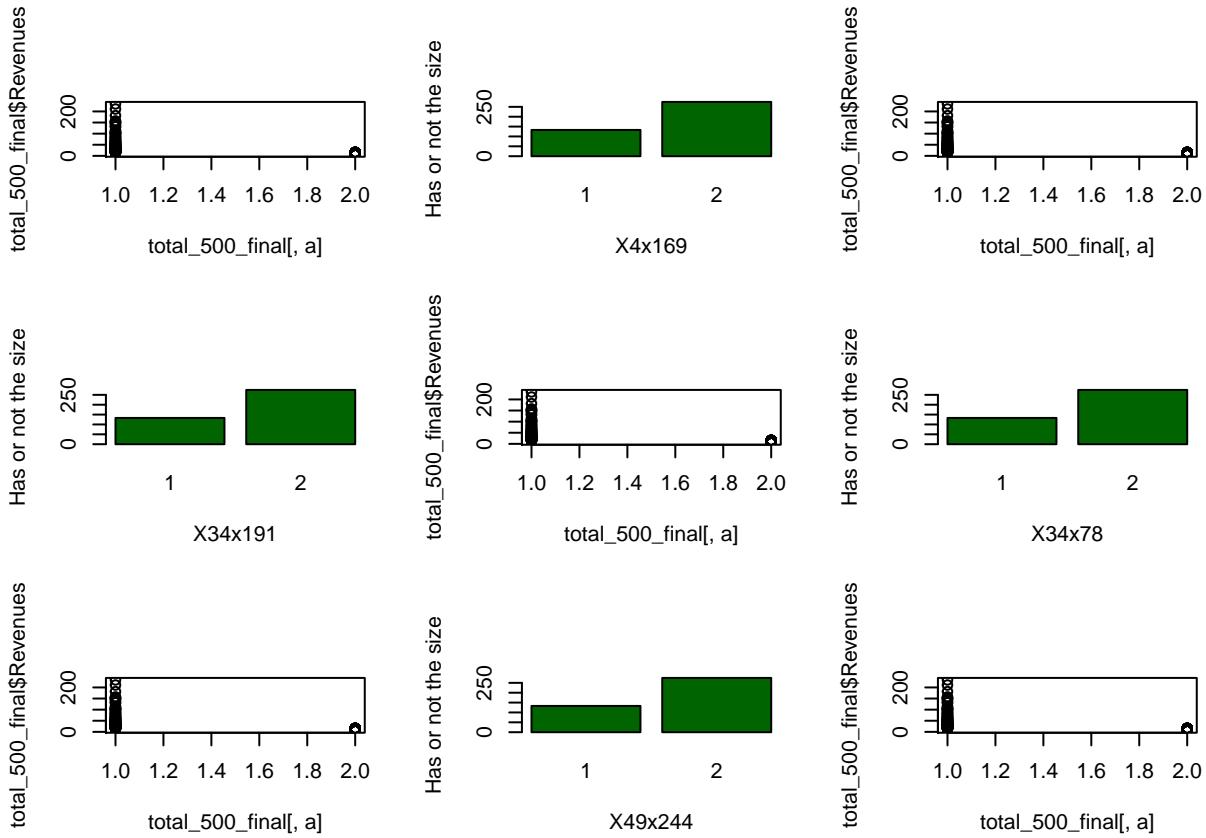


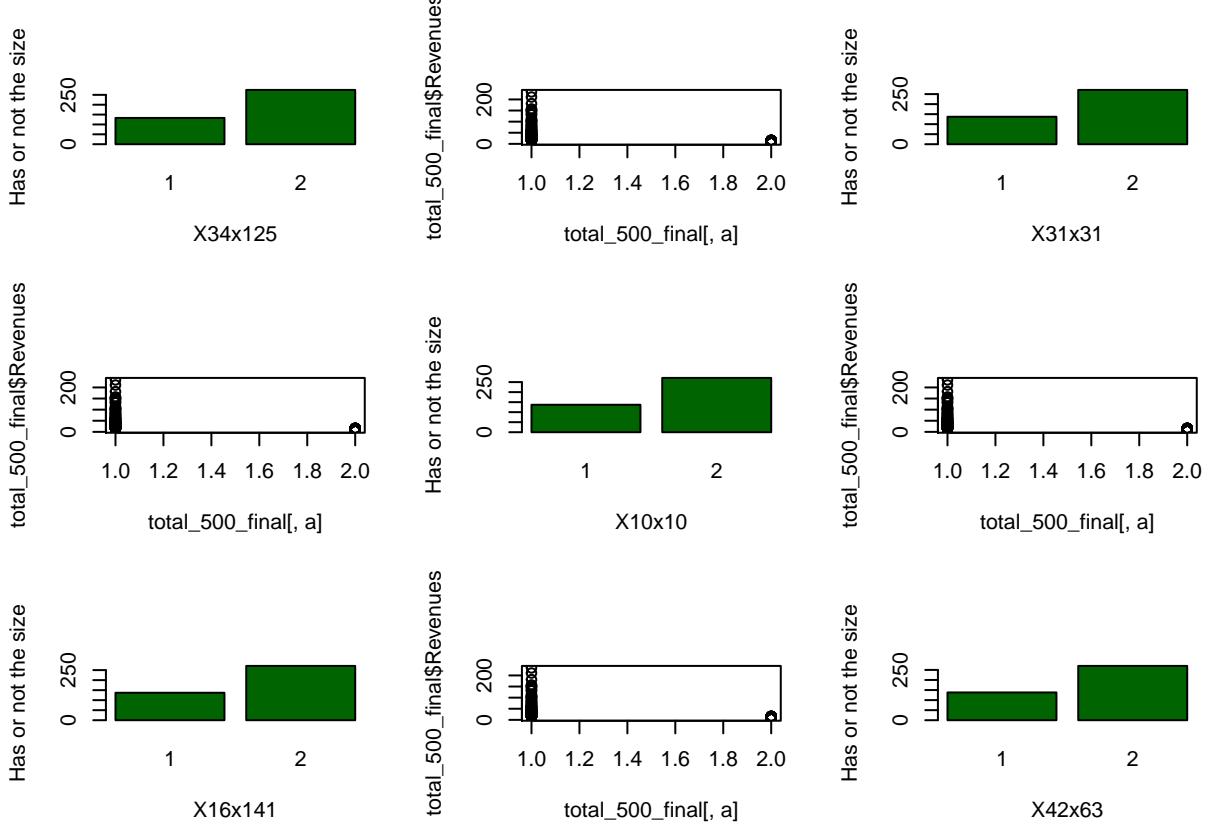


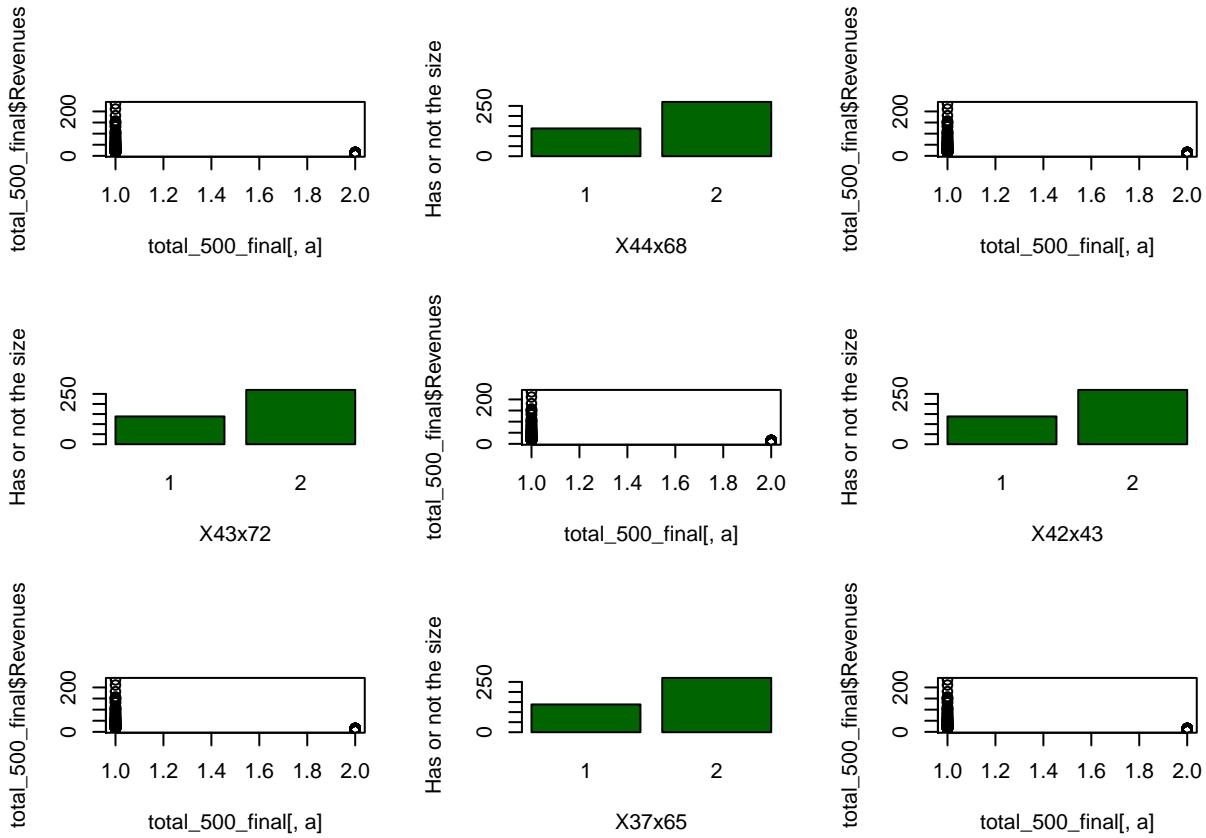


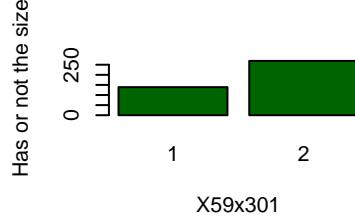
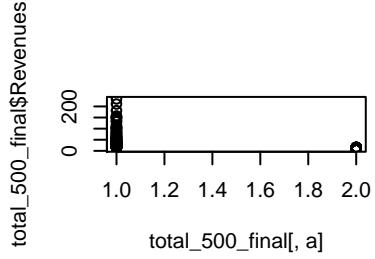
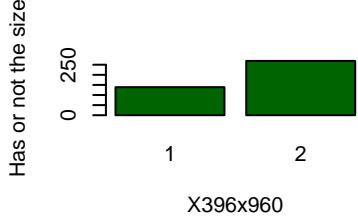
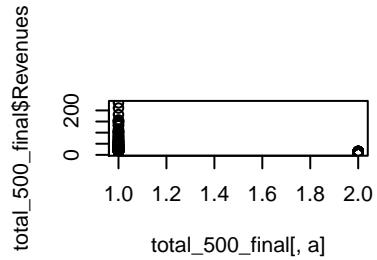
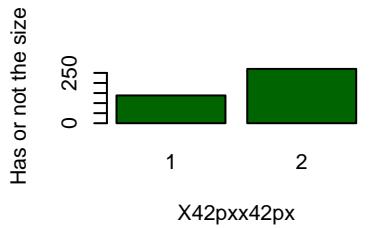
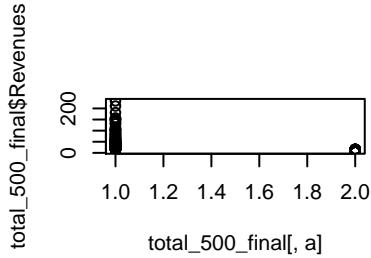
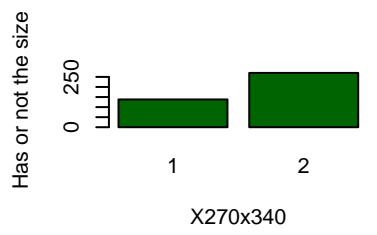
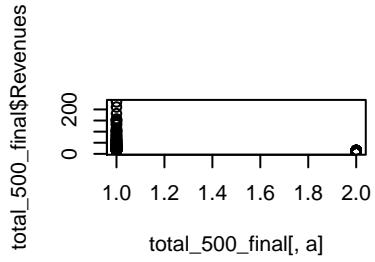
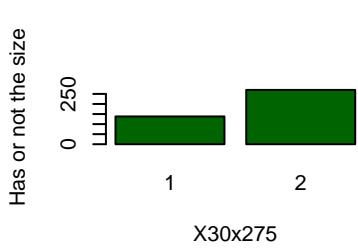


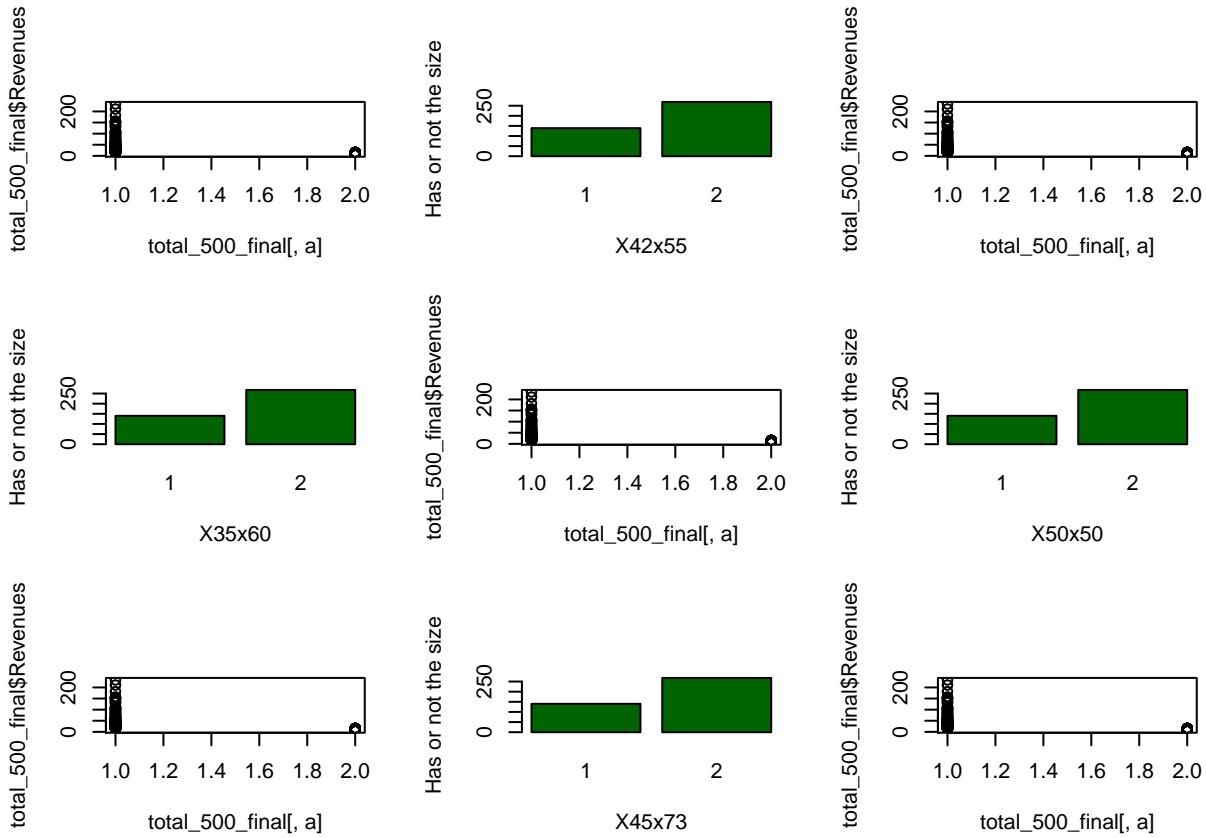


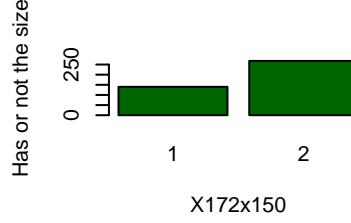
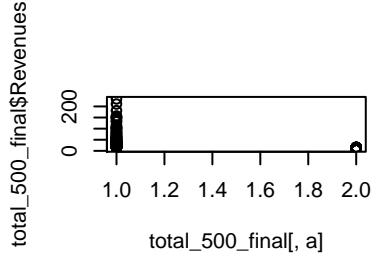
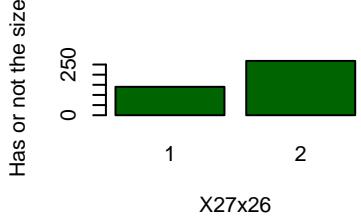
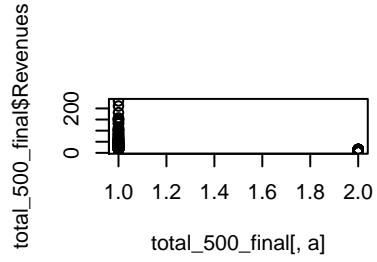
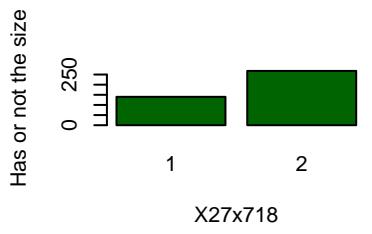
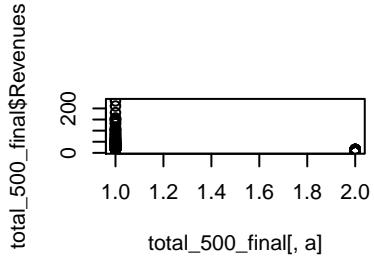
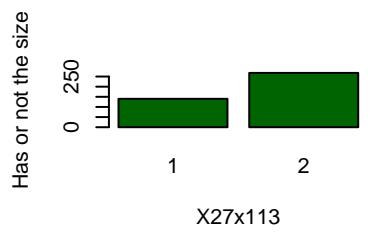
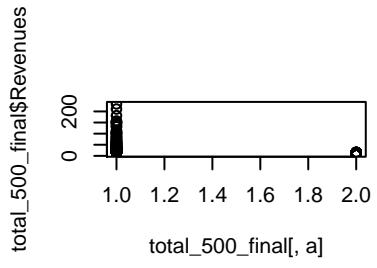
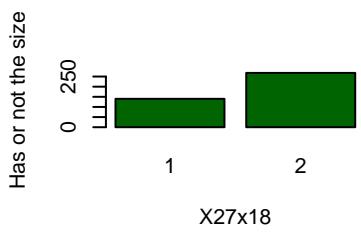


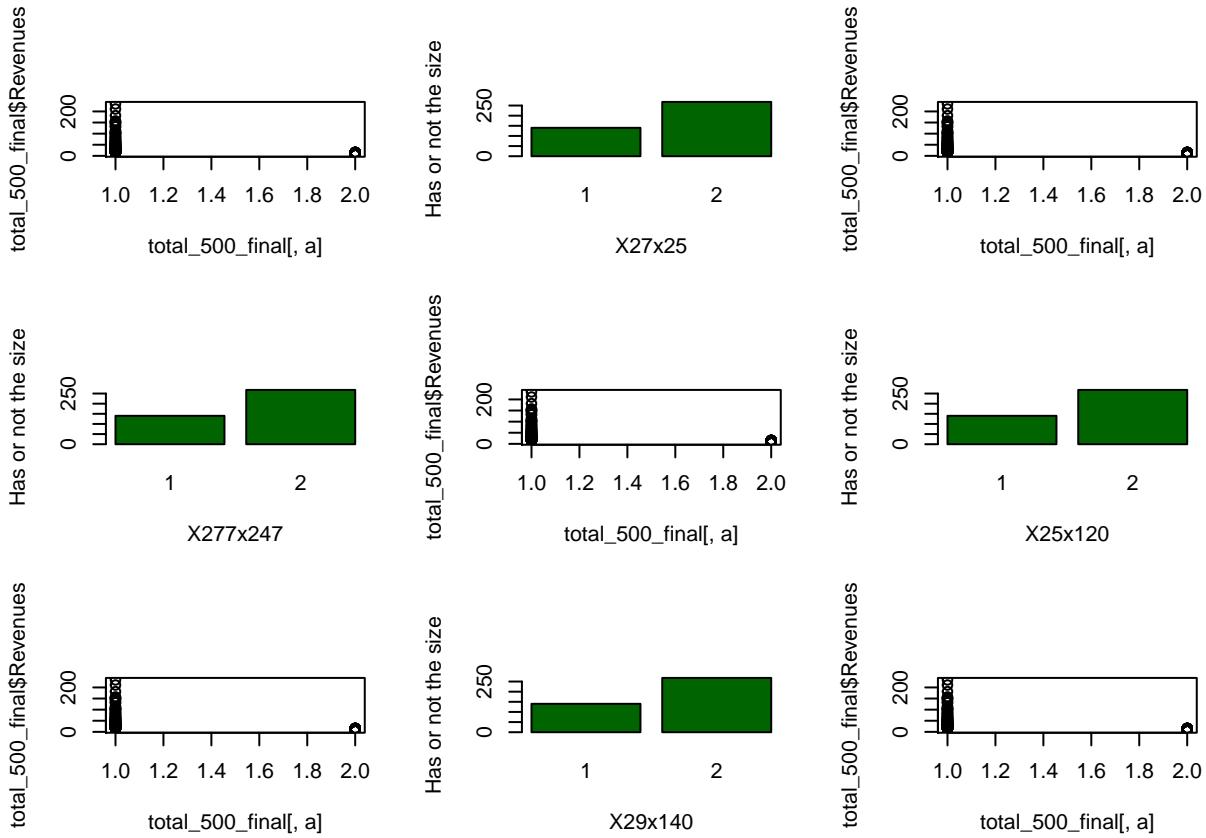


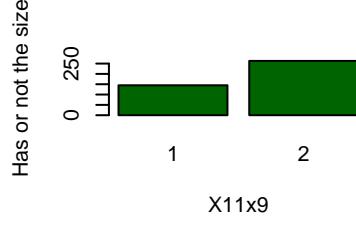
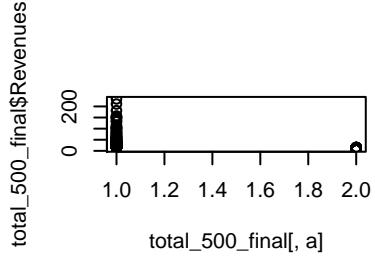
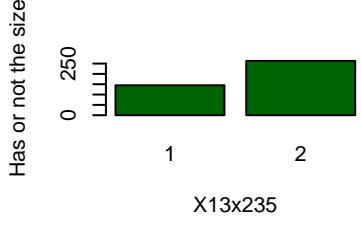
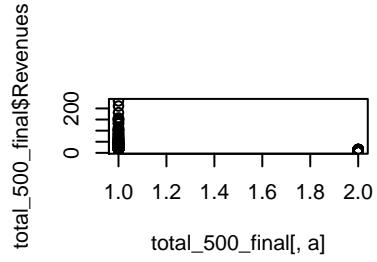
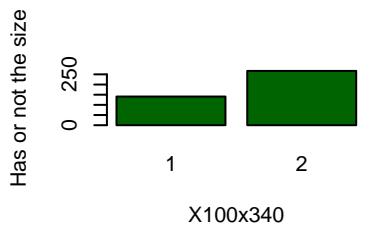
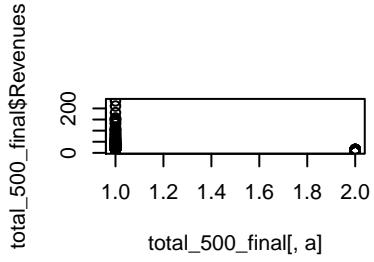
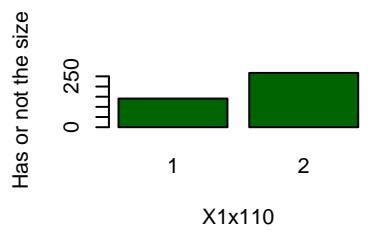
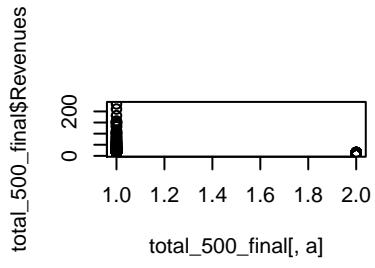
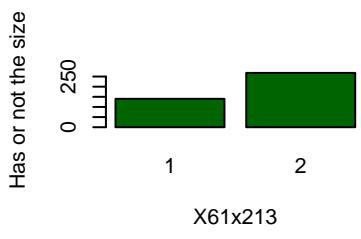


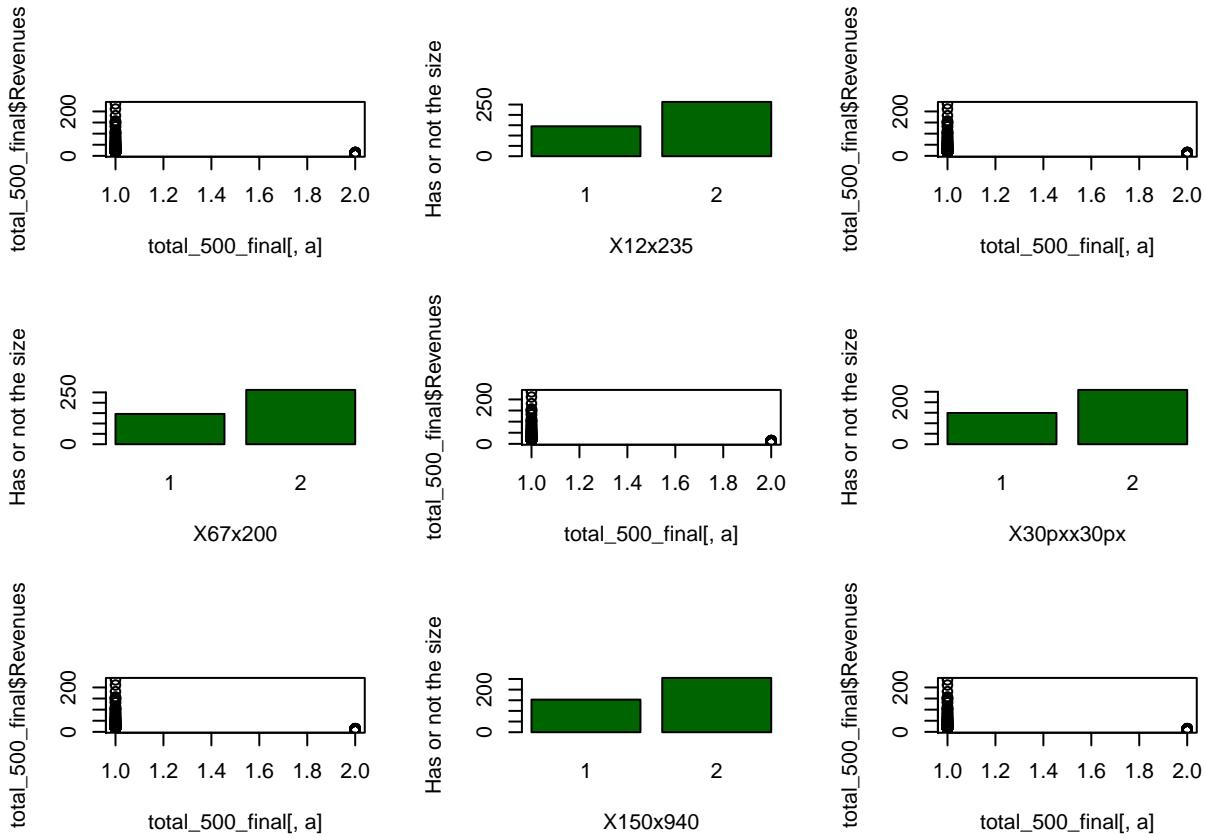


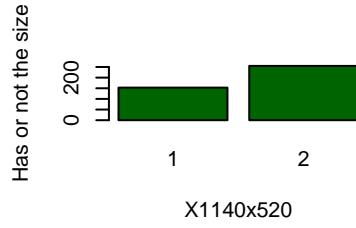
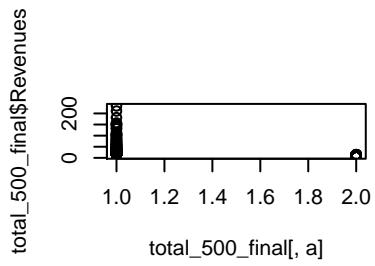
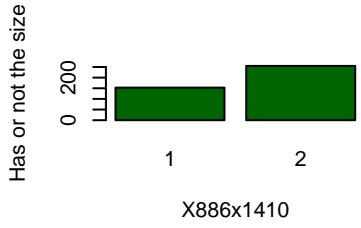
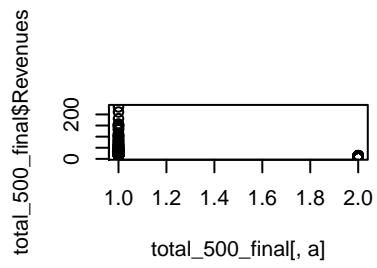
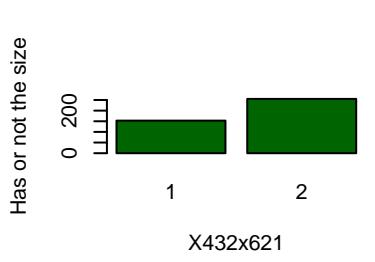
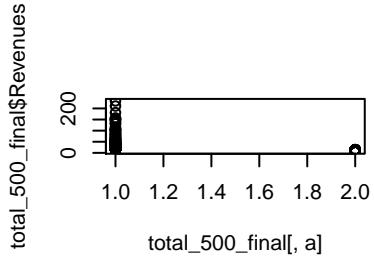
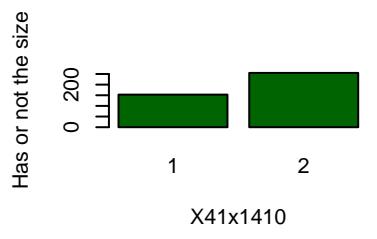
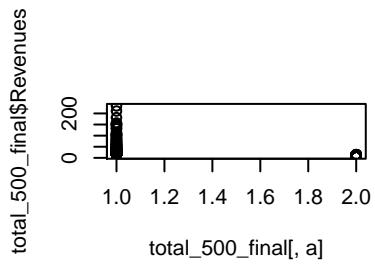
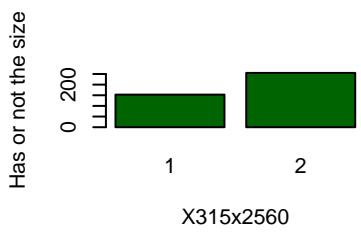


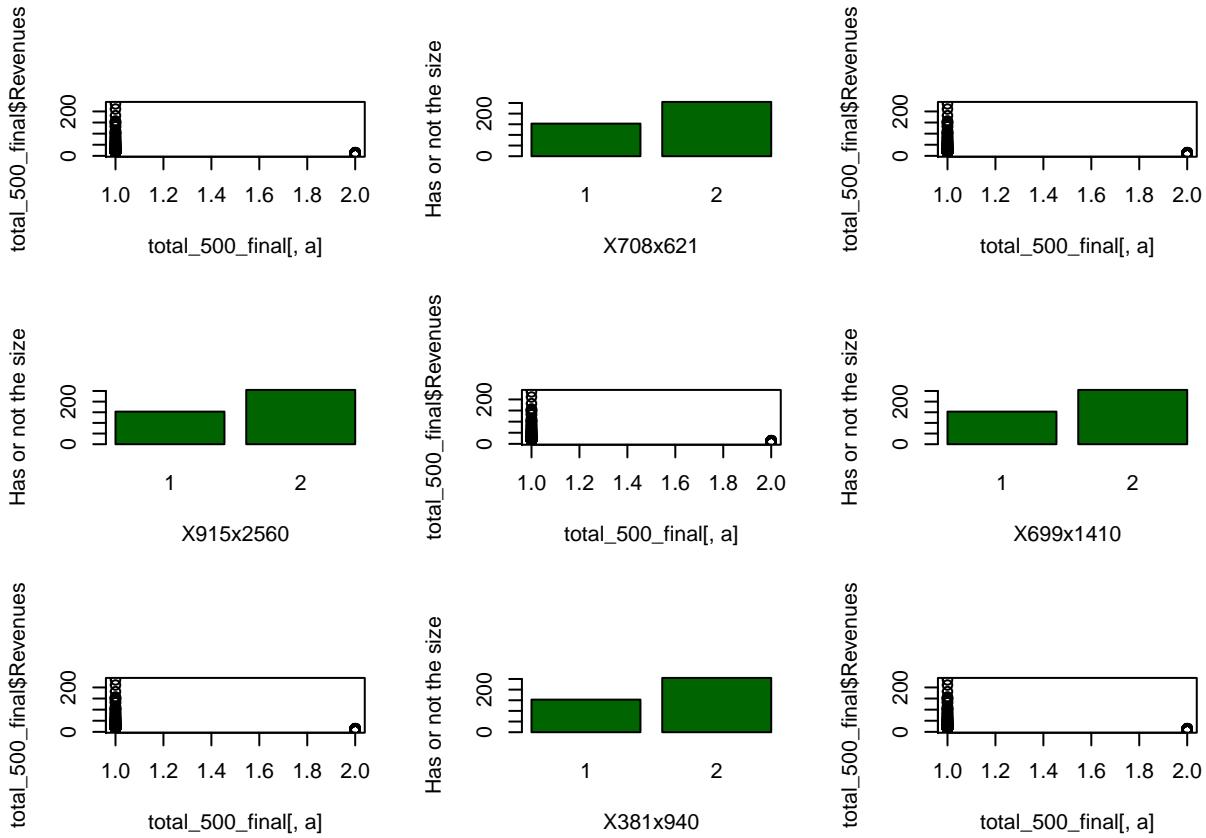


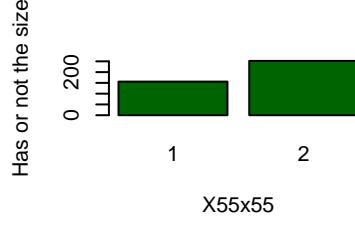
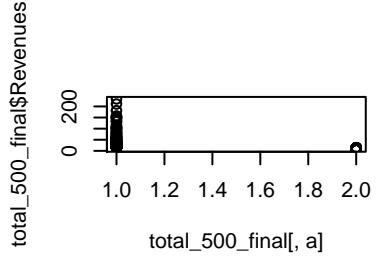
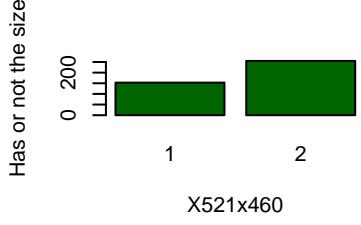
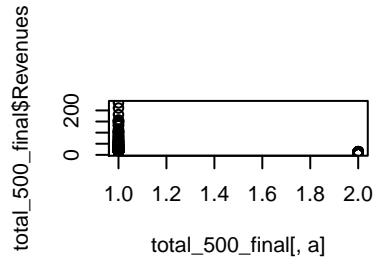
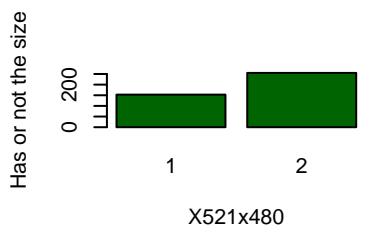
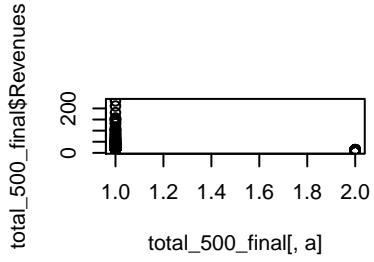
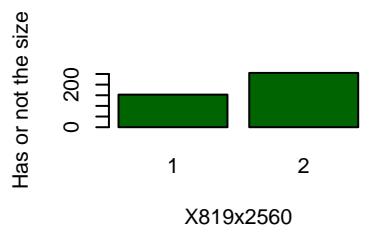
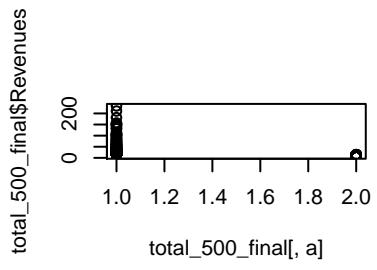
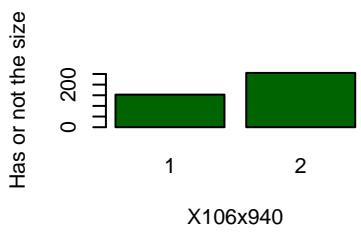


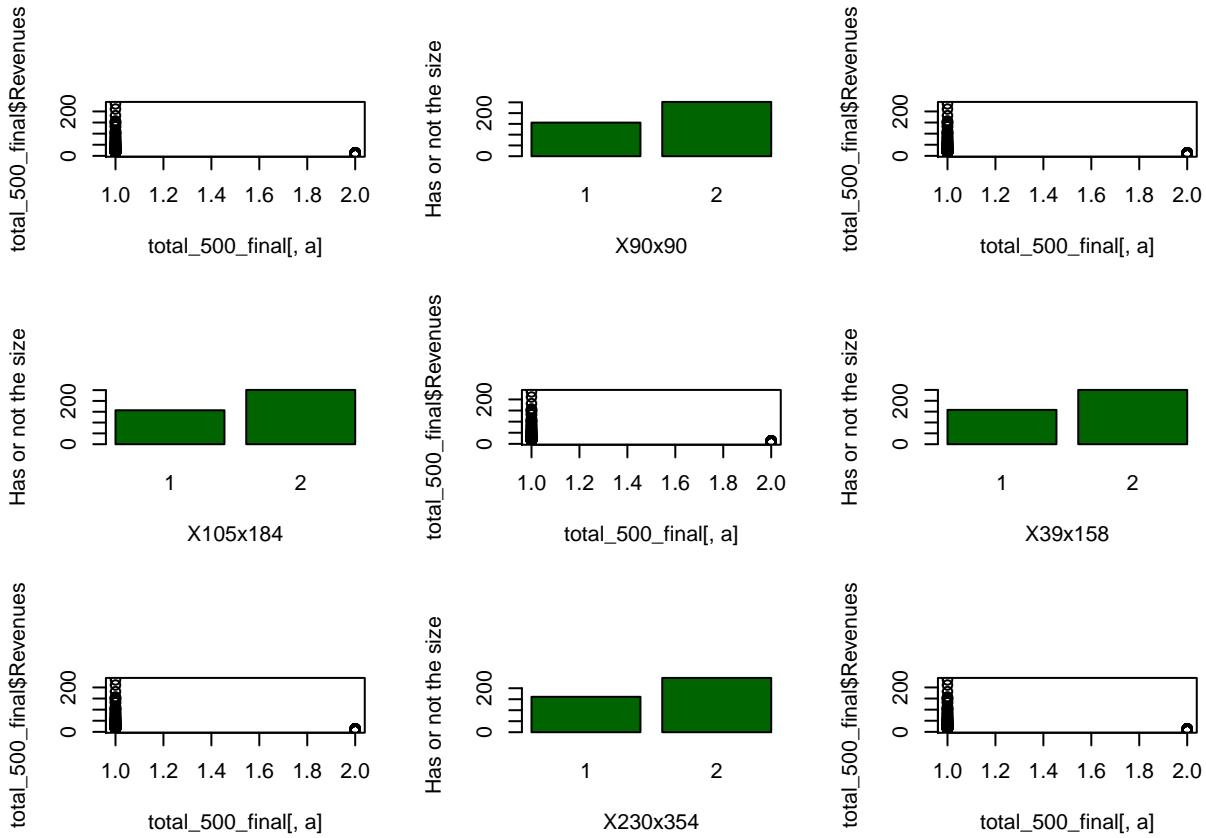


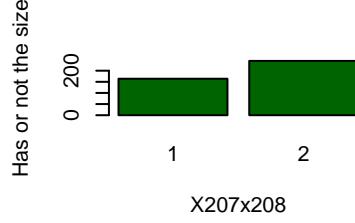
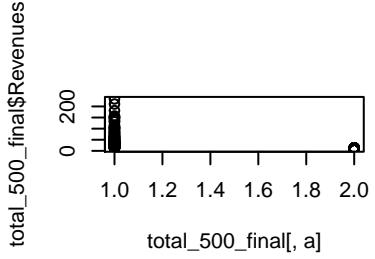
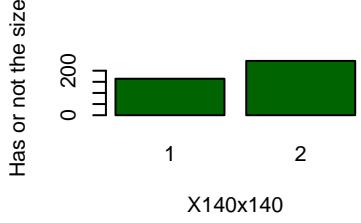
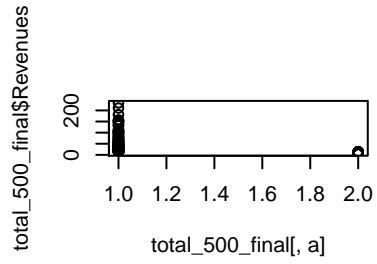
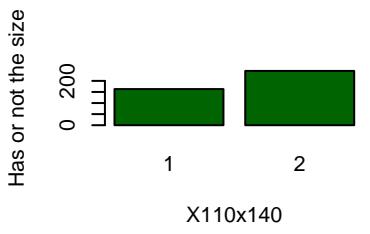
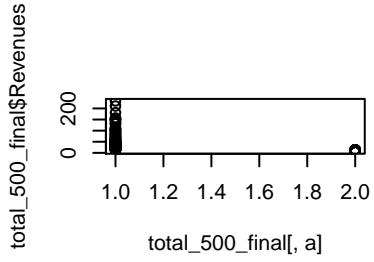
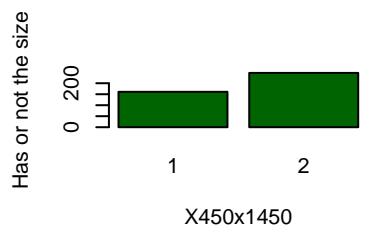
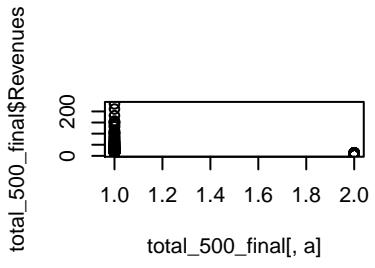
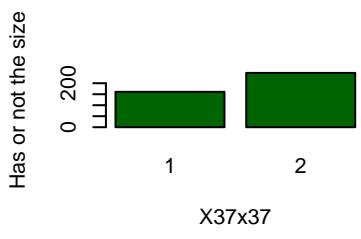


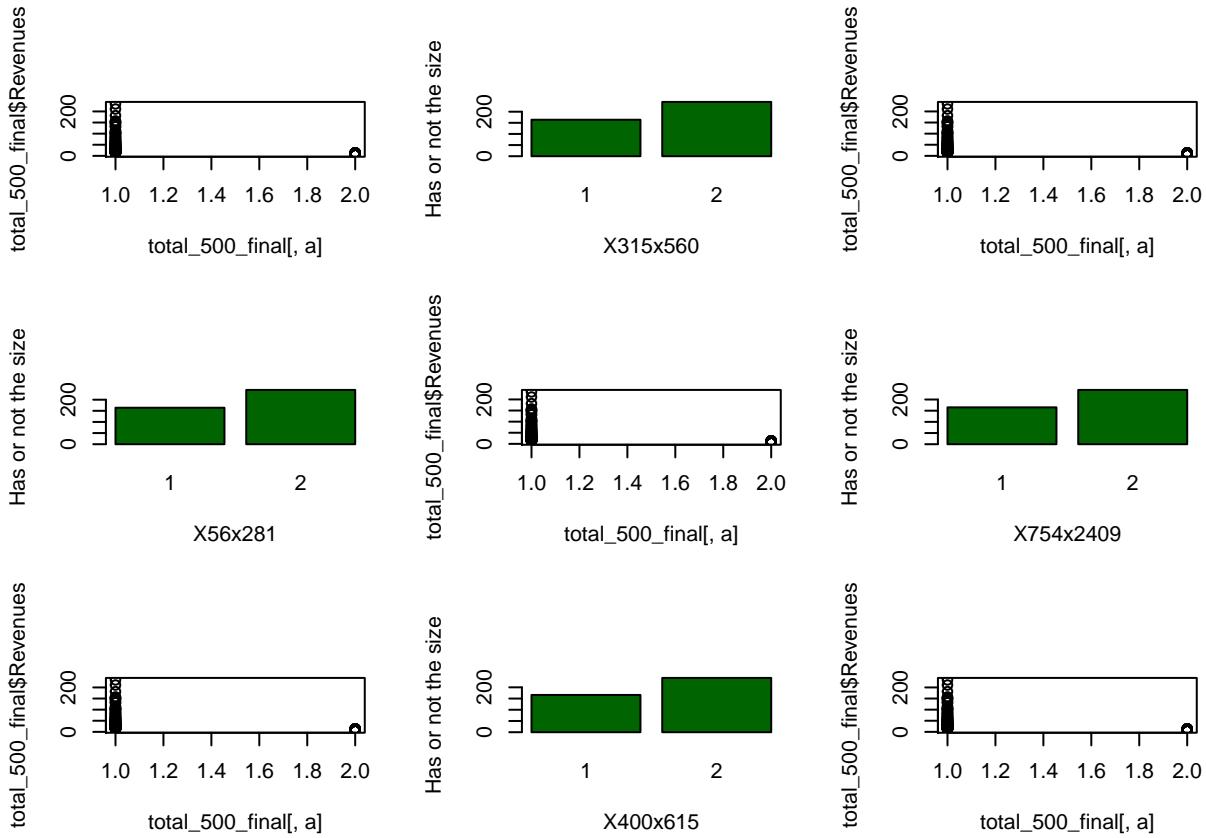


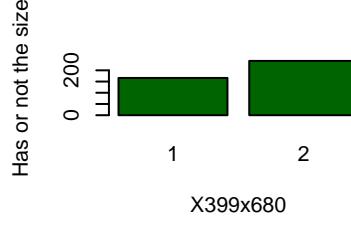
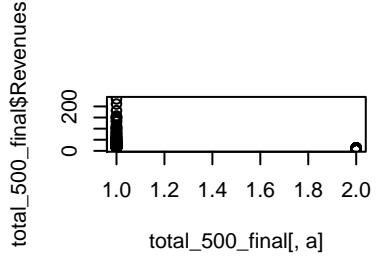
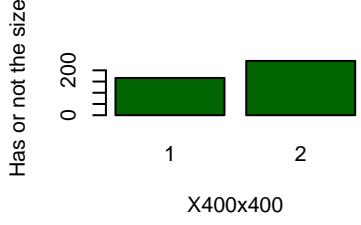
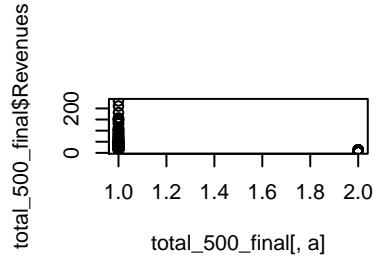
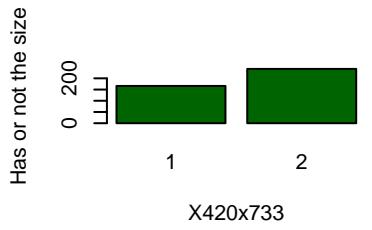
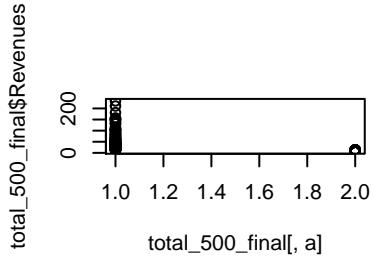
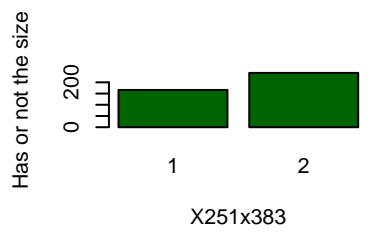
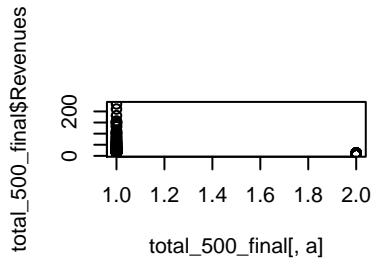
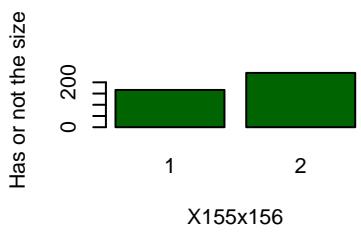


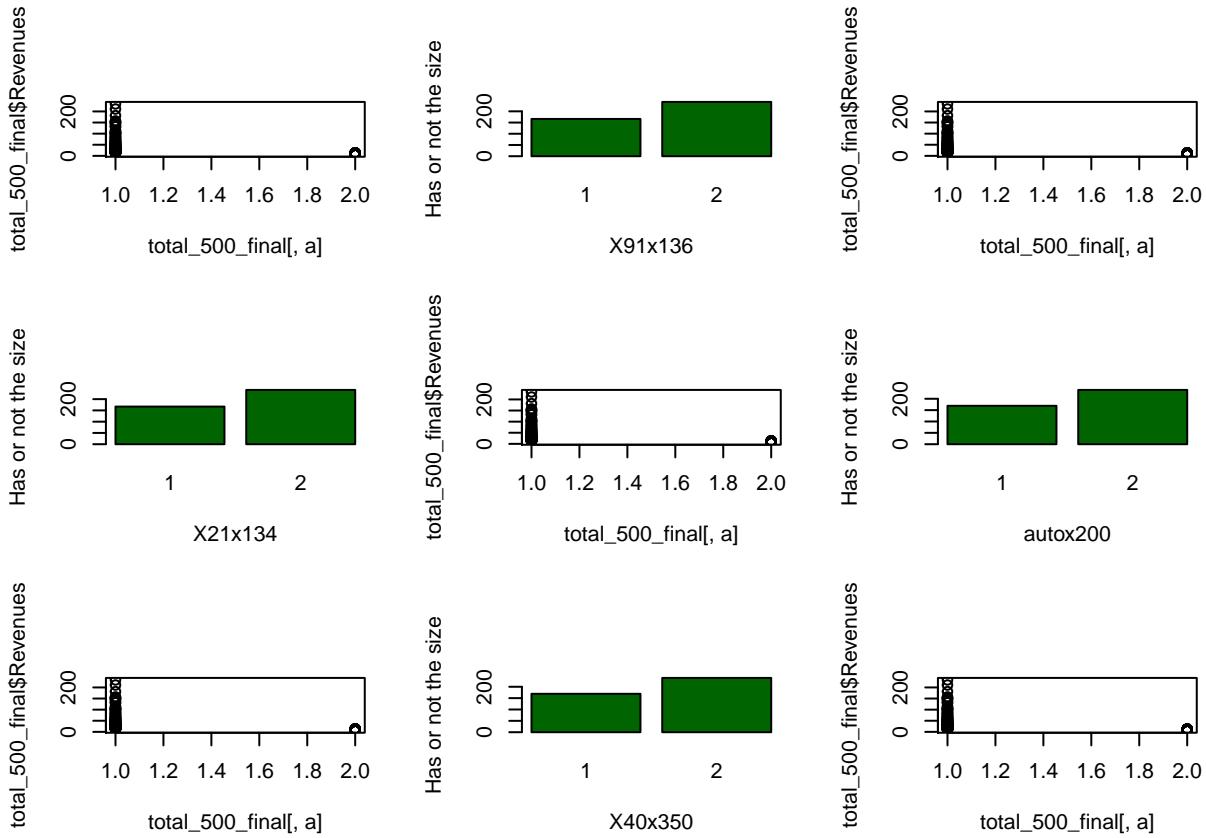


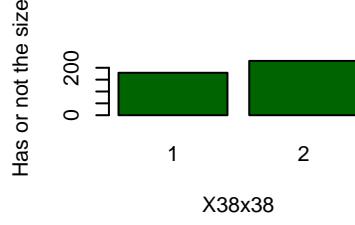
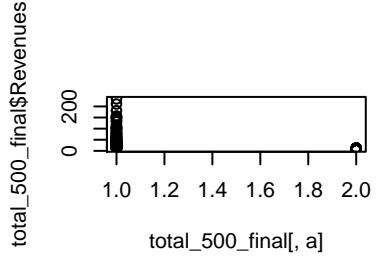
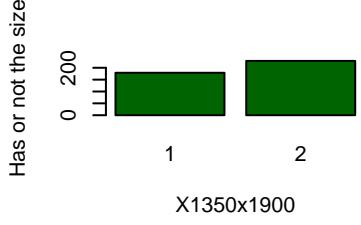
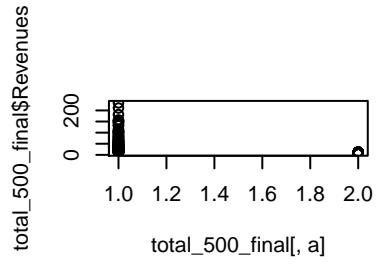
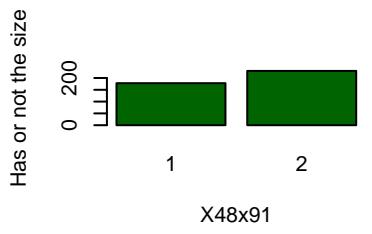
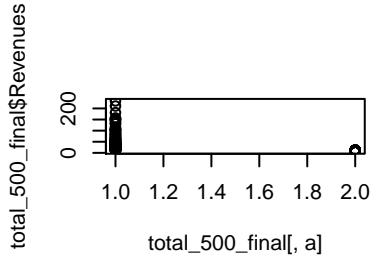
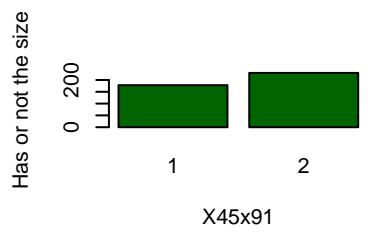
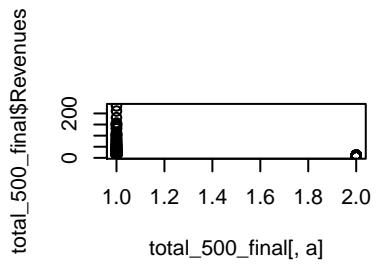
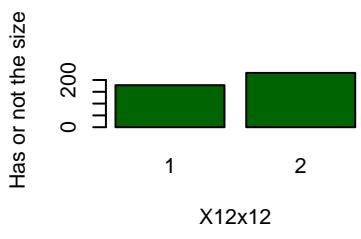


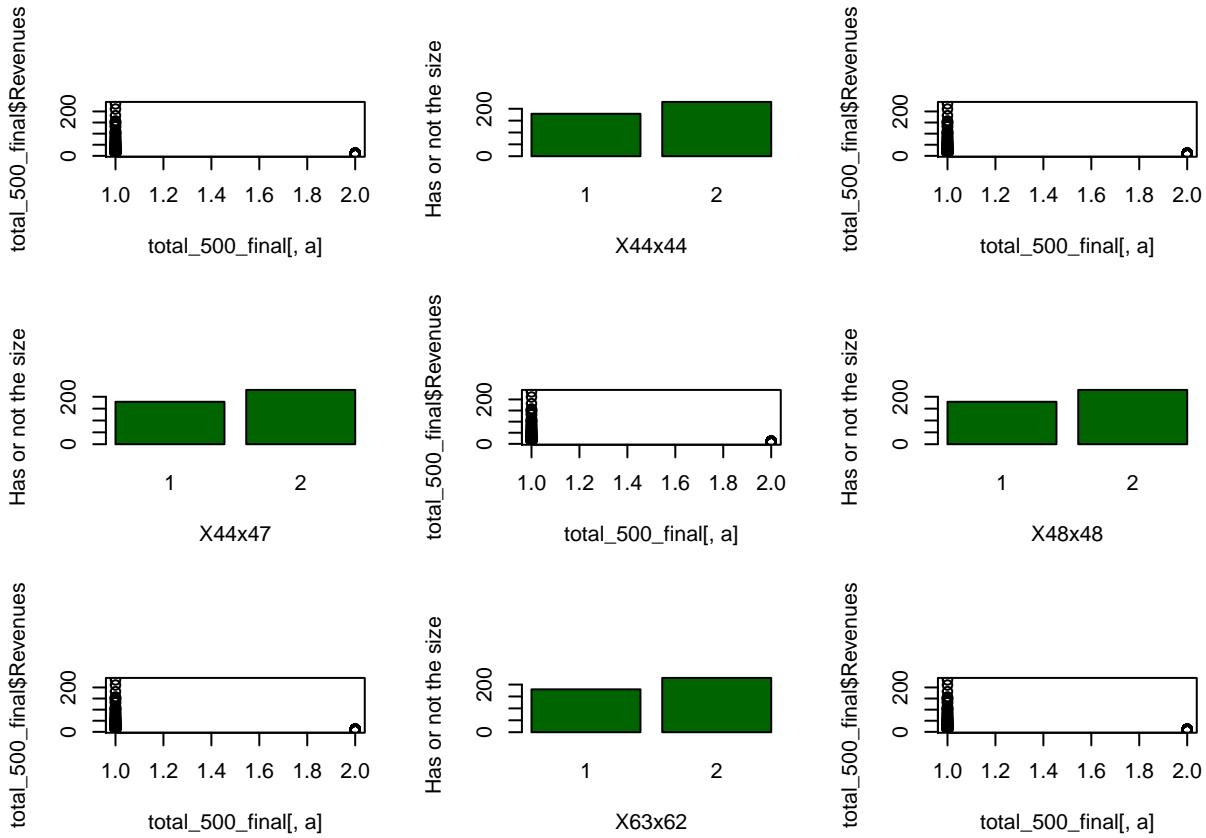


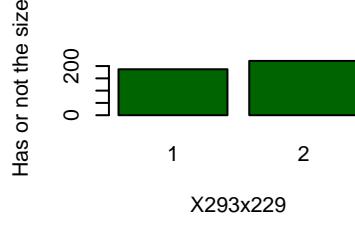
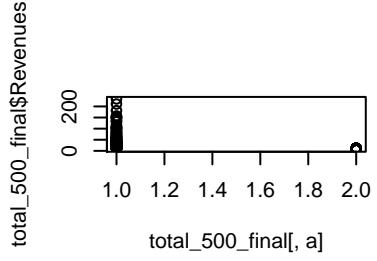
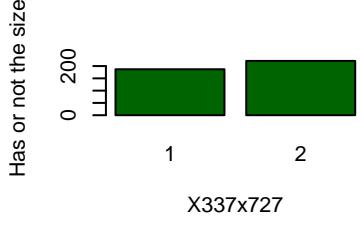
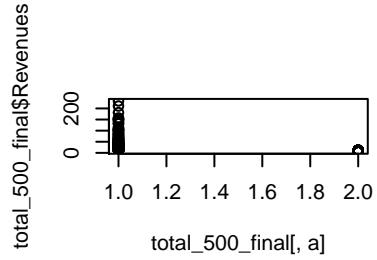
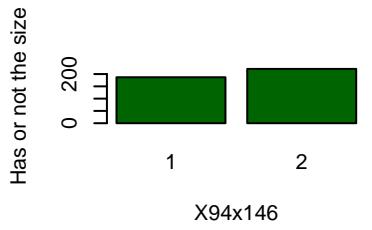
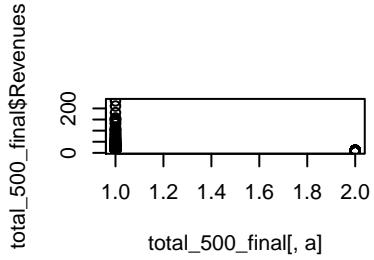
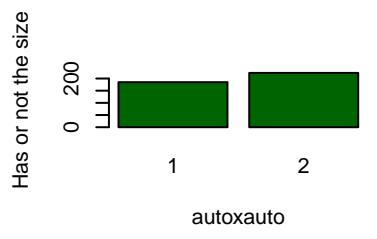
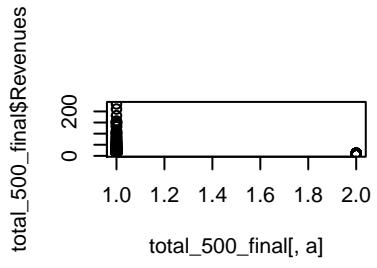
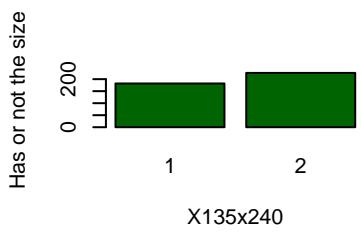


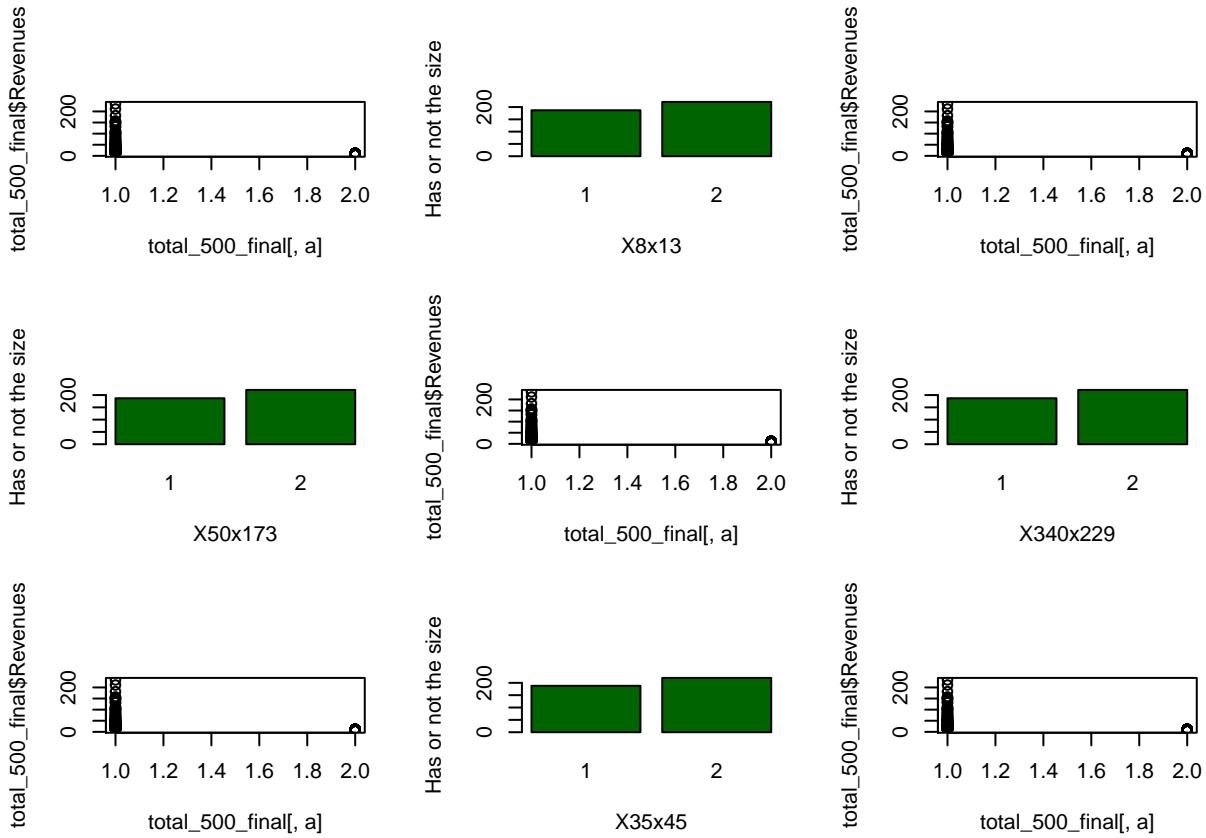


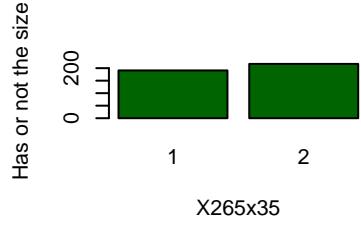
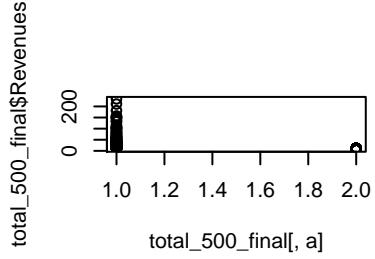
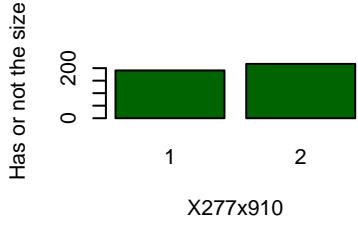
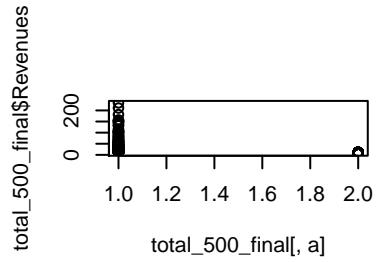
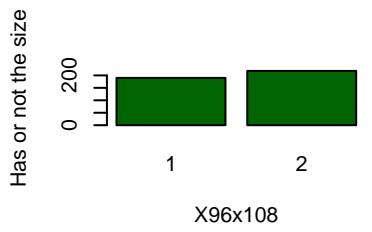
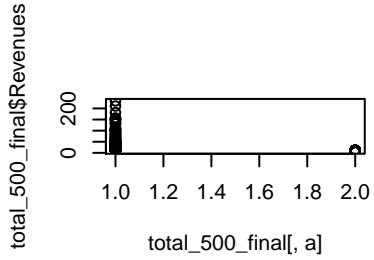
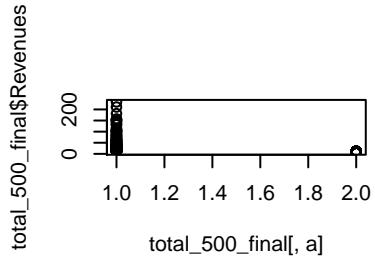
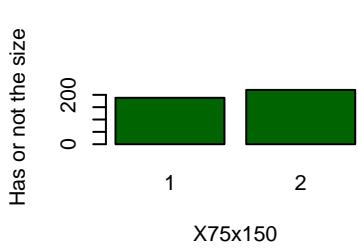


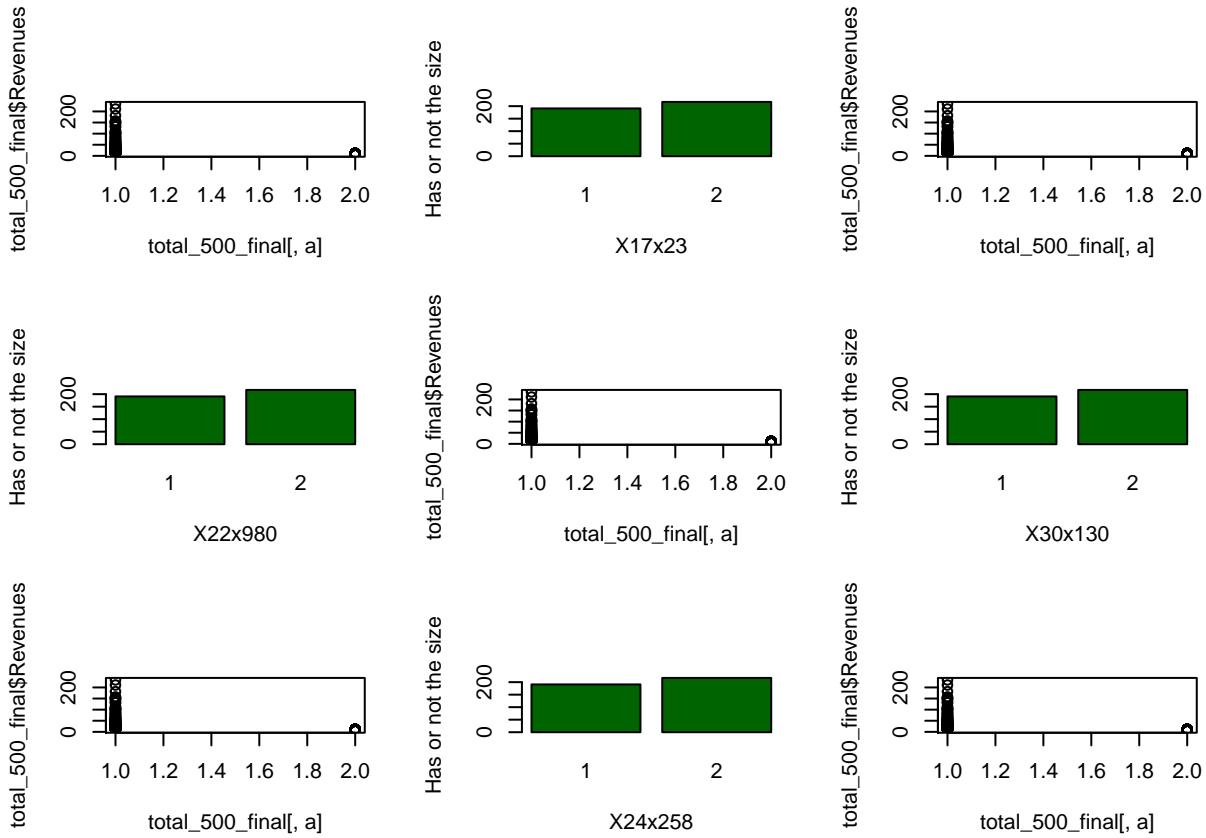


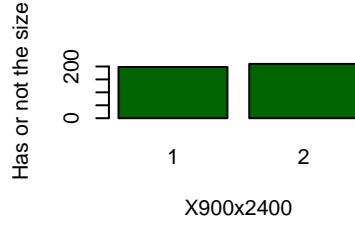
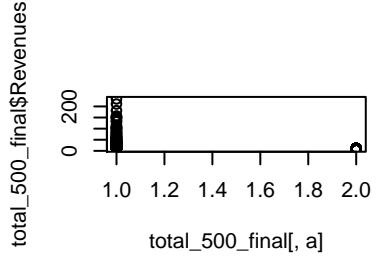
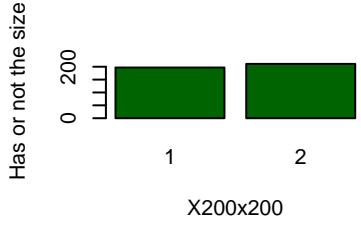
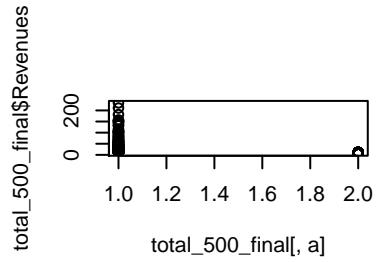
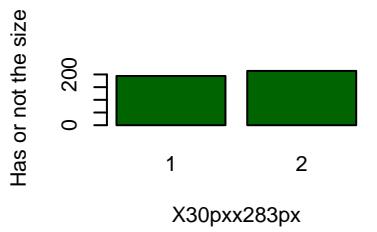
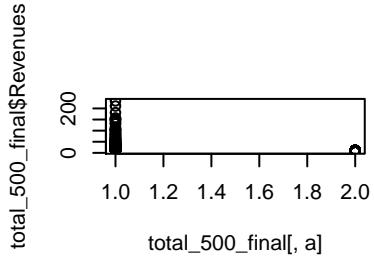
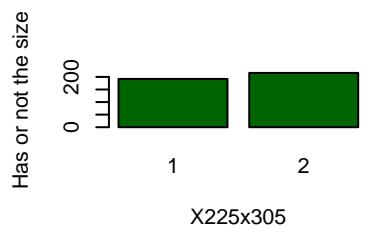
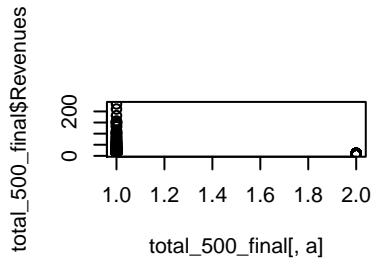
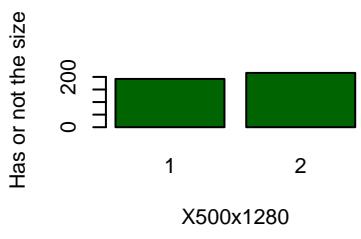


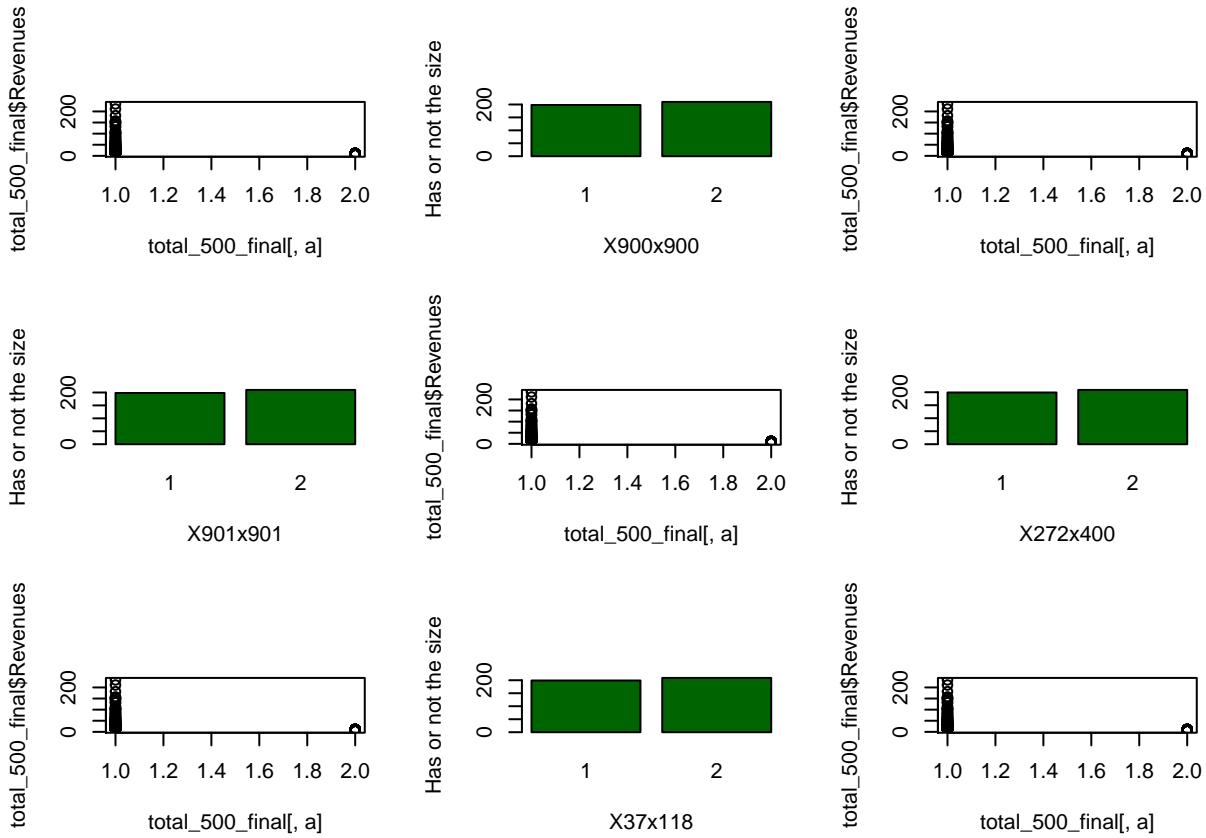


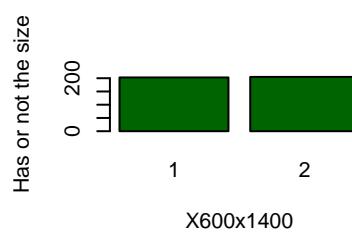
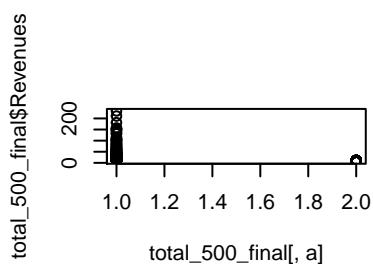
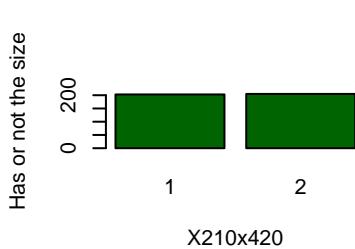












```
#By checking the above plots we can see that the 24 first sizes do appear to have some differentiation
par(mfrow=c(3,3))
keep = c()
for(i in 1:24){
  a = true_existing[i]
  keep = union (keep, c(a))}
keep

## [1] 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46
## [24] 47

#As we can see they are the variables from 24 to 47 and these are the only sizes we are going to keep
total_500_final <- total_500_final[,-c(48:715)]

#Also we remove the other Fortune 500 variables since they will interfere in the outcome of the model analysis
total_500_final$Market_Value <- NULL
total_500_final$Assets <- NULL
total_500_final$Ranking <- NULL
total_500_final$Total_SH_Equity <- NULL
total_500_final$The_page_opened <- NULL

summary(total_500_final)

##      Revenues      non.document.error number_of_errors number_of_warning 
## Min.   : 5.130   Min.   :0.0000    Min.   : 0.00   Min.   : 0.000  
## 1st Qu.: 7.047   1st Qu.:0.0000    1st Qu.: 0.00   1st Qu.: 0.000  
## Median :11.118   Median :0.0000    Median :13.00   Median : 3.000  
## Mean   :22.244   Mean   :0.2451    Mean   :37.36   Mean   : 8.669
```

```

## 3rd Qu.: 20.858   3rd Qu.:0.0000    3rd Qu.: 37.00   3rd Qu.:  9.000
## Max.   :233.715   Max.   :1.0000    Max.   :995.00   Max.   :214.000
##   facebook      instagram     linkedin      pinterest
## Min.   :0.0000    Min.   :0.000    Min.   :0.00000  Min.   :0.00000
## 1st Qu.:0.0000    1st Qu.:0.000    1st Qu.:0.00000  1st Qu.:0.00000
## Median :1.0000    Median :0.000    Median :1.00000  Median :0.00000
## Mean   :0.6471    Mean   :0.223    Mean   :0.5711   Mean   :0.09804
## 3rd Qu.:1.0000    3rd Qu.:0.000    3rd Qu.:1.00000  3rd Qu.:0.00000
## Max.   :1.0000    Max.   :1.000    Max.   :1.00000  Max.   :1.00000
##   twitter       youtube     Flesh_Mesaure   Readability
## Min.   :0.0000    Min.   :0.00000  Min.   :-3422.40  Min.   :1.000
## 1st Qu.:0.0000    1st Qu.:0.00000  1st Qu.: 34.85   1st Qu.:5.000
## Median :1.0000    Median :1.00000  Median : 45.55   Median :6.000
## Mean   :0.6863    Mean   :0.5833   Mean   : 35.37   Mean   :5.517
## 3rd Qu.:1.0000    3rd Qu.:1.00000  3rd Qu.: 55.40   3rd Qu.:6.000
## Max.   :1.0000    Max.   :1.00000  Max.   :121.20   Max.   :7.000
##   Sentences     Unique.words    Words      external
## Min.   : 1.0      Min.   : 0.0    Min.   : 1.0    Min.   : 0.00
## 1st Qu.: 73.5    1st Qu.: 59.0   1st Qu.: 287.0  1st Qu.: 2.00
## Median :139.5    Median :113.0   Median : 504.5   Median : 5.00
## Mean   :178.9    Mean   :152.3   Mean   :692.1    Mean   :18.00
## 3rd Qu.:243.5    3rd Qu.:191.5   3rd Qu.: 905.5  3rd Qu.:13.25
## Max.   :1350.0   Max.   :1910.0  Max.   :8306.0   Max.   :545.00
##   internal      total.links    X15x75     X8x15
## Min.   : 0.0      Min.   : 0.0    Min.   :1.000   Min.   :1.000
## 1st Qu.: 73.0    1st Qu.: 81.0   1st Qu.:2.000   1st Qu.:2.000
## Median :117.0    Median :134.5   Median :2.000   Median :2.000
## Mean   :154.9    Mean   :172.9   Mean   :1.998   Mean   :1.995
## 3rd Qu.:183.0    3rd Qu.:212.2   3rd Qu.:2.000   3rd Qu.:2.000
## Max.   :1254.0   Max.   :1255.0  Max.   :2.000   Max.   :2.000
##   X44x556      X1x1        X800x1200  autox100.
## Min.   :1.000    Min.   :1.000   Min.   :1.00    Min.   :1.000
## 1st Qu.:2.000   1st Qu.:2.000   1st Qu.:2.00   1st Qu.:2.000
## Median :2.000    Median :2.000   Median :2.00   Median :2.000
## Mean   :1.993    Mean   :1.993   Mean   :1.99    Mean   :1.985
## 3rd Qu.:2.000   3rd Qu.:2.000   3rd Qu.:2.00   3rd Qu.:2.000
## Max.   :2.000    Max.   :2.000   Max.   :2.00   Max.   :2.000
##   X24pxx133px   X21pxx173px  X46x214     X49x49
## Min.   :1.000    Min.   :1.000   Min.   :1.00    Min.   :1.00
## 1st Qu.:2.000   1st Qu.:2.000   1st Qu.:2.00   1st Qu.:2.00
## Median :2.000    Median :2.000   Median :2.00   Median :2.00
## Mean   :1.983    Mean   :1.983   Mean   :1.98    Mean   :1.98
## 3rd Qu.:2.000   3rd Qu.:2.000   3rd Qu.:2.00   3rd Qu.:2.00
## Max.   :2.000    Max.   :2.000   Max.   :2.00   Max.   :2.00
##   X50x45       X400x300    X292pxx292px  X200pxx200px
## Min.   :1.000    Min.   :1.000   Min.   :1.000   Min.   :1.000
## 1st Qu.:2.000   1st Qu.:2.000   1st Qu.:2.000  1st Qu.:2.000
## Median :2.000    Median :2.000   Median :2.000  Median :2.000
## Mean   :1.975    Mean   :1.973   Mean   :1.968   Mean   :1.968
## 3rd Qu.:2.000   3rd Qu.:2.000   3rd Qu.:2.000  3rd Qu.:2.000
## Max.   :2.000    Max.   :2.000   Max.   :2.000   Max.   :2.000
##   X1279pxx984px X300pxx1500px X29x29      X115x223
## Min.   :1.000    Min.   :1.000   Min.   :1.000   Min.   :1.000
## 1st Qu.:2.000   1st Qu.:2.000   1st Qu.:2.000  1st Qu.:2.000

```

```

## Median :2.000  Median :2.000  Median :2.000  Median :2.000
## Mean   :1.968  Mean   :1.968  Mean   :1.961  Mean   :1.951
## 3rd Qu.:2.000 3rd Qu.:2.000 3rd Qu.:2.000 3rd Qu.:2.000
## Max.   :2.000  Max.   :2.000  Max.   :2.000  Max.   :2.000
## X160x233      X300x993      X41x192      X28x221
## Min.   :1.000  Min.   :1.000  Min.   :1.000  Min.   :1.000
## 1st Qu.:2.000 1st Qu.:2.000 1st Qu.:2.000 1st Qu.:2.000
## Median :2.000  Median :2.000  Median :2.000  Median :2.000
## Mean   :1.951  Mean   :1.951  Mean   :1.951  Mean   :1.951
## 3rd Qu.:2.000 3rd Qu.:2.000 3rd Qu.:2.000 3rd Qu.:2.000
## Max.   :2.000  Max.   :2.000  Max.   :2.000  Max.   :2.000
## X15x12        X60x60        .bmp         .dib
## Min.   :1.000  Min.   :1.000  Min.   :0.00000  Min.   :0.0000
## 1st Qu.:2.000 1st Qu.:2.000 1st Qu.:0.00000 1st Qu.:0.0000
## Median :2.000  Median :2.000  Median :0.00000  Median :0.0000
## Mean   :1.949  Mean   :1.946  Mean   :0.06863  Mean   :0.1838
## 3rd Qu.:2.000 3rd Qu.:2.000 3rd Qu.:0.00000 3rd Qu.:0.0000
## Max.   :2.000  Max.   :2.000  Max.   :23.00000  Max.   :35.0000
## .gif          .jpe          .jpeg        .jpg
## Min.   : 0.000  Min.   : 0.000  Min.   : 0.000  Min.   : 0.00
## 1st Qu.: 0.000  1st Qu.: 0.000  1st Qu.: 0.000  1st Qu.: 2.00
## Median : 1.000  Median : 0.000  Median : 0.000  Median : 8.50
## Mean   : 4.081  Mean   : 2.863  Mean   : 2.821  Mean   : 18.09
## 3rd Qu.: 3.000  3rd Qu.: 0.000  3rd Qu.: 0.000  3rd Qu.: 15.25
## Max.   :143.000  Max.   :968.000  Max.   :968.000  Max.   :363.00
## .png          .tif          .tiff       total.images
## Min.   : 0.00  Min.   : 0.000  Min.   :0.00000  Min.   : 0.00
## 1st Qu.: 3.00  1st Qu.: 0.000  1st Qu.:0.00000  1st Qu.: 13.00
## Median : 8.00  Median : 0.000  Median :0.00000  Median : 24.00
## Mean   :15.51  Mean   : 4.211  Mean   :0.01471  Mean   : 47.84
## 3rd Qu.:18.00  3rd Qu.: 3.000  3rd Qu.:0.00000  3rd Qu.: 43.25
## Max.   :304.00  Max.   :301.000  Max.   :2.00000  Max.   :2162.00
## loading.time
## Min.   :0.00000
## 1st Qu.:0.08125
## Median :0.27400
## Mean   :0.36143
## 3rd Qu.:0.50850
## Max.   :4.06800

names(total_500_final)

## [1] "Revenues"           "non.document.error" "number_of_errors"
## [4] "number_of_warning"   "facebook"           "instagram"
## [7] "linkedin"            "pinterest"          "twitter"
## [10] "youtube"             "Flesh_Mesaure"     "Readability"
## [13] "Sentences"           "Unique.words"       "Words"
## [16] "external"            "internal"           "total.links"
## [19] "X15x75"              "X8x15"              "X44x556"
## [22] "X1x1"                "X800x1200"          "autox100."
## [25] "X24pxx133px"         "X21pxx173px"        "X46x214"
## [28] "X49x49"               "X50x45"              "X400x300"
## [31] "X292pxx292px"        "X200pxx200px"        "X1279pxx984px"
## [34] "X300pxx1500px"        "X29x29"              "X115x223"
## [37] "X160x233"             "X300x993"            "X41x192"

```

```

## [40] "X28x221"           "X15x12"           "X60x60"
## [43] ".bmp"               ".dib"              ".gif"
## [46] ".jpe"               ".jpeg"             ".jpg"
## [49] ".png"               ".tif"              ".tiff"
## [52] "total.images"        "loading.time"

total_500_final$X15x12<- gsub("1","0", total_500_final$X15x12)
total_500_final$X15x12 <- gsub("2", "1", total_500_final$X15x12 )

total_500_final$X60x60<- gsub("1","0", total_500_final$X60x60)
total_500_final$X60x60 <- gsub("2", "1", total_500_final$X60x60 )

total_500_final$X15x75<- gsub("1","0", total_500_final$X15x75)
total_500_final$X15x75 <- gsub("2", "1", total_500_final$X15x75 )

total_500_final$X28x221<- gsub("1","0", total_500_final$X28x221)
total_500_final$X28x221 <- gsub("2", "1", total_500_final$X28x221 )

total_500_final$X41x192 <- gsub("1","0", total_500_final$X41x192 )
total_500_final$X41x192 <- gsub("2", "1", total_500_final$X41x192 )

total_500_final$X300x993 <- gsub("1","0", total_500_final$X300x993 )
total_500_final$X300x993 <- gsub("2", "1", total_500_final$X300x993 )

total_500_final$X160x233 <- gsub("1","0", total_500_final$X160x233 )
total_500_final$X160x233 <- gsub("2", "1", total_500_final$X160x233 )

total_500_final$X29x29 <- gsub("1","0", total_500_final$X29x29 )
total_500_final$X29x29 <- gsub("2", "1", total_500_final$X29x29 )

total_500_final$X300pxx1500px <- gsub("1","0", total_500_final$X300pxx1500px )
total_500_final$X300pxx1500px <- gsub("2", "1", total_500_final$X300pxx1500px )

total_500_final$X200pxx200px<- gsub("1","0", total_500_final$X200pxx200px )
total_500_final$X200pxx200px <- gsub("2", "1", total_500_final$X200pxx200px )

total_500_final$X292pxx292px <- gsub("1","0", total_500_final$X292pxx292px )
total_500_final$X292pxx292px <- gsub("2", "1", total_500_final$X292pxx292px )

total_500_final$X400x300 <- gsub("1","0", total_500_final$X400x300 )
total_500_final$X400x300 <- gsub("2", "1", total_500_final$X400x300 )

total_500_final$X115x223 <- gsub("1","0", total_500_final$X115x223 )
total_500_final$X115x223 <- gsub("2", "1", total_500_final$X115x223 )

total_500_final$X1279pxx984px <- gsub("1","0", total_500_final$X1279pxx984px )
total_500_final$X1279pxx984px<- gsub("2", "1", total_500_final$X1279pxx984px )

total_500_final$X8x15 <- gsub("1","0", total_500_final$X8x15 )
total_500_final$X8x15 <- gsub("2", "1", total_500_final$X8x15 )

total_500_final$X44x556 <- gsub("1","0", total_500_final$X44x556 )
total_500_final$X44x556 <- gsub("2", "1", total_500_final$X44x556 )

```

```

total_500_final$X1x1 <- gsub("1","0", total_500_final$X1x1 )
total_500_final$X1x1 <- gsub("2", "1", total_500_final$X1x1 )

total_500_final$autox100. <- gsub("1","0", total_500_final$autox100. )
total_500_final$autox100. <- gsub("2", "1", total_500_final$autox100. )
colnames(total_500_final)[24] <- "X100x100"

total_500_final$X800x1200 <- gsub("1","0", total_500_final$X800x1200 )
total_500_final$X800x1200 <- gsub("2", "1", total_500_final$X800x1200 )

total_500_final$X24pxx133px <- gsub("1","0", total_500_final$X24pxx133px )
total_500_final$X24pxx133px <- gsub("2", "1", total_500_final$X24pxx133px )

total_500_final$X21pxx173px <- gsub("1","0", total_500_final$X21pxx173px )
total_500_final$X21pxx173px <- gsub("2", "1", total_500_final$X21pxx173px )

total_500_final$X46x214 <- gsub("1","0", total_500_final$X46x214)
total_500_final$X46x214 <- gsub("2", "1", total_500_final$X46x214 )

total_500_final$X49x49 <- gsub("1","0", total_500_final$X49x49)
total_500_final$X49x49 <- gsub("2", "1", total_500_final$X49x49 )

total_500_final$X50x45 <- gsub("1","0", total_500_final$X50x45)
total_500_final$X50x45 <- gsub("2", "1", total_500_final$X50x45 )

for(i in 19:42){
  total_500_final[,i] <- as.numeric(total_500_final[,i])}

#We split the set to training and test set
library(caret)
set.seed(20)
sampling_vector <- createDataPartition(total_500_final$Revenues, p = 0.70, list = FALSE)
total_500_final_train <- total_500_final[sampling_vector,]
total_500_final_test <- total_500_final[-sampling_vector,]

#We will try to create a regression model to see which of the variables of the websites play the most important role
#We create the empty lm model
model_null = lm(Revenues~1,data=total_500_final_train)
summary(model_null)

## 
## Call:
## lm(formula = Revenues ~ 1, data = total_500_final_train)
## 
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -18.070 -16.119 -12.152 -2.342 210.515 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept)  23.200     2.023   11.47  <2e-16 ***  
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 
## 
## Residual standard error: 34.33 on 287 degrees of freedom

```

```

#####
#LASSO and Logistic Regression models
library(glmnet)

## Warning: package 'glmnet' was built under R version 3.3.3
## Loading required package: Matrix
## Loading required package: foreach
## Loaded glmnet 2.0-5
#We create a full model for the variable Ranking
full <- lm(Revenues~., data=total_500_final_train)
summary(full)

##
## Call:
## lm(formula = Revenues ~ ., data = total_500_final_train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -20.227  -7.486  -3.338   1.303   66.316
##
## Coefficients: (13 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.266e+02  1.618e+01 14.004 < 2e-16 ***
## non.document.error 5.657e-01  2.195e+00  0.258 0.796868
## number_of_errors  1.230e-02  1.113e-02  1.105 0.270122
## number_of_warning 4.159e-03  4.276e-02  0.097 0.922592
## facebook        -4.503e+00  2.796e+00 -1.611 0.108480
## instagram        1.358e+00  2.551e+00  0.532 0.595127
## linkedin         -3.492e-01  2.435e+00 -0.143 0.886095
## pinterest        -1.613e+00  3.598e+00 -0.448 0.654262
## twitter          2.911e+00  2.727e+00  1.068 0.286712
## youtube          5.036e+00  2.221e+00  2.267 0.024250 *
## Flesh_Mesaure   -3.099e-02  3.103e-02 -0.999 0.318996
## Readability      6.973e-01  9.704e-01  0.719 0.473101
## Sentences         2.104e-03  1.216e-02  0.173 0.862755
## Unique.words     -2.273e-02  2.089e-02 -1.088 0.277518
## Words            5.093e-03  5.053e-03  1.008 0.314457
## external         -6.190e-04  1.927e-02 -0.032 0.974397
## internal         -8.858e-04  1.083e-02 -0.082 0.934908
## total.links      NA        NA        NA        NA
## X15x75          -1.762e+01  2.070e+01 -0.851 0.395495
## X8x15           -3.169e+01  2.095e+01 -1.513 0.131677
## X44x556         -1.722e+01  2.104e+01 -0.819 0.413786
## X1x1             NA        NA        NA        NA
## X800x1200        -1.293e+01  1.815e+01 -0.712 0.477033
## X100x100         -3.007e+01  5.554e+01 -0.541 0.588661
## X24pxx133px     2.467e+01  5.588e+01  0.441 0.659272
## X21pxx173px     NA        NA        NA        NA
## X46x214          -9.818e+00  1.819e+01 -0.540 0.589877
## X49x49           NA        NA        NA        NA
## X50x45           -7.650e+00  1.819e+01 -0.421 0.674433
## X400x300         -2.693e+00  2.133e+01 -0.126 0.899623

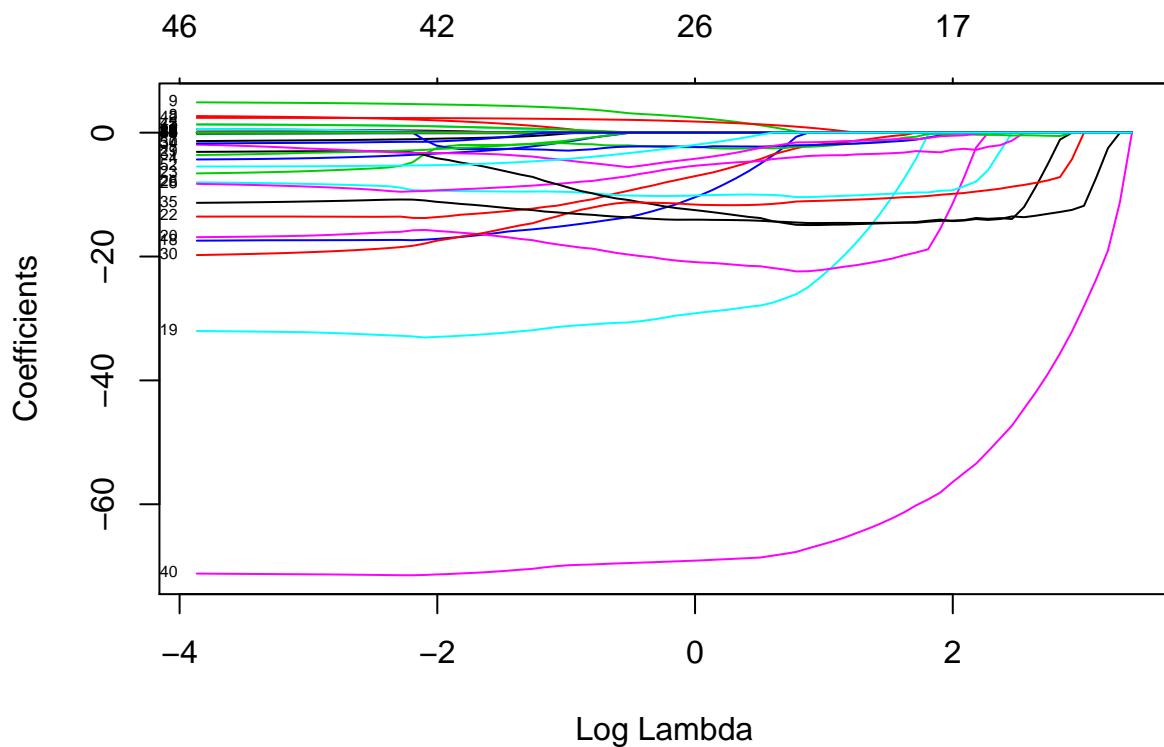
```

```

## X292pxx292px      -2.495e+01  1.887e+01  -1.322 0.187254
## X200pxx200px       NA          NA          NA          NA
## X1279pxx984px      NA          NA          NA          NA
## X300pxx1500px      NA          NA          NA          NA
## X29x29            -1.477e+00  1.823e+01  -0.081 0.935472
## X115x223           -1.150e+01  2.099e+01  -0.548 0.584217
## X160x233             NA          NA          NA          NA
## X300x993             NA          NA          NA          NA
## X41x192              NA          NA          NA          NA
## X28x221              NA          NA          NA          NA
## X15x12            -7.116e+01  1.501e+01  -4.741 3.59e-06 ***
## X60x60               NA          NA          NA          NA
## .bmp                2.396e+00  6.357e-01   3.770 0.000204 ***
## .dib                1.431e+00  1.315e+00   1.088 0.277552
## .gif                -1.133e-01  8.477e-02  -1.337 0.182535
## .jpe                1.133e-01  4.324e+00   0.026 0.979123
## .jpeg               3.710e-01  4.367e+00   0.085 0.932364
## .jpg                4.584e-03  3.071e-02   0.149 0.881438
## .png                -1.891e-02  4.015e-02  -0.471 0.637976
## .tif                3.002e-03  4.911e-02   0.061 0.951302
## .tiff               -1.127e+00  7.240e+00  -0.156 0.876394
## total.images          NA          NA          NA          NA
## loading.time        -5.314e+00  2.578e+00  -2.061 0.040332 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.38 on 248 degrees of freedom
## Multiple R-squared:  0.8484, Adjusted R-squared:  0.8246
## F-statistic: 35.59 on 39 and 248 DF,  p-value: < 2.2e-16
x <- model.matrix(full) [,-1]
dim(x)

## [1] 288  52
lasso <- glmnet (x, total_500_final_train$Revenues)
par(mfrow=c(1,1),no.readonly = TRUE)
plot(lasso, xvar='lambda', label=T)

```



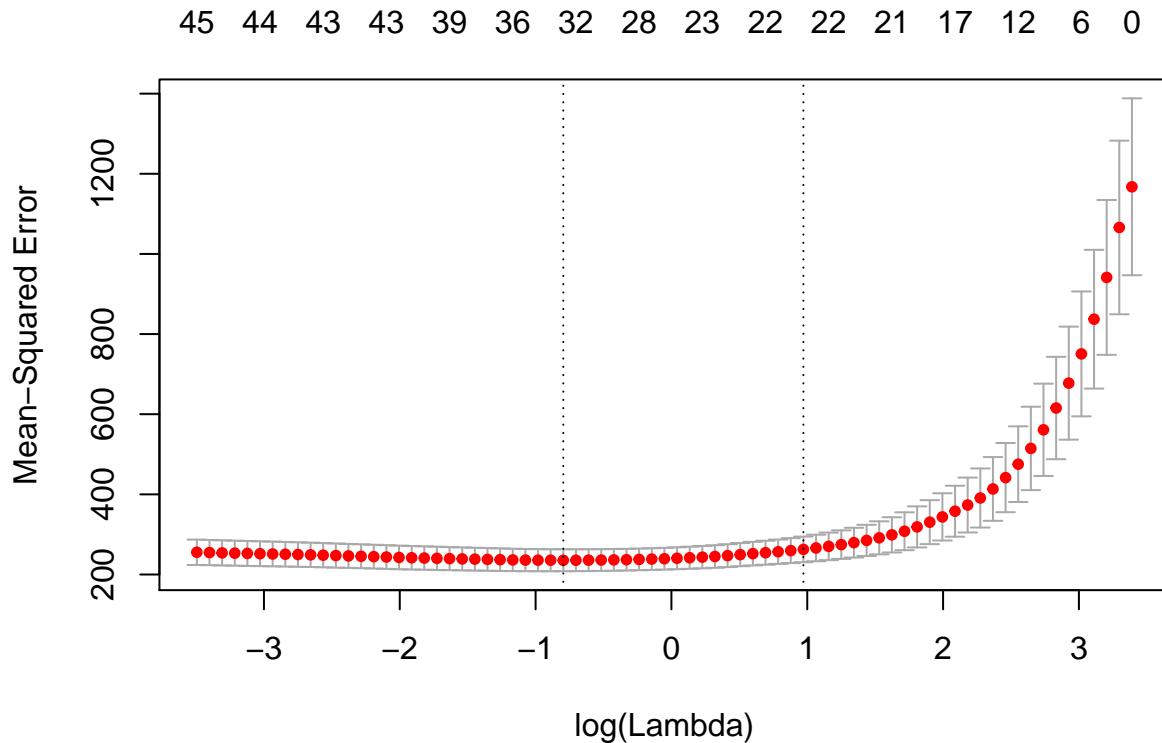
```

lassob <- cv.glmnet(x, total_500_final_train$Revenues)
lassob$lambda.min

## [1] 0.4511413
lassob$lambda.1se

## [1] 2.642344
plot(lassob)

```



```
#We see the coefficients for lambda min
lasso <- coef(lassob, s="lambda.min")
lasso
```

```
## 53 x 1 sparse Matrix of class "dgCMatrix"
##                               1
## (Intercept)      2.230833e+02
## non.document.error .
## number_of_errors 6.149259e-03
## number_of_warning .
## facebook        -8.053246e-01
## instagram       .
## linkedin         .
## pinterest       .
## twitter          9.824382e-02
## youtube          3.711612e+00
## Flesh_Mesaure   -1.720979e-02
## Readability      .
## Sentences         2.566489e-03
## Unique.words     .
## Words             .
## external          .
## internal          .
## total.links      .
## X15x75           -1.448135e+01
## X8x15            -3.095291e+01
## X44x556          -1.873621e+01
```

```

## X1x1          -5.102733e-11
## X800x1200    -1.075697e+01
## X100x100     -1.877745e+00
## X24pxx133px  -2.682910e+00
## X21pxx173px .
## X46x214      -9.778805e+00
## X49x49       .
## X50x45       -7.787668e+00
## X400x300     -9.798054e+00
## X292pxx292px -1.209080e+01
## X200pxx200px -7.655717e-01
## X1279pxx984px -2.020626e-03
## X300pxx1500px -2.665792e-12
## X29x29       -4.932214e+00
## X115x223     -1.326506e+01
## X160x233     -2.505006e-10
## X300x993     -3.322224e-14
## X41x192      -9.966671e-14
## X28x221      -1.162778e-13
## X15x12       -6.972276e+01
## X60x60       -7.618330e-11
## .bmp          2.148840e+00
## .dib          3.222502e-01
## .gif          -7.060778e-02
## .jpe          .
## .jpeg         .
## .jpg          1.022970e-02
## .png          .
## .tif          .
## .tiff         .
## total.images .
## loading.time -3.897970e+00

dim(blasso)

## [1] 53  1

zblasso <- blasso[-1] * apply(x, 2, sd)
zbolt <- coef(full)[-1] * apply(x, 2, sd)
azbolt <- abs(zbolt)
sum(azbolt)

## [1] NA

#since the sum is NA that means we have to subtract some variables
# in order to find which variables to subtract we run the coefficients and we see which of them has NA
coef(full)

##          (Intercept) non.document.error number_of_errors
## 2.265809e+02      5.656880e-01      1.230253e-02
## number_of_warning      facebook      instagram
## 4.158870e-03     -4.503280e+00      1.357670e+00
## linkedin            pinterest      twitter
## -3.491824e-01     -1.613471e+00      2.911065e+00
## youtube             Flesh_Mesaure Readability
## 5.036021e+00     -3.098642e-02      6.972808e-01
## Sentences           Unique.words      Words

```

```

##      2.104257e-03      -2.273419e-02      5.093208e-03
##      external           internal          total.links
##      -6.190402e-04      -8.857651e-04          NA
##      X15x75              X8x15            X44x556
##      -1.761698e+01      -3.169113e+01      -1.722275e+01
##      X1x1                X800x1200        X100x100
##      NA                  -1.292930e+01      -3.007324e+01
##      X24pxx133px        X21pxx173px        X46x214
##      2.466914e+01          NA          -9.817905e+00
##      X49x49              X50x45          X400x300
##      NA                  -7.650074e+00      -2.692975e+00
##      X292pxx292px        X200pxx200px        X1279pxx984px
##      -2.494753e+01          NA          NA
##      X300pxx1500px        X29x29        X115x223
##      NA                  -1.477068e+00      -1.150271e+01
##      X160x233              X300x993        X41x192
##      NA                  NA          NA
##      X28x221              X15x12          X60x60
##      NA                  -7.115613e+01          NA
##      .bmp                 .dib          .gif
##      2.396454e+00          1.430784e+00      -1.133168e-01
##      .jpe                 .jpeg          .jpg
##      1.132600e-01          3.710030e-01      4.584279e-03
##      .png                 .tif          .tiff
##      -1.891363e-02          3.002462e-03      -1.127253e+00
##      total.images          loading.time
##      NA                  -5.314492e+00

```

#Now we create a new model with only the variables with coef different from NA

```

full_2 <- lm(Revenues~. - total.images - total.links - X1x1 - X21pxx173px - X49x49 - X200pxx200px - X1279pxx984px - X300pxx1500px -
summary(full_2)

```

```

##
## Call:
## lm(formula = Revenues ~ . - total.images - total.links - X1x1 -
##     X21pxx173px - X49x49 - X200pxx200px - X1279pxx984px - X300pxx1500px -
##     X160x233 - X300x993 - X41x192 - X28x221 - X60x60, data = total_500_final_train)
##
## Residuals:
##    Min      1Q  Median      3Q      Max
## -20.227  -7.486  -3.338   1.303   66.316
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)             2.266e+02  1.618e+01 14.004 < 2e-16 ***
## non.document.error  5.657e-01  2.195e+00  0.258 0.796868
## number_of_errors   1.230e-02  1.113e-02  1.105 0.270122
## number_of_warning  4.159e-03  4.276e-02  0.097 0.922592
## facebook            -4.503e+00  2.796e+00 -1.611 0.108480
## instagram           1.358e+00  2.551e+00  0.532 0.595127
## linkedin            -3.492e-01  2.435e+00 -0.143 0.886095
## pinterest          -1.613e+00  3.598e+00 -0.448 0.654262
## twitter             2.911e+00  2.727e+00  1.068 0.286712
## youtube             5.036e+00  2.221e+00  2.267 0.024250 *
## Flesh_Mesaure      -3.099e-02  3.103e-02 -0.999 0.318996

```

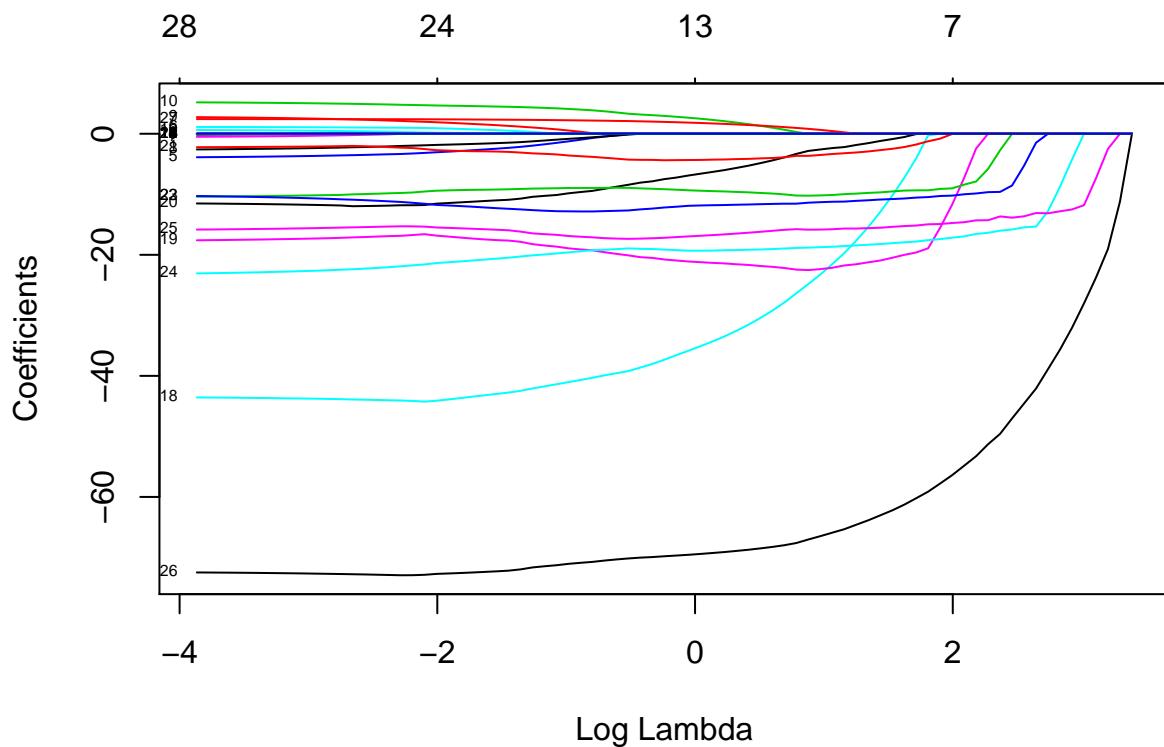
```

## Readability      6.973e-01 9.704e-01 0.719 0.473101
## Sentences        2.104e-03 1.216e-02 0.173 0.862755
## Unique.words    -2.273e-02 2.089e-02 -1.088 0.277518
## Words           5.093e-03 5.053e-03 1.008 0.314457
## external        -6.190e-04 1.927e-02 -0.032 0.974397
## internal        -8.858e-04 1.083e-02 -0.082 0.934908
## X15x75          -1.762e+01 2.070e+01 -0.851 0.395495
## X8x15            -3.169e+01 2.095e+01 -1.513 0.131677
## X44x556          -1.722e+01 2.104e+01 -0.819 0.413786
## X8000x1200       -1.293e+01 1.815e+01 -0.712 0.477033
## X100x100         -3.007e+01 5.554e+01 -0.541 0.588661
## X24pxx133px      2.467e+01 5.588e+01 0.441 0.659272
## X46x214          -9.818e+00 1.819e+01 -0.540 0.589877
## X50x45            -7.650e+00 1.819e+01 -0.421 0.674433
## X400x300          -2.693e+00 2.133e+01 -0.126 0.899623
## X292pxx292px     -2.495e+01 1.887e+01 -1.322 0.187254
## X29x29            -1.477e+00 1.823e+01 -0.081 0.935472
## X115x223          -1.150e+01 2.099e+01 -0.548 0.584217
## X15x12             -7.116e+01 1.501e+01 -4.741 3.59e-06 ***
## .bmp               2.396e+00 6.357e-01 3.770 0.000204 ***
## .dib               1.431e+00 1.315e+00 1.088 0.277552
## .gif               -1.133e-01 8.477e-02 -1.337 0.182535
## .jpe               1.133e-01 4.324e+00 0.026 0.979123
## .jpeg              3.710e-01 4.367e+00 0.085 0.932364
## .jpg               4.584e-03 3.071e-02 0.149 0.881438
## .png               -1.891e-02 4.015e-02 -0.471 0.637976
## .tif               3.002e-03 4.911e-02 0.061 0.951302
## .tiff              -1.127e+00 7.240e+00 -0.156 0.876394
## loading.time       -5.314e+00 2.578e+00 -2.061 0.040332 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.38 on 248 degrees of freedom
## Multiple R-squared:  0.8484, Adjusted R-squared:  0.8246
## F-statistic: 35.59 on 39 and 248 DF,  p-value: < 2.2e-16
x <- model.matrix(full_2) [,-c(18,22,28,26,34,32,33,42,37,38,39,40,52)]
dim(x)

## [1] 288 29
lasso <- glmnet (x, total_500_final_train$Revenues)

plot(lasso, xvar='lambda', label=T)

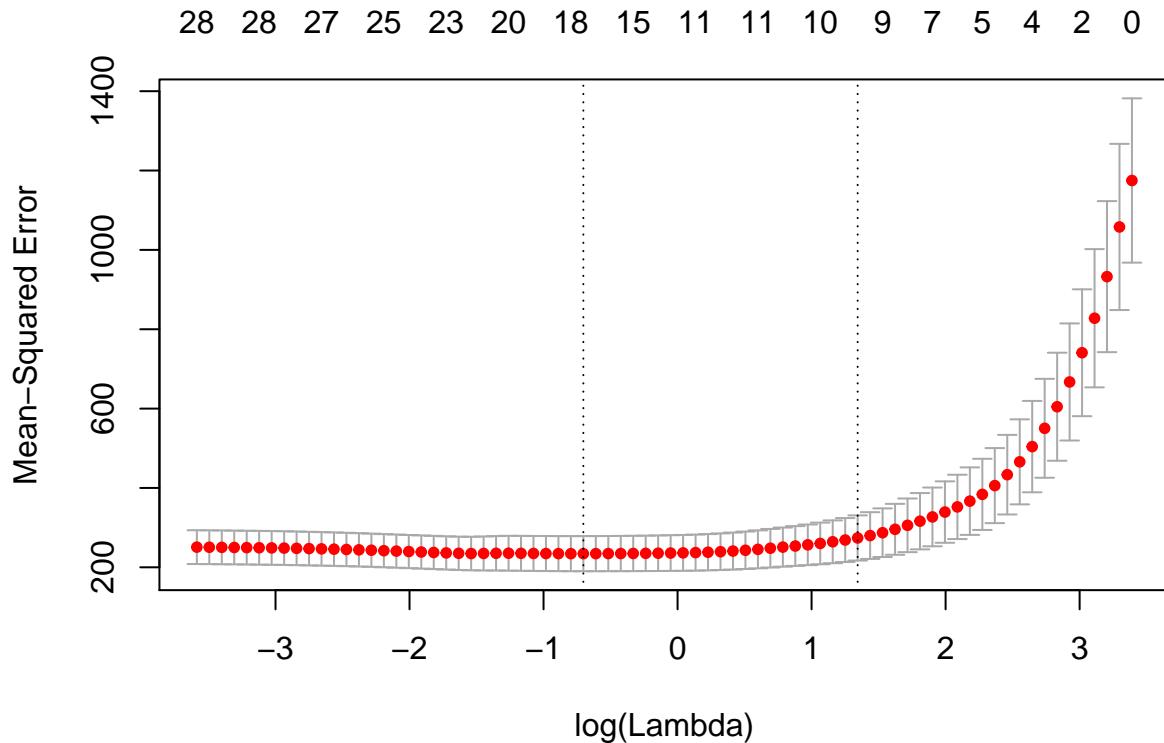
```



```
lassob <- cv.glmnet(x, total_500_final_train$Revenues)
lassob$lambda.min

## [1] 0.495127
lassob$lambda.1se

## [1] 3.833588
plot(lassob)
```



```

#coefiecinets for lammda min
lasso <- coef(lassob, s="lambda.min")
lasso

## 30 x 1 sparse Matrix of class "dgCMatrix"
##                               1
## (Intercept)      215.050596035
## (Intercept)      .
## non.document.error   .
## number_of_errors    0.004545402
## number_of_warning   .
## facebook        -0.465821955
## instagram       .
## linkedin        .
## pinterest       -0.442067577
## twitter         .
## youtube          3.714394348
## Flesh_Mesaure   -0.012455651
## Readability      .
## Sentences         0.001245264
## Unique.words     .
## Words            .
## external          .
## internal          .
## X8x15           -39.910397889
## X44x556          -19.570357469
## X800x1200         -9.091373726

```

```

## X24pxx133px      -3.964585608
## X46x214          -8.939435748
## X50x45           -12.779830022
## X292pxx292px     -19.158990860
## X115x223         -17.248924356
## X15x12            -70.512070793
## .bmp              2.152463769
## .jpeg              .
## .jpg               0.008247367

dim(blasso)

## [1] 30  1

zblasso <- blasso[-1] * apply(x, 2, sd)
zbolt <- coef(full_2)[-1] * apply(x, 2, sd)

## Warning in coef(full_2)[-1] * apply(x, 2, sd): longer object length is not
## a multiple of shorter object length
azbolt <- abs(zbolt)
sum(azbolt)

## [1] 5546.237

s <- sum(abs(zblasso))/sum(abs(azbolt))
s

## [1] 0.007532779

full_3 <- lm(Revenues~1 +number_of_errors +facebook +pinterest +youtube+ Flesh_Mesaure +Sentences +X8x15
summary(full_3)

##
## Call:
## lm(formula = Revenues ~ 1 + number_of_errors + facebook + pinterest +
##      youtube + Flesh_Mesaure + Sentences + X8x15 + X44x556 + X800x1200 +
##      X24pxx133px + X46x214 + X50x45 + X292pxx292px + X115x223 +
##      X15x12 + .bmp + .jpg, data = total_500_final_train)
##
## Residuals:
##       Min     1Q   Median     3Q    Max 
## -14.444 -7.482 -3.728  1.968 67.468 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 220.708931  10.142389  21.761 < 2e-16 ***
## number_of_errors  0.009297  0.009860  0.943 0.346575    
## facebook      -2.590609  2.142461  -1.209 0.227654    
## pinterest     -1.966700  3.155150  -0.623 0.533594    
## youtube        5.721687  2.009439  2.847 0.004746 **  
## Flesh_Mesaure -0.028894  0.024264  -1.191 0.234779    
## Sentences       0.003712  0.005628  0.660 0.510092    
## X8x15          -43.592759 17.439657  -2.500 0.013026 *  
## X44x556        -16.356982 20.137030  -0.812 0.417344    
## X800x1200      -12.902443 16.481121  -0.783 0.434394    
## X24pxx133px     -2.075224 16.689391  -0.124 0.901135    
## X46x214        -9.506257 17.481911  -0.544 0.587044

```

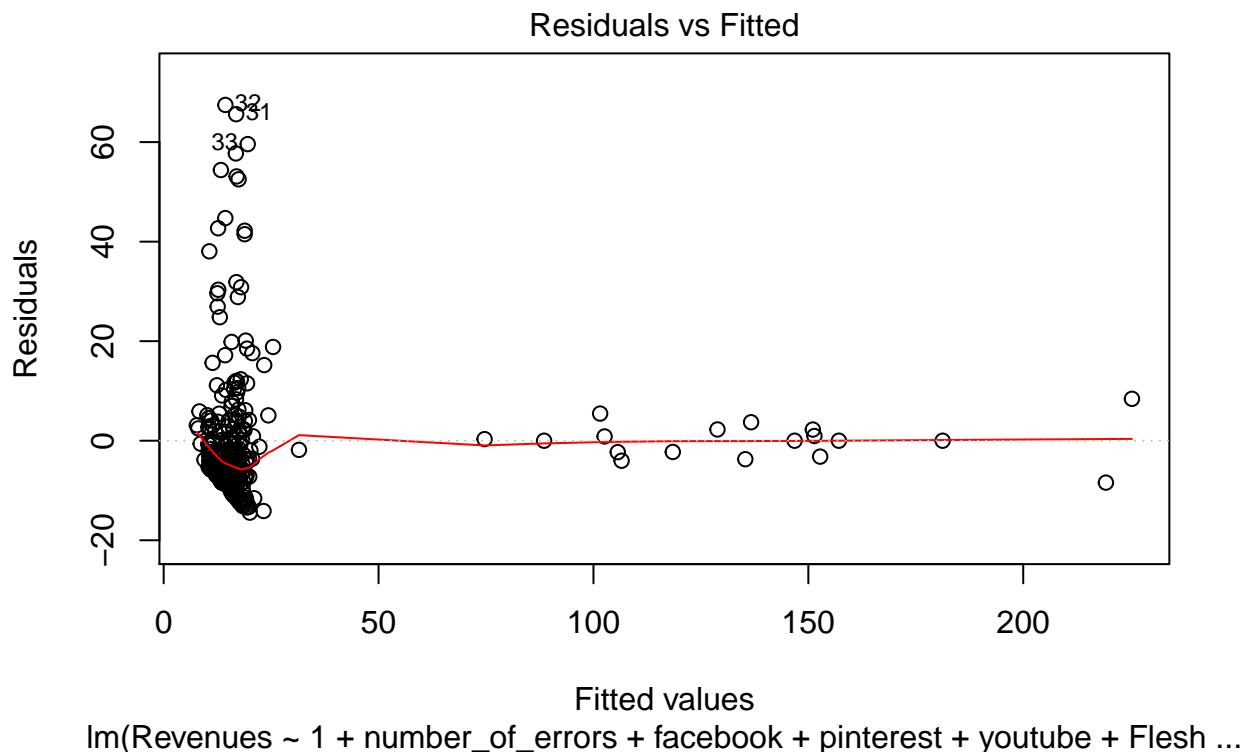
```

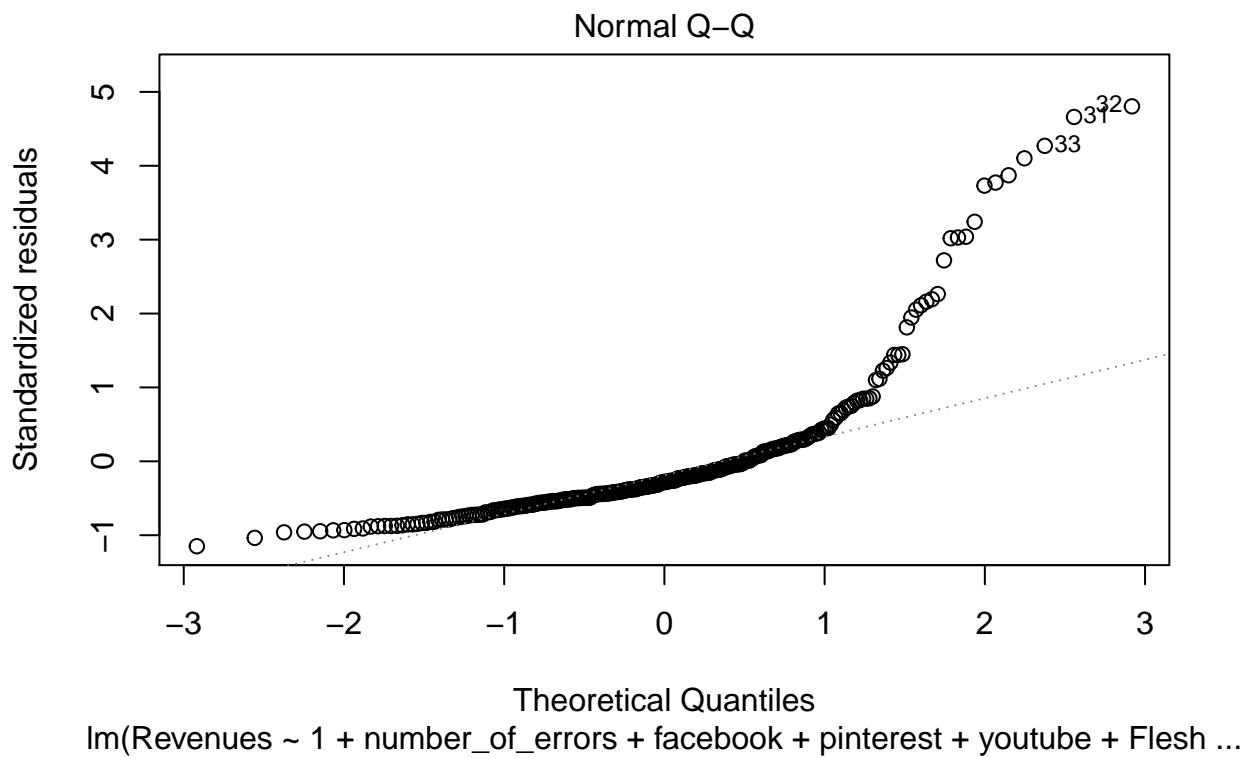
## X50x45      -13.137735 14.219494 -0.924 0.356351
## X292pxx292px -20.460731 12.440363 -1.645 0.101194
## X115x223    -15.298995 16.102120 -0.950 0.342899
## X15x12       -73.307867 14.451593 -5.073 7.3e-07 ***
## .bmp          2.435106  0.619522  3.931 0.000108 ***
## .jpg          0.015724  0.020551  0.765 0.444875
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.14 on 270 degrees of freedom
## Multiple R-squared:  0.8405, Adjusted R-squared:  0.8304
## F-statistic: 83.67 on 17 and 270 DF,  p-value: < 2.2e-16
ad_r_sq_f3 <- summary(full_3)$adj.r.squared
aic_f3 <- AIC(full_3)

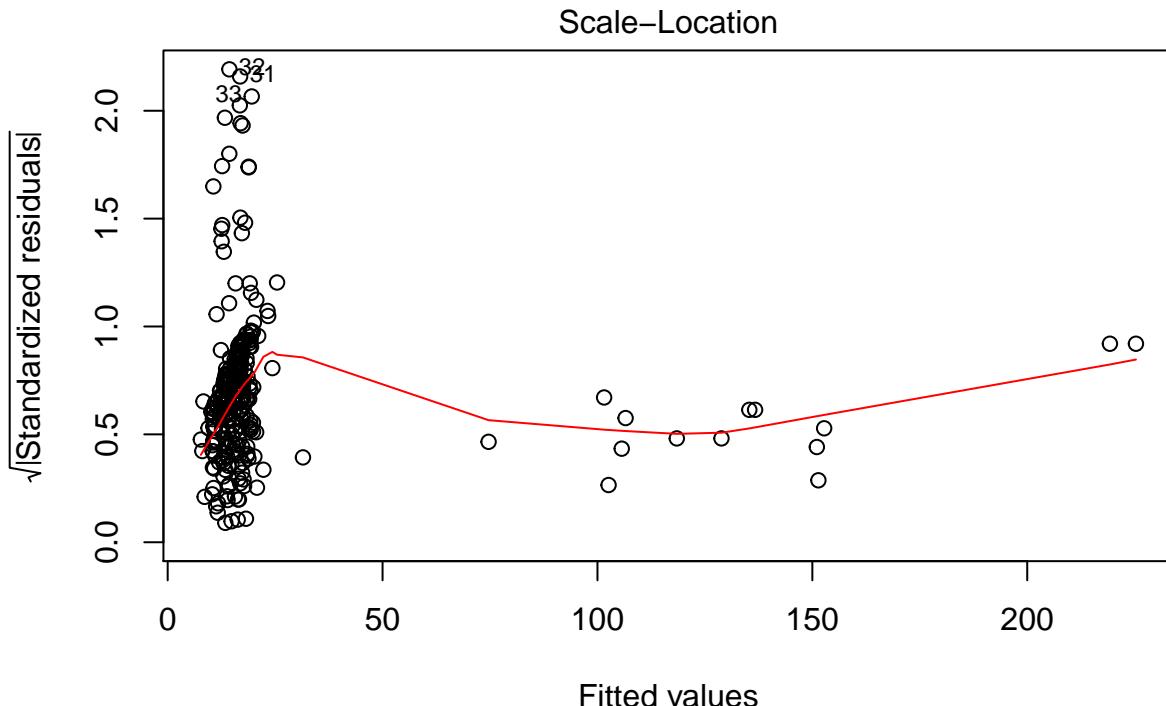
plot(full_3,which=1:3)

```

Warning: not plotting observations with leverage one:
3, 4, 8, 17







lm(Revenues ~ 1 + number_of_errors + facebook + pinterest + youtube + Flesh ...

```
#####
blassob <- coef(lassob, s="lambda.1se")
blassob
```

```
## 30 x 1 sparse Matrix of class "dgCMatrix"
##                                     1
## (Intercept)      176.417391
## (Intercept)      .
## non.document.error   .
## number_of_errors   .
## number_of_warning   .
## facebook   .
## instagram   .
## linkedin   .
## pinterest   .
## twitter   .
## youtube   .
## Flesh_Mesaure   .
## Readability   .
## Sentences   .
## Unique.words   .
## Words   .
## external   .
## internal   .
## X8x15      -15.418145
## X44x556     -21.240420
## X800x1200    -1.553671
```

```

## X24pxx133px      -2.926734
## X46x214          -9.772246
## X50x45           -10.987677
## X292pxx292px     -18.334032
## X115x223         -15.592779
## X15x12            -63.865600
## .bmp              .
## .jpeg              .
## .jpg               .

zblassob <- blassob[-1] * apply(x, 2, sd)
zboltb <- coef(full_2)[-1] * apply(x, 2, sd)

## Warning in coef(full_2)[-1] * apply(x, 2, sd): longer object length is not
## a multiple of shorter object length
s <- sum(abs(zblassob))/sum(abs(zboltb))
s

## [1] 0.005413772

#The model based on the lasso method by taking the lambda.1se is the null model only with the intercept

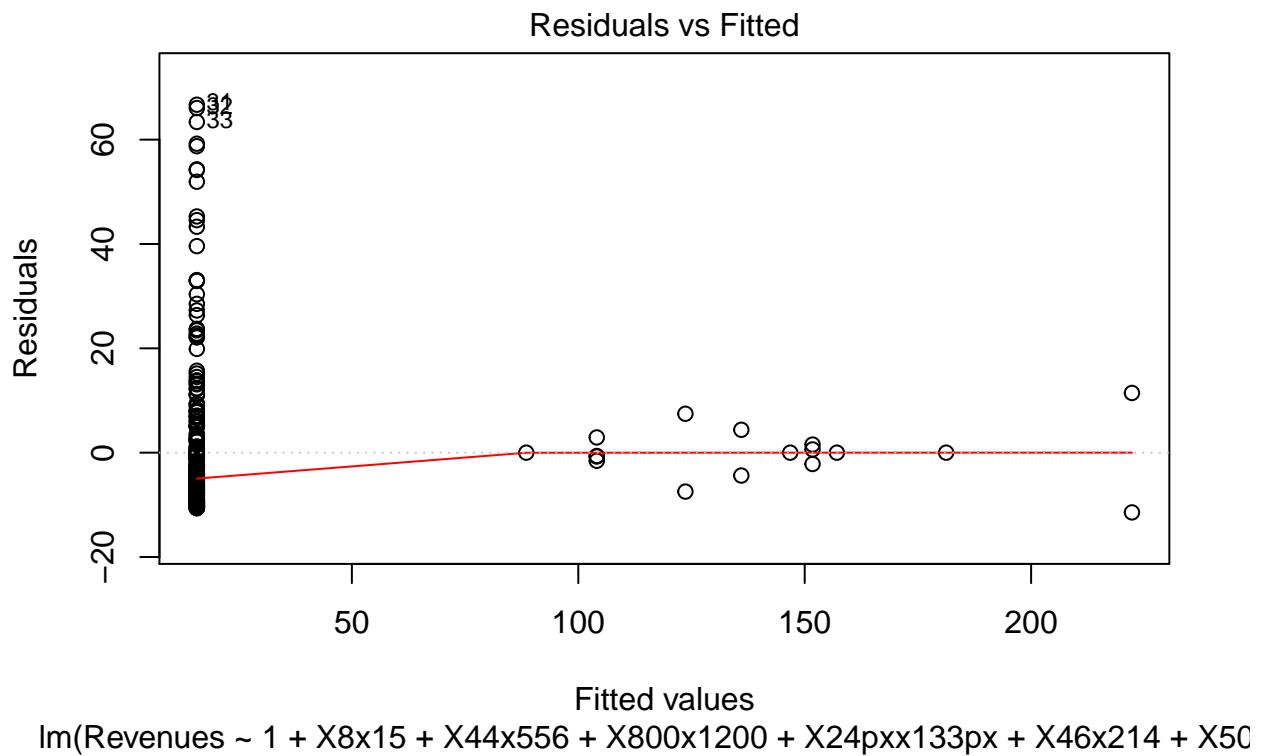
full_4 <- lm(Revenues~1 +X8x15 +X44x556 +X800x1200 +X24pxx133px +X46x214 +X50x45 +X292pxx292px +X115x223 +X15x12, data = total_500_final_train)
summary(full_4)

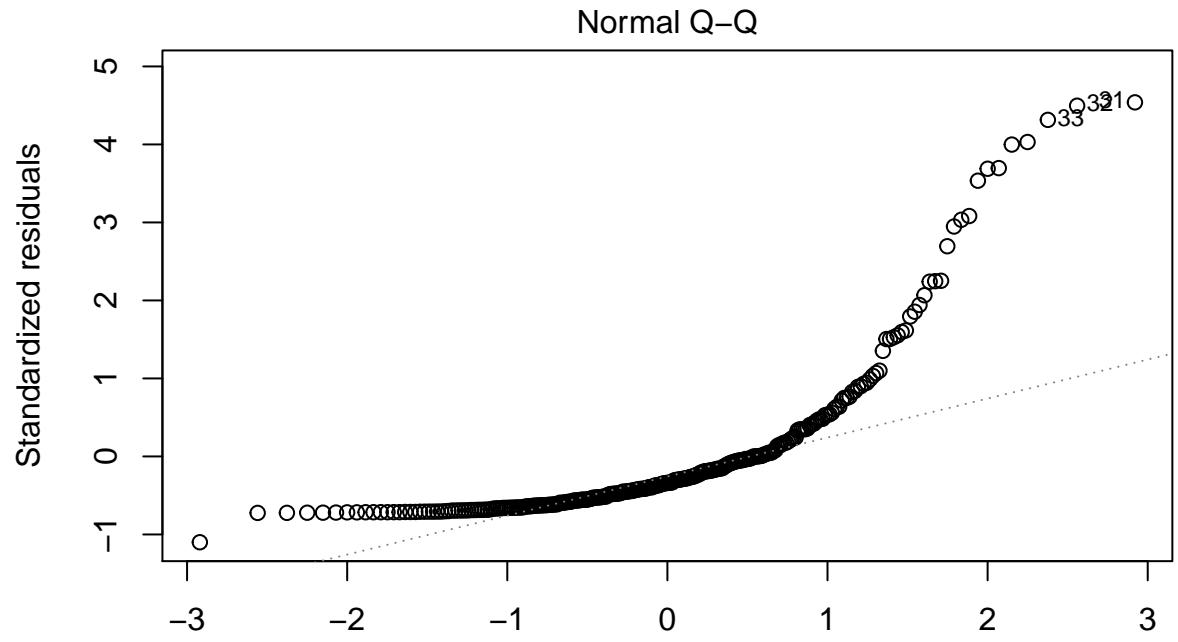
##
## Call:
## lm(formula = Revenues ~ 1 + X8x15 + X44x556 + X800x1200 + X24pxx133px +
##       X46x214 + X50x45 + X292pxx292px + X115x223 + X15x12, data = total_500_final_train)
##
## Residuals:
##    Min      1Q  Median      3Q      Max
## -11.447  -8.735  -4.926   1.100  66.695
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 222.268    10.408  21.356 < 2e-16 ***
## X8x15      -41.027    18.027  -2.276  0.0236 *
## X44x556     -24.134    20.816  -1.159  0.2473
## X800x1200    -5.372    16.996  -0.316  0.7522
## X24pxx133px   -4.934    16.996  -0.290  0.7718
## X46x214      -10.796   18.027  -0.599  0.5497
## X50x45       -12.346   14.719  -0.839  0.4023
## X292pxx292px  -19.575   12.747  -1.536  0.1258
## X115x223     -15.565   16.456  -0.946  0.3450
## X15x12        -72.753   14.746  -4.934  1.39e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14.72 on 278 degrees of freedom
## Multiple R-squared:  0.822, Adjusted R-squared:  0.8162
## F-statistic: 142.6 on 9 and 278 DF,  p-value: < 2.2e-16
ad_r_sq_f4 <- summary(full_4)$adj.r.squared
aic_f4 <- AIC(full_4)

```

```
plot(full_4,which=1:3)

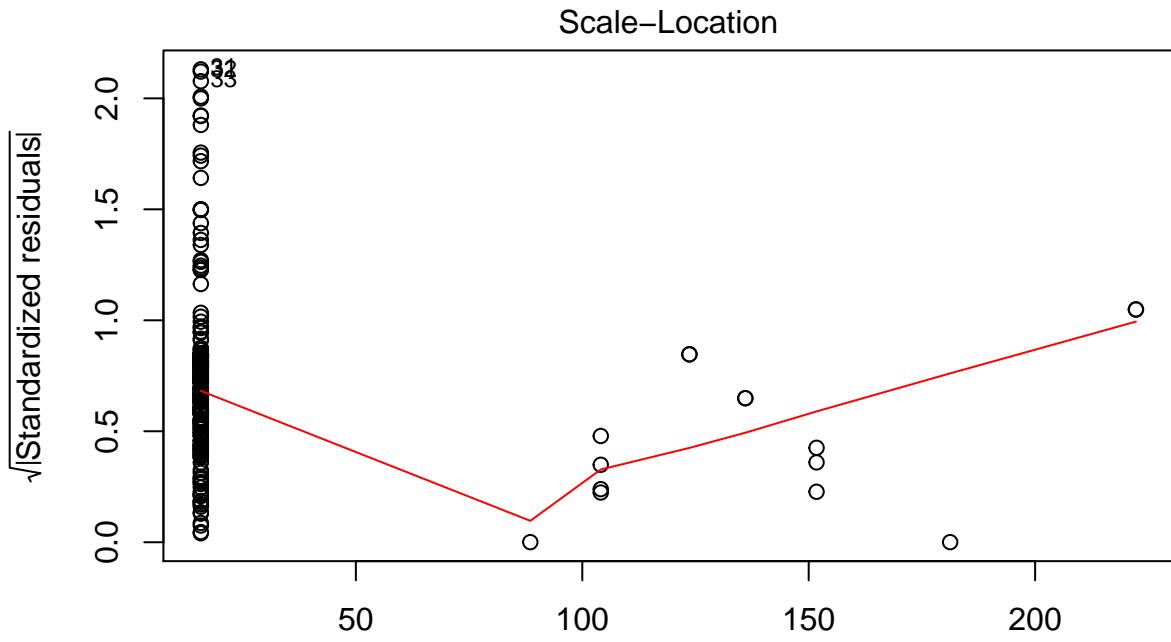
## Warning: not plotting observations with leverage one:
##      4, 8
```





Theoretical Quantiles

Im(Revenues ~ 1 + X8x15 + X44x556 + X800x1200 + X24pxx133px + X46x214 + X50



Im(Revenues ~ 1 + X8x15 + X44x556 + X800x1200 + X24pxx133px + X46x214 + X50

```
#####
#####
```

We use the "both" method to compare the full_3 model with the null model to see how many variables are
model_a <- step(model_null, scope = list(lower = model_null, upper=full_2), direction = "both")

```
## Start: AIC=2037.77
## Revenues ~ 1
##
##          Df Sum of Sq    RSS    AIC
## + X15x12      1   253735  84542 1640.4
## + X115x223    1   251468  86810 1648.0
## + X29x29      1   245941  92337 1665.8
## + X292pxx292px  1   229850 108428 1712.1
## + X400x300    1   222075 116203 1732.0
## + X50x45      1   210336 127942 1759.8
## + X46x214      1   184923 153355 1811.9
## + X24pxx133px  1   169418 168860 1839.7
## + X100x100     1   153400 184878 1865.8
## + X800x1200    1   120731 217547 1912.6
## + X44x556      1   104196 234082 1933.7
## + X8x15        1    79810 258468 1962.3
## + X15x75        1   44471 293807 1999.2
## + .jpe         1   14156 324122 2027.5
## + .jpeg        1   13535 324743 2028.0
## + .dib         1    7274 331004 2033.5
## + loading.time  1    5444 332834 2035.1
## + number_of_warning 1   4240 334037 2036.1
```

```

## + non.document.error 1      3575 334703 2036.7
## + youtube             1      3338 334940 2036.9
## + .bmp                1      2634 335644 2037.5
## <none>                 338278 2037.8
## + number_of_errors    1      1708 336570 2038.3
## + pinterest           1      1505 336773 2038.5
## + internal             1      907 337371 2039.0
## + .gif                 1      666 337612 2039.2
## + instagram            1      646 337632 2039.2
## + linkedin             1      570 337708 2039.3
## + Readability          1      563 337715 2039.3
## + Unique.words         1      552 337726 2039.3
## + Words                1      201 338077 2039.6
## + Sentences             1      189 338089 2039.6
## + external              1      131 338147 2039.7
## + .png                 1      119 338159 2039.7
## + .tiff                1      107 338171 2039.7
## + twitter               1      106 338172 2039.7
## + Flesh_Mesaure        1      54 338224 2039.7
## + .tif                 1      9 338269 2039.8
## + facebook              1      7 338271 2039.8
## + .jpg                 1      4 338274 2039.8
##
## Step: AIC=1640.43
## Revenues ~ X15x12
##
##                                     Df Sum of Sq   RSS   AIC
## + X44x556                      1   16295 68247 1580.8
## + X800x1200                     1   15259 69283 1585.1
## + X46x214                      1   14900 69642 1586.6
## + X24pxx133px                   1   14711 69832 1587.4
## + X8x15                         1   14710 69833 1587.4
## + X100x100                      1   14623 69919 1587.7
## + X50x45                        1   13965 70577 1590.4
## + X400x300                      1   13532 71011 1592.2
## + X292pxx292px                  1   11756 72787 1599.3
## + X15x75                        1   8994 75548 1610.0
## + X29x29                        1   4835 79708 1625.5
## + .bmp                           1   3641 80901 1629.8
## + X115x223                      1   3006 81536 1632.0
## + youtube                        1   1650 82892 1636.8
## + loading.time                   1   908 83635 1639.3
## + .dib                           1   710 83832 1640.0
## + .gif                           1   599 83943 1640.4
## <none>                          84542 1640.4
## + pinterest                     1   513 84030 1640.7
## + Readability                    1   400 84142 1641.1
## + instagram                      1   372 84171 1641.2
## + linkedin                       1   332 84211 1641.3
## + .jpg                           1   257 84285 1641.5
## + Flesh_Mesaure                  1   243 84299 1641.6
## + number_of_errors                1   219 84324 1641.7
## + number_of_warning               1   174 84368 1641.8
## + .jpe                           1   158 84384 1641.9

```

```

## + .jpeg          1      152  84390 1641.9
## + external      1      87   84455 1642.1
## + facebook      1      76   84466 1642.2
## + non.document.error 1      41   84502 1642.3
## + .png           1      31   84511 1642.3
## + twitter        1      28   84514 1642.3
## + Sentences      1      18   84525 1642.4
## + Words          1      12   84530 1642.4
## + .tif           1      11   84531 1642.4
## + Unique.words   1      2    84541 1642.4
## + internal       1      1    84542 1642.4
## + .tiff          1      0    84542 1642.4
## - X15x12         1    253735 338278 2037.8
##
## Step: AIC=1580.76
## Revenues ~ X15x12 + X44x556
##
##                                     Df Sum of Sq   RSS   AIC
## + X400x300                  1      5982 62265 1556.3
## + X50x45                     1      5558 62689 1558.3
## + X292pxx292px               1      5424 62823 1558.9
## + X46x214                   1      4647 63600 1562.5
## + X24pxx133px                1      3699 64548 1566.7
## + .bmp                       1      3641 64606 1567.0
## + X100x100                  1      2757 65490 1570.9
## + X29x29                     1      2367 65880 1572.6
## + X115x223                  1      1626 66621 1575.8
## + youtube                    1      1604 66643 1575.9
## + X8x15                      1      1122 67125 1578.0
## + X800x1200                 1      952  67295 1578.7
## + loading.time                1      950  67297 1578.7
## + X15x75                     1      947  67300 1578.7
## + .jpe                       1      730  67517 1579.7
## + .jpeg                      1      691  67556 1579.8
## <none>                      68247 1580.8
## + .gif                       1      414  67833 1581.0
## + .jpg                       1      394  67853 1581.1
## + number_of_errors            1      348  67899 1581.3
## + twitter                     1      286  67961 1581.5
## + linkedin                   1      254  67993 1581.7
## + instagram                  1      252  67995 1581.7
## + external                    1      189  68058 1582.0
## + Flesh_Mesaure              1      186  68061 1582.0
## + number_of_warning           1      168  68079 1582.0
## + Readability                 1      162  68085 1582.1
## + .dib                        1      144  68103 1582.2
## + pinterest                  1      121  68126 1582.2
## + Sentences                   1      99   68148 1582.3
## + Words                       1      37   68210 1582.6
## + .tif                        1      37   68210 1582.6
## + non.document.error          1      19   68228 1582.7
## + internal                    1      10   68237 1582.7
## + .png                        1      8    68239 1582.7
## + facebook                    1      7    68240 1582.7

```

```

## + Unique.words      1      2  68245 1582.8
## + .tiff            1      0  68247 1582.8
## - X44x556          1    16295  84542 1640.4
## - X15x12           1   165835 234082 1933.7
##
## Step: AIC=1556.34
## Revenues ~ X15x12 + X44x556 + X400x300
##
##                                     Df Sum of Sq    RSS     AIC
## + .bmp                  1    3641  58624 1541.0
## + youtube              1    1284  60981 1552.3
## + X8x15                1    1122  61143 1553.1
## + X15x75               1     947  61318 1553.9
## + loading.time         1     856  61409 1554.4
## + X46x214              1     571  61694 1555.7
## + X24pxx133px         1     486  61778 1556.1
## + .jpg                 1     459  61806 1556.2
## <none>                  62265 1556.3
## + .gif                 1     416  61849 1556.4
## + X100x100             1     386  61879 1556.5
## + X115x223             1     270  61995 1557.1
## + number_of_errors     1     244  62021 1557.2
## + X50x45               1     229  62036 1557.3
## + X292pxx292px        1     193  62072 1557.5
## + X29x29               1     191  62074 1557.5
## + twitter              1     172  62093 1557.5
## + X800x1200            1     160  62105 1557.6
## + Flesh_Mesaure        1     136  62129 1557.7
## + Sentences             1     115  62150 1557.8
## + Readability           1     112  62153 1557.8
## + .jpe                 1     100  62165 1557.9
## + .dib                 1     100  62165 1557.9
## + .jpeg                1      97  62168 1557.9
## + external              1      97  62168 1557.9
## + linkedin              1      79  62186 1558.0
## + pinterest            1      74  62191 1558.0
## + .tif                 1      68  62197 1558.0
## + instagram             1      67  62198 1558.0
## + Words                 1      32  62233 1558.2
## + internal              1      17  62247 1558.3
## + Unique.words          1       6  62259 1558.3
## + number_of_warning     1       2  62263 1558.3
## + facebook              1       1  62264 1558.3
## + non.document.error   1       1  62264 1558.3
## + .tiff                 1       0  62265 1558.3
## + .png                  1       0  62265 1558.3
## - X400x300              1    5982  68247 1580.8
## - X44x556               1    8746  71011 1592.2
## - X15x12                1   45193 107457 1711.5
##
## Step: AIC=1540.99
## Revenues ~ X15x12 + X44x556 + X400x300 + .bmp
##
##                                     Df Sum of Sq    RSS     AIC

```

```

## + X8x15          1    1122  57501 1537.4
## + youtube       1    1035  57588 1537.9
## + X15x75         1     947  57677 1538.3
## + loading.time   1     742  57882 1539.3
## + X46x214        1     571  58053 1540.2
## + X24pxx133px   1     486  58137 1540.6
## + .jpg           1     429  58195 1540.9
## <none>           58624 1541.0
## + .gif           1     401  58223 1541.0
## + X100x100       1     386  58237 1541.1
## + twitter         1     321  58302 1541.4
## + X115x223       1     270  58354 1541.7
## + number_of_errors 1     232  58391 1541.8
## + X50x45          1     229  58394 1541.9
## + Sentences        1     209  58414 1542.0
## + X292pxx292px   1     193  58430 1542.0
## + X29x29          1     191  58432 1542.0
## + linkedin         1     167  58457 1542.2
## + X800x1200       1     160  58464 1542.2
## + Flesh_Mesaure   1     141  58482 1542.3
## + .jpe            1     106  58517 1542.5
## + .jpeg           1     102  58521 1542.5
## + .dib            1      90  58534 1542.5
## + instagram        1      87  58536 1542.6
## + Words            1      80  58544 1542.6
## + Readability      1      79  58545 1542.6
## + pinterest        1      61  58562 1542.7
## + external          1      57  58567 1542.7
## + internal          1      52  58571 1542.7
## + .tif             1      26  58597 1542.9
## + facebook          1      11  58613 1542.9
## + number_of_warning 1       3  58620 1543.0
## + Unique.words      1       1  58623 1543.0
## + non.document.error 1       1  58623 1543.0
## + .png              1       0  58623 1543.0
## + .tiff             1       0  58624 1543.0
## - .bmp              1    3641  62265 1556.3
## - X400x300          1    5982  64606 1567.0
## - X44x556           1    8746  67369 1579.0
## - X15x12            1   45466 104089 1704.3
##
## Step: AIC=1537.42
## Revenues ~ X15x12 + X44x556 + X400x300 + .bmp + X8x15
##
##                               Df Sum of Sq    RSS    AIC
## + youtube                   1   1146  56355 1533.6
## + loading.time               1    601  56900 1536.4
## + X46x214                   1    571  56931 1536.5
## + X24pxx133px               1    486  57015 1537.0
## + twitter                    1    464  57037 1537.1
## + .jpg                       1    438  57063 1537.2
## + .gif                       1    407  57094 1537.4
## <none>                      57501 1537.4
## + X100x100                  1    386  57115 1537.5

```

```

## + number_of_errors      1    273  57229 1538.0
## + X115x223             1    270  57232 1538.1
## + X15x75                1    262  57239 1538.1
## + Sentences              1    258  57243 1538.1
## + X50x45                1    229  57272 1538.3
## + linkedin               1    212  57289 1538.4
## + X292pxx292px          1    193  57308 1538.5
## + X29x29                 1    191  57310 1538.5
## + Flesh_Mesaure          1    180  57322 1538.5
## + instagram              1    179  57323 1538.5
## + X800x1200              1    160  57342 1538.6
## + Readability             1    129  57372 1538.8
## + .jpe                    1    106  57395 1538.9
## + .jpeg                   1    102  57399 1538.9
## + Words                   1    101  57400 1538.9
## + .dib                    1     90  57412 1539.0
## + pinterest              1     61  57440 1539.1
## + internal                1     58  57443 1539.1
## + external                1     57  57444 1539.1
## + facebook                1     46  57456 1539.2
## + number_of_warning        1     38  57464 1539.2
## + .tif                     1     26  57475 1539.3
## + Unique.words             1     11  57490 1539.4
## + non.document.error       1      1  57501 1539.4
## + .png                     1      0  57501 1539.4
## + .tiff                    1      0  57501 1539.4
## - X8x15                   1    1122 58624 1541.0
## - X44x556                  1    1150 58651 1541.1
## - .bmp                      1    3641 61143 1553.1
## - X400x300                  1    5982 63484 1563.9
## - X15x12                   1    45466 102967 1703.2
##
## Step: AIC=1533.62
## Revenues ~ X15x12 + X44x556 + X400x300 + .bmp + X8x15 + youtube
##
##                               Df Sum of Sq   RSS   AIC
## + loading.time            1    617  55738 1532.5
## + X46x214                  1    536  55819 1532.9
## + X24pxx133px              1    486  55869 1533.1
## + X100x100                 1    416  55939 1533.5
## <none>                      56355 1533.6
## + .gif                      1    389  55966 1533.6
## + X115x223                 1    349  56006 1533.8
## + X29x29                    1    304  56051 1534.1
## + number_of_errors           1    302  56053 1534.1
## + .jpg                      1    284  56071 1534.2
## + X50x45                    1    264  56091 1534.3
## + X292pxx292px              1    261  56094 1534.3
## + X800x1200                 1    255  56100 1534.3
## + Flesh_Mesaure             1    210  56145 1534.5
## + pinterest                 1    200  56155 1534.6
## + X15x75                     1    179  56176 1534.7
## + Sentences                  1    152  56203 1534.8
## + facebook                   1    138  56216 1534.9

```

```

## + .dib          1    104  56251 1535.1
## + Readability  1     97  56258 1535.1
## + .jpe          1     74  56281 1535.2
## + .jpeg         1     72  56283 1535.2
## + Words         1     62  56293 1535.3
## + twitter       1     37  56318 1535.4
## + external      1     34  56321 1535.5
## + number_of_warning 1     32  56323 1535.5
## + .tif          1     24  56331 1535.5
## + internal      1     15  56340 1535.5
## + instagram     1     10  56345 1535.6
## + .tiff         1      5  56350 1535.6
## + .png          1      2  56353 1535.6
## + Unique.words  1      1  56353 1535.6
## + linkedin      1      0  56354 1535.6
## + non.document.error 1      0  56355 1535.6
## - X44x556       1   1085  57440 1537.1
## - youtube       1   1146  57501 1537.4
## - X8x15          1   1233  57588 1537.9
## - .bmp           1   3380  59735 1548.4
## - X400x300      1   5677  62032 1559.3
## - X15x12         1   45753 102108 1702.8
##
## Step: AIC=1532.45
## Revenues ~ X15x12 + X44x556 + X400x300 + .bmp + X8x15 + youtube +
##           loading.time
##
##                               Df Sum of Sq   RSS   AIC
## + X46x214            1     597 55141 1531.3
## + X24pxx133px        1     518 55220 1531.8
## + X100x100           1     446 55292 1532.1
## + number_of_errors    1     395 55343 1532.4
## <none>                  55738 1532.5
## + X115x223          1     384 55354 1532.5
## + .gif                1     377 55361 1532.5
## + X29x29             1     357 55381 1532.6
## + .jpg                1     353 55385 1532.6
## + X800x1200          1     332 55406 1532.7
## + Sentences           1     287 55451 1533.0
## + X50x45              1     282 55456 1533.0
## + Flesh_Mesaure       1     253 55485 1533.1
## + X292pxx292px        1     241 55497 1533.2
## + pinterest          1     193 55545 1533.5
## + Words                1     189 55549 1533.5
## + X15x75              1     180 55558 1533.5
## + .dib                1     177 55561 1533.5
## - loading.time         1     617 56355 1533.6
## + facebook             1     139 55599 1533.7
## + external              1     115 55623 1533.9
## + Readability           1      93 55645 1534.0
## + .jpe                 1      76 55662 1534.1
## + .jpeg                1      75 55663 1534.1
## + internal              1      74 55664 1534.1
## + twitter               1      65 55673 1534.1

```

```

## + .tif          1      54  55684 1534.2
## + Unique.words 1      50  55688 1534.2
## + .png          1      32  55706 1534.3
## + number_of_warning 1      25  55713 1534.3
## + instagram     1      16  55722 1534.4
## + .tiff         1       8  55730 1534.4
## + linkedin      1       7  55731 1534.4
## + non.document.error 1       0  55738 1534.5
## - X8x15         1    1085  56823 1536.0
## - youtube       1    1162  56900 1536.4
## - X44x556       1    1199  56937 1536.6
## - .bmp          1    3278  59016 1546.9
## - X400x300      1    5596  61334 1558.0
## - X15x12        1   45266 101004 1701.7
##
## Step: AIC=1531.35
## Revenues ~ X15x12 + X44x556 + X400x300 + .bmp + X8x15 + youtube +
##           loading.time + X46x214
##
##                               Df Sum of Sq   RSS   AIC
## + .gif                 1     405 54736 1531.2
## + number_of_errors     1     405 54736 1531.2
## + X115x223            1     385 54756 1531.3
## <none>                55141 1531.3
## + X29x29              1     358 54783 1531.5
## + .jpg                1     335 54806 1531.6
## + Flesh_Mesaure       1     284 54857 1531.9
## + Sentences            1     273 54868 1531.9
## + X292pxx292px        1     239 54902 1532.1
## + X15x75              1     181 54960 1532.4
## + Words                1     174 54967 1532.4
## + facebook             1     172 54969 1532.5
## - X46x214             1     597 55738 1532.5
## + .dib                1     167 54974 1532.5
## + pinterest           1     165 54976 1532.5
## + Readability          1     148 54993 1532.6
## + X800x1200            1     134 55007 1532.7
## - loading.time          1     678 55819 1532.9
## + external              1      87 55054 1532.9
## + internal              1      68 55073 1533.0
## + X100x100              1      65 55076 1533.0
## + Unique.words           1      53 55087 1533.1
## + .tif                  1      51 55090 1533.1
## + twitter               1      47 55094 1533.1
## + X24pxx133px           1      37 55104 1533.2
## + X50x45                1      31 55110 1533.2
## + number_of_warning      1      21 55120 1533.2
## + instagram              1      15 55126 1533.3
## - X44x556              1    760 55901 1533.3
## + .png                  1      11 55130 1533.3
## + .tiff                 1       8 55133 1533.3
## + linkedin              1       7 55134 1533.3
## + .jpe                  1       6 55135 1533.3
## + .jpeg                 1       5 55136 1533.3

```

```

## + non.document.error 1 0 55141 1533.3
## - X8x15 1 1076 56217 1534.9
## - youtube 1 1126 56267 1535.2
## - X400x300 1 1712 56853 1538.2
## - .bmp 1 3277 58418 1546.0
## - X15x12 1 45240 100381 1701.9
##
## Step: AIC=1531.23
## Revenues ~ X15x12 + X44x556 + X400x300 + .bmp + X8x15 + youtube +
## loading.time + X46x214 + .gif
##
## Df Sum of Sq RSS AIC
## + number_of_errors 1 457 54279 1530.8
## <none> 54736 1531.2
## + X115x223 1 366 54369 1531.3
## - .gif 1 405 55141 1531.3
## + X29x29 1 355 54381 1531.3
## + .jpg 1 353 54383 1531.4
## + Flesh_Mesaure 1 293 54443 1531.7
## + X292pxx292px 1 237 54499 1532.0
## + Sentences 1 223 54513 1532.0
## + facebook 1 200 54536 1532.2
## + .dib 1 183 54553 1532.3
## + X15x75 1 182 54554 1532.3
## + Readability 1 154 54582 1532.4
## + Words 1 147 54589 1532.5
## - X46x214 1 625 55361 1532.5
## + pinterest 1 126 54610 1532.6
## + X800x1200 1 117 54619 1532.6
## - loading.time 1 666 55402 1532.7
## + external 1 93 54642 1532.7
## + internal 1 63 54673 1532.9
## + .tif 1 59 54677 1532.9
## + X24pxx133px 1 45 54691 1533.0
## + X100x100 1 41 54695 1533.0
## + Unique.words 1 40 54696 1533.0
## - X44x556 1 726 55462 1533.0
## + .png 1 36 54700 1533.0
## + .jpe 1 30 54706 1533.1
## + .jpeg 1 29 54707 1533.1
## + X50x45 1 28 54708 1533.1
## + twitter 1 22 54714 1533.1
## + number_of_warning 1 20 54716 1533.1
## + instagram 1 20 54716 1533.1
## + .tiff 1 6 54730 1533.2
## + non.document.error 1 6 54730 1533.2
## + linkedin 1 0 54736 1533.2
## - X8x15 1 1083 55819 1534.9
## - youtube 1 1107 55843 1535.0
## - X400x300 1 1683 56419 1538.0
## - .bmp 1 3266 58002 1545.9
## - X15x12 1 45298 100033 1702.9
##
## Step: AIC=1530.81

```

```

## Revenues ~ X15x12 + X44x556 + X400x300 + .bmp + X8x15 + youtube +
##      loading.time + X46x214 + .gif + number_of_errors
##
##                                     Df Sum of Sq   RSS     AIC
## <none>                               54279 1530.8
## + X115x223                           1    360 53919 1530.9
## + Flesh_Mesaure                      1    353 53926 1530.9
## + X29x29                            1    334 53945 1531.0
## - number_of_errors                   1    457 54736 1531.2
## - .gif                                1    458 54736 1531.2
## + facebook                           1    284 53995 1531.3
## + Readability                        1    219 54060 1531.7
## + X292pxx292px                      1    195 54084 1531.8
## + X15x75                            1    178 54101 1531.9
## + .jpg                                1    176 54102 1531.9
## + .dib                                1    165 54114 1531.9
## + pinterest                          1    138 54141 1532.1
## + X800x1200                          1    134 54144 1532.1
## - X46x214                            1    638 54916 1532.2
## + Sentences                           1     83 54196 1532.4
## + Words                               1     73 54206 1532.4
## + non.document.error                 1     63 54216 1532.5
## + X100x100                           1     62 54216 1532.5
## + X50x45                             1     49 54230 1532.5
## + X24pxx133px                       1     37 54242 1532.6
## - X44x556                            1    726 55004 1532.6
## + external                            1     22 54257 1532.7
## + twitter                            1     15 54263 1532.7
## + .png                                1     14 54264 1532.7
## + Unique.words                        1     13 54265 1532.7
## + .tif                                1     10 54269 1532.8
## + instagram                           1      6 54273 1532.8
## + .tiff                               1      5 54274 1532.8
## + .jpe                                1      4 54275 1532.8
## + .jpeg                               1      3 54275 1532.8
## + linkedin                            1      3 54276 1532.8
## + number_of_warning                   1      0 54278 1532.8
## + internal                            1      0 54278 1532.8
## - loading.time                        1    771 55050 1532.9
## - X8x15                              1   1126 55405 1534.7
## - youtube                            1   1142 55421 1534.8
## - X400x300                           1   1610 55888 1537.2
## - .bmp                                1   3237 57516 1545.5
## - X15x12                            1  45164 99442 1703.2

summary(model_a)

##
## Call:
## lm(formula = Revenues ~ X15x12 + X44x556 + X400x300 + .bmp +
##      X8x15 + youtube + loading.time + X46x214 + .gif + number_of_errors,
##      data = total_500_final_train)
##
## Residuals:
##   Min    1Q  Median    3Q   Max

```

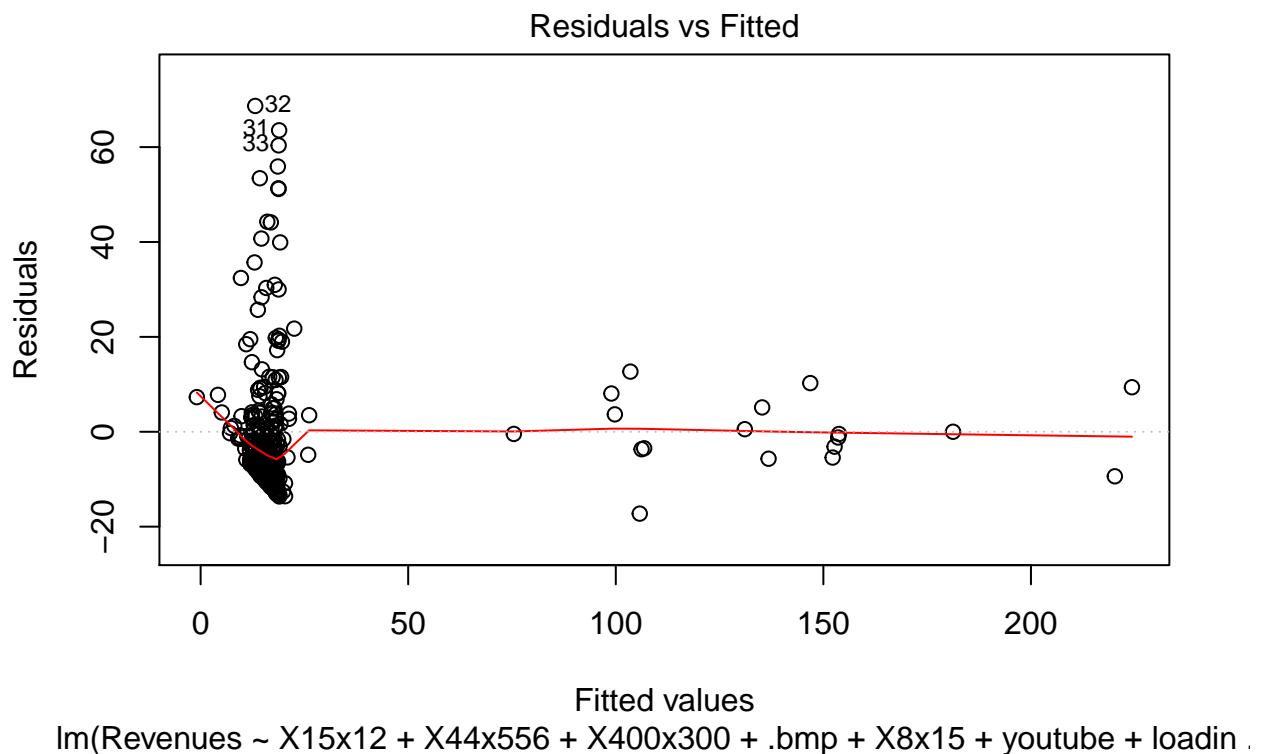
```

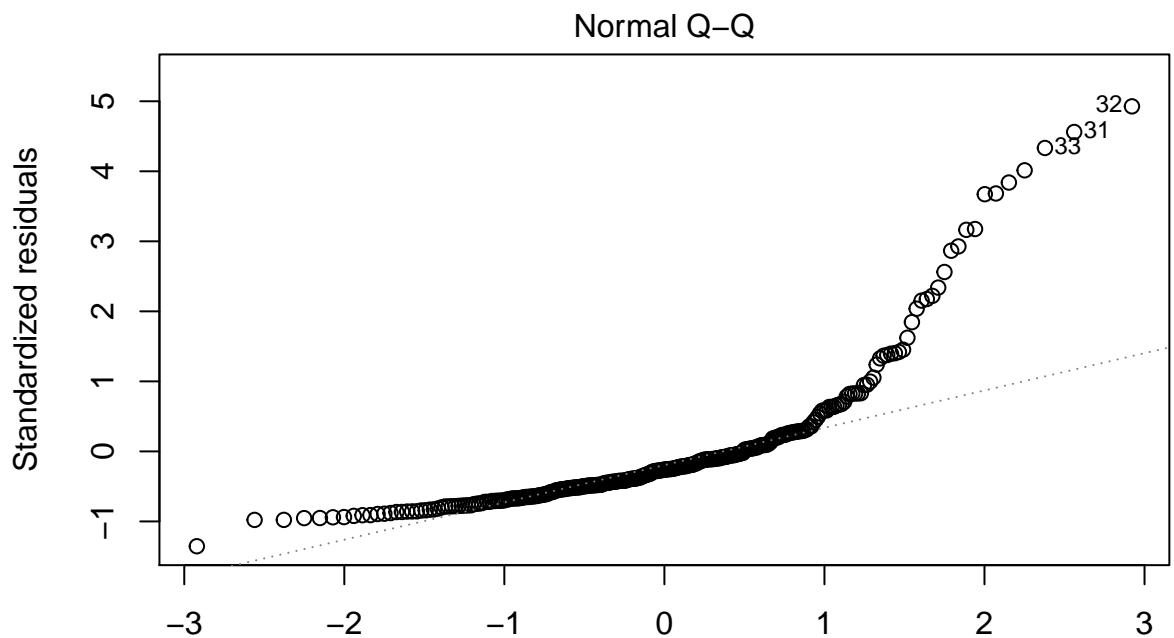
## -17.254 -7.695 -3.529  2.137 68.673
##
## Coefficients:
##                               Estimate Std. Error t value Pr(>|t|)
## (Intercept)            220.273625   9.934832 22.172 < 2e-16 ***
## X15x12                -87.818989   5.784528 -15.182 < 2e-16 ***
## X44x556                -29.578085  15.369922  -1.924  0.05533 .
## X400x300                -28.443747  9.923508  -2.866  0.00447 **
## .bmp                   2.470545   0.607833   4.065 6.28e-05 ***
## X8x15                  -41.311818  17.232973  -2.397  0.01718 *
## youtube                 4.058630   1.681080   2.414  0.01641 *
## loading.time              -4.490750  2.264113  -1.983  0.04830 *
## X46x214                 -18.477039 10.243269  -1.804  0.07235 .
## .gif                   -0.120050   0.078556  -1.528  0.12760
## number_of_errors        0.013675   0.008952   1.528  0.12774
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14 on 277 degrees of freedom
## Multiple R-squared:  0.8395, Adjusted R-squared:  0.8338
## F-statistic: 144.9 on 10 and 277 DF,  p-value: < 2.2e-16
ad_r_sq_ma <- summary(model_a)$adj.r.squared
aic_ma <- AIC(model_a)

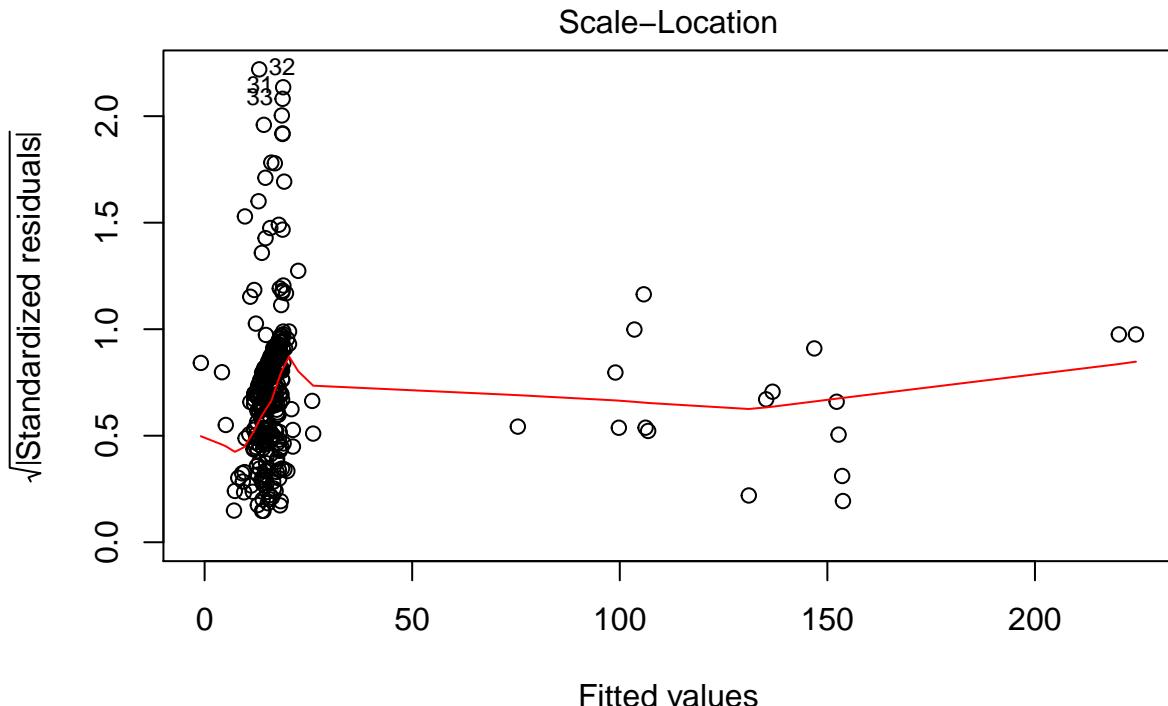
plot(model_a,which=1:3)

## Warning: not plotting observations with leverage one:
##      3

```







```
lm(Revenues ~ X15x12 + X44x556 + X400x300 + .bmp + X8x15 + youtube + loadin .
```

```
#####
#We compare the Adjusted R squares of the models and also the AIC of the models we created to find the
```

```
ad_r_sq_f3
```

```
## [1] 0.830413
```

```
ad_r_sq_f4
```

```
## [1] 0.8161957
```

```
ad_r_sq_ma #BEST
```

```
## [1] 0.8337517
```

#The best Adkusted R square is the one in model a (the closer to 1 the better)

```
aic_f3
```

```
## [1] 2362.474
```

```
aic_f4
```

```
## [1] 2378.069
```

```
aic_ma #Best
```

```
## [1] 2350.119
```

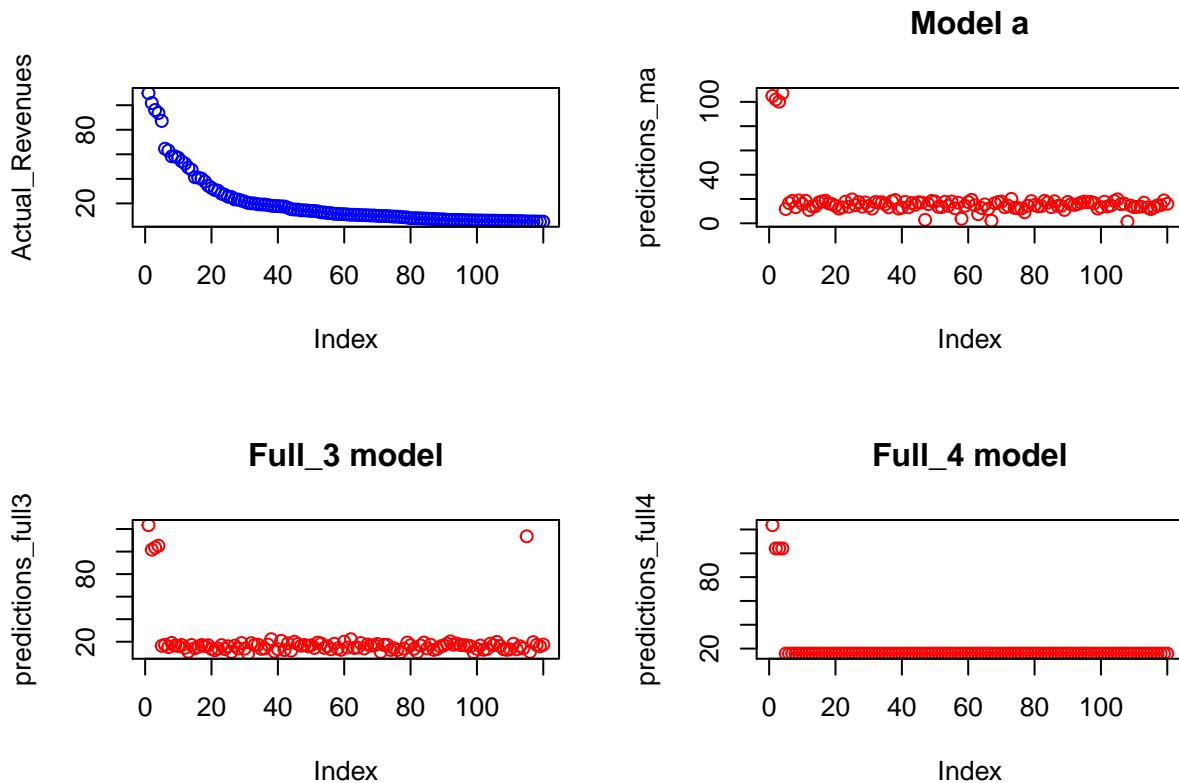
#The best AIC and the best Adjusted R square is for model ma

```
#####
par(mfrow=c(2,2))
```

```

Actual_Revenues<- total_500_final_test$Revenues
plot (Actual_Revenues, col = "blue")
#####
predictions_ma <- predict(model_a,total_500_final_test)
plot (predictions_ma, col = "Red",main = "Model a")
#####
predictions_full3 <- predict(full_3,total_500_final_test)
plot (predictions_full3, col = "Red",main = "Full_3 model")
#####
predictions_full4 <- predict(full_4,total_500_final_test)
plot (predictions_full4, col = "Red",main = "Full_4 model")

```



```

#####
#From the plots above we can see that the actual Revenues have a more smooth way of leveling up except .
#The prediction model that is more smooth is the model a which has as we said before the best Adjusted R^2
names(total_500_final_train)
```

```

## [1] "Revenues"           "non.document.error" "number_of_errors"
## [4] "number_of_warning"   "facebook"          "instagram"
## [7] "linkedin"           "pinterest"         "twitter"
## [10] "youtube"            "Flesh_Mesaure"    "Readability"
## [13] "Sentences"          "Unique.words"     "Words"
## [16] "external"           "internal"         "total.links"
## [19] "X15x75"             "X8x15"            "X44x556"
## [22] "X1x1"               "X800x1200"       "X100x100"
## [25] "X24pxx133px"      "X21pxx173px"     "X46x214"
```

```

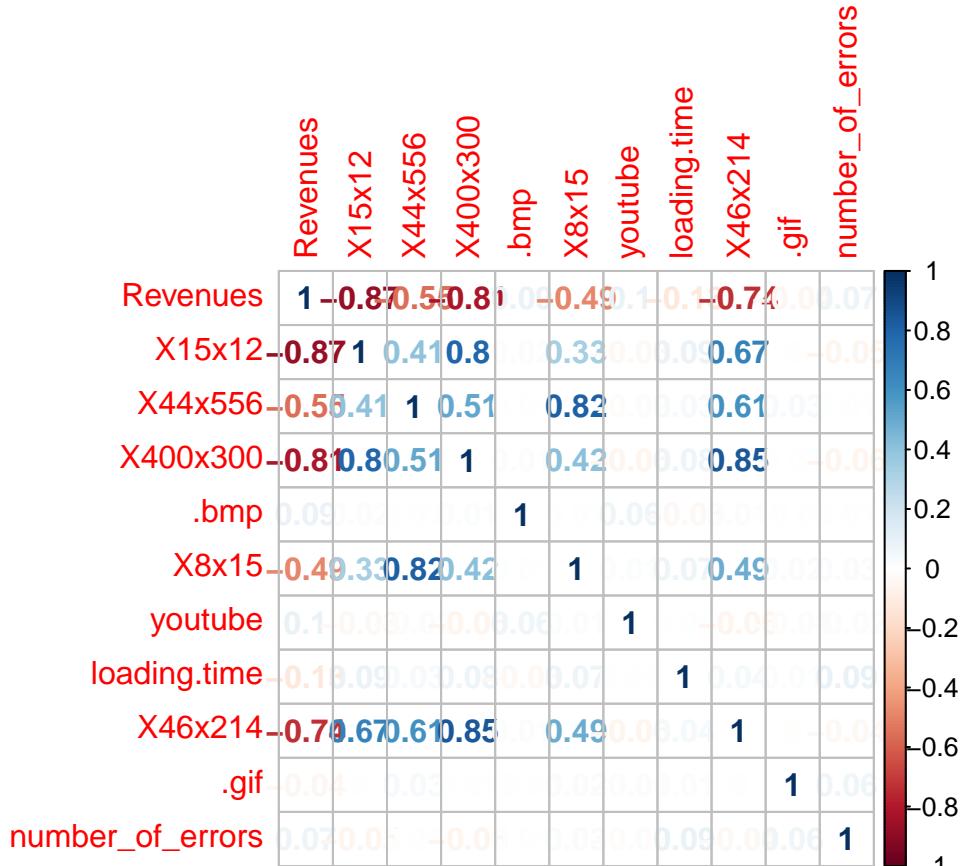
## [28] "X49x49"          "X50x45"           "X400x300"
## [31] "X292pxx292px"    "X200pxx200px"    "X1279pxx984px"
## [34] "X300pxx1500px"   "X29x29"          "X115x223"
## [37] "X160x233"        "X300x993"       "X41x192"
## [40] "X28x221"         "X15x12"          "X60x60"
## [43] ".bmp"             ".dib"            ".gif"
## [46] ".jpe"             ".jpeg"           ".jpg"
## [49] ".png"             ".tif"            ".tiff"
## [52] "total.images"     "loading.time"

par(mfrow=c(1,1))
total_500_final_reg <- total_500_final_train[,c(1,41,21,30,43,20,10,53,27,45,3)]
cor(total_500_final_reg)

##                                Revenues      X15x12      X44x556      X400x300
## Revenues                 1.00000000 -0.866071521 -0.554995003 -0.810238376
## X15x12                  -0.866071521  1.000000000  0.409636250  0.795640062
## X44x556                  -0.554995000  0.409636250  1.000000000  0.514851212
## X400x300                 -0.810238388  0.795640062  0.514851212  1.000000000
## .bmp                      0.088236868  0.017894718  0.007330325  0.014237755
## X8x15                     -0.48572735  0.333881354  0.815067889  0.419638691
## youtube                   0.09933975 -0.034098547 -0.018043770 -0.059515452
## loading.time               -0.12686096  0.086897868  0.030631302  0.079596684
## X46x214                   -0.73936498  0.674879385  0.606976979  0.848221975
## .gif                      -0.04436402  0.002618064  0.030679347  0.008558755
## number_of_errors           0.07106309 -0.052737836  0.005969412 -0.058134366
##                               .bmp          X8x15      youtube loading.time
## Revenues                  0.088236861 -0.485727351  0.09933975 -0.12686096
## X15x12                    0.017894718  0.333881354 -0.03409855  0.08689787
## X44x556                   0.007330325  0.815067889 -0.01804377  0.03063130
## X400x300                  0.014237755  0.419638691 -0.05951545  0.07959668
## .bmp                      1.000000000  0.005974713  0.06081623 -0.03202094
## X8x15                     0.005974713  1.000000000  0.01353035  0.07305785
## youtube                   0.060816230  0.013530346  1.000000000  0.00711365
## loading.time                -0.032020942 0.073057846  0.00711365  1.000000000
## X46x214                   0.012076776  0.494727445 -0.05826545  0.04051586
## .gif                      -0.006020021 0.022359056 -0.01420548  0.01202806
## number_of_errors           0.005372103  0.026722692 -0.01898525  0.09290871
##                               X46x214      .gif number_of_errors
## Revenues                  -0.739364984 -0.044364026  0.071063094
## X15x12                     0.674879385  0.0026180642 -0.052737836
## X44x556                   0.606976978  0.0306793475  0.005969412
## X400x300                  0.8482219754 0.0085587551 -0.058134366
## .bmp                      0.0120767763 -0.0060200206  0.005372103
## X8x15                     0.4947274449 0.0223590563  0.026722692
## youtube                   -0.0582654543 -0.0142054765 -0.018985254
## loading.time                0.0405158629 0.0120280632  0.092908709
## X46x214                   1.0000000000 -0.0009473611 -0.038355173
## .gif                      -0.0009473611 1.0000000000  0.062343582
## number_of_errors           -0.0383551731 0.0623435816  1.0000000000

corrplot(cor(total_500_final_reg),method="number")

```



#We can see here that the variable x8x15 has a very high correlation with the variable x44x556 and also

#So we can try creating a new model excluding the 2 variables that are correlated from each pair to see
full_5 <- lm(Revenues~1 +X44x556 +X400x300 + .bmp + youtube +loading.time + X46x214 + .gif + number_of_

```

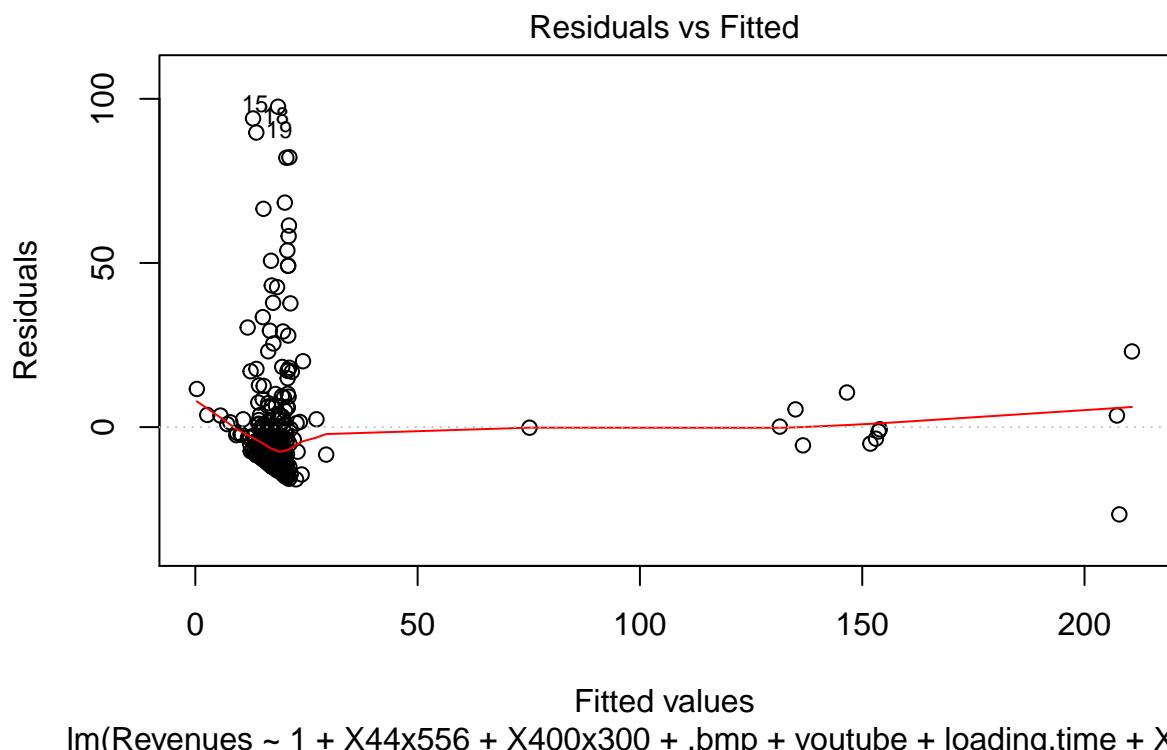
## Call:
## lm(formula = Revenues ~ 1 + X44x556 + X400x300 + .bmp + youtube +
##     loading.time + X46x214 + .gif + number_of_errors, data = total_500_final_train)
##
## Residuals:
##      Min    1Q   Median    3Q   Max 
## -26.584 -9.736 -5.298  0.999 97.593 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 207.33952  11.08956 18.697 < 2e-16 ***
## X44x556    -57.05604  13.88428 -4.109 5.22e-05 ***
## X400x300   -113.81359  11.08164 -10.270 < 2e-16 ***
## .bmp        2.37017   0.82413   2.876  0.00434 ** 
## youtube     3.34344   2.27623   1.469  0.14300  
## loading.time -6.30802   3.05772  -2.063  0.04004 *  
## X46x214    -18.96451  13.88890  -1.365  0.17321  
## .gif       -0.11152   0.10651  -1.047  0.29600  
## number_of_errors 0.01489   0.01213   1.227  0.22072  
## ---
```

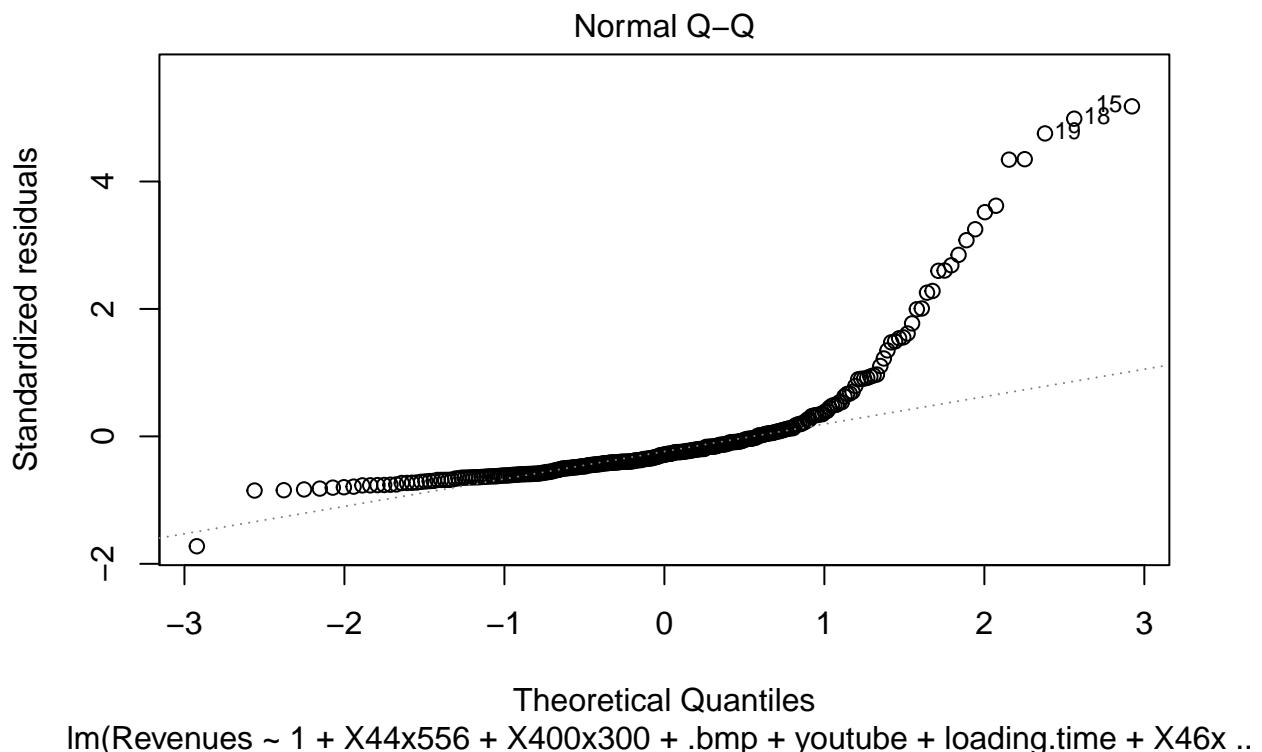
```

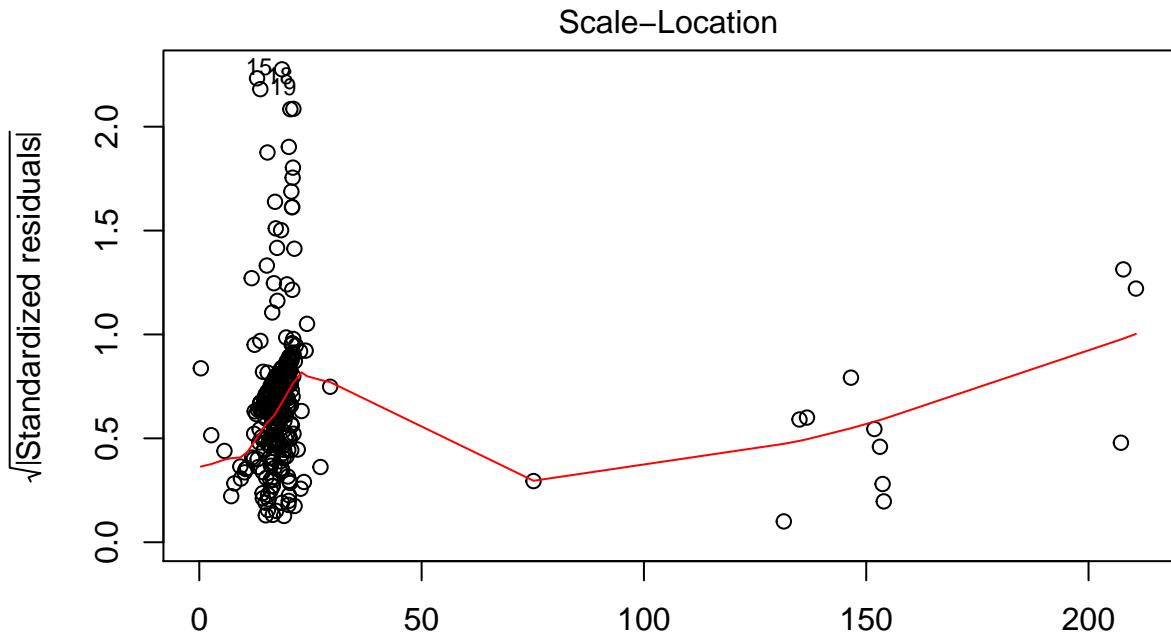
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 18.98 on 279 degrees of freedom
## Multiple R-squared:  0.7029, Adjusted R-squared:  0.6943
## F-statistic:  82.5 on 8 and 279 DF,  p-value: < 2.2e-16
adj_r_square_full5 <- summary(full_5)$adj.r.squared
aic_full5 <- AIC(full_5)

```

#We create the 2 basic plots so as to be able to explain the regression model







Fitted values

lm(Revenues ~ 1 + X44x556 + X400x300 + .bmp + youtube + loading.time + X46x ..

```
ad_r_sq_ma
```

```
## [1] 0.8337517
```

```
adj_r_square_full5
```

```
## [1] 0.6943474
```

```
aic_ma
```

```
## [1] 2350.119
```

```
aic_full5
```

```
## [1] 2523.573
```

#The adjusted R square and the aic are a little worse than before

```
#####
#####
```

#Clustering

#Kmeans clustering

#Based on those results we will try to cluster the companies based on the results of the regression

```
set.seed(220)
```

```
clusters <- hclust(dist(total_500_final_reg[, 1]))
```

```
plot(clusters)
```

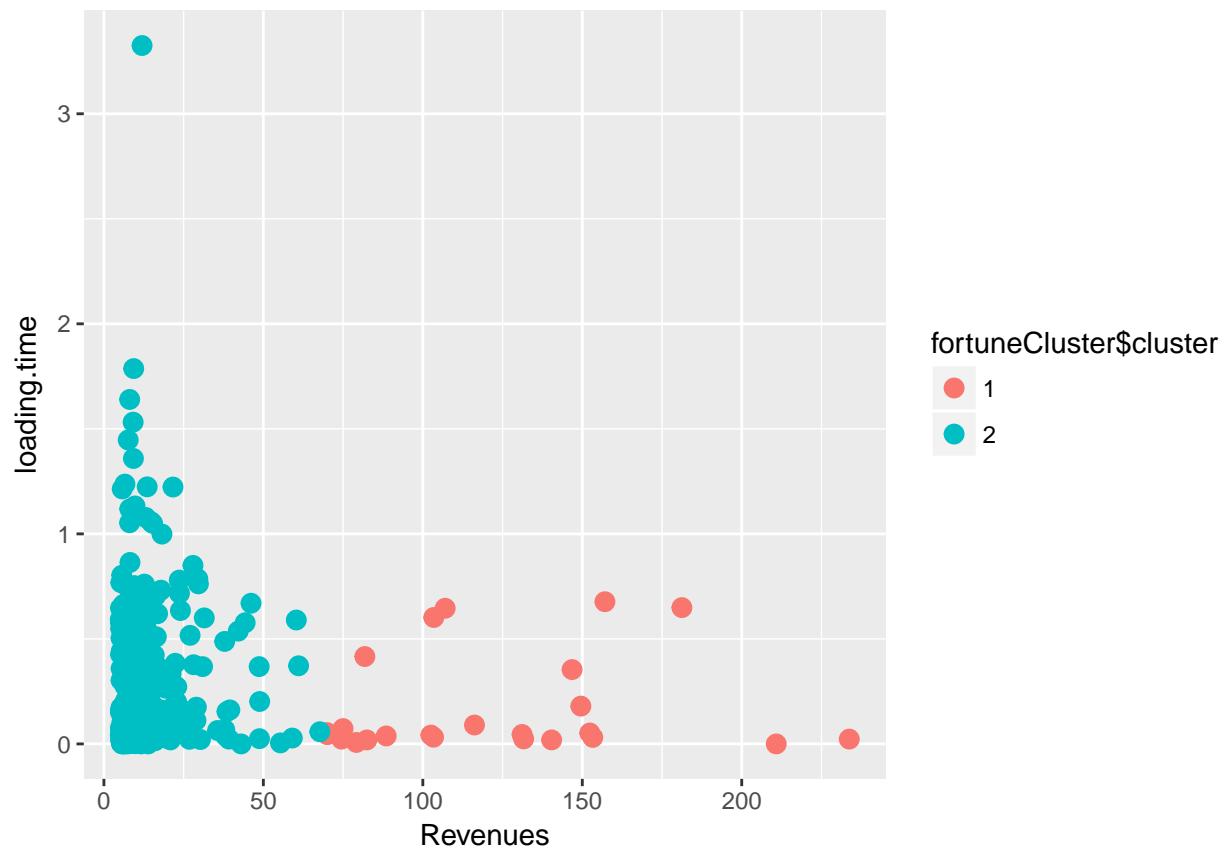
Cluster Dendrogram

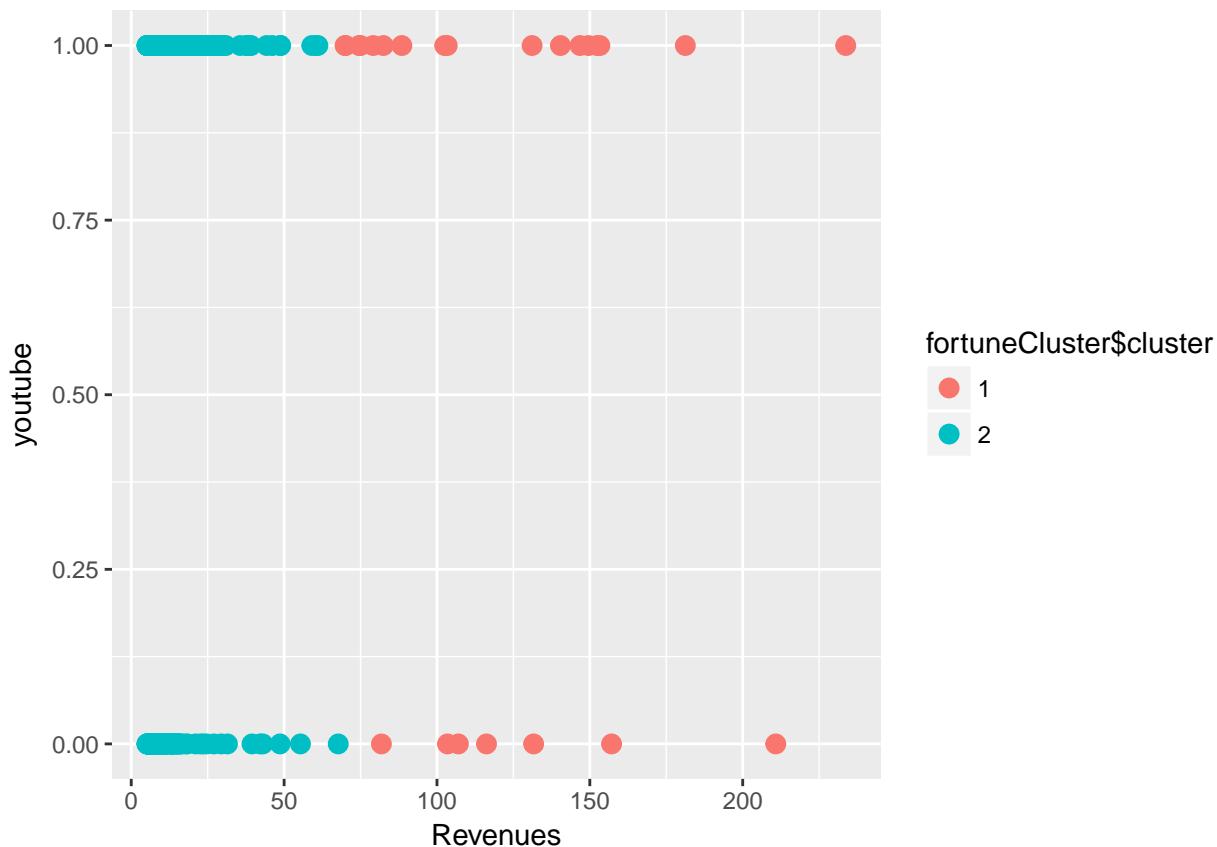


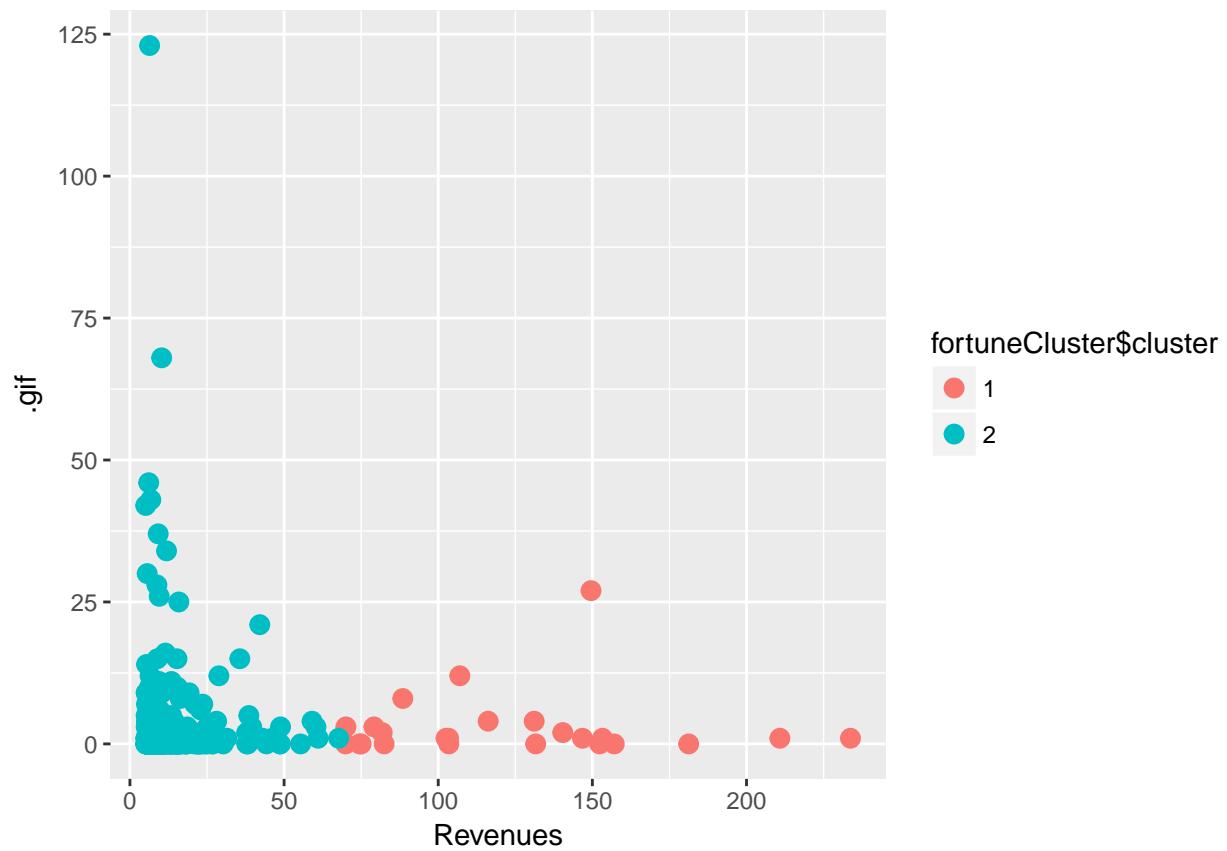
```
dist(total_500_final_reg[, 1])  
hclust (*, "complete")
```

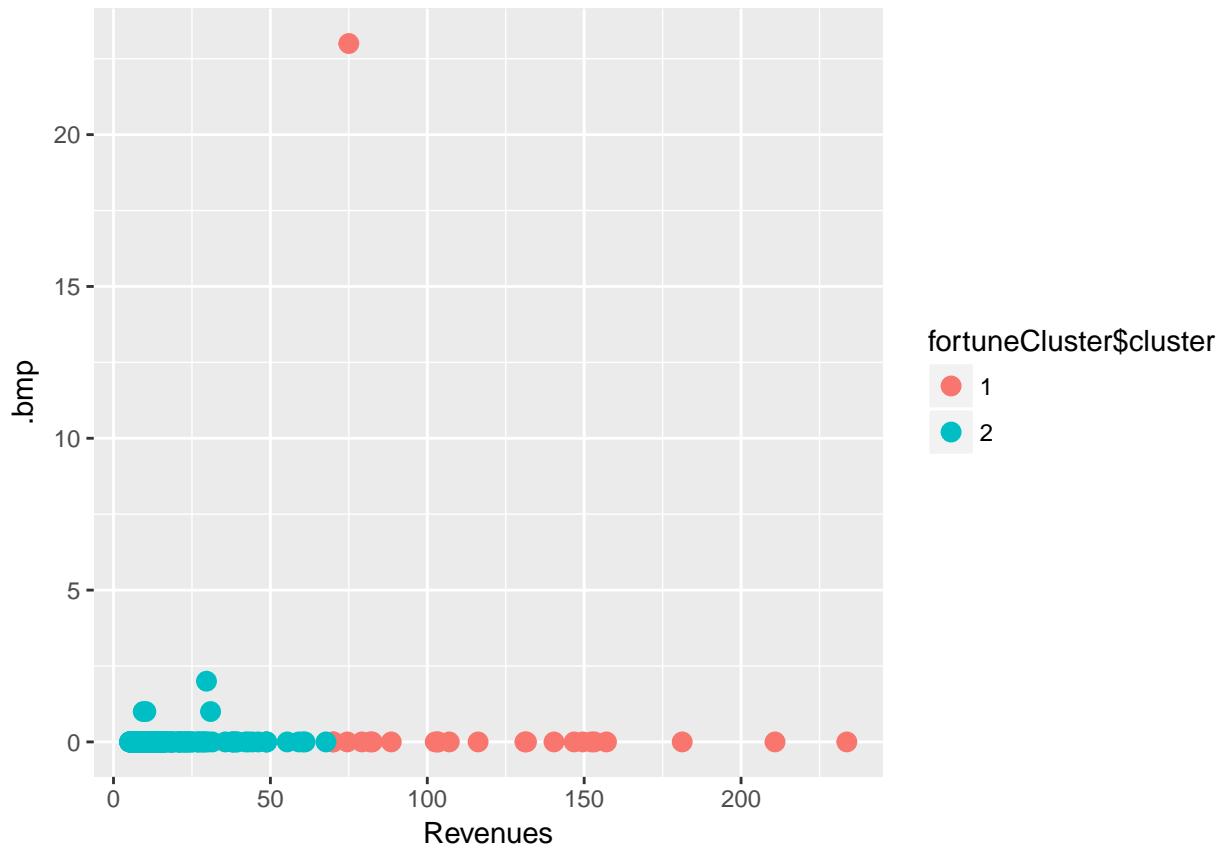
```
fortuneCluster <- kmeans(total_500_final_reg[, 1], 2, iter.max = 500, nstart = 1)
cluster <- table(fortuneCluster$cluster)
fortuneCluster$cluster <- as.factor(fortuneCluster$cluster)

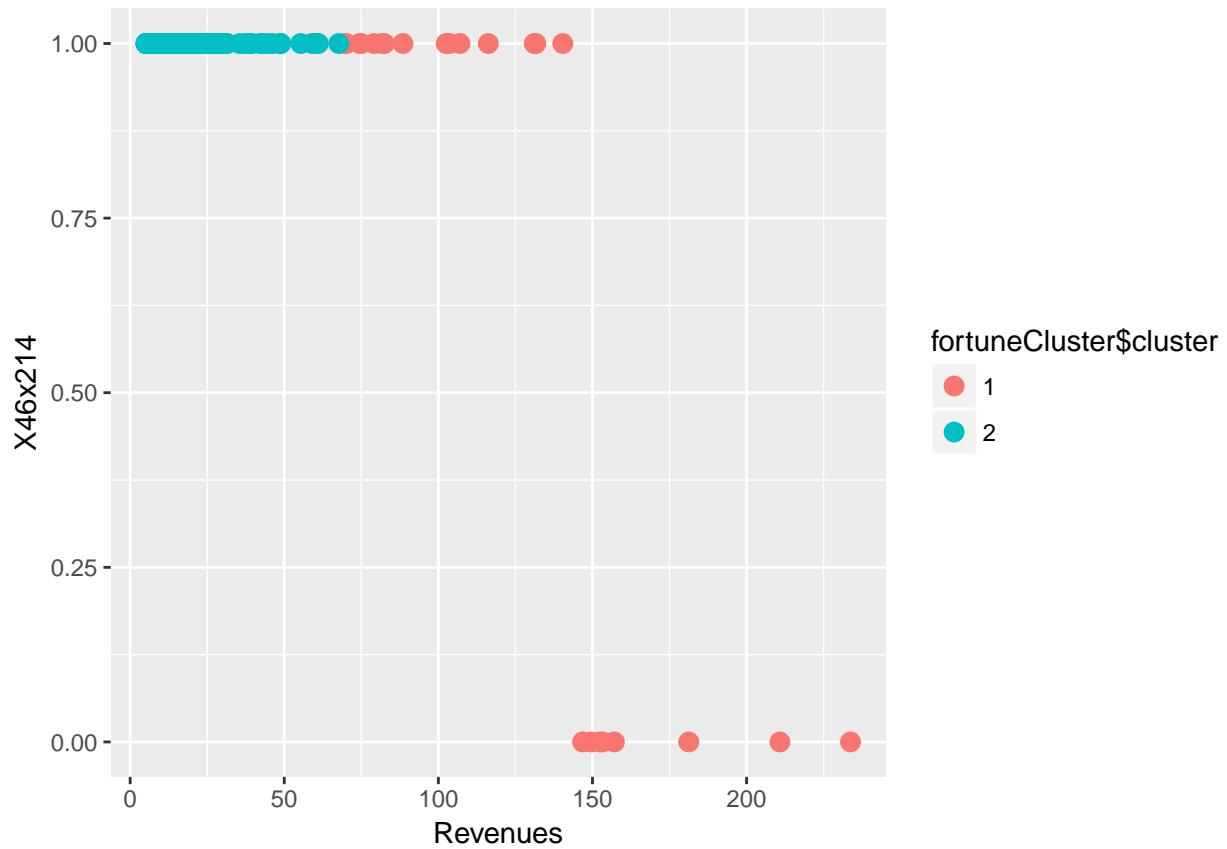
ggplot(total_500_final_reg, aes(Revenues, loading.time, color = fortuneCluster$cluster)) + geom_point(s
```

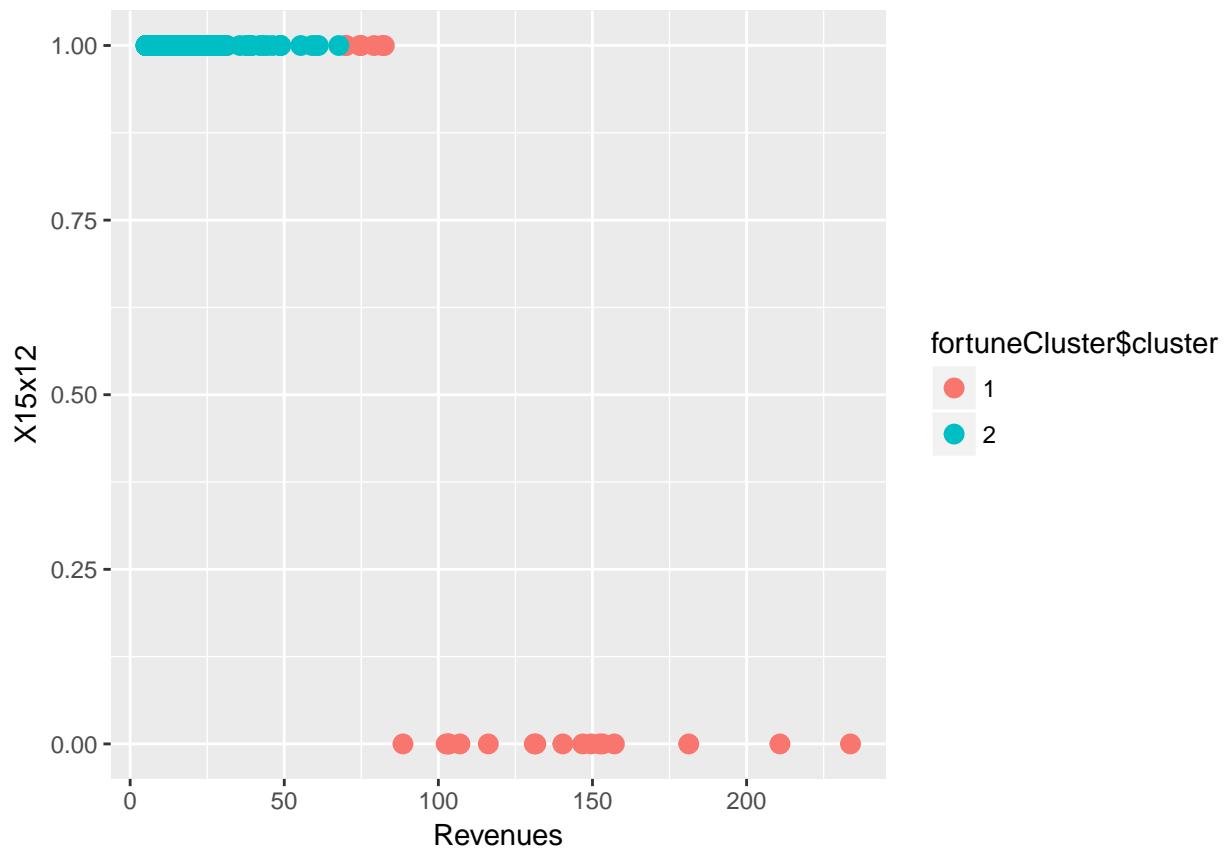


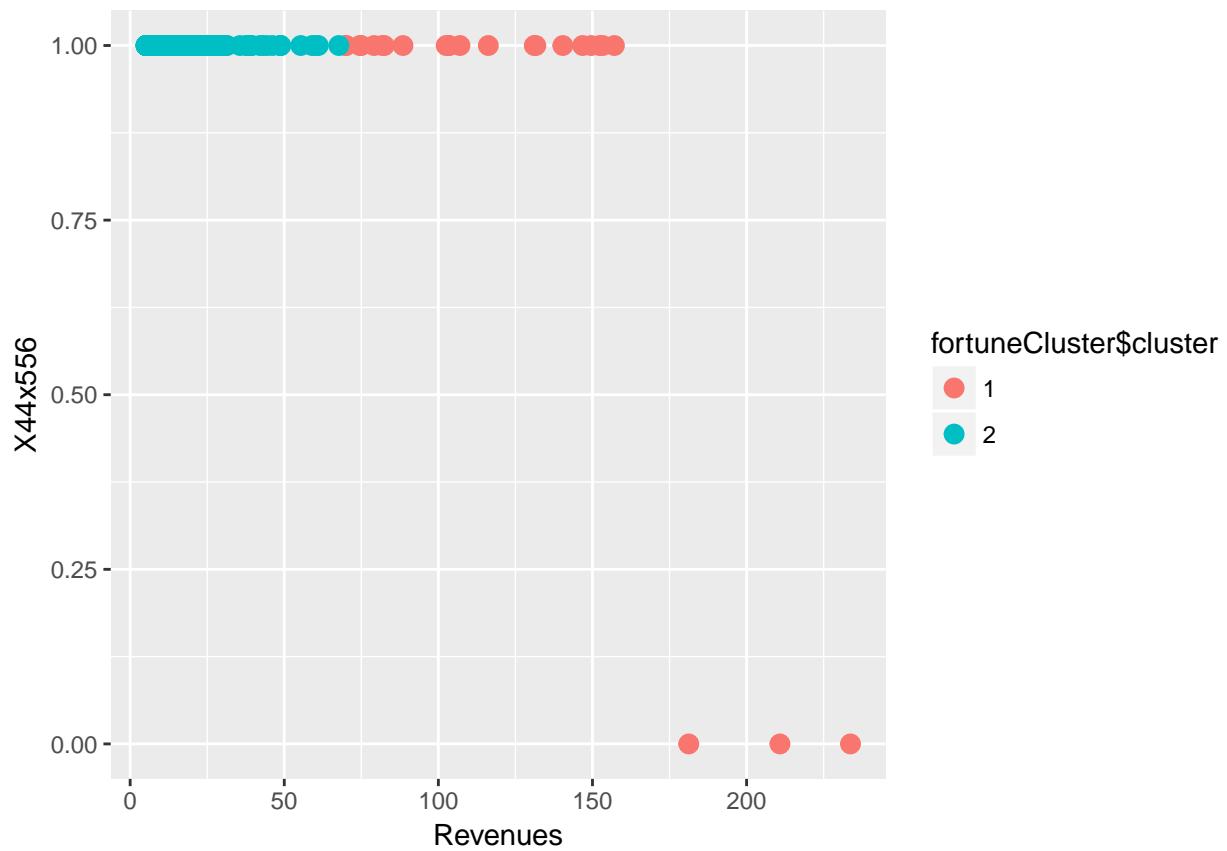


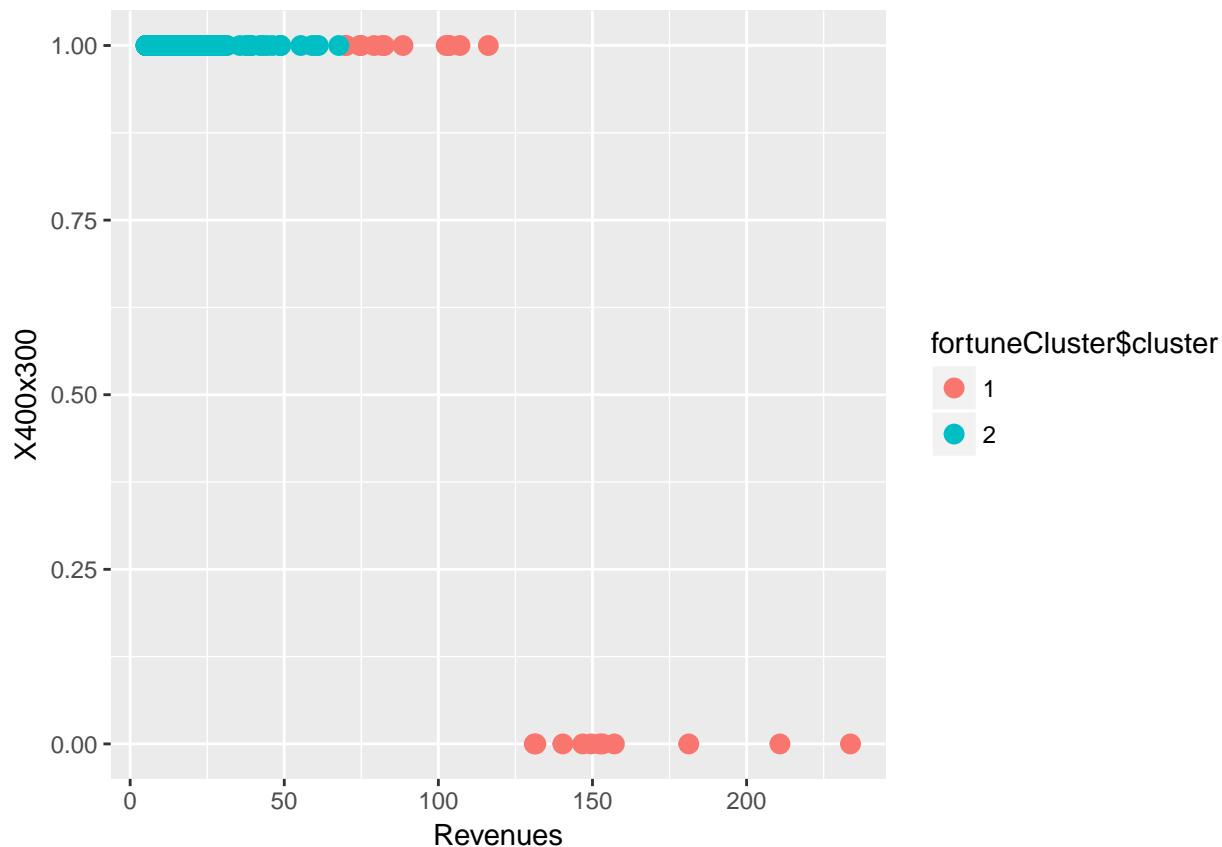




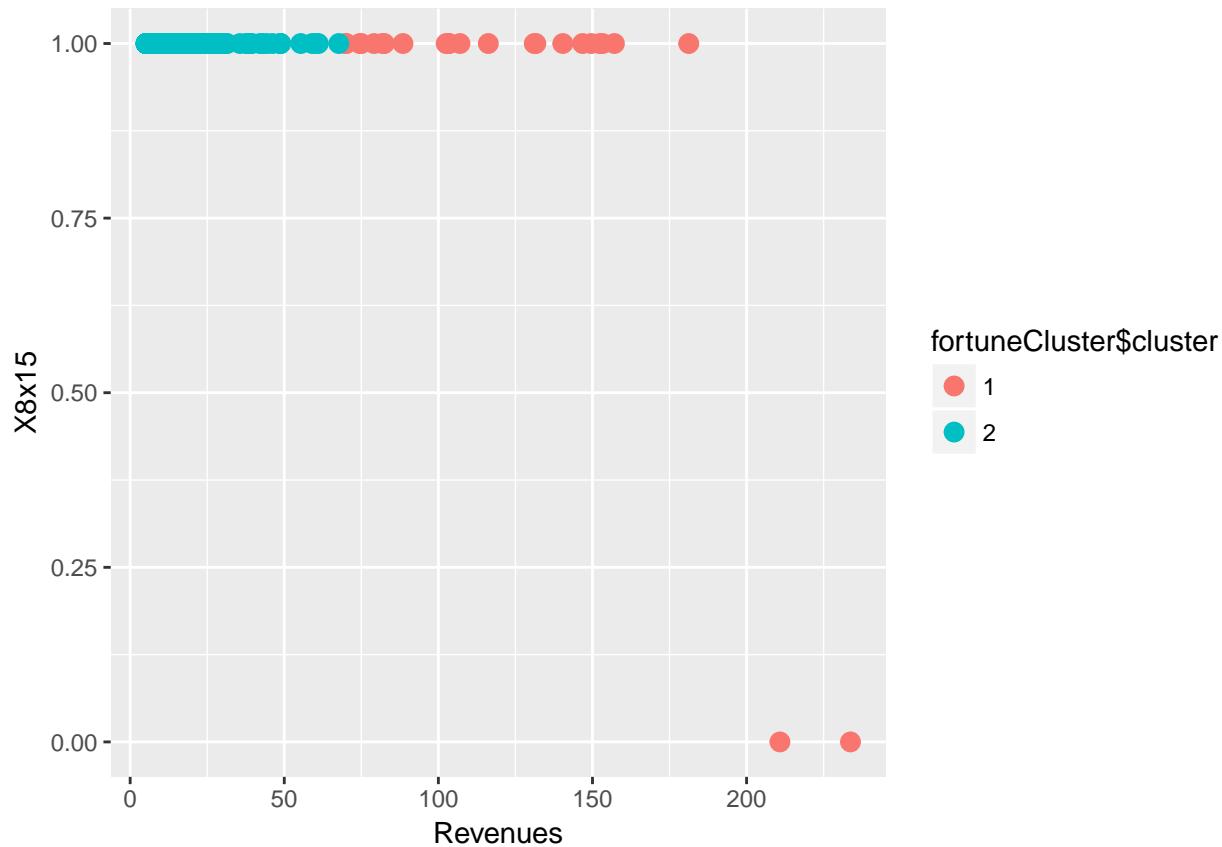








```
ggplot(total_500_final_reg, aes(Revenues, X8x15, color = fortuneCluster$cluster)) + geom_point(size=3)
```



#From the clustering we can see that the variables do indeed devide the most high revenues from the small ones.

```

## Call:
## lm(formula = Revenues ~ X15x12 + X44x556 + X400x300 + .bmp +
##     X8x15 + youtube + loading.time + X46x214 + .gif + number_of_errors,
##     data = total_500_final_train)
##
## Residuals:
##      Min    1Q   Median    3Q   Max 
## -17.254 -7.695 -3.529  2.137 68.673 
## 
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 220.273625  9.934832 22.172 < 2e-16 ***
## X15x12      -87.818989  5.784528 -15.182 < 2e-16 ***
## X44x556     -29.578085 15.369922 -1.924  0.05533 .  
## X400x300    -28.443747  9.923508 -2.866  0.00447 ** 
## .bmp         2.470545  0.607833  4.065 6.28e-05 ***
## X8x15       -41.311818 17.232973 -2.397  0.01718 *  
## youtube      4.058630  1.681080  2.414  0.01641 *  
## loading.time -4.490750  2.264113 -1.983  0.04830 *  
## X46x214     -18.477039 10.243269 -1.804  0.07235 .  
## .gif        -0.120050  0.078556 -1.528  0.12760  
## number_of_errors  0.013675  0.008952  1.528  0.12774 
## 
```

```

## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14 on 277 degrees of freedom
## Multiple R-squared: 0.8395, Adjusted R-squared: 0.8338
## F-statistic: 144.9 on 10 and 277 DF, p-value: < 2.2e-16

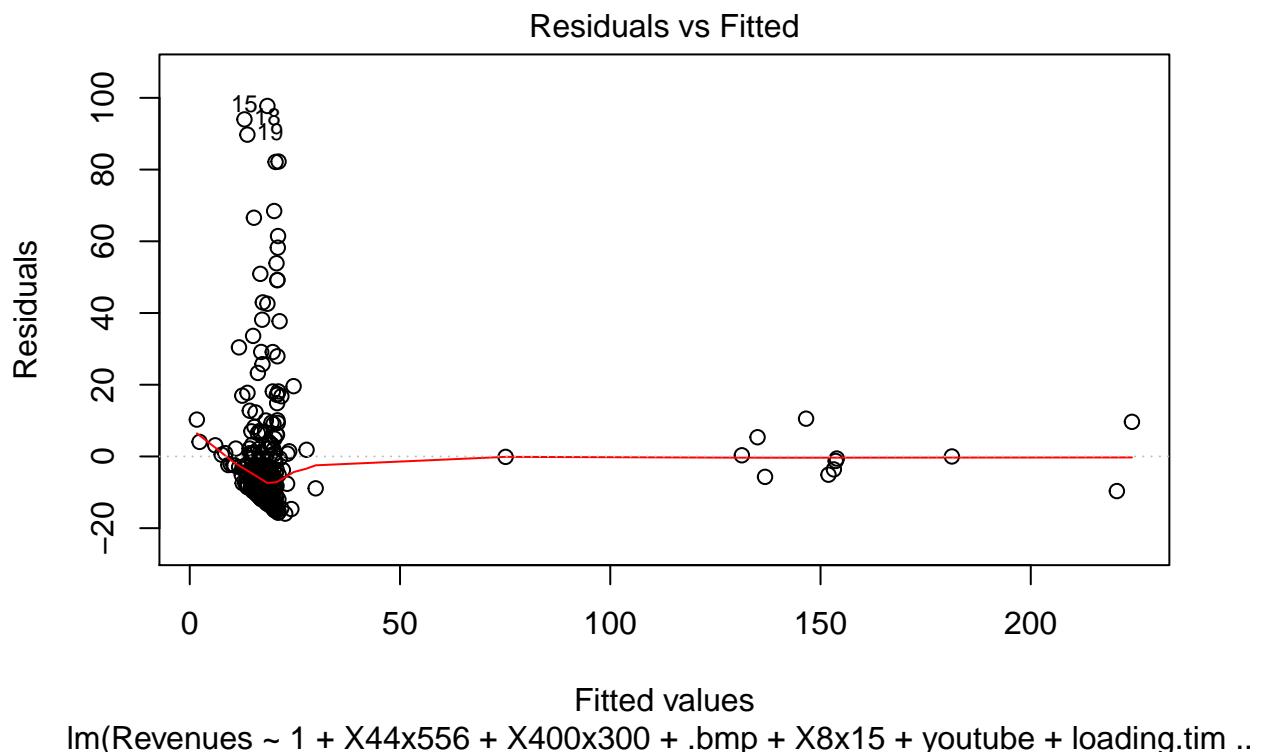
#We can see from the model that the basic variable that effect a companys ranking
#is whether or not it has an image in size X15x12
#We will try to make a model that we will not take into consideration this variable at all just in order
full_6 <- lm(Revenues~1 +X44x556 +X400x300 + .bmp + X8x15+ youtube +loading.time + X46x214 + .gif + num
summary(full_6)

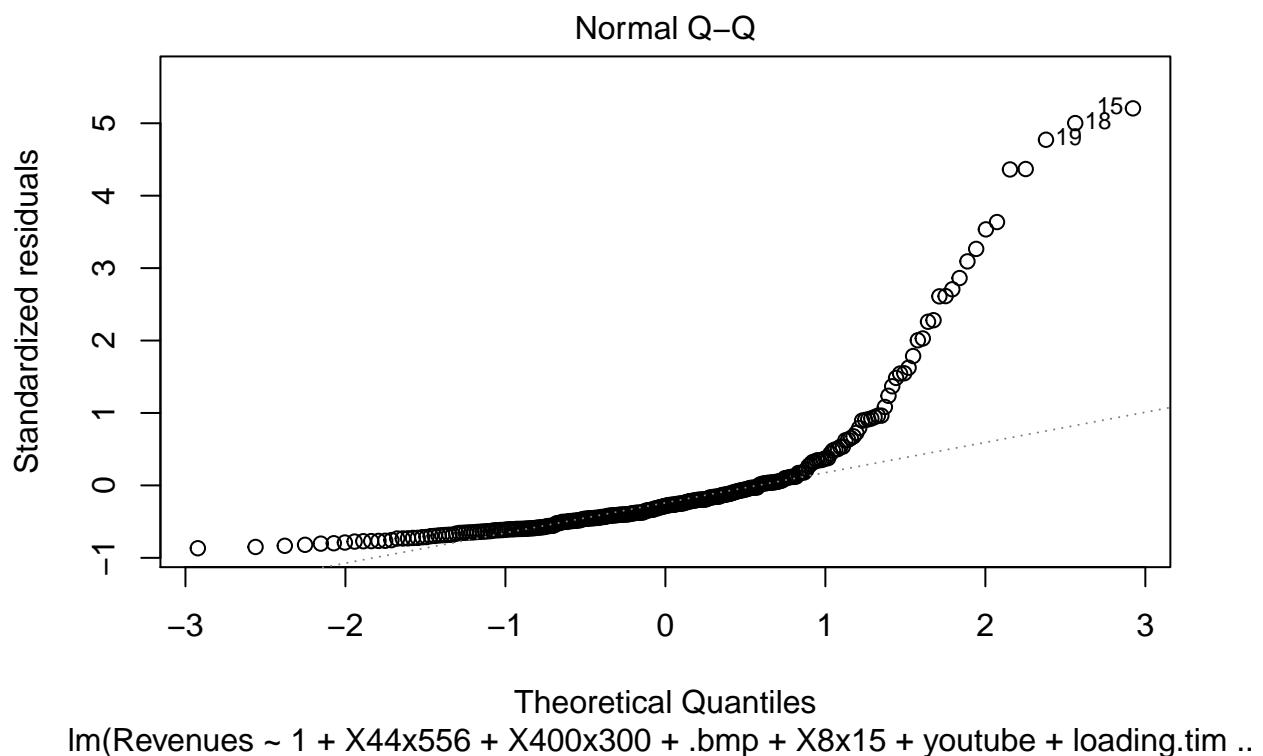
##
## Call:
## lm(formula = Revenues ~ 1 + X44x556 + X400x300 + .bmp + X8x15 +
##     youtube + loading.time + X46x214 + .gif + number_of_errors,
##     data = total_500_final_train)
##
## Residuals:
##      Min        1Q    Median        3Q       Max
## -15.941   -9.536   -5.194    0.764   97.740
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 220.52558  13.42298 16.429 < 2e-16 ***
## X44x556     -30.25871  20.76628 -1.457  0.14622
## X400x300    -113.89510  11.04234 -10.314 < 2e-16 ***
## .bmp          2.36934   0.82120  2.885  0.00422 **
## X8x15        -40.28865  23.28336 -1.730  0.08468 .
## youtube      3.53481   2.27084  1.557  0.12070
## loading.time -5.88761   3.05653 -1.926  0.05509 .
## X46x214     -18.85123  13.83968 -1.362  0.17426
## .gif          -0.11274   0.10614 -1.062  0.28906
## number_of_errors 0.01555   0.01209  1.285  0.19971
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 18.91 on 278 degrees of freedom
## Multiple R-squared: 0.706, Adjusted R-squared: 0.6965
## F-statistic: 74.19 on 9 and 278 DF, p-value: < 2.2e-16

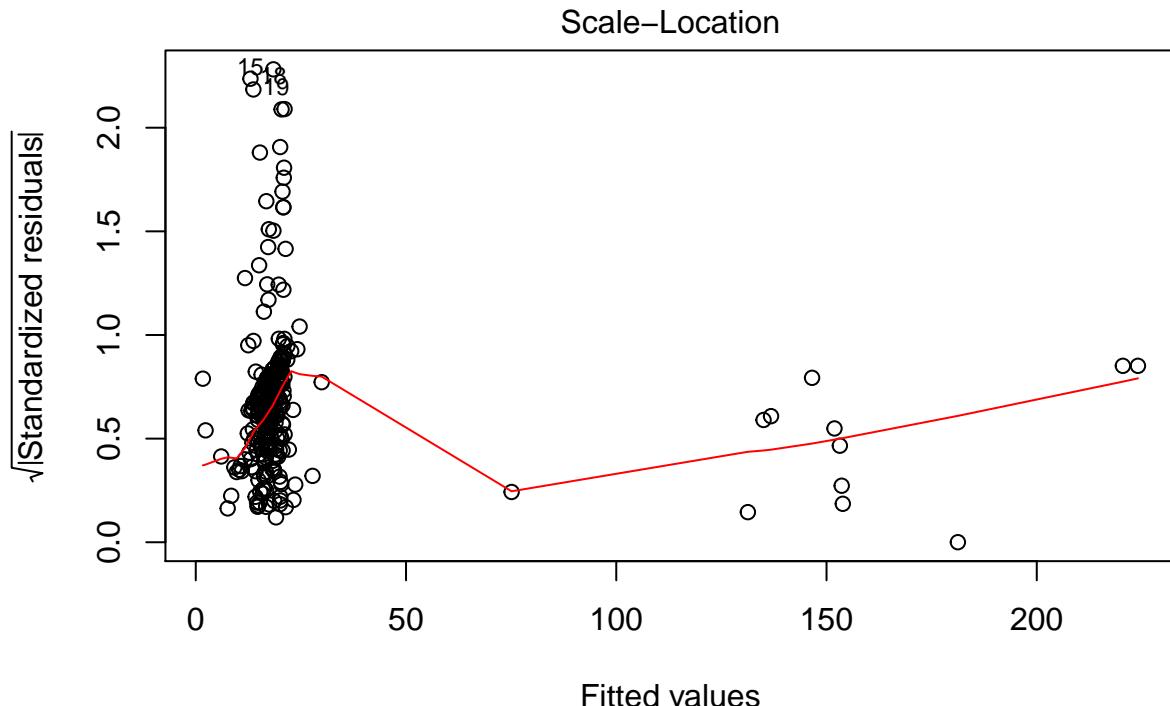
adj_r_square_full6 <- summary(full_6)$adj.r.squared
aic_full6 <- AIC(full_6)

#We create the 2 basic plots so as to be able to explain the regression model
plot(full_6,which=1:3)

```

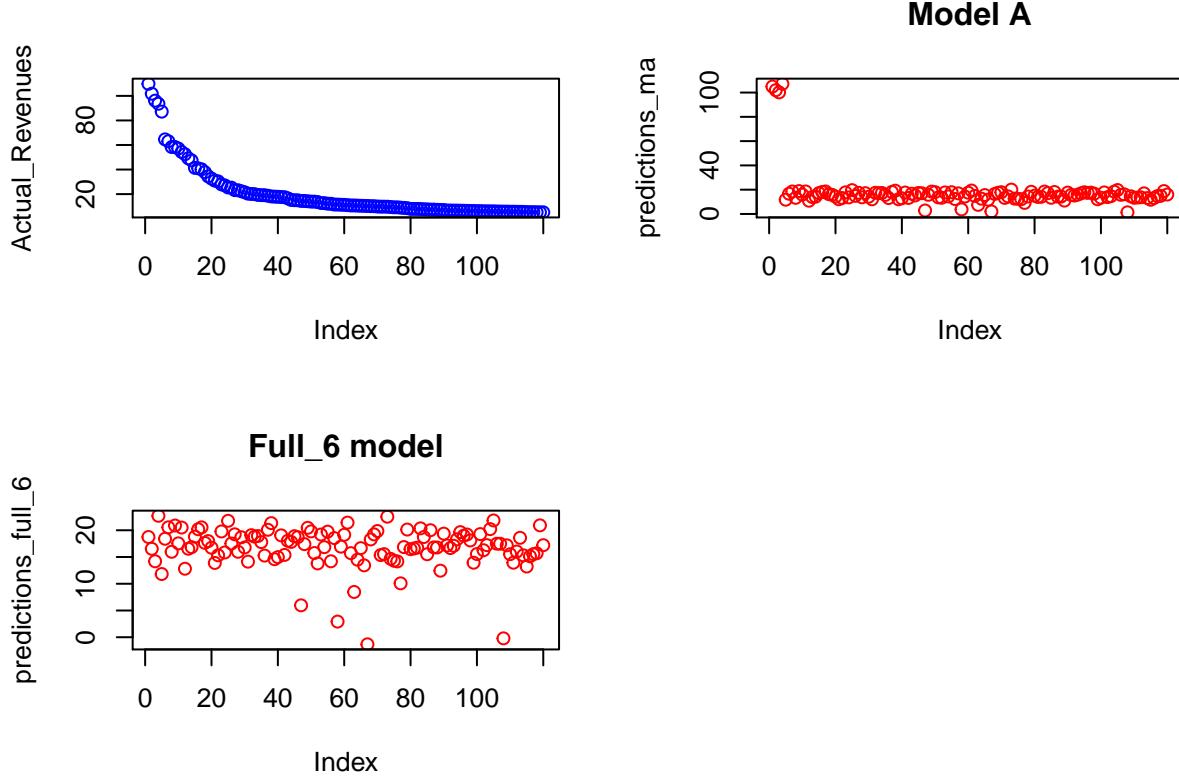






```
predictions_ma <- predict(model_a, total_500_final_test)
Actual_Revenues<- total_500_final_test$Revenues
```

```
par(mfrow=c(2,2))
plot (Actual_Revenues, col = "blue")
plot (predictions_ma, col = "Red",main = "Model A")
#####
predictions_full_6 <- predict(full_6,total_500_final_test)
plot (predictions_full_6, col = "Red",main = "Full_6 model")
#####
```



```
#We can see that here the prediction of the new model is not as good as the previous one so now that we
summary(model_a)
```

```
##
## Call:
## lm(formula = Revenues ~ X15x12 + X44x556 + X400x300 + .bmp +
##     X8x15 + youtube + loading.time + X46x214 + .gif + number_of_errors,
##     data = total_500_final_train)
##
## Residuals:
##      Min      1Q      Median      3Q      Max 
## -17.254  -7.695  -3.529   2.137  68.673 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 220.273625  9.934832 22.172 < 2e-16 ***
## X15x12      -87.818989  5.784528 -15.182 < 2e-16 ***
## X44x556     -29.578085 15.369922 -1.924  0.05533 .
## X400x300    -28.443747  9.923508 -2.866  0.00447 ** 
## .bmp         2.470545  0.607833  4.065 6.28e-05 ***
## X8x15       -41.311818 17.232973 -2.397  0.01718 *  
## youtube      4.058630  1.681080  2.414  0.01641 *  
## loading.time -4.490750  2.264113 -1.983  0.04830 *  
## X46x214     -18.477039 10.243269 -1.804  0.07235 .  
## .gif        -0.120050  0.078556 -1.528  0.12760  
## number_of_errors  0.013675  0.008952  1.528  0.12774
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 14 on 277 degrees of freedom
## Multiple R-squared:  0.8395, Adjusted R-squared:  0.8338
## F-statistic: 144.9 on 10 and 277 DF,  p-value: < 2.2e-16
```