

```

# math189 final

setwd("C:\\Users\\hsw11\\Desktop")
library(data.table)
mydata<-read("translate.csv")
pindata<-read("pindata.csv")

L = 1000 # number of samples
l = 914 # number of pindata

#####
## deal prob
data.prob = mydata$deal_prob$ability$prob
hist(data.prob,labels=T,ylm=c(0,800))
summary(data.prob)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 0.0000 0.0000 0.0000 0.1390 0.1679 1.0000
hist(data.prob, labels = T, ylim = c(0,700))
hist(data.prob[which(data.prob<0.1)], labels = T, ylim = c(0,640))

#####
## Price
data.Price.0 = mydata$price
ind.throw = c(which(data.Price.0==0),which(is.na(data.Price.0)))ind.throw
data.Price = data.Price.0[-ind.throw]
data.Price

length(data.Price) # 932 informative points
Price.row = seq(1,L,1)[-ind.throw] # price informative rows

hist(data.Price,breaks=30)
hist(data.Price[which(data.Price<0.5e07)],breaks=30)
hist(data.Price[which(data.Price<100000)],breaks=30, labels = T)

summary(data.Price)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 1 500 1500 342611 10000 36000000

#may be a GOF test to see this is exponential

#####
## date
data.Date = mydata$activation_date
data.Date

length(data.Date) # 1000
string.date = numeric(l)
for(i in 1:l){
  date = data.Date[i]
  # split the word and only take the last two numbers as the date of March
  string.date[i] = as.numeric(strsplit(date, "-")[, 2])
}
string.date
table(string.date)
# Month day 15 16 17 18 19 20 21 22 23 24 25 26 27 28
# Count 46 62 68 65 68 75 56 77 69 64 68 74 75 73
plot.monthday ~ barplot(table(string.date), ylim = c(0,90),
  main = "Barplot of activation date (month day)")
text(plot.monthday, 5, labels = as.matrix(table(string.date)))

week.day = weekdays(as.Date(data.Date))
week.day

length(week.day) # 1000
table(week.day)
# Friday Monday Saturday Sunday Thursday Tuesday Wednesday
# 132 150 153 142 151 129 143
week.table = data.frame(Mon=132, Tue=129, Wed=143, Thur=151, Fri=132,
  Sat=150, Sun=142) # reorder the data
plot.weekday ~ barplot(as.matrix(week.table), ylim = c(0,160),
  main = "Barplot of activation date (week day)")
text(plot.weekday, 5, labels = as.matrix(week.table))

## list of the activation_date is weekend or not
data.isWeekend=rep(1,1000)
for(i in 1:l){
  if (week.day[i]=="Saturday" || week.day[i]=="Sunday"){
    data.isWeekend[i] = 1
  }
  else{
    data.isWeekend[i] = 0
  }
}
data.isWeekend
table(data.isWeekend)
# 705 weekdays(Mon-Fri)
# 293 weekends
isweekend.table = data.frame(weekdays=705, weekends=293)

plot.isweekend ~ barplot(as.matrix(isweekend.table), ylim = c(0,710),
  main = "Barplot of activation date (weekdays/weekends)")
text(plot.isweekend, 60, labels = as.matrix(isweekend.table))

#####
## User type
data.user = mydata$user_type
data.user
table(data.user)
# Company Private Shop
# 271 681 48
plot.usertype ~ barplot(table(data.user), ylim = c(0,690),
  main = "Barplot of user type")
text(plot.usertype, 20, labels = as.matrix(table(data.user)))

#####
## with or without image
data.image = mydata$image
data.withimage=rep(0,1000)
for(i in 1:l){
  if (data.image[i] == ""){
    data.withimage[i] = 0
  }
  else{
    data.withimage[i] = 1
  }
}
data.withimage
table(data.withimage)

withimage.table = data.frame(with=914, without=86)
plot.withimage ~ barplot(as.matrix(withimage.table), ylim = c(0,920),
  main = "Barplot of using of image")
text(plot.withimage, 40, labels = as.matrix(withimage.table))

#####
## Text
#--- Word count function
wordcount <- function(x){
  apply(gregexpr("\\b\\w+\\b", str, perl=TRUE), function(x) sum(x)) + 1
}
#--- End

## description
data.description = mydata$en_desc
# word counts
data.des.WordCount = rep(0,L)
for(i in 1:l){
  # data.des.WordCount[i]= length(gregexpr("\\w+", data.description[i]))
  if (data.description[i] != "missing"){
    data.des.WordCount[i]= wordcount(data.description[i])
  }
}
data.des.WordCount
summary(data.des.WordCount)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 0.00 0.000 17.00 21.97 36.00 384.00
hist(data.des.WordCount,labels=T,ylim=c(0,900))
hist(data.des.WordCount[which(data.des.WordCount<50)],labels=T,ylim=c(0,900))

# Capital count
data.des.CapCount = rep(0,L)
for(i in 1:l){
  data.des.CapCount[i]= apply(gregexpr("\\b[A-Z][2,]\\b", data.description[i]),
    function(x) length(x[x>0]))
}
data.des.CapCount
summary(data.des.CapCount)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 0.000 0.000 0.000 0.905 0.000 39.000
hist(data.des.CapCount, labels = T, ylim = c(0,980))
hist(data.des.CapCount[which(data.des.CapCount<5)], labels = T, ylim = c(0,980))

# Digit count
data.des.DigitCount = rep(0,L)
for(i in 1:l){
  Digit<-gregexpr("^0-9","",data.description[i])
  # data.des.DigitCount[i]= length(gregexpr("\\w+", Digitx))
  if (length("","",Digitx) != ""){
    data.des.DigitCount[i]= wordcount(Digitx)
  }
}
data.des.DigitCount
summary(data.des.DigitCount)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 0.000 0.000 1.000 1.151 2.000 50.000
hist(data.des.DigitCount, labels = T, ylim = c(0,920))
hist(data.des.DigitCount[which(data.des.DigitCount<5)], labels = T, ylim = c(0,450))

#####
## Title
data.title = mydata$en_title
# word counts
data.ti.WordCount = rep(0,L)
for(i in 1:l){
  # data.ti.WordCount[i]= length(gregexpr("\\w+", data.title[i]))
  if (data.title[i] != "missing"){
    data.ti.WordCount[i]= wordcount(data.title[i])
  }
}
data.ti.WordCount
summary(data.ti.WordCount)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 1 1 3 12
ti.WordCount.table = table(data.ti.WordCount)
plot.ti.WordCount ~ barplot(tl.WordCount.table, main = "Barplot of word count for title")
text(plot.ti.WordCount, 4, labels = ti.WordCount.table)

# Capital count
data.ti.CapCount = rep(0,L)
for(i in 1:l){
  data.ti.CapCount[i]= apply(gregexpr("\\b[A-Z][2,]\\b", data.title[i]),
    function(x) length(x[x>0]))
}
data.ti.CapCount
summary(data.ti.CapCount)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 0.000 0.000 0.000 0.091 0.000 3.000
ti.CapCount.table = table(data.ti.CapCount)
plot.ti.CapCount ~ barplot(tl.CapCount.table, main = "Barplot of capital count for title")
text(plot.ti.CapCount, 40, labels = ti.CapCount.table)

# Digit count
data.ti.DigitCount = rep(0,L)
for(i in 1:l){
  Digit<-gregexpr("^0-9","",data.title[i])
  # data.ti.DigitCount[i]= length(gregexpr("\\w+", Digitx))
  if (length("","",Digitx) != ""){
    data.ti.DigitCount[i]= wordcount(Digitx)
  }
}
data.ti.DigitCount
summary(data.ti.DigitCount)
# Min. 1st Qu. Median Mean 3rd Qu. Max.
# 0.000 0.000 0.000 0.603 1.000 6.000
ti.DigitCount.table = table(data.ti.DigitCount)
plot.ti.DigitCount ~ barplot(tl.DigitCount.table, main = "Barplot of digit count for title")
text(plot.ti.DigitCount, 40, labels = ti.DigitCount.table)

#####
## Region and city
data.region = mydata$en_region

```

[illegible]

```
length(p.blurnee.100)(p.blurnee.100<0.000001) #142 deal prob not 0
length(p.blurnee.100)(p.blurnee.100>=0.000001) #87 deal prob not 0

matrx3<-matrix(c(length(p.blurnee.25)(p.blurnee.25<0.000001),
length(p.blurnee.25)(p.blurnee.25>=0.000001),
length(p.blurnee.50)(p.blurnee.50<0.000001),
length(p.blurnee.50)(p.blurnee.50>=0.000001),
length(p.blurnee.75)(p.blurnee.75<0.000001),
length(p.blurnee.75)(p.blurnee.75>=0.000001),
length(p.blurnee.100)(p.blurnee.100<0.000001),
length(p.blurnee.100)(p.blurnee.100>=0.000001)),
ncol=2,byrow=T)
T3<-as.data.frame(matrx3);T3
names(T3)<-c("deal prob is 0", "deal prob not 0");T3
rownames(T3)<-c("p<=0.79,0.9", ">0.79,0.9 and <=0.91,0.9", ">0.91,0.9 and <=0.95,0.9", ">0.95,0.9")

# deal prob is 0 deal prob not 0
# <=0.79,0.9 155 74
# >0.79,0.9 and <=0.91,0.9 155 74
# >0.91,0.9 and <=0.95,0.9 147 81
# >0.95,0.9 142 87

chiq.test(T3)
# X-squared = 2.448, df = 3, p-value = 0.4848
# FAIL TO REJECT

#####
# dollness
data.dollness<-importa$14$dollness;data.dollness
summary(data.dollness)
length(data.dollness)(data.dollness=0) #401 dollness=0 #313 dollness not 0

index.0.doll<-which(data.dollness=0);index.0.doll
prob.doll.0<-deal.prob$14(index.0.doll);length(prob.doll.0)
prob.doll.non.0<-deal.prob$14(-index.0.doll);length(prob.doll.non.0)

length(prob.doll.0)(prob.doll.0=0) #among 401 dollness=0, 394 deal prob =0
length(prob.doll.0)(prob.doll.0>0) # 207 deal prob not 0

length(prob.doll.non.0)(prob.doll.non.0=0) #among 313 dollness not 0, 205 deal prob =0
length(prob.doll.non.0)(prob.doll.non.0>0) # 108 Deal prob not 0

matrx4<-matrix(c(length(prob.doll.0)(prob.doll.0=0)),
length(prob.doll.0)(prob.doll.0>0)),
length(prob.doll.non.0)(prob.doll.non.0=0)),
length(prob.doll.non.0)(prob.doll.non.0>0)),
ncol=2,byrow=T)

matrx4
T4<-as.data.frame(matrx4);T4
names(T4)<-c("deal prob is 0", "deal prob not 0");T4
rownames(T4)<-c("dollness=0", "dollness larger than 0");T4

# dollness=0 deal prob is 0 deal prob not 0
# dollness larger than 0 394 207
# dollness larger than 0 205 108
chiq.test(T4)
#X-squared = 2.4722=30, df = 1, p-value = 1

#####
#####
# image size
data.imp.size<-importa$14$image_size;data.imp.size
summary(data.imp.size)
quantile(data.imp.size,seq(0.1,0.9,0.1))

# 10% 20% 30% 40% 50% 60% 70% 80% 90%
# 20094.6 24590.4 28770.8 32937.0 37025.5 40999.4 44986.2 50786.4 55066.3
index.size.10<-which(data.imp.size>20094.6);length(index.size.10) #92 imp size <=20094.6
p.imp.size.10<-deal.prob$14(index.size.10)

length(p.imp.size.10)(p.imp.size.10<0.000001) #61 deal prob 0
length(p.imp.size.10)(p.imp.size.10>=0.000001) #31 deal prob not 0

index.size.20<-which(data.imp.size>20094.6 & data.imp.size<=24590.4);length(index.size.20) #91 imp size >20094.6 & <=24590.4
p.imp.size.20<-data.imp.size[index.size.20]

length(p.imp.size.20)(p.imp.size.20<0.000001) #56 deal prob 0
length(p.imp.size.20)(p.imp.size.20>=0.000001) #35 deal prob not 0

index.size.30<-which(data.imp.size>24590.4 & data.imp.size<=28770.8);length(index.size.30) #91 imp size >24590.4 & <=28770.8
p.imp.size.30<-data.imp.size[index.size.30]
p.imp.size.30<-deal.prob$14(index.size.30)

length(p.imp.size.30)(p.imp.size.30<0.000001) #63 deal prob 0
length(p.imp.size.30)(p.imp.size.30>=0.000001) #28 deal prob not 0

index.size.40<-which(data.imp.size>28770.8 & data.imp.size<=32937.0);length(index.size.40) #93 imp size >28770.8 & <=32937.0
p.imp.size.40<-data.imp.size[index.size.40]
p.imp.size.40<-deal.prob$14(index.size.40)

length(p.imp.size.40)(p.imp.size.40<0.000001) #57 deal prob 0
length(p.imp.size.40)(p.imp.size.40>=0.000001) #36 deal prob not 0

index.size.50<-which(data.imp.size>32937.0 & data.imp.size<=37025.5);length(index.size.50) #90 imp
p.imp.size.50<-data.imp.size[index.size.50]
p.imp.size.50<-deal.prob$14(index.size.50)

length(p.imp.size.50)(p.imp.size.50<0.000001) #64 deal prob 0
length(p.imp.size.50)(p.imp.size.50>=0.000001) #26 deal prob not 0

index.size.60<-which(data.imp.size>37025.5 & data.imp.size<=40999.4);length(index.size.60) #91 imp
p.imp.size.60<-data.imp.size[index.size.60]
p.imp.size.60<-deal.prob$14(index.size.60)

length(p.imp.size.60)(p.imp.size.60<0.000001) #61 deal prob 0
length(p.imp.size.60)(p.imp.size.60>=0.000001) #30 deal prob not 0

index.size.70<-which(data.imp.size>40999.4 & data.imp.size<=44986.2);length(index.size.70) #92 imp
p.imp.size.70<-data.imp.size[index.size.70]
p.imp.size.70<-deal.prob$14(index.size.70)

length(p.imp.size.70)(p.imp.size.70<0.000001) #60 deal prob 0
length(p.imp.size.70)(p.imp.size.70>=0.000001) #32 deal prob not 0

index.size.80<-which(data.imp.size>44986.2 & data.imp.size<=50786.4);length(index.size.80) #91 imp
p.imp.size.80<-data.imp.size[index.size.80]
p.imp.size.80<-deal.prob$14(index.size.80)

length(p.imp.size.80)(p.imp.size.80<0.000001) #59 deal prob 0
length(p.imp.size.80)(p.imp.size.80>=0.000001) #32 deal prob not 0

index.size.90<-which(data.imp.size>50786.4 & data.imp.size<=55066.3);length(index.size.90) #91 imp
p.imp.size.90<-data.imp.size[index.size.90]
p.imp.size.90<-deal.prob$14(index.size.90)

length(p.imp.size.90)(p.imp.size.90<0.000001) #62 deal prob 0
length(p.imp.size.90)(p.imp.size.90>=0.000001) #29 deal prob not 0

index.size.100<-which(data.imp.size>55066.3);length(index.size.100) #92 imp
p.imp.size.100<-data.imp.size[index.size.100]
p.imp.size.100<-deal.prob$14(index.size.100)

length(p.imp.size.100)(p.imp.size.100<0.000001) #56 deal prob 0
length(p.imp.size.100)(p.imp.size.100>=0.000001) #36 deal prob not 0

matrx5<-matrix(c(61,31,56,35,63,28,57,36,64,26,61,30,60,32,59,32,62,29,56,36),ncol=2,byrow=T)
matrx5
T5<-as.data.frame(matrx5);T5
names(T5)<-c("deal prob is 0", "deal prob not 0");T5
rownames(T5)<-c("p<=0.1 quantile", ">0.1,0.2 quantile", ">0.2,0.3 quantile",
">0.3,0.4 quantile", ">0.4,0.5 quantile", ">0.5,0.6 quantile",
">0.6,0.7 quantile", ">0.7,0.8 quantile", ">0.8,0.9 quantile",
">0.9,1 quantile")

T5
chiq.test(T5)
# X-squared = 4.4713, df = 9, p-value = 0.8777

#####
#####
# image width
data.imp.width<-importa$14$width
summary(data.imp.width)
#Min. 1st Qu. Median Mean 3rd Qu. Max.
#108.0 365.0 365.0 415.8 480.0 640.0

quantile(data.imp.width,seq(0.1,0.9,0.1))

table(data.imp.width)
#108 154 161 185 215 222 245 245 250 255 259 266 267 270 271 274 275 279 281 284 287 288
#1 1 1 1 1 1 1 1 1 1 1 1 1 1 96 2 1 1 1 1 1 1 7

#289 295 298 300 301 303 306 313 315 317 318 319 320 322 323 334 343 350 355 356 357 358 359
#2 1 1 1 1 2 2 2 2 2 2 2 3 38 1 1 1 2 3 2 5 2 6 8

#560 380 385 392 394 400 401 402 409 415 427 430 432 439 441 442 443 444 450 461 454 456 458
#325 1 1 1 2 2 1 2 1 2 1 1 2 1 1 1 1 1 1 1 1 2 1 1 1

#470 472 477 480 481 482 483 485 486 493 494 498 506 507 521 527 531 532 533 537 538 540 541
#1 1 1 206 4 8 1 1 1 1 3 1 1 1 2 1 1 1 1 1 1 1 2 1 15 6

#642 543 549 557 559 564 569 573 580 589 597 600 602 603 608 628 624 638 640
#4 2 1 3 1 3 1 2 1 1 1 3 1 1 1 1 1 1 1 6 80

index.imp.width.first<-which(data.imp.width<300); length(index.imp.width.first) #129 width <=300
index.imp.width.second<-which(data.imp.width>=400 & data.imp.width<=500); length(index.imp.width.second) #392 width >=400 and <=500
index.imp.width.third<-which(data.imp.width>=600 & data.imp.width<=800); length(index.imp.width.third) #248 width >=600 and <=800
index.imp.width.fourth<-which(data.imp.width>=900); length(index.imp.width.fourth) #144 width >=900

p.imp.width.first<-deal.prob$14(index.imp.width.first); length(p.imp.width.first)
p.imp.width.second<-deal.prob$14(index.imp.width.second); length(p.imp.width.second)
p.imp.width.third<-deal.prob$14(index.imp.width.third); length(p.imp.width.third)
p.imp.width.fourth<-deal.prob$14(index.imp.width.fourth); length(p.imp.width.fourth)

length(p.imp.width.first)(p.imp.width.first=0) #among 129 width<=300, 91 deal prob =0
length(p.imp.width.first)(p.imp.width.first>0) # 38 not 0

length(p.imp.width.second)(p.imp.width.second=0) #among 392, 293 is zero
length(p.imp.width.second)(p.imp.width.second>0) # 99 not 0

length(p.imp.width.third)(p.imp.width.third=0) #among 248, 137 deal prob =0
length(p.imp.width.third)(p.imp.width.third>0) # 112 not 0

length(p.imp.width.fourth)(p.imp.width.fourth=0) #among 144, 78 deal prob =0
length(p.imp.width.fourth)(p.imp.width.fourth>0) # 66 not 0

matrx6<-matrix(c(length(p.imp.width.first)(p.imp.width.first=0)),
length(p.imp.width.second)(p.imp.width.second=0)),
length(p.imp.width.third)(p.imp.width.third=0)),
length(p.imp.width.fourth)(p.imp.width.fourth=0)),
length(p.imp.width.fourth)(p.imp.width.fourth>0)),byrow=T,ncol=2)

matrx6
T6<-as.data.frame(matrx6);T6
names(T6)<-c("deal prob is 0", "deal prob not 0");T6
rownames(T6)<-c("width <=300", "width >=300 and <=400", "width >=400 and <=500", "width >=500");T6

# width <=300 91 38
# width >=300 and <=400 293 99
# width >=400 and <=500 137 112
# width >=500 78 66

chiq.test(T6)
# X-squared = 36.582, df = 3, p-value = 5.66e-08
# X-squared = 36.582, df = 3, p-value = 5.66e-08

#####
#####
data.imp.height<-importa$14$height
summary(data.imp.height)
table(data.imp.height)

data.imp.height
# 165 179 184 187 206 213 215 225 250 261 268 296 305 312 313 314 319 320 322 340 349 350 359 360 361 362 366
# 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 2 1 1 1 3 1 1 1 2 7 414 1 1 1

# 371 374 377 379 381 383 386 399 400 403 406 408 413 414 417 420 424 428 435 437 441 447 450 452 454 456 460
# 1 2 1 1 1 1 0 1 2 1 1 2 2 1 1 1 1 1 1 1 2 1 1 2 1 1 1

# 461 479 480
# 1 2 428

index.imp.height.first<-which(data.imp.height<=400); length(index.imp.height.first) #462 height<=400
index.imp.height.second<-which(data.imp.height>=500); length(index.imp.height.second) #425 height>=500

p.imp.height.first<-deal.prob$14(index.imp.height.first); length(p.imp.height.first)
p.imp.height.second<-deal.prob$14(index.imp.height.second);length(p.imp.height.second)

length(p.imp.height.first)(p.imp.height.first=0) #among 462 height <=400, 266 deal prob =0
```

```
length(p.imp.height.first[p.imp.height.first==0]) # 196 not 0

length(p.imp.height.second[p.imp.height.second==0]) #among 452 height >400, 332 deal prob =0
length(p.imp.height.second[p.imp.height.second==0]) # 119 not 0

matriz<-matrix(c(length(p.imp.height.first[p.imp.height.first==0]),
length(p.imp.height.first[p.imp.height.first==1]),
length(p.imp.height.second[p.imp.height.second==0]),
length(p.imp.height.second[p.imp.height.second==1])),nrow=2,byrow=T)

matriz?

T<-as.data.frame(matriz);?T
names(T)<-c("deal prob is 0", "deal prob not 0");?T
rownames(T)<-c("height <=400", "height>400");?T
chisq.test(T)

#          deal prob is 0    deal prob not 0
# height <=400             266             196
# height>400               333             119

summary(lm(data.imp.height~data.imp.width))
#####
# Certain type of dimension of image

index.width.270<-which(data.imp.width==270);length(index.width.270) #96 pic with width 270
width.270.height<-data.imp.height[index.width.270];width.270.height
length(width.270.height,width.270,height==480) #all 96 pics are 270*480

index.width.360<-which(data.imp.width==360);length(index.width.360) #325 pic with width 360
width.360.height<-data.imp.height[index.width.360];width.360.height
summary(width.360.height)
length(width.360.height,width.360,height==480) # 258 pics are 360*480
index.360.480<-which(width.360.height==480);length(index.360.480)

p.360.480<-deal.prob94[index.360.480];length(p.360.480) # 258 imp dim 360*480
p.270.480<-deal.prob94[index.width.270];length(p.270.480) # 96 imp dim 270*480

length(p.360.480[p.360.480==0]) #among 258 imp 360*480, 171 deal prob 0
length(p.360.480[p.360.480==1]) #97 deal prob not 0

length(p.270.480[p.270.480==0]) #among 96 imp 270*480, 67 deal prob 0
length(p.270.480[p.270.480==1]) #29 deal prob not 0

matriz<-matrix(c(length(p.360.480[p.360.480==0]),
length(p.360.480[p.360.480==1]),
length(p.270.480[p.270.480==0]),
length(p.270.480[p.270.480==1])),nrow=2,byrow=T)

matriz?
T<-as.data.frame(matriz);?T
names(T)<-c("deal prob is 0", "deal prob not 0");?T
rownames(T)<-c("<=360*480", ">270*480");?T
#          deal prob is 0    deal prob not 0
# <=360*480             171             87
# >270*480              67             29
# 360*480               112             94
# 270*480               67             29

chisq.test(T)
# X-squared = 0.24862, df = 1, p-value = 0.618

#####
# width 480

index.width.480<-which(data.imp.width==480);length(index.width.480) #206 pic with width 480
width.480.height<-data.imp.height[index.width.480];length(width.480.height==360) # all 206 pics dim 480*360
p.480.360<-deal.prob94[index.width.480]
length(p.480.360[p.480.360==0]) #among 206 480*360 pics, 112 deal prob 0
length(p.480.360[p.480.360==1]) # 94 not 0

#####
# width 640

index.width.640<-which(data.imp.width==640);length(index.width.640) #80 pic with width 640
width.640.height<-data.imp.height[index.width.640]
table(width.640.height) #mainly 640*360
# 265 296 359 360
# 1 1 3 71

index.640.360<-which(width.640.height==360);length(index.640.360)

p.640.360<-deal.prob94[index.640.360];length(p.640.360) #71 imp dim 640*360
length(p.640.360[p.640.360==0]) #74 deal prob 0
length(p.640.360[p.640.360==1]) # 17 deal prob not 0

matriz<-matrix(c(length(p.360.480[p.360.480==0]),
length(p.360.480[p.360.480==1]),
length(p.270.480[p.270.480==0]),
length(p.270.480[p.270.480==1]),
length(p.480.360[p.480.360==0]),
length(p.480.360[p.480.360==1]),
length(p.640.360[p.640.360==0]),
length(p.640.360[p.640.360==1])),nrow=2,byrow=T)

matriz?
T<-as.data.frame(matriz);?T
names(T)<-c("deal prob is 0", "deal prob not 0");?T
rownames(T)<-c("<=360*480", ">270*480", "<=640*360", ">480*360");?T
#          deal prob is 0    deal prob not 0
#          360*480             171             87
#          270*480             67             29
#          480*360            112             94
#          640*360             54             17

chisq.test(T)
# X-squared = 14.756, df = 3, p-value = 0.002937

#####
# region

data.whiteness<-lm(data$R[whiteness];data.whiteness)
summary(data.whiteness)
length(data.whiteness[data.whiteness==0]) #599 whiteness=0 #315 whiteness not 0

index.0.white<-which(data.whiteness==0);index.0.white
prob.white.0<-deal.prob94[index.0.white];length(prob.white.0)
prob.white.non.0<-deal.prob94[-index.0.white];length(prob.white.non.0)

length(prob.white.0[prob.white.0==0]) #among 599 whiteness=0, 415 deal prob =0
length(prob.white.0[prob.white.0==1]) # 184 deal prob not 0

length(prob.white.non.0[prob.white.non.0==0]) #among 315 whiteness not 0, 184 deal prob =0
length(prob.white.non.0[prob.white.non.0==1]) # 131 deal prob not 0

matriz<-matrix(c(length(prob.white.0[prob.white.0==0]),
length(prob.white.0[prob.white.0==1]),
length(prob.white.non.0[prob.white.non.0==0]),
length(prob.white.non.0[prob.white.non.0==1])),
nrow=2,byrow=T)

matriz?
T<-as.data.frame(matriz);?T
names(T)<-c("deal prob is 0", "deal prob not 0");?T
rownames(T)<-c("whiteness=0", "whiteness larger than 0");?T
#          deal prob is 0    deal prob not 0
# whiteness=0              415             184
# whiteness larger than 0  184             131

chisq.test(T)
# X-squared = 10.323, df = 1, p-value = 0.001314

#####
# region

table(data.region)

#Altai region (Siberia) Bashkortostan (Volga)
#25 50
#Belgorod region (Central) Chelyabinsk region (Ural)
#18 53
#Irkutsk region (Siberia) Kaliningrad region (Northwest)
#30 23
#Kemerovo Region (Siberia) Khanty-Mansiysk Autonomous Okrug (Ural)
#21 16
#Krasnodar region (South) Krasnoyarsk region (Siberia)
#104 39
#Nizhny Novgorod Region (Volga) Novosibirsk region (Siberia)
#48 46
#Omsk Region (Siberia) Orenburg region (Volga)
#37 25
#Perm Region (Volga) Rostov region (South)
#16 58
#Samara Region (Volga) Saratov region (Volga)
#48 35
#Stavropol region (North Caucasus) Sverdlovsk region (Ural)
#14 57
#Tatarstan (Volga) Tula region (Central)
#45 18
#Tyumen region (Ural) Udmurtia (Volga)
#23 16
#Vladimir region (Central) Volgograd region (South)
#16 27
#Voronezh region (Central) Yarovavl region (Central)
#23 23

### Siberia\Altai region\ Irkutsk region\ Kemerovo Region\ Krasnoyarsk region\
# Novosibirsk region\ Omsk Region\

### Volga: Bashkortostan\ Nizhny Novgorod Region\ Orenburg region\
# Perm Region\ Samara Region\ Saratov region\ Tatarstan\ Udmurtia

### Central: Belgorod region\ Tula region\ Vladimir region\ Voronezh region\
# Yarovavl region

### Ural: Chelyabinsk region\ Khanty-Mansiysk Autonomous Okrug\ Sverdlovsk region
# Tyumen region

### Northwest: Kaliningrad region\

### South: Krasnodar region\ Rostov region\ Volgograd region

### North Caucasus: Stavropol region

length(table(data.region)) #28 regions in 6 Federal District
wydataFederal.District<-rep(NA,1000)

for (i in 1:1000)
{
  if (data.region[i]=="Altai region" | data.region[i]=="Irkutsk region"
| data.region[i]=="Kemerovo Region" | data.region[i]=="Krasnoyarsk region"
| data.region[i]=="Novosibirsk region" | data.region[i]=="Omsk Region")
{
  wydataFederal.District[i]<"Siberia"
}

  else if (data.region[i]=="Bashkortostan" | data.region[i]=="Nizhny Novgorod Region"
| data.region[i]=="Orenburg region" | data.region[i]=="Perm Region"
| data.region[i]=="Samara Region" | data.region[i]=="Saratov region"
| data.region[i]=="Tatarstan" | data.region[i]=="Udmurtia")
{
  wydataFederal.District[i]<"Volga"
}

  else if (data.region[i]=="Belgorod region" | data.region[i]=="Tula region"
| data.region[i]=="Vladimir region" | data.region[i]=="Voronezh region"
| data.region[i]=="Yarovavl region")
{
  wydataFederal.District[i]<"Central"
}

  else if (data.region[i]=="Chelyabinsk region" | data.region[i]=="Khanty-Mansiysk Autonomous Okrug"
| data.region[i]=="Sverdlovsk region" | data.region[i]=="Tyumen region")
{
  wydataFederal.District[i]<"Ural"
}

  else if (data.region[i]=="Kaliningrad region" )
{
  wydataFederal.District[i]<"Northwest"
}

  else if (data.region[i]=="Krasnodar region" | data.region[i]=="Rostov region"
| data.region[i]=="Volgograd region")
{
  wydataFederal.District[i]<"South"
}

  else if (data.region[i]=="Stavropol region")
{
  wydataFederal.District[i]<"North Caucasus"
}
}

Fed.6<-wydataFederal.District

index.Volga<-which(Fed.6=="Volga");length(index.Volga) #289 in Volga Federal District
index.Ural <-which(Fed.6=="Ural"); length(index.Ural) #153 in Ural Federal District
index.Siberia <-which(Fed.6=="Siberia"); length(index.Siberia) #438 in Siberia Federal District
index.Central <-which(Fed.6=="Central"); length(index.Central) #106 in Central Federal District
index.South <-which(Fed.6=="South"); length(index.South) #497 in South Federal District
index.Northwest <-which(Fed.6=="Northwest"); length(index.Northwest) #23 in Northwest Federal District
index.North.Caucasus <-which(Fed.6=="North Caucasus"); length(index.North.Caucasus) #34 in North Caucasus Federal District

deal.prob1000<-wydataDeal_probability
```



```

data.title.wordCount<-mydata2$data.t1.WordCount
quantile(data.title.wordCount,seq(0.2,0.8,0.2))

index.title.word.count.20<-which(data.title.wordCount<=0)length(index.title.word.count.20)
index.title.word.count.40<-which(data.title.wordCount<=4)length(index.title.word.count.40)
index.title.word.count.60<-which(data.title.wordCount<=4)length(index.title.word.count.60)
index.title.word.count.80<-which(data.title.wordCount<=4)length(index.title.word.count.80)
index.title.word.count.100<-which(data.title.wordCount<=4)length(index.title.word.count.100)

p.title.word.count.20<-deal.prob1000(index.title.word.count.20)
p.title.word.count.40<-deal.prob1000(index.title.word.count.40)
p.title.word.count.60<-deal.prob1000(index.title.word.count.60)
p.title.word.count.80<-deal.prob1000(index.title.word.count.80)
p.title.word.count.100<-deal.prob1000(index.title.word.count.100)

length(p.title.word.count.20) ##333
PWT20.0<-length(p.title.word.count.20[p.title.word.count.20==0])PWT20.0 #239
PWT20.1<-length(p.title.word.count.20[p.title.word.count.20!=0])PWT20.1 #54

length(p.title.word.count.40) ##199
PWT40.0<-length(p.title.word.count.40[p.title.word.count.40==0])PWT40.0 #119
PWT40.1<-length(p.title.word.count.40[p.title.word.count.40!=0])PWT40.1 #60

length(p.title.word.count.60) ##110
PWT60.0<-length(p.title.word.count.60[p.title.word.count.60==0])PWT60.0 #98
PWT60.1<-length(p.title.word.count.60[p.title.word.count.60!=0])PWT60.1 #52

length(p.title.word.count.80) ##147
PWT80.0<-length(p.title.word.count.80[p.title.word.count.80==0])PWT80.0 #63
PWT80.1<-length(p.title.word.count.80[p.title.word.count.80!=0])PWT80.1 #54

length(p.title.word.count.100) ##171
PWT100.0<-length(p.title.word.count.100[p.title.word.count.100==0])PWT100.0 #78
PWT100.1<-length(p.title.word.count.100[p.title.word.count.100!=0])PWT100.1 #93

matrix1<-matrix(c(PWT20.0,PWT40.0,PWT60.0,PWT80.0,PWT100.0,PWT20.0,PWT40.0,PWT60.0,PWT80.0,PWT100.0,1,1,
                    ncol=4,byrow=T);matrix1)
T1<-as.data.frame(matrix1);T1
names(T1)<-c("deal prob is 0", "deal prob not 0"); T1
rownames(T1)<-c("title word count <=4", "title word count >= 4 <=3", "title word count >= 4 <=4", "title word count >= 4 <=6", "title word count >=6");T1

#
# title word count <=2      deal prob is 0      deal prob not 0
# title word count >= 2 <=3      239      54
# title word count >= 2 <=4      119      60
# title word count >= 2 <=4      98      52
# title word count >= 4 <=6      93      54
# title word count >=6      78      93

chiSq.test(T1)

setwd("C:\\Users\\hsaw1\\Desktop")
library(data.table)

mydata<-read("translate.csv")
mydata<-read("picdata_mw_1000.csv")
newdata<-read("newdata.csv")
data.prob = newdata$data.prob

Boostarp_CT_Method <- function(data1,data2,n){
  u = mean(data1)
  N = length(data1)
  simulation<-rep(NA,n)
  for (i in 1:n)
  {
    btample = sample(data2,N,replace = F)
    simulation[i]= mean(btample)
  }
  lower = quantile(simulation,0.05,na.rm = T)
  upper = quantile(simulation,0.95,na.rm = T)
  inside = NA
  if(u >=lower is u <= upper){
    inside = T
  }
  else{
    inside = F
  }
  return(inside)
}

Boostarp_Derr_Method <- function(data1,data2,n){
  u = length(which(data1 == 0))
  N = length(data1)
  simulation<-rep(NA,n)
  for (i in 1:n)
  {
    btample = sample(data2,N,replace = F)
    simulation[i]= length(which(btample == 0))
  }
  lower = quantile(simulation,0.05,na.rm = T)
  upper = quantile(simulation,0.95,na.rm = T)
  inside = NA
  if(u >=lower is u <= upper){
    inside = T
  }
  else{
    inside = F
  }
  return(inside)
}

###
### price
###
summary(newdata$data.Price.0[is.na(newdata$data.Price.0)=F])
{
  category = 10
  quant = quantile(data.price,seq(0,1,by = 1/category))
  list = rep(0,category)
  for (i in 1:category){
    cate_prob = data.prob[which(data.price> quant[i] & data.price <= quant[i+1])]
    list[i] = Boostarp_CT_Method(cate_prob,data.prob,500)
  }
  list
} # not rejecting

###
### weekend date on deal_pro
###
index.weekend = which(newdata$data.IsWeekend == 1)
weekend_pro = data.prob[index.weekend]
Boostarp_CT_Method(weekend_pro,data.prob,500) #not reject

# mean(weekend_pro)
# mean(data.prob)

index.weekday = which(newdata$data.IsWeekend == 0)
weekday_pro = data.prob[index.weekend]
Boostarp_CT_Method(weekday_pro,data.prob[-index.weekend],500) #not reject

###
### User Type
###
index.type = which(newdata$data.UserType == 0)
type1_prob = data.prob[index.type]
Boostarp_CT_Method(type1_prob,data.prob,500) #not reject
index.type = which(newdata$data.UserType == 1)
type2_prob = data.prob[index.type]
Boostarp_CT_Method(type1_prob,data.prob,500) #not reject
index.type = which(newdata$data.UserType == 2)
type2_prob = data.prob[index.type]
Boostarp_CT_Method(type1_prob,data.prob,500) #not reject

###
### description
###
summary(newdata$data.dea.WordCount)
wordcount = newdata$data.dea.WordCount
{
  category = 4
  quant = quantile(wordcount,seq(0,1,by = 1/category));quant
  list = rep(0,category)
  for (i in 1:category){
    if (i == 1){
      cate_prob = data.prob[which(wordcount <= quant[i+1])]
    }
    else{
      cate_prob = data.prob[which(wordcount> quant[i] & wordcount <= quant[i+1])]
    }
    list[i] = Boostarp_CT_Method(cate_prob,data.prob,500)
  }
  list
} # rejecting only on the lower 25 %
index.word = which(wordcount> quant[1] & wordcount <= quant[2])
mean(data.prob[index.word])
mean(data.prob) # the lower 25 % (description word count smaller than 7) has a smaller deal prob

### digit count
summary(newdata$data.dea.DigitCount)
digitcount = newdata$data.dea.DigitCount
index.digit = which(digitcount == 0)
digit0_prob = data.prob[index.digit]
Boostarp_CT_Method(type0_prob,data.prob,500) #not reject
index.digit = which(digitcount > 0 & digitcount<= 1)
digit1_prob = data.prob[index.digit]
Boostarp_CT_Method(digit1_prob,data.prob,500) #not reject
index.digit = which(digitcount > 1)
digit2_prob = data.prob[index.digit]
Boostarp_CT_Method(digit2_prob,data.prob,500) #not reject

### capital count
summary(newdata$data.dea.CapacCount)
capacount = newdata$data.dea.CapacCount
index.caga = which(capacount == 0)
capal_prob = data.prob[index.caga]
Boostarp_CT_Method(capal_prob,data.prob,500) # reject
mean(capal_prob)
mean(data.prob) # the description without capalock has a smaller deal prob
index.caga = which(capacount > 0)
capal_prob = data.prob[index.caga]
Boostarp_CT_Method(capal_prob,data.prob,500) # reject
mean(capal_prob)
mean(data.prob) # the description with capalock has a bigger deal prob

###
### title
###
summary(newdata$data.t1.WordCount)
wordcount = newdata$data.t1.WordCount
{
  category = 4
  quant = quantile(wordcount,seq(0,1,by = 1/category));quant
  list = rep(0,category)
  for (i in 1:category){
    if (i == 1){
      cate_prob = data.prob[which(wordcount <= quant[i+1])]
    }
    else{
      cate_prob = data.prob[which(wordcount> quant[i] & wordcount <= quant[i+1])]
    }
    list[i] = Boostarp_CT_Method(cate_prob,data.prob,500)
  }
  list
} # not rejecting

### digit count
summary(newdata$data.t1.DigitCount)
digitcount = newdata$data.t1.DigitCount
index.digit = which(digitcount == 0)
digit0_prob = data.prob[index.digit]
Boostarp_CT_Method(type0_prob,data.prob,500) #not reject
index.digit = which(digitcount > 0)
digit2_prob = data.prob[index.digit]
Boostarp_CT_Method(digit2_prob,data.prob,500) # reject
mean(digit2_prob)
mean(data.prob) # the title with digits has a bigger deal prob

```

```

### capital count
summary(picdata$capital)
capcount = newdata$capital
index.caps = which(capcount == 0)
caps_prob = data.prob[index.caps]
Boostarp_CT_Method(caps_prob, data, prob, 500) # reject
mean(caps_prob)
# the title without caplock has a smaller deal prob

index.caps = which(capcount > 0)
caps_prob = data.prob[index.caps]
Boostarp_CT_Method(caps_prob, data, prob, 500) # reject
mean(caps_prob)
# the description with caplock has a bigger deal prob
mean(data_prob)

### pic.klar on deal_prob
###
summary(picdata$blurness)
blurness = picdata$blurness
{
  category = 5
  quant = quantile(blurness, seq(0,1,by = 1/category), na.rm = T) : quant
  list = rep(0, category)
  for (i in 1:category){
    if (i == 1){
      cate_prob = data_prob[which(blurness <= quant[i+1])]
    }
    else{
      cate_prob = data_prob[which(blurness> quant[i] & blurness <= quant[i+1])]
    }
  }
  list[i] = Boostarp_CT_Method(cate_prob, data, prob, 500)
}
list
} # not rejecting

### dullness
###
summary(picdata$dullness)
dullness = picdata$dullness
category = 2
quant = quantile(dullness, seq(0,1,by = 1/category), na.rm = T) : quant
list = rep(0, category)
for (i in 1:category){
  if (i == 1){
    cate_prob = data_prob[which(dullness <= quant[i+1])]
  }
  else{
    cate_prob = data_prob[which(dullness> quant[i] & dullness <= quant[i+1])]
  }
  list[i] = Boostarp_CT_Method(cate_prob, data, prob, 500)
}
list
} # not rejecting

### whiteness
###
summary(picdata$whiteness)
whiteness = picdata$whiteness
category = 2
quant = quantile(whiteness, seq(0,1,by = 1/category), na.rm = T) : quant
list = rep(0, category)
for (i in 1:category){
  if (i == 1){
    cate_prob = data_prob[which(whiteness <= quant[i+1])]
  }
  else{
    cate_prob = data_prob[which(whiteness> quant[i] & whiteness <= quant[i+1])]
  }
  list[i] = Boostarp_CT_Method(cate_prob, data, prob, 500)
}
list
} # not rejecting

### size
###
summary(picdata$image_size)
size = picdata$image_size
category = 4
quant = quantile(size, seq(0,1,by = 1/category), na.rm = T) : quant
list = rep(0, category)
for (i in 1:category){
  if (i == 1){
    cate_prob = data_prob[which(size <= quant[i+1])]
  }
  else{
    cate_prob = data_prob[which(size> quant[i] & size <= quant[i+1])]
  }
  list[i] = Boostarp_CT_Method(cate_prob, data, prob, 500)
}
list
} # rejecting on the (size bigger than 24931.5)
mean(data_prob[which(size > quant[i+1])])
mean(data_prob) #the bigger pictures (upper 25% bigger than 47860) have a smaller deal prob

###
### width
###
summary(picdata$width)
width = picdata$width
{
  category = 3
  quant = quantile(width, seq(0,1,by = 1/category), na.rm = T) : quant
  list = rep(0, category)
  for (i in 1:category){
    if (i == 1){
      cate_prob = data_prob[which(width <= quant[i+1])]
    }
    else{
      cate_prob = data_prob[which(width> quant[i] & width <= quant[i+1])]
    }
  }
  list[i] = Boostarp_CT_Method(cate_prob, data, prob, 500)
}
list
} # rejecting on all width

###
### height
###
summary(picdata$height)
height = picdata$height
{
  category = 2
  quant = quantile(height, seq(0,1,by = 1/category), na.rm = T) : quant
  list = rep(0, category)
  for (i in 1:category){
    if (i == 1){
      cate_prob = data_prob[which(height <= quant[i+1])]
    }
    else{
      cate_prob = data_prob[which(height> quant[i] & height <= quant[i+1])]
    }
  }
  list[i] = Boostarp_CT_Method(cate_prob, data, prob, 500)
}
list
} # rejecting on all height

###
### temp_size
###
summary(picdata$temp_size)
temp_size = picdata$temp_size
index.size = which(temp_size == "(360, 480)")
size1_prob = data_prob[index.size]
Boostarp_CT_Method(size1_prob, data, prob, 500) # not reject
index.size = which(temp_size == "(270, 480)")
size1_prob = data_prob[index.size]
Boostarp_CT_Method(size1_prob, data, prob, 500) # not reject
index.size = which(temp_size == "(480, 360)")
size2_prob = data_prob[index.size]
Boostarp_CT_Method(size2_prob, data, prob, 500) # reject
Boostarp_Merc_Method(size2_prob, data, prob, 500)
mean(size1_prob)
mean(data_prob) # higher deal pro mean in this category
length(which(size2_prob == 0)) / length(size2_prob) # higher deal pro non zero rate in this category

index.size = which(temp_size == "(640, 360)")
size3_prob = data_prob[index.size]
Boostarp_CT_Method(size3_prob, data, prob, 500) # not reject
Boostarp_Merc_Method(size3_prob, data, prob, 500)
length(which(size3_prob == 0)) / length(size3_prob)

setwd("C:\\Users\\hsw1\\Desktop")

library(data.table)
mydata<-fread("translate.csv")
picdata<-fread("picdata_merc_1000.csv")
newdata<-fread("newdata.csv")

data_prob = newdata$capital
data_region = mydata$caption
data_price = mydata$price

index.throw = which(is.na(data_price))~T | data_price == 0)
data.lprice = log(data_price[-index.throw])

table(mydata$owner_id)
table(mydata$owner_type)

Boostarp_CT_Method <- function(data1, data2, n){
  u = mean(data1)
  N = length(data1)
  simulation<-rep(0,n)
  for (i in 1:n)
  {
    btsample = sample(data2, N, replace = F)
    simulation[i] = mean(btsample)
  }
  lower = quantile(simulation, 0.05, na.rm = T)
  upper = quantile(simulation, 0.95, na.rm = T)
  inside = NA
  if(u >= lower && u <= upper){
    inside = T
  }
  else{
    inside = F
  }
  return(inside)
}

Boostarp_Merc_Method <- function(data1, data2, n){
  u = length(which(data1 == 0))
  N = length(data1)
  simulation<-rep(0,n)
  for (i in 1:n)
  {
    btsample = sample(data2, N, replace = F)
    simulation[i] = length(which(btsample == 0))
  }
  lower = quantile(simulation, 0.05, na.rm = T)
  upper = quantile(simulation, 0.95, na.rm = T)
  inside = NA
  if(u >= lower && u <= upper){
    inside = T
  }
  else{
    inside = F
  }
  return(inside)
}

###
### with/ without image
###
index.image = which(is.na(picdata$image))~T)
image0_prob = data_prob[index.image]
image1_prob = data_prob[-index.image]

```

```

#bootstarp_C1_Method(image0,prob,data,prob,500) #not reject
#bootstarp_Zero_Method(image0,prob,data,prob,500) #reject

length(which(image0,prob!=0))/length(image0,prob) #0.6744186
length(which(data,prob!=0))/length(data,prob) #0.373

#bootstarp_C1_Method(image1,prob,data,prob,500) #not reject
#bootstarp_Zero_Method(image1,prob,data,prob,500) #reject

length(which(image1,prob!=0))/length(image1,prob) #0.3446389
length(which(data,prob !=0))/length(data,prob) #0.373

#kruskal_type
#bootstarp_C1_Method(image0,prob,image1,prob,500) #not reject
#bootstarp_Zero_Method(image0,prob,image1,prob,500) #reject

length(which(image0,prob!=0))/length(image0,prob) #0.6744186
length(which(image1,prob!=0))/length(image1,prob) #0.3446389
```

```
##### log price #####
```

```
###
### weekday
###

index.day0 = which(newdata$data.isWeekend == 0 )
index.day1 = which(newdata$data.isWeekend == 1 )

day0.price = data.price[index.day0]
day1.price = data.price[index.day1]

index.na0 = which(is.na(day0.price)==T | day0.price == 0)
index.na1 = which(is.na(day1.price)==T | day1.price == 0)

day0.lprice = log(day0.price[index.na0])
day1.lprice = log(day1.price[index.na1])

#bootstarp_C1_Method(day0.lprice,data.lprice,500) # not reject
#bootstarp_C1_Method(day1.lprice,data.lprice,500) # not reject

#kruskal_type
#bootstarp_C1_Method(day0.lprice,day0.lprice,500) # reject

mean(day0.lprice) # not weekend mean 8.198072
mean(day1.lprice) # weekend mean 7.88289
mean(data.lprice)
```

```
###
### description word count
###

summary(newdata$dev.WordCount)
wordcount = newdata$dev.WordCount
category = 4
quant = quantile(wordcount,seq(0,1,by = 1/category))quant
# 0 250 500 750 1000
# 1 2 3 5 12

cate.price0 = data.price[which(wordcount <= quant[2])]
cate.price1 = data.price[which(wordcount> quant[2] & wordcount <= quant[3])]
cate.price2 = data.price[which(wordcount> quant[3] & wordcount <= quant[4])]
cate.price3 = data.price[which(wordcount > quant[4])]

cate.lprice0 = log(cate.price0[which(is.na(cate.price0)==T | cate.price0 == 0)])
cate.lprice1 = log(cate.price1[which(is.na(cate.price1)==T | cate.price1 == 0)])
cate.lprice2 = log(cate.price2[which(is.na(cate.price2)==T | cate.price2 == 0)])
cate.lprice3 = log(cate.price3[which(is.na(cate.price3)==T | cate.price3 == 0)])

length(cate.lprice0)+length(cate.lprice1)+length(cate.lprice2)+length(cate.lprice3) #932
length(data.lprice) #932

#bootstarp_C1_Method(cate.lprice0,data.lprice,500) # reject
#bootstarp_C1_Method(cate.lprice1,data.lprice,500) # reject
#bootstarp_C1_Method(cate.lprice2,data.lprice,500) # not reject
#bootstarp_C1_Method(cate.lprice3,data.lprice,500) # reject

mean(cate.lprice0) #6.887436
mean(cate.lprice1) #7.226705
mean(cate.lprice2) #8.22525
mean(cate.lprice3) #9.816935
mean(data.lprice) #8.104735
```

```
###
### title word count
###

summary(newdata$t1.WordCount)
wordcount = newdata$t1.WordCount
category = 4
quant = quantile(wordcount,seq(0,1,by = 1/category))quant
# 0 250 500 750 1000
# 1 2 3 5 12

cate.price0 = data.price[which(wordcount <= quant[2])]
cate.price1 = data.price[which(wordcount> quant[2] & wordcount <= quant[3])]
cate.price2 = data.price[which(wordcount> quant[3] & wordcount <= quant[4])]
cate.price3 = data.price[which(wordcount > quant[4])]

cate.lprice0 = log(cate.price0[which(is.na(cate.price0)==T | cate.price0 == 0)])
cate.lprice1 = log(cate.price1[which(is.na(cate.price1)==T | cate.price1 == 0)])
cate.lprice2 = log(cate.price2[which(is.na(cate.price2)==T | cate.price2 == 0)])
cate.lprice3 = log(cate.price3[which(is.na(cate.price3)==T | cate.price3 == 0)])

length(cate.lprice0)+length(cate.lprice1)+length(cate.lprice2)+length(cate.lprice3) #932
length(data.lprice) #932

#bootstarp_C1_Method(cate.lprice0,data.lprice,500) # reject
#bootstarp_C1_Method(cate.lprice1,data.lprice,500) # not reject
#bootstarp_C1_Method(cate.lprice2,data.lprice,500) # reject
#bootstarp_C1_Method(cate.lprice3,data.lprice,500) # reject

mean(cate.lprice0) #6.944468
mean(cate.lprice1) #8.397611
mean(cate.lprice2) #7.687226
mean(cate.lprice3) #9.862613
mean(data.lprice) #8.104735
```

```
xydata<-fread("translate.csv")
picdata<-fread("picdata_mort_1000.csv")
newdata<-fread("newdata.csv")
```

```
data,prob
hist(data,prob)
mean(data,prob)
mean,prob = mean(data,prob)

nonzero,prob = data,prob[which(data,prob!= 0)]
hist(nonzero,prob)

# by MDE
theta = 1/ mean,prob
sim.map = rep(1000, rate = theta)
hist(sim.map)
hist(data,prob)

expquantile<-rep(NA,19)
for (i in 1:19)
{
  expquantile[i]<-qexp(0.05*i,rate=theta)
}
expquantile

chnumber<-rep(NA,20)
for (i in 1:20){
  chnumber[i]<-length(data,prob[data,prob==expquantile[i]])
  for (i in 2:19)
  {
    chnumber[i]<-length(data,prob[data,prob==expquantile[i-1] & data,prob==expquantile[i]])
  }
  chnumber[20]<-length(data,prob[data,prob==expquantile[19]])
  chnumber
  chiazq.test(chnumber, p = rep(0.05,20)) # reject

## non zero probe
nonzero,prob
hist(nonzero,prob)
mean(nonzero,prob)

uniquequatile<-rep(NA,19)
for (i in 1:19)
{
  uniquequatile[i]<-qunif(0.05*i,0,1)
}
uniquequatile

chnumber<-rep(NA,20)
chnumber[1]<-length(nonzero,prob[nonzero,prob==uniquequatile[1]])
for (i in 2:19)
{
  chnumber[i]<-length(nonzero,prob[nonzero,prob==uniquequatile[i-1] & nonzero,prob==uniquequatile[i]])
}
chnumber[20]<-length(nonzero,prob[nonzero,prob==uniquequatile[19]])
chnumber

chiazq.test(chnumber, p = rep(0.05,20)) # reject

sim.dis <- function(u1,u2,v1,v2,p,n){
  value= rep(0,n)
  for(i in 1:n){
    luck = runif(1,0,1)
    if ( luck > p ){
      value [i] = rnorm(1,u2,v2)
    }
    else{
      value [i] = rnorm(1,u1,v1)
    }
    if ( value [i] < 0 ){
      value [i] = 0
    }
    else if (value [i] > 1){
      value [i] = 1
    }
  }
  return(value)
}

hist(sim.dis(0.18,0.1,0.8,0.08,0.65,1000))
hist(nonzero,prob)
```

```
### comparison in 0 dealprob, small dealprob and large dealprob
# price%\ user type%\ whitesas%\ image width
```

```

#bootstarp_Zero_Method <- function(data1,data2,n){
  u = length(which(data1 == 0))
  N = length(data1)
  simulation<-rep(0,n)
  for (i in 1:n)
  {
    btsample = sample(data2,N,replace = F)
    simulation[i]= length(which(btsample == 0))
  }
  lower = quantile(simulation,0.05,na.rm = T)
  upper = quantile(simulation,0.95,na.rm = T)
  inside = NA
  if(n >lower && u <= upper){
    inside = T
  }
  else{
    inside = F
  }
  return(inside)
}

#bootstarp_Small_Method <- function(data1,data2,n){
  u = length(which(data1 > 0 & data1 < 0.5))
  N = length(data1)
  simulation<-rep(0,n)
  for (i in 1:n)
  {
    btsample = sample(data2,N,replace = F)
    simulation[i]= length(which(btsample > 0 & btsample < 0.5))
  }
  lower = quantile(simulation,0.05,na.rm = T)
  upper = quantile(simulation,0.95,na.rm = T)
  inside = NA
```



```

if(u >= lower && u <= upper){
  inside = T
}
else{
  inside = F
}
return(inside)
}

Roostarp_Big_Method <- function(data1,data2,n){
  u = length(which(data1 >=0.5))
  N = length(data1)
  simulation<-rep(0,n)
  for (i in 1:n){
    {
      btsample = sample(data1,N,replace = T)
      simulation[i]= length(which(btsample >= 0.5))
    }
  }
  lower = quantile(simulation,0.05,na.rm = T)
  upper = quantile(simulation,0.95,na.rm = T)
  inside = NA
  if(u >=lower && u <= upper){
    inside = T
  }
  else{
    inside = F
  }
  return(inside)
}

#----> dealprob  zero, small and big deal probe
# probab user type% whiteneas% image width
data.prob
prob.4 = data.prob[which(data.prob==0)]
length(prob.0) #627

prob.small = data.prob[which(data.prob > 0 & data.prob < 0.5)]
length(prob.small) #257

prob.big = data.prob[which(data.prob >=0.5)]
length(prob.big) #164

percentrate <-function(data){
  a=length(data[which(data == 0)]/length(data)
  b=length(data[which(data >0 & data <= 0.5)]/length(data)
  c=length(data[which(data >=0.5)]/length(data)
  return(c(a,b,c,"zero",small,big + 0.027, 0.257,0.116*))
}

###
### price
###
summary(data.price)

category = 4
quant = quantile(data.price,seq(0,1,by = 1/category))quant
# 0 250 500 750 1000
# 1 2 3 5 12

prob.price0 = data.prob[which(data.price <= quant[2])]
prob.price1 = data.prob[which(data.price> quant[2] & data.price <= quant[3])]
prob.price2 = data.prob[which(data.price> quant[3] & data.price <= quant[4])]
prob.price3 = data.prob[which(data.price > quant[4])]

percentrate(prob.price0) #0.4346782 0.2437938 0.1002460
percentrate(prob.price1) #0.4391763 0.2577320 0.1030208
percentrate(prob.price2) #0.6553812 0.2511211 0.1448978
percentrate(prob.price3) #0.4184211 0.2476438 0.1149301

Roostarp_Zero_Method(prob.price0,data.prob,500) # not reject
Roostarp_Big_Method(prob.price0,data.prob,500) # not reject

Roostarp_Small_Method(prob.price1,data.prob,500) # not reject
Roostarp_Big_Method(prob.price1,data.prob,500) # not reject

Roostarp_Zero_Method(prob.price2,data.prob,500) # not reject
Roostarp_Small_Method(prob.price2,data.prob,500) # not reject
Roostarp_Big_Method(prob.price2,data.prob,500) # not reject

Roostarp_Zero_Method(prob.price3,data.prob,500) # not reject
Roostarp_Small_Method(prob.price3,data.prob,500) # not reject
Roostarp_Big_Method(prob.price3,data.prob,500) # not reject

###
### user type
###
index.type = which(newdata$Data.UserType == 0)
type0_prob = data.prob[index.type0]
Roostarp_Zero_Method(type0_prob,data.prob,500) #not reject
Roostarp_Small_Method(type0_prob,data.prob,500) #not reject
Roostarp_Big_Method(type0_prob,data.prob,500) #not reject

index.type = which(newdata$Data.UserType == 1)
type1_prob = data.prob[index.type1]
Roostarp_Zero_Method(type1_prob,data.prob,500) #reject
Roostarp_Small_Method(type1_prob,data.prob,500) #reject
Roostarp_Big_Method(type1_prob,data.prob,500) #not reject

index.type = which(newdata$Data.UserType == 2)
type2_prob = data.prob[index.type2]
Roostarp_Zero_Method(type2_prob,data.prob,500) # reject
Roostarp_Small_Method(type2_prob,data.prob,500) # reject
Roostarp_Big_Method(type2_prob,data.prob,500) # reject

percentrate(type0_prob)
percentrate(type1_prob)
percentrate(type2_prob)

###
### whiteneas
###
index.type = which(pindex$whiteneas == 0)
whitel_prob = data.prob[index.type]
percentrate(whitel_prob)
Roostarp_Zero_Method(whitel_prob,data.prob,500) # reject
Roostarp_Small_Method(whitel_prob,data.prob,500) # reject
Roostarp_Big_Method(whitel_prob,data.prob,500) # not reject

index.type = which(pindex$whiteneas > 0)
whitel_prob = data.prob[index.type]
percentrate(whitel_prob)
Roostarp_Zero_Method(whitel_prob,data.prob,500) # reject
Roostarp_Small_Method(whitel_prob,data.prob,500) # reject
Roostarp_Big_Method(whitel_prob,data.prob,500) # not reject

###
### tempsize
###
summary(pindex$temp_size)
temp_size=pindex$temp_size

index.size = which(temp_size == "(360, 480)")
size0_prob = data.prob[index.size0]
percentrate(size0_prob)
Roostarp_Zero_Method(size0_prob,data.prob,500) # reject
Roostarp_Small_Method(size0_prob,data.prob,500) # reject
Roostarp_Big_Method(size0_prob,data.prob,500) # not reject

index.size = which(temp_size == "(270, 480)")
size1_prob = data.prob[index.size1]
percentrate(size1_prob)
Roostarp_Zero_Method(size1_prob,data.prob,500) # not reject
Roostarp_Small_Method(size1_prob,data.prob,500) # not reject
Roostarp_Big_Method(size1_prob,data.prob,500) # not reject

index.size = which(temp_size == "(480, 360)")
size2_prob = data.prob[index.size2]
percentrate(size2_prob)
Roostarp_Zero_Method(size2_prob,data.prob,500) # reject
Roostarp_Small_Method(size2_prob,data.prob,500) # not reject
Roostarp_Big_Method(size2_prob,data.prob,500) # reject

index.size = which(temp_size == "(640, 360)")
size3_prob = data.prob[index.size3]
percentrate(size3_prob)
Roostarp_Zero_Method(size3_prob,data.prob,500) # reject
Roostarp_Small_Method(size3_prob,data.prob,500) # reject
Roostarp_Big_Method(size3_prob,data.prob,500) # reject

newd1<-"C:\\Oncro\\ham1\\Desktop"
newdata<-fread("newdata.csv")
newdata$tic<-newdata[,-index.no.image,]
mydata$with.image<-rep(1,1000)
for (i in 1:1000)
{
  if (mydata$image[i]==F){
    {
      mydata$with.image[i]<=0
    }
  }
  length(index.no.image) #86 of 1000 ad no image
  mean(mydata$deal_probability[index.no.image])#mean(mydata$deal_probability[-index.no.image])
  length(which(mydata$deal_probability[index.no.image]==0)) #28/86 no img deal prob 0
  length(which(mydata$deal_probability[-index.no.image]==0)) #599/914 with img deal prob 0
  length(mydata$with.image[mydata$with.image==0])

mydata$with.deal<-rep(0,1000)
for (i in 1:1000)
{
  if (mydata$deal_probability[i]>0)
  {
    mydata$with.deal[i]<=1
  }
  length(mydata$with.deal[mydata$with.deal==1])

mydata$is.private<-rep(0,1000)
for (i in 1:1000)
{
  if (mydata$user_type[i]==Private*)
  {
    mydata$is.private[i]<=1
  }
}

mydata$is.company<-rep(0,1000)
for (i in 1:1000)
{
  if (mydata$user_type[i]==Company*)
  {
    mydata$is.company[i]<=1
  }
}

mydata$is.360x480<-rep(0,1000)
for (i in 1:914)
{
  if (pindex$temp_size[i]==(360, 480)*)
  {
    mydata$is.360x480[i]<=1
  }
}

mydata$is.270x480<-rep(0,1000)
for (i in 1:914)
{
  if (pindex$temp_size[i]==(270, 480)*)
  {
    mydata$is.270x480[i]<=1
  }
}

mydata$is.480x360<-rep(0,1000)
for (i in 1:914)
{
  if (pindex$temp_size[i]==(480, 360)*)
  {
    mydata$is.480x360[i]<=1
  }
}

glm.1<-glm(mydata$with.deal~mydata$price+mydata$with.image+newdata$Data.t0Weekend+newdata$Data.des.WordCount
+newdata$Data.des.CapeCount+newdata$Data.des.DigitCount+newdata$Data.t1.CapeCount+newdata$Data.t1.DigitCount
+newdata$Data.t1.WordCount+mydata$is.private+mydata$is.private.family ~ binomial(link=logit))
summary(glm.1)

# get rid of newdata$Data.t1.WordCount -3.388e-03 3.917e-02 -0.086 0.9311

glm.2<-glm(mydata$with.deal~mydata$price+mydata$with.image+newdata$Data.t0Weekend+newdata$Data.des.WordCount
+newdata$Data.des.CapeCount+newdata$Data.des.DigitCount+newdata$Data.t1.CapeCount+newdata$Data.t1.DigitCount
+mydata$is.company+mydata$is.private, family = binomial(link=logit))

```

[illegible]

```

    return(sav.out)
  }

###
### function of splitting deal probability into groups for categorical variables
###
###
anovaFun = function(variable, n){
  category = n
  quant = quantile(variable, seq(0,1,by = 1/category))
  cate.prob = list()
  for (i in 1:category){
    cate.prob[i] = list(data.prob[which(variable== quant[i] & variable <= quant[i+1])])
  }
  names(cate.prob) = 1:category
  sav.out = sav[value=ind, stack(cate.prob)]
  return(sav.out)
}

###
### function of splitting deal probability into groups for numerical variables
###
###
small.outcat = function(variable){
  var.fac = as.factor(variable)
  category = length(levels(var.fac))
  cate.prob = list()
  j = 1
  for (i in levels(var.fac)){
    cate.prob[j] = list(small.data.prob[var.fac == i])
    j = j+1
  }
  names(cate.prob) = levels(var.fac)
  sav.out = sav[value=ind, stack(cate.prob)]
  return(sav.out)
}

###
### function of splitting deal probability into groups for numerical variables
###
###
big.outcat = function(variable){
  var.fac = as.factor(variable)
  category = length(levels(var.fac))
  cate.prob = list()
  j = 1
  for (i in levels(var.fac)){
    cate.prob[j] = list(small.data.prob[var.fac == i])
    j = j+1
  }
  names(cate.prob) = levels(var.fac)
  sav.out = sav[value=ind, stack(cate.prob)]
  return(sav.out)
}

###
### function of splitting deal probability into groups for categorical variables
###
###
small.outnum = function(variable, n){
  category = n
  quant = quantile(variable, seq(0,1,by = 1/category))
  cate.prob = list()
  for (i in 1:category){
    cate.prob[i] = list(small.data.prob[which(variable== quant[i] & variable <= quant[i+1])])
  }
  names(cate.prob) = 1:category
  sav.out = sav[value=ind, stack(cate.prob)]
  return(sav.out)
}

###
### function of splitting deal probability into groups for categorical variables
###
###
big.outnum = function(variable, n){
  category = n
  quant = quantile(variable, seq(0,1,by = 1/category))
  cate.prob = list()
  for (i in 1:category){
    cate.prob[i] = list(big.data.prob[which(variable== quant[i] & variable <= quant[i+1])])
  }
  names(cate.prob) = 1:category
  sav.out = sav[value=ind, stack(cate.prob)]
  return(sav.out)
}

###
### price
###
###
# check data
summary(newdata$Data.Price.0)
data.price = newdata$Data.Price.0[is.na(newdata$Data.Price.0)==F]
cat(data.price, 5) # split into groups
price.sav = anovaFun(data.price, 5)
summary(price.sav) #ANOVA
TukeyHSD(price.sav) #Tukey's Honest Significance Differences
plot(TukeyHSD(price.sav), las = 2) #plot
# none of these reject

summary(small.tukey.newdata$Data.Price.0)
small.data.price = small.tukey.newdata$Data.Price.0[is.na(small.tukey.newdata$Data.Price.0)==F]
cat(small.data.price, 5) # split into groups
small.price.sav = anovaFun(small.data.price, 5)
summary(small.price.sav) #ANOVA
TukeyHSD(small.price.sav) #Tukey's Honest Significance Differences
plot(TukeyHSD(small.price.sav), las = 2) #plot
# none of these reject

###
### month day on deal_pro
###
###
# check data
data.monthday = newdata$Data.monthday[is.na(newdata$Data.monthday)==F]
cat(data.monthday, 7)
monthday.sav = anovaFun(data.monthday, 7) # convert into ANOVA table
summary(monthday.sav) #ANOVA
TukeyHSD(monthday.sav) #Tukey's Honest Significance Differences
plot(TukeyHSD(monthday.sav), las = 2) #plot
# none of these reject

###
### week day on deal_pro
###
###
# check data
week.day = subset(weekdays(as.Date(mydata$activation_date)),1,3)
data.week.day = week.day[is.na(week.day)==F]
weekday.sav = anovaFun(data.week.day) # convert into ANOVA table
summary(weekday.sav) #ANOVA
TukeyHSD(weekday.sav) #Tukey's Honest Significance Differences
plot(TukeyHSD(weekday.sav), las = 2) #plot
# none of these reject

###
### weekend data on deal_pro
###
###
# check data
data.isweekend = newdata$Data.isWeekend[is.na(newdata$Data.isWeekend)==F]
isweekend.sav = anovaFun(data.isweekend) # convert into ANOVA table
summary(isweekend.sav) #ANOVA
TukeyHSD(isweekend.sav) #Tukey's Honest Significance Differences
plot(TukeyHSD(isweekend.sav), las = 2) #plot
# not reject

###
### user type
###
###
# check data
data.usertype = newdata$Data.UserType[is.na(newdata$Data.UserType)==F]
usertype.sav = anovaFun(data.usertype) # convert into ANOVA table
summary(usertype.sav) #ANOVA
TukeyHSD(usertype.sav) #Tukey's Honest Significance Differences
plot(TukeyHSD(usertype.sav), las = 2) #plot
# none of these reject

```