

# Tarefa1

Lucas G Nachtigall

2022-09-14

```
## Warning: package 'gridExtra' was built under R version 4.2.1
```

```
## Warning: package 'GGally' was built under R version 4.2.1
```

1. Carregue a base de dados e mostre a estrutura do dataset (`str()`). Comente sobre o número de amostras, de variáveis (e seu tipo). O arquivo do dataset não pode ser modificado de forma alguma. A leitura deverá ser realizada de tal maneira qualquer característica dos dados.

```
rm(list=ls())
#setwd("./Downloads")
dados = read.csv2("Dry_Bean_Dataset.csv", header=T)
str(dados)

## 'data.frame': 13611 obs. of 17 variables:
## $ Area : int 28395 28734 29380 30008 30140 30279 30477 30519 30685 30834 ...
## $ Perimeter : num 610 638 624 646 620 ...
## $ MajorAxisLength: num 208 201 213 211 202 ...
## $ MinorAxisLength: num 174 183 176 183 190 ...
## $ AspectRatio : num 1.2 1.1 1.21 1.15 1.06 ...
## $ Eccentricity : num 0.55 0.412 0.563 0.499 0.334 ...
## $ ConvexArea : int 28715 29172 29690 30724 30417 30600 30970 30847 31044 31120 ...
## $ EquivDiameter : num 190 191 193 195 196 ...
## $ Extent : num 0.764 0.784 0.778 0.783 0.773 ...
## $ Solidity : num 0.989 0.985 0.99 0.977 0.991 ...
## $ roundness : num 0.958 0.887 0.948 0.904 0.985 ...
## $ Compactness : num 0.913 0.954 0.909 0.928 0.971 ...
## $ ShapeFactor1 : num 0.00733 0.00698 0.00724 0.00702 0.0067 ...
## $ ShapeFactor2 : num 0.00315 0.00356 0.00305 0.00321 0.00366 ...
## $ ShapeFactor3 : num 0.834 0.91 0.826 0.862 0.942 ...
## $ ShapeFactor4 : num 0.999 0.998 0.999 0.994 0.999 ...
## $ Class : chr "SEKER" "SEKER" "SEKER" "SEKER" ...
```

O dataset possui 17 colunas: 16 colunas numéricas e 1 categórica. A variável categórica é o rótulo da classe.

2. Altere a variável do tipo do feijão (Class) para um factor. Utilizando um comando mostre como estimar o número de classes existentes.

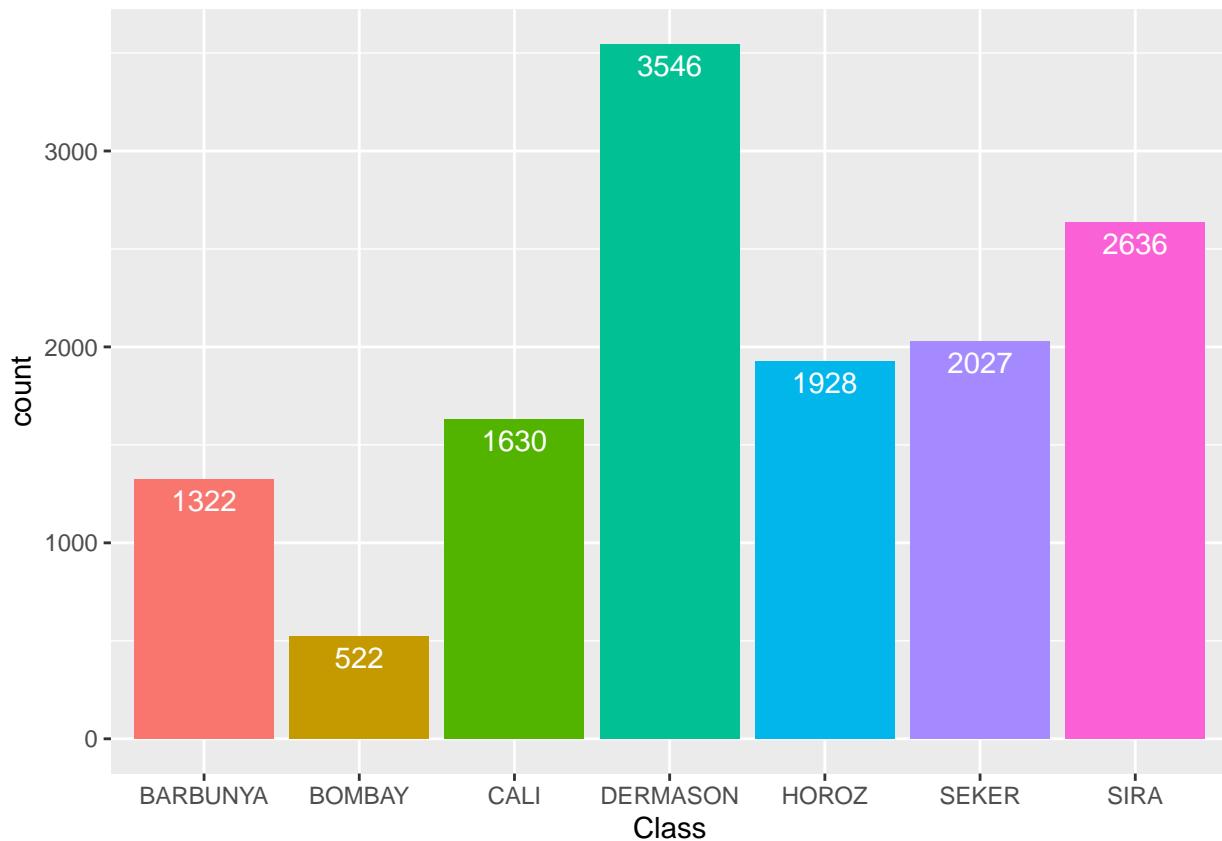
```
dados$Class=factor(dados$Class)
str(dados$Class)
```

```
##  Factor w/ 7 levels "BARBUNYA","BOMBAY",...: 6 6 6 6 6 6 6 6 6 ...
```

O database possui 7 classes.

3. Quantas amostras existem por classe? Use um gráfico de barras para ilustrar as quantidades.

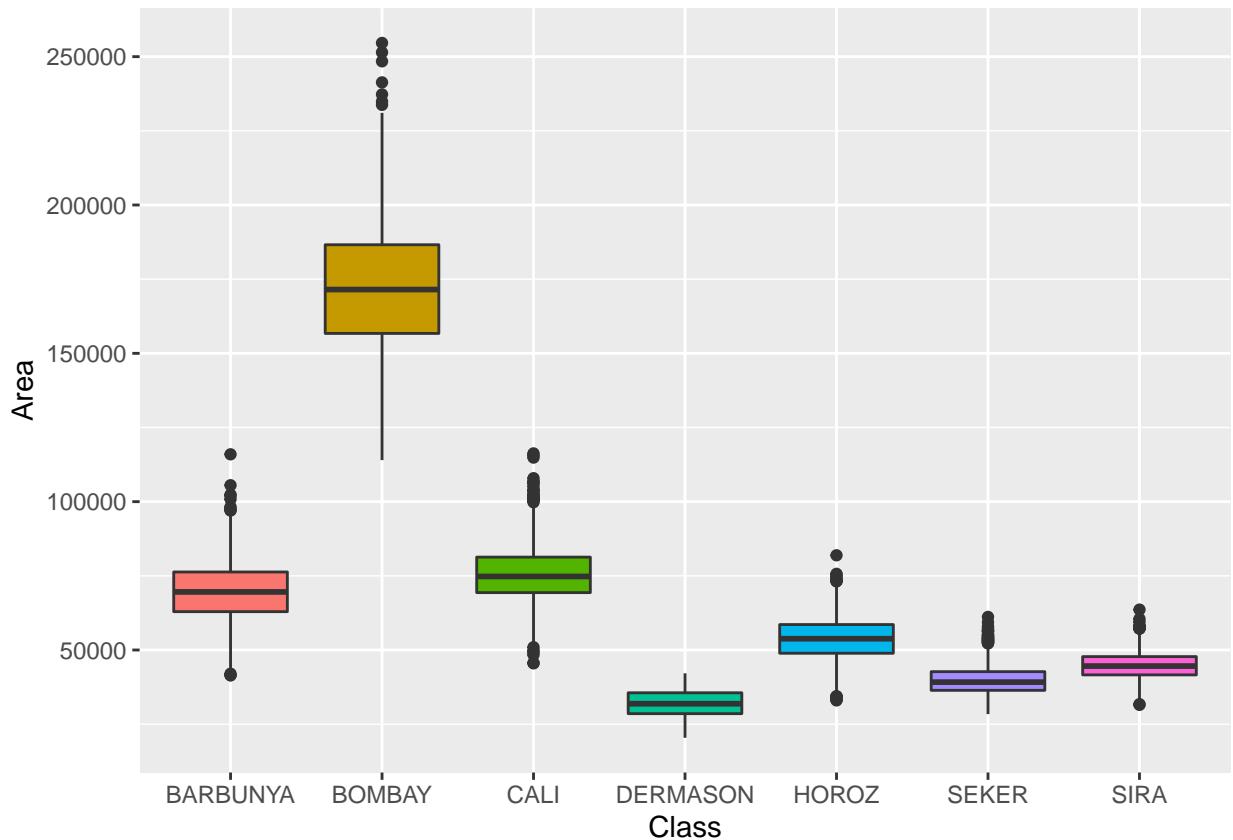
```
ggplot(dados,aes(Class,fill=Class)) +
  geom_bar() + geom_text(aes(label = ..count..), stat = "count", vjust = 1.5, colour = "white") + theme(...)
```

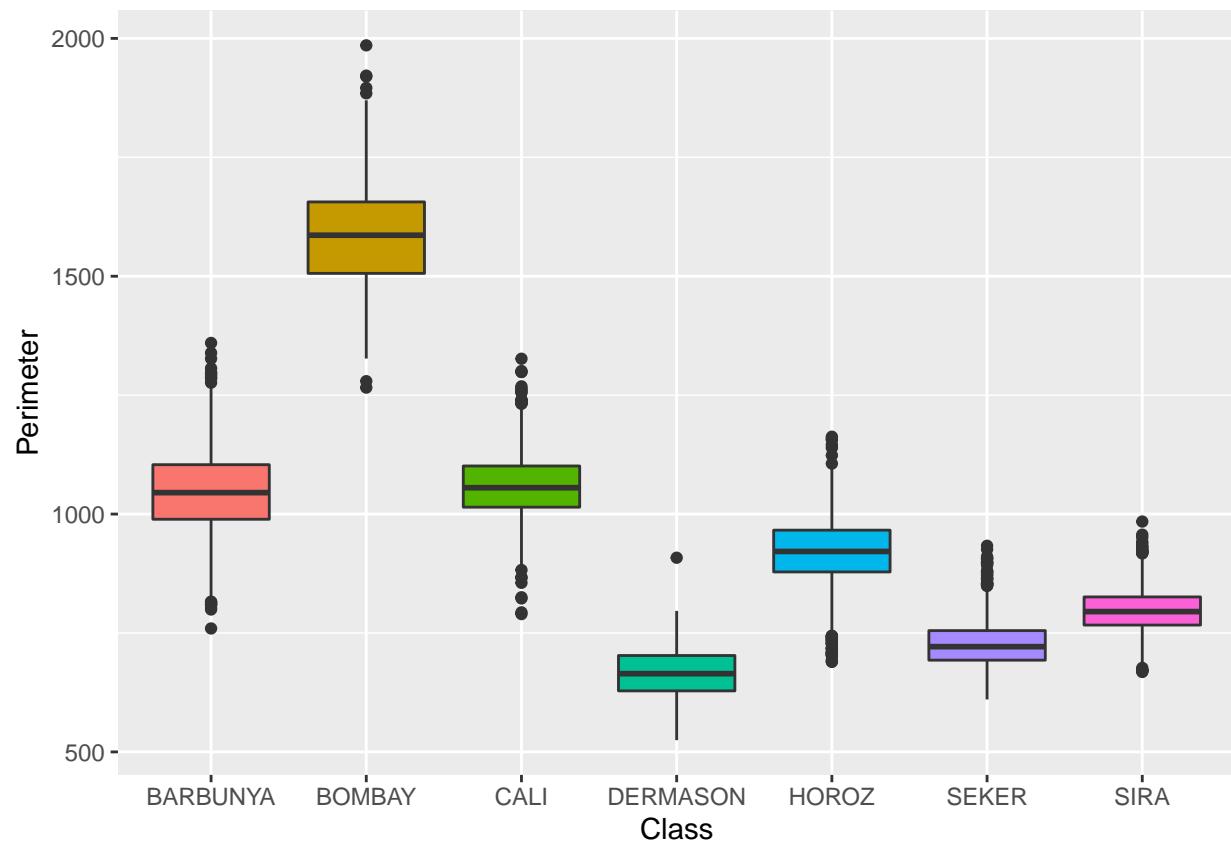


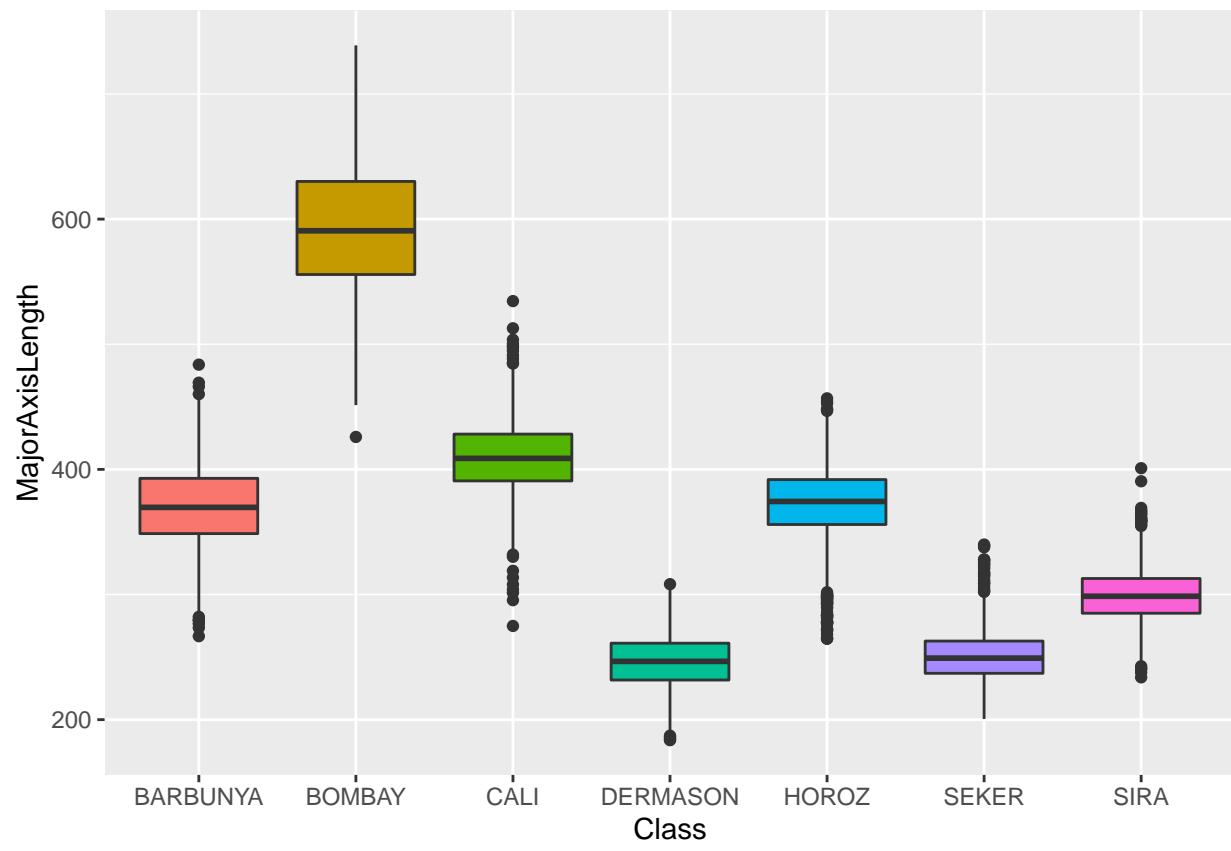
Das 7 classes (Barbunya, Bombay, Cali, Dermason, Horoz, Seker e Sira) existentes, sendo que as classes possuem o número de amostras divergentes. Isto é, o dataset é desbalanceado.

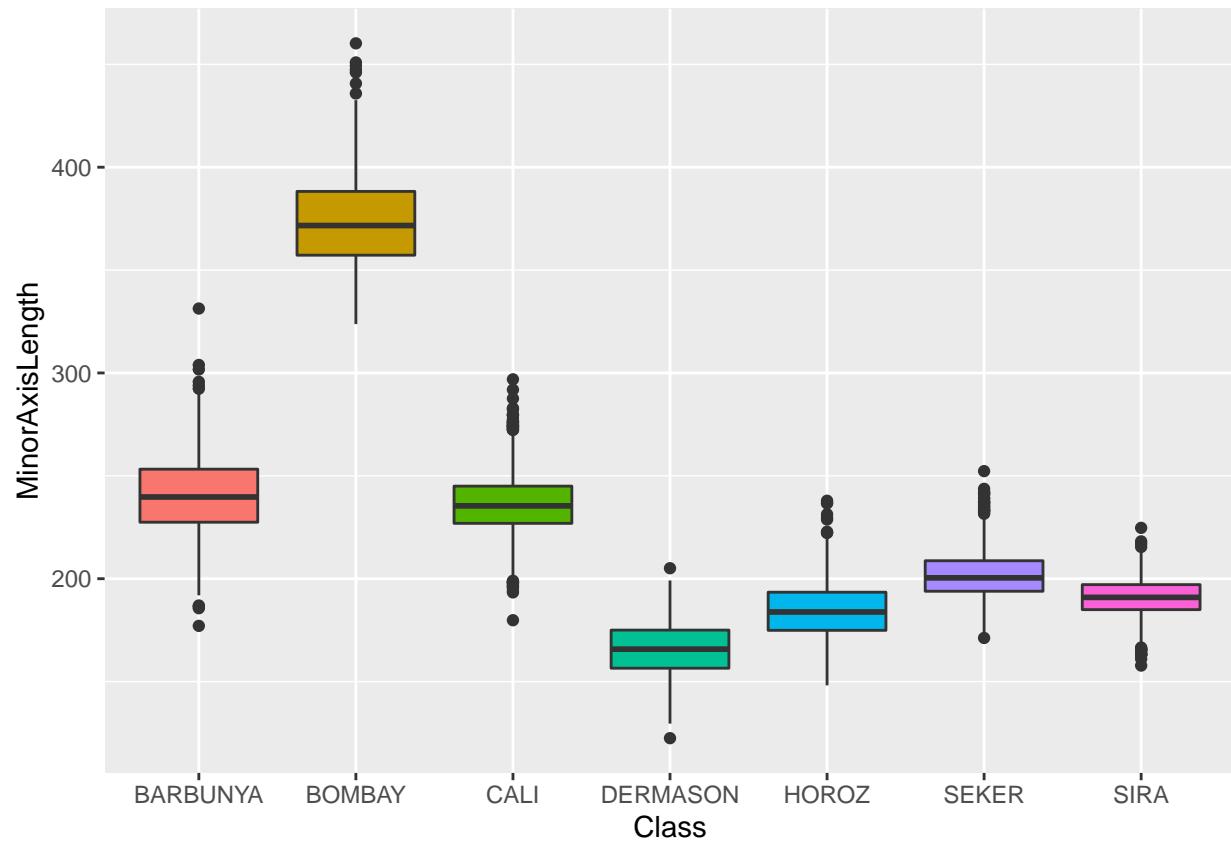
4. Alguma variável apresenta outliers? Tem alguma que não apresenta? Se as amostras com outliers fossem removidas, reduziria em quanto o número de amostras? Alguma classe sofreria uma redução maior do que a outra? Crie um boxplot por variável (boxplot()) para auxiliar na explicação. Não remova os outliers para as próximas etapas!

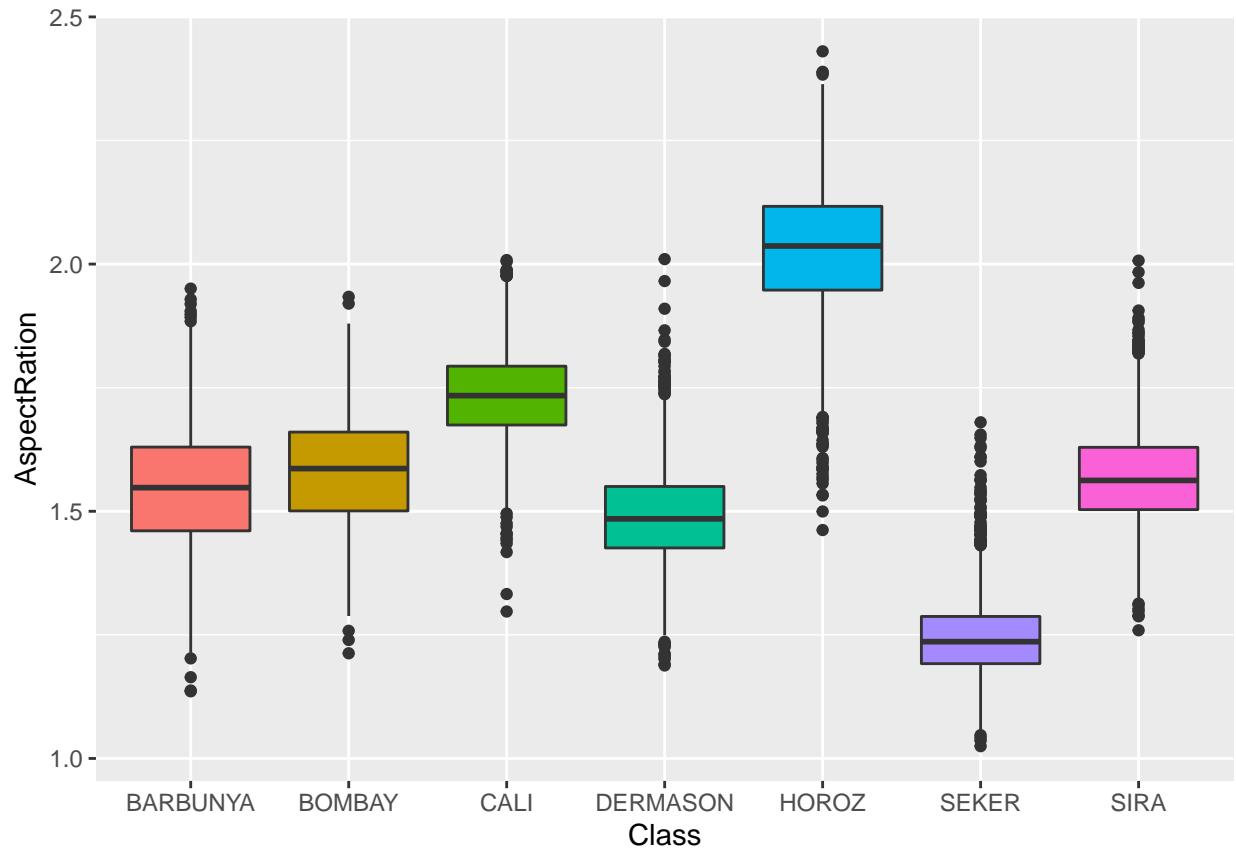
```
p = list()
for(i in 1:16){
  p[[i]] = ggplot(dados, aes_string(x="Class", y=names(dados)[i], fill="Class")) + geom_boxplot() + theme_minimal()
  do.call(grid.arrange, c(p[i], ncol=1))
}
```

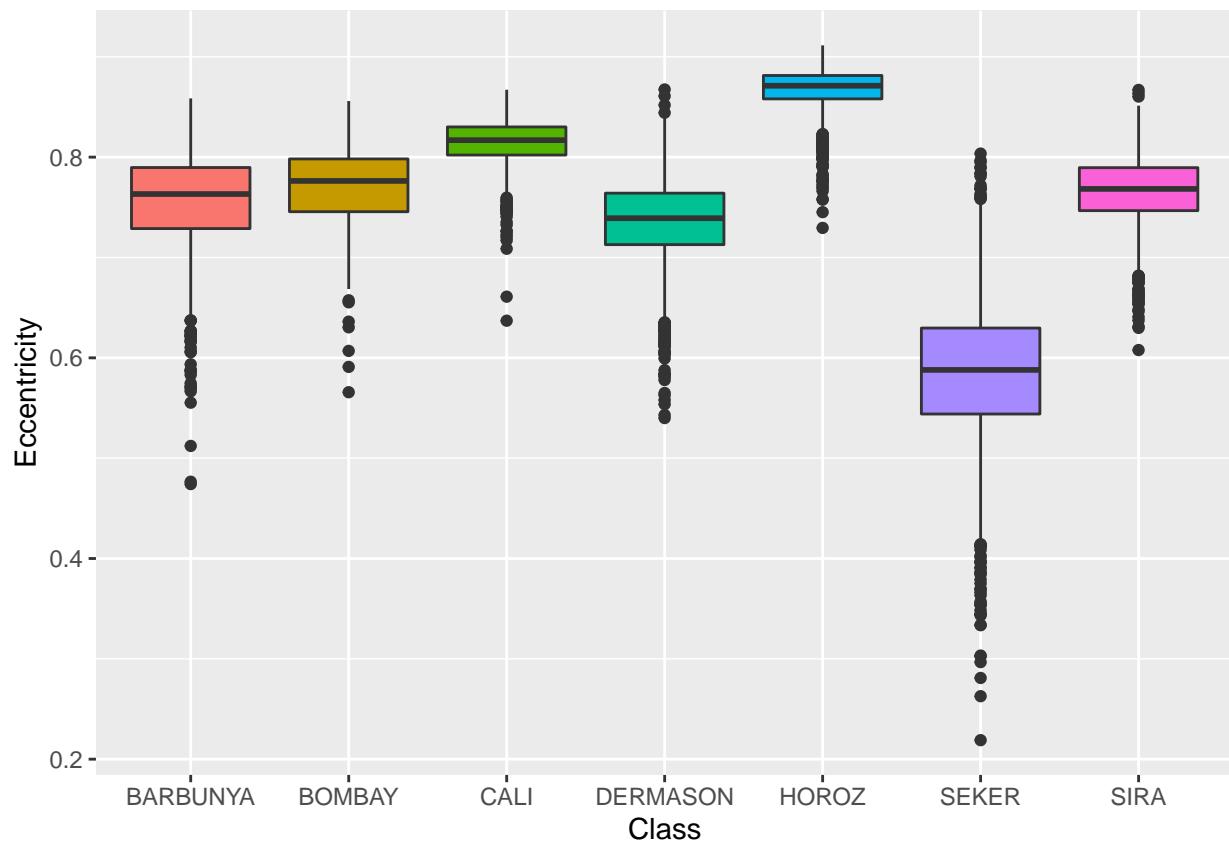


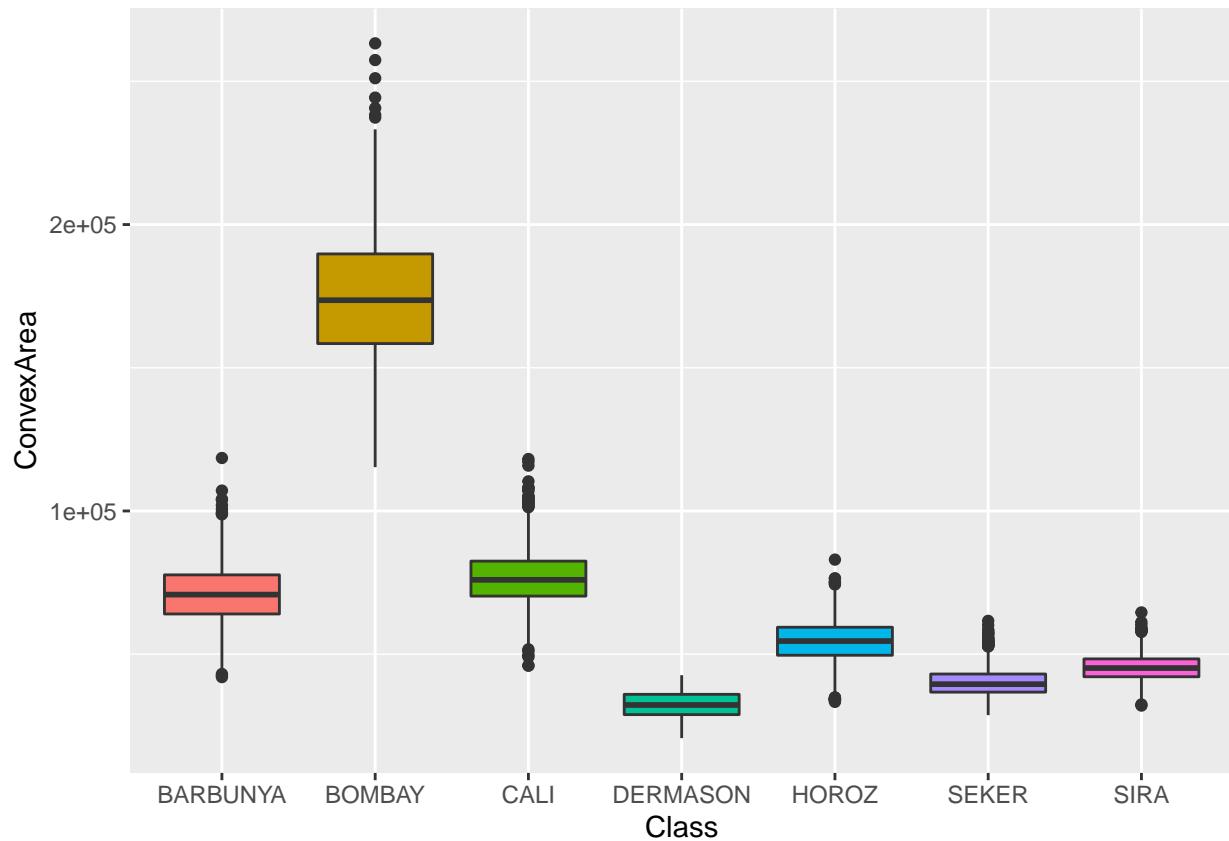


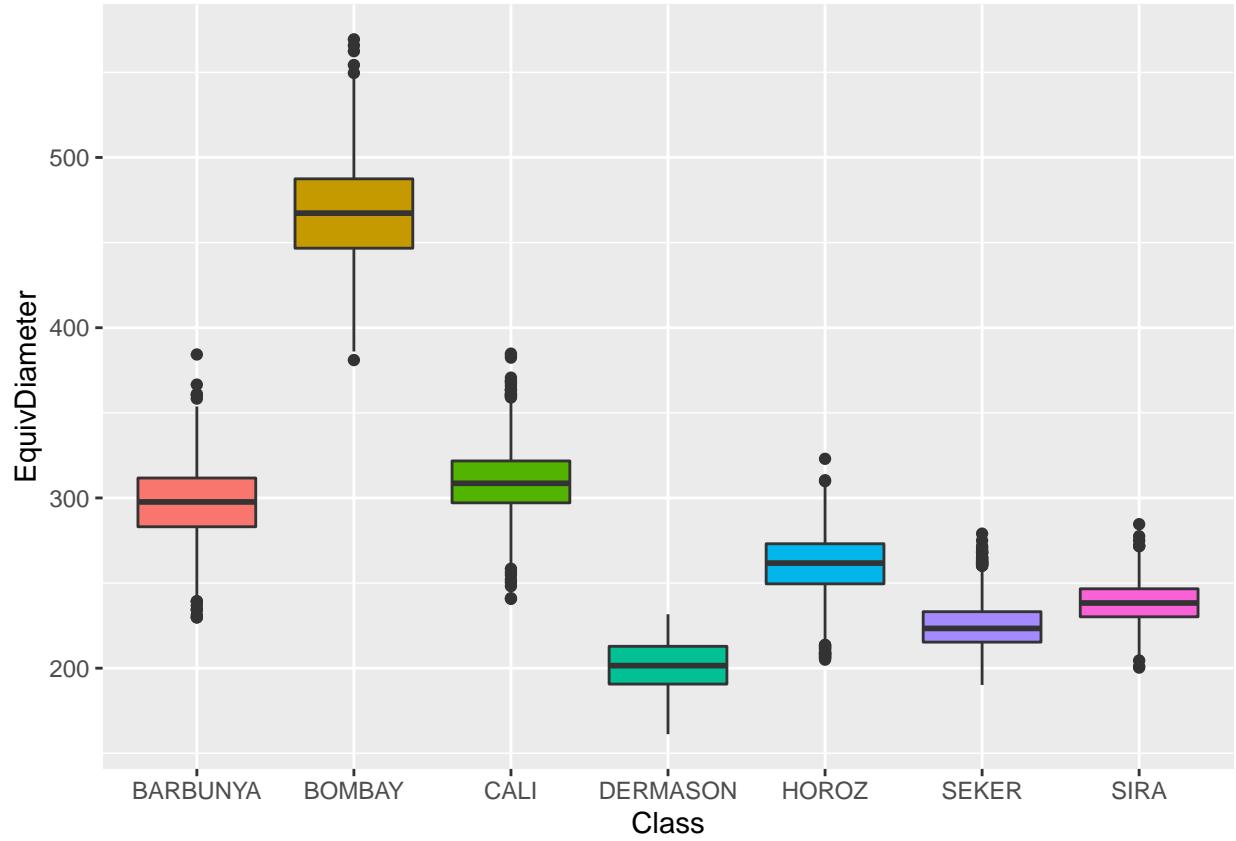


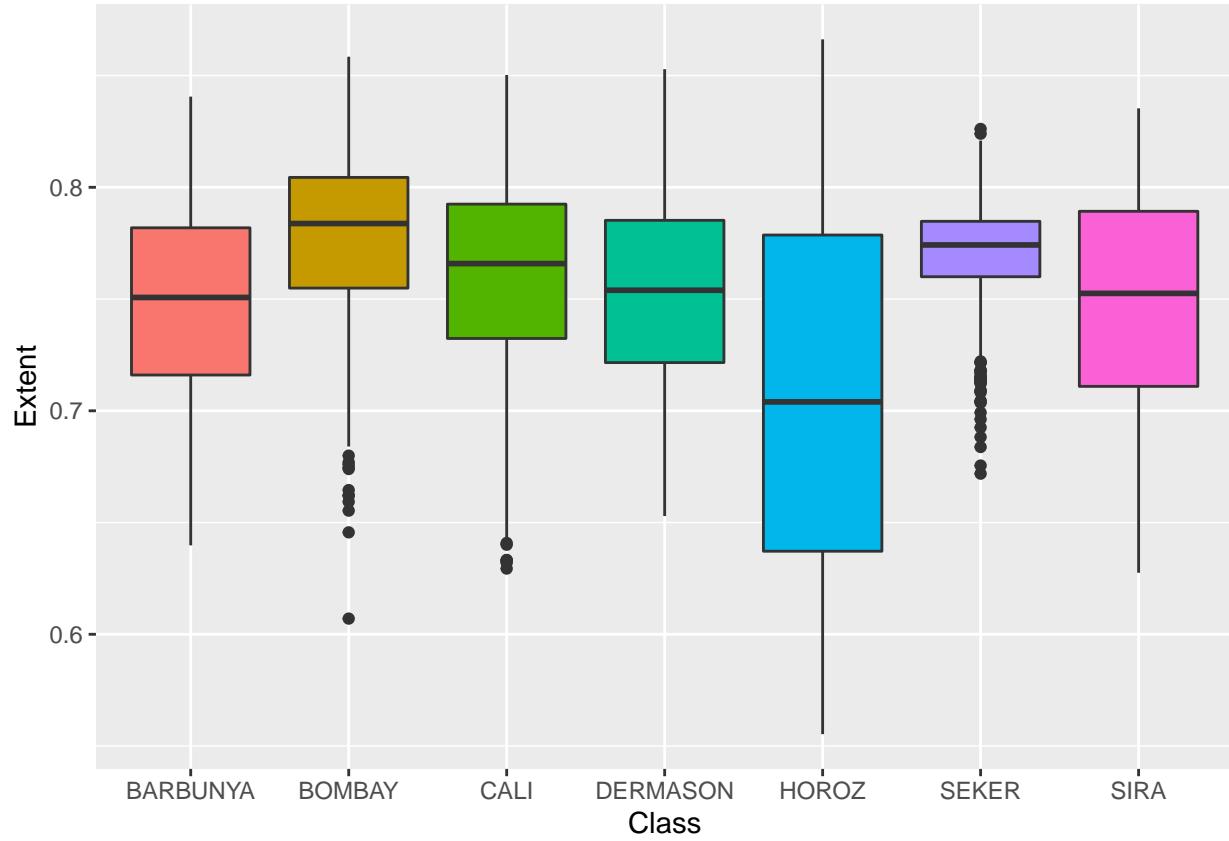


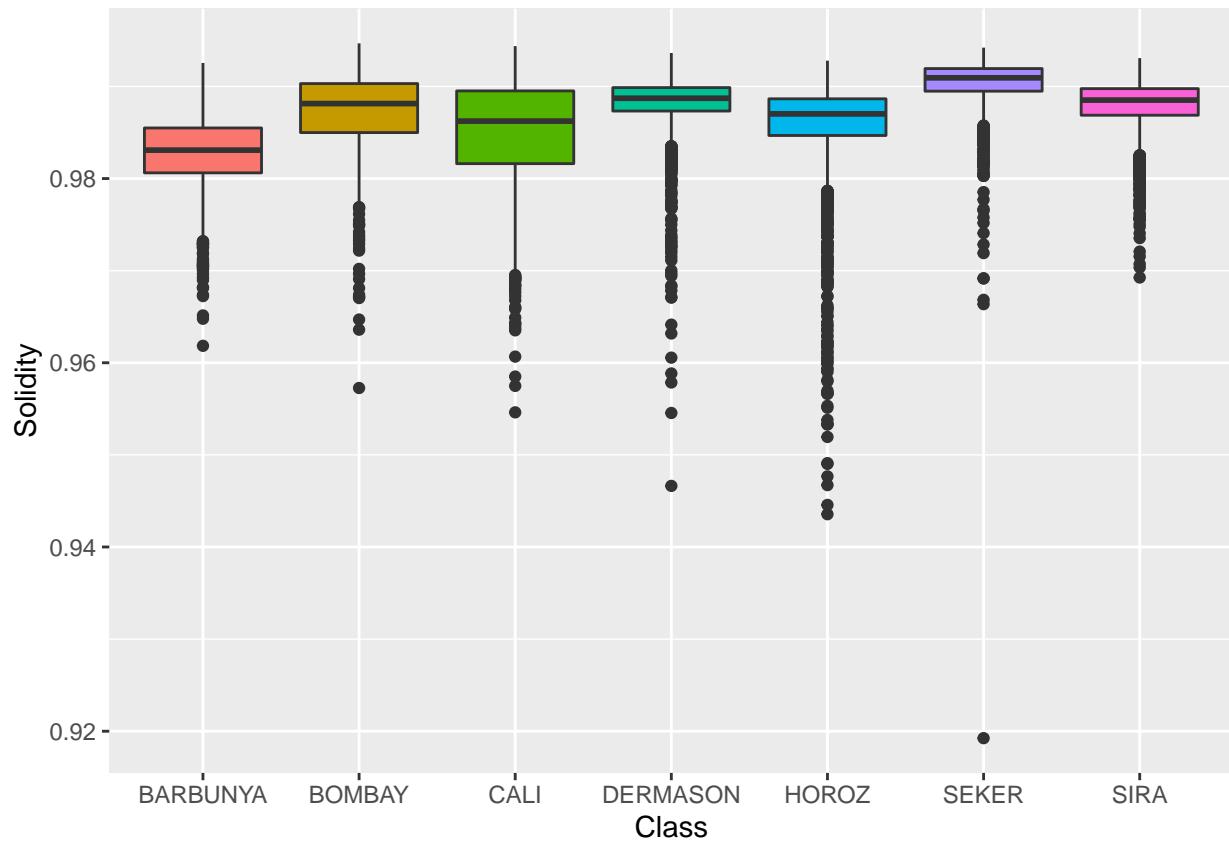


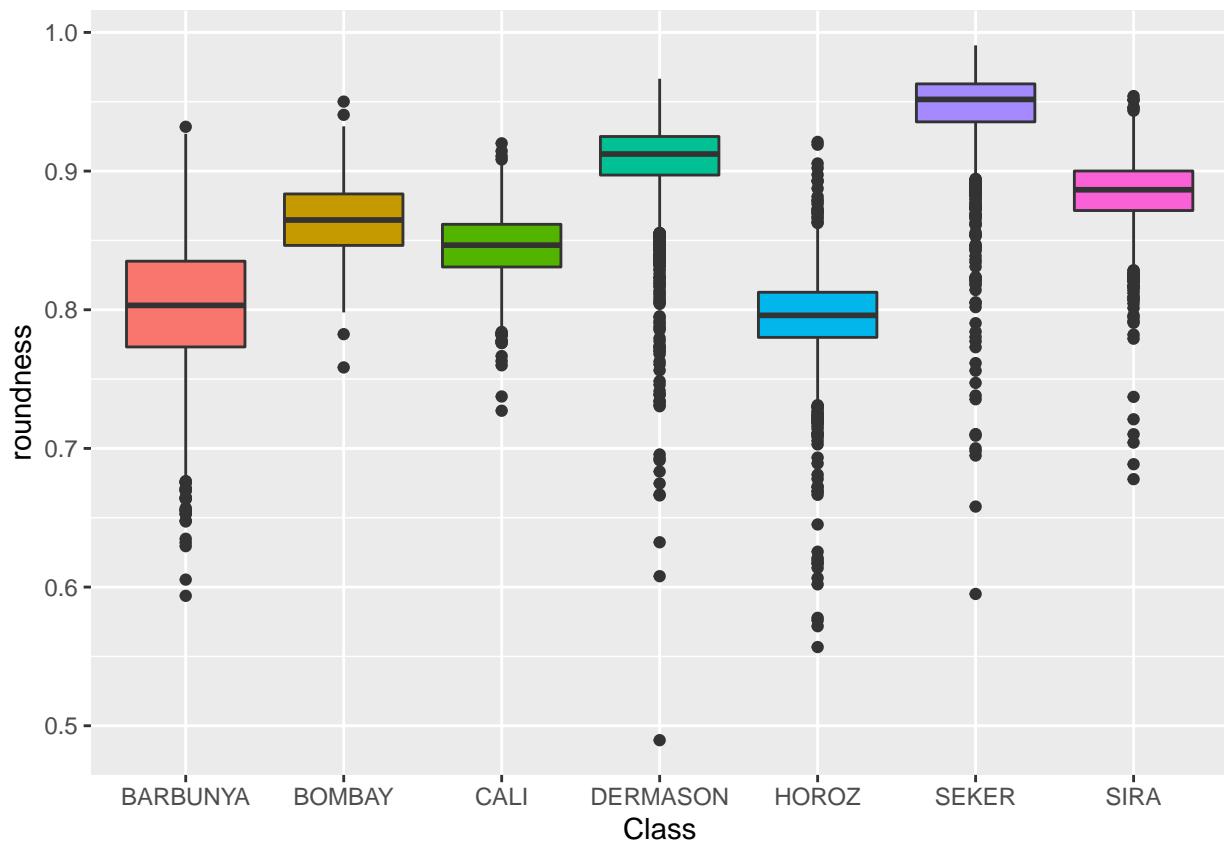


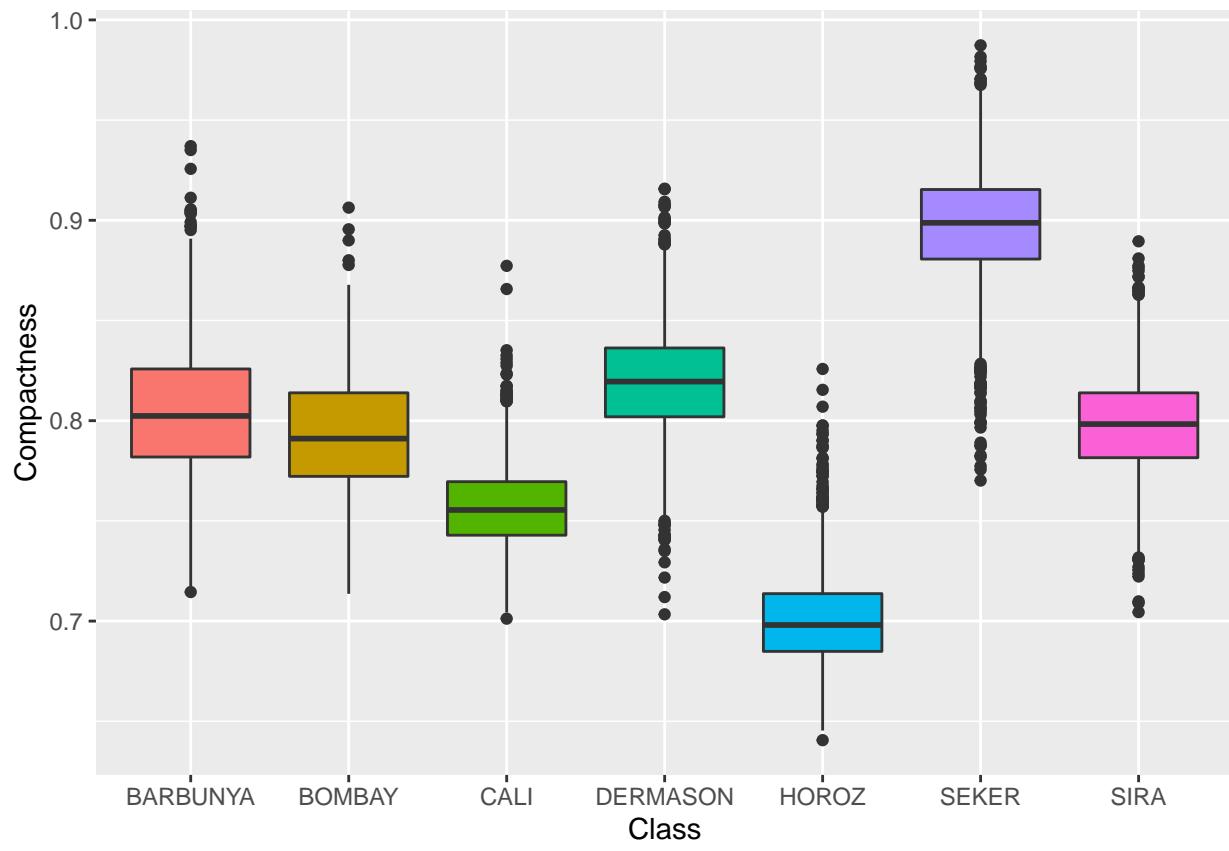


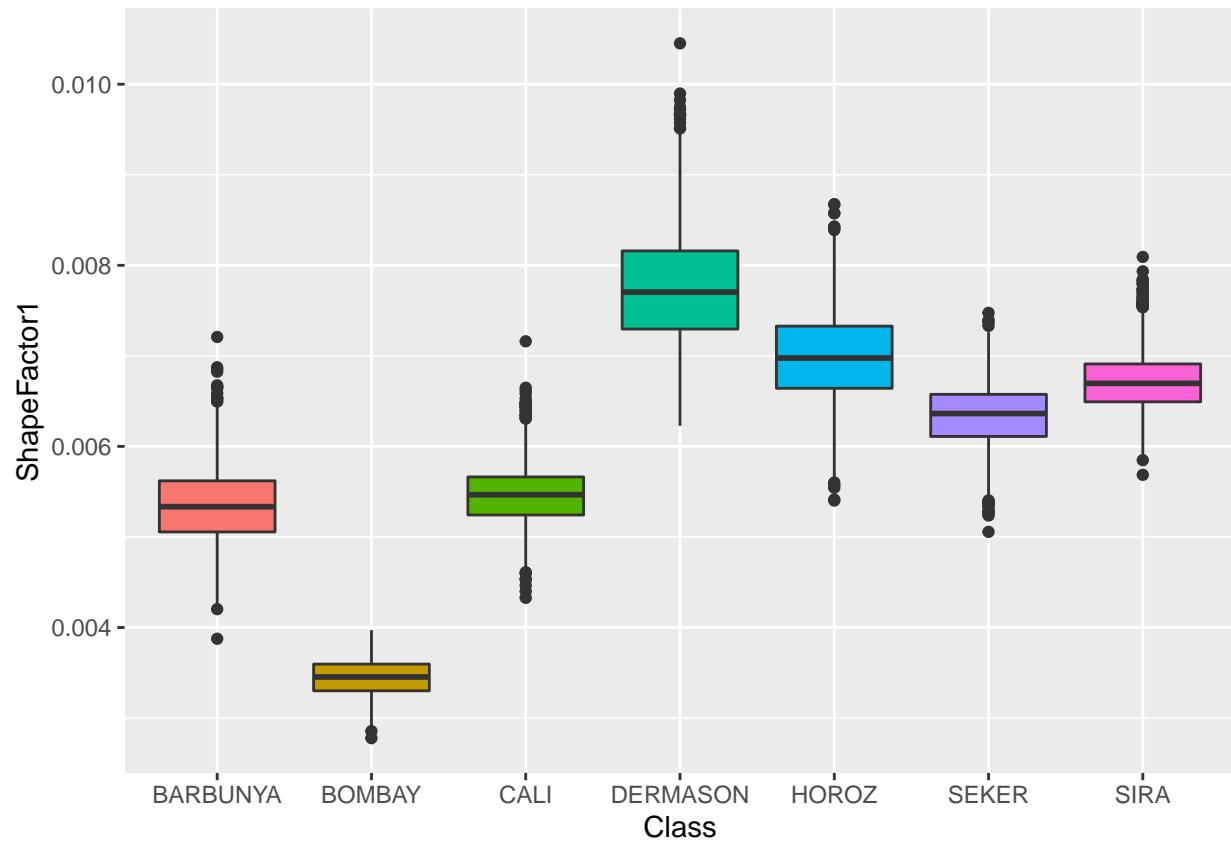


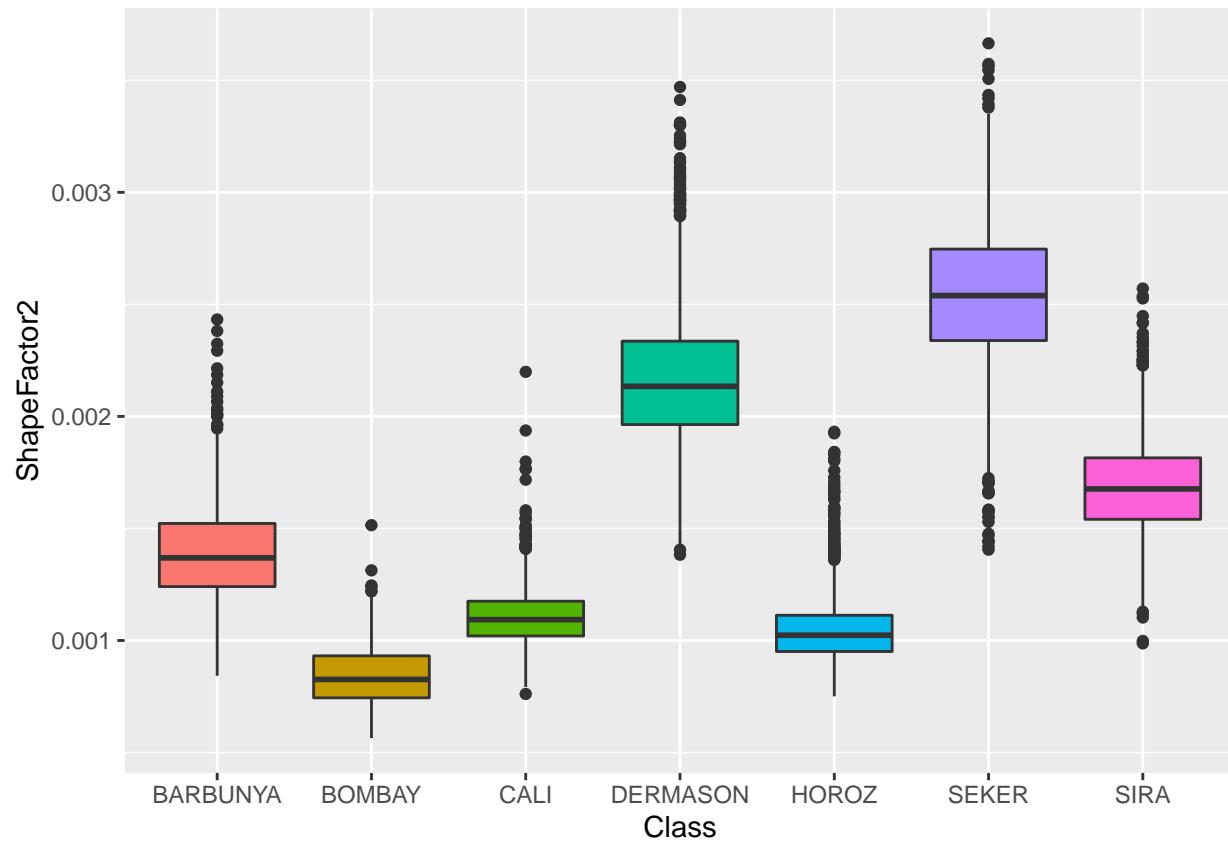


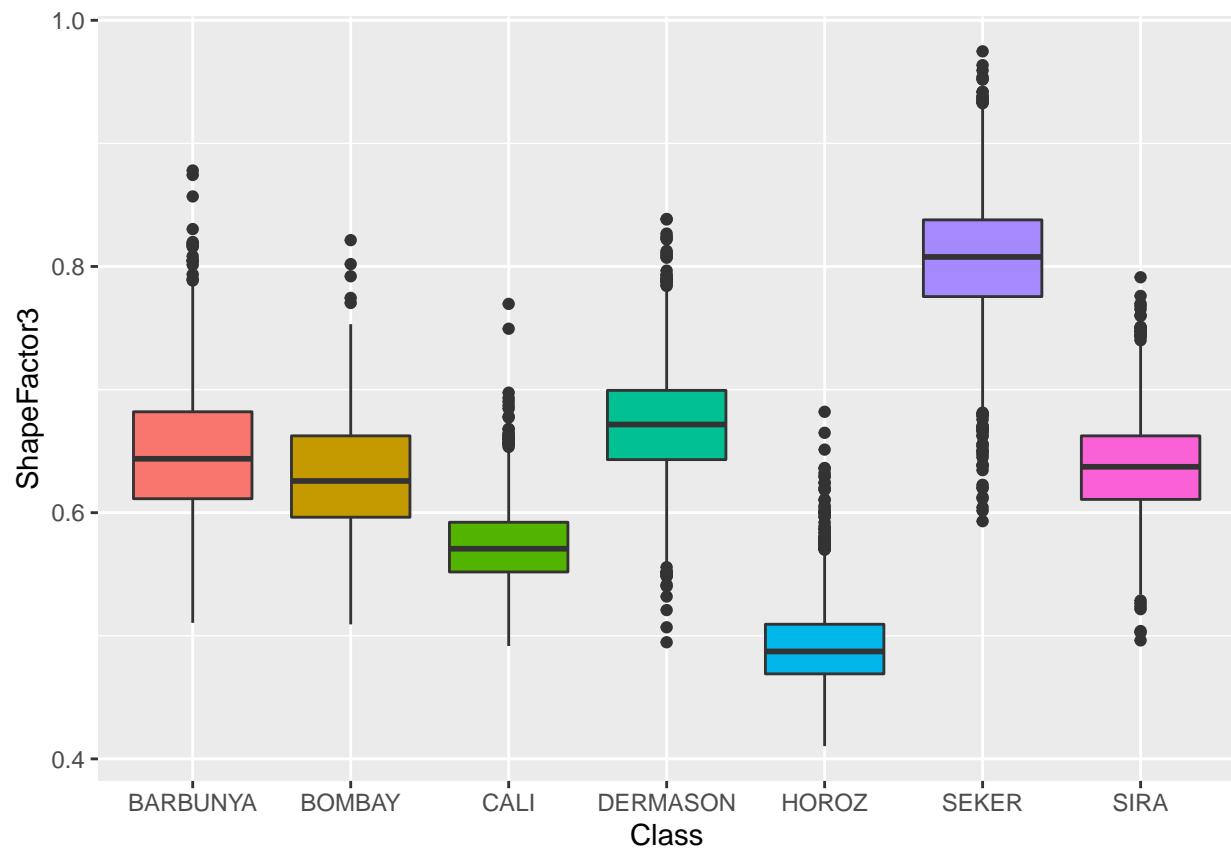


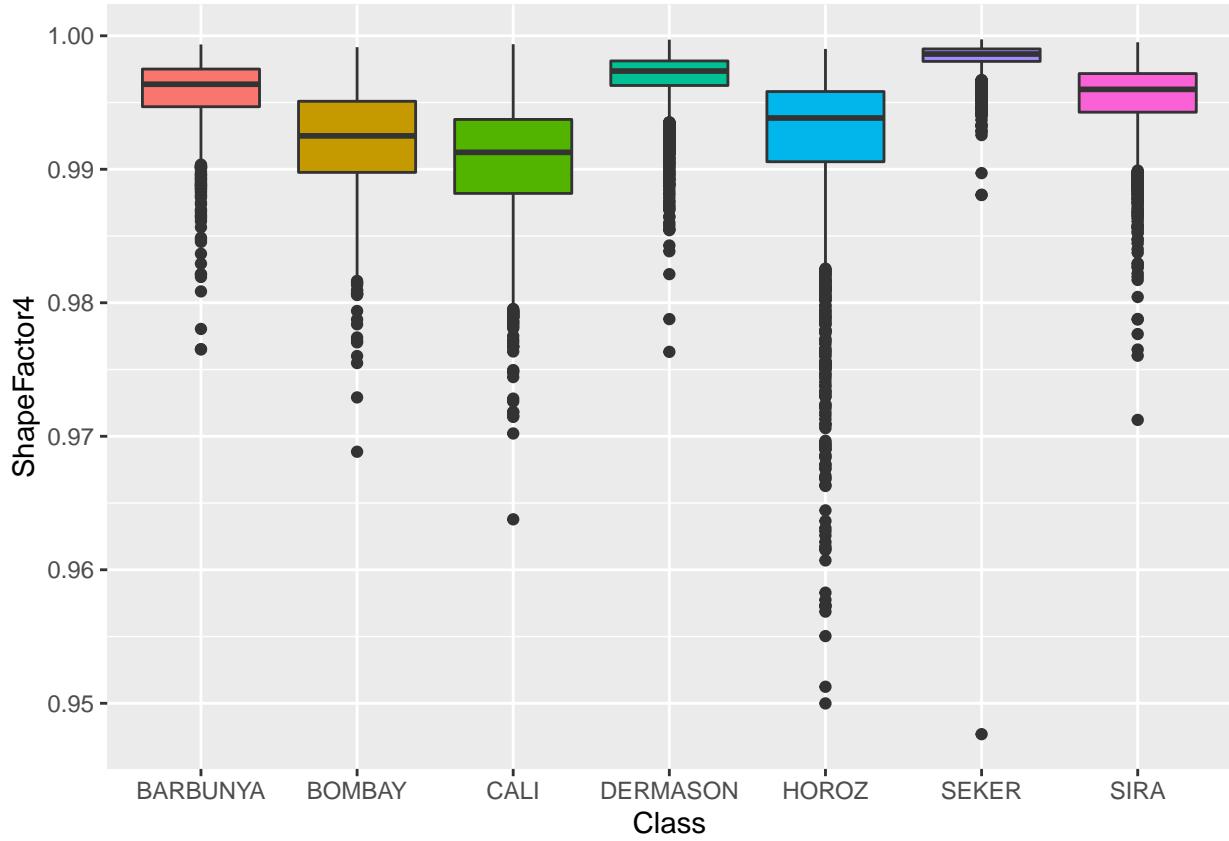












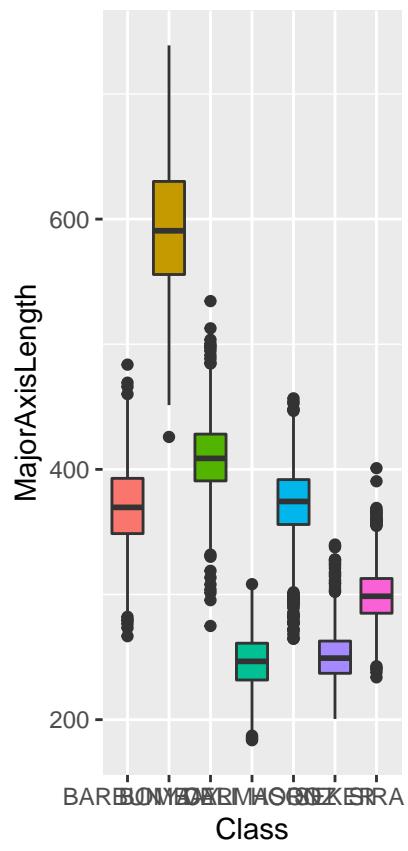
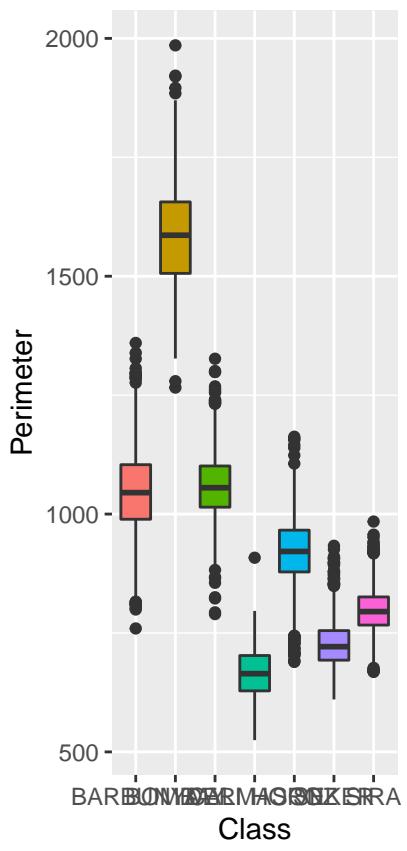
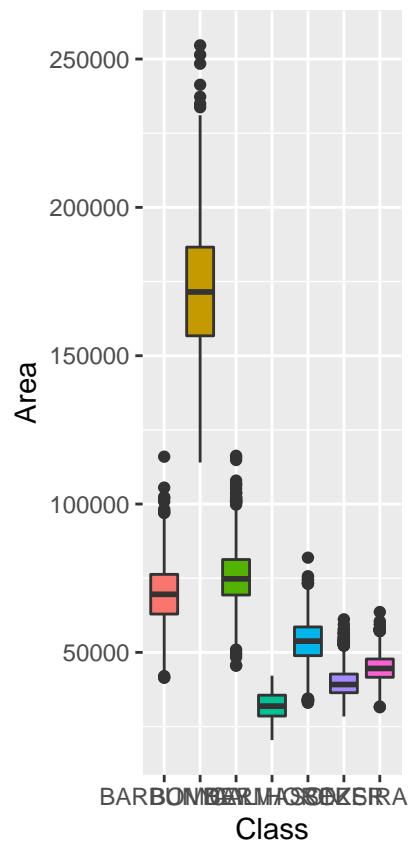
```
colnames(dados)
```

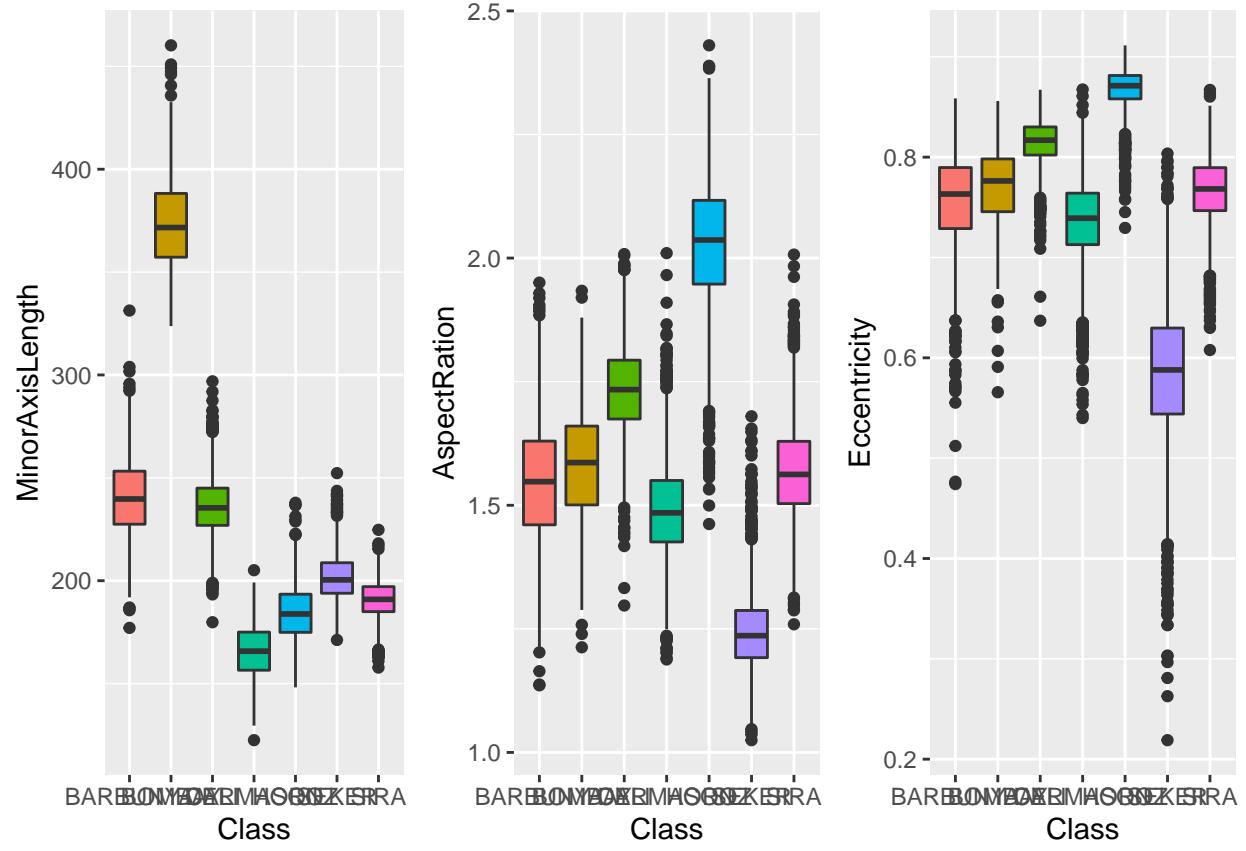
```
## [1] "Area"           "Perimeter"       "MajorAxisLength" "MinorAxisLength"
## [5] "AspectRation"   "Eccentricity"     "ConvexArea"      "EquivDiameter"
## [9] "Extent"          "Solidity"         "roundness"       "Compactness"
## [13] "ShapeFactor1"    "ShapeFactor2"     "ShapeFactor3"    "ShapeFactor4"
## [17] "Class"
```

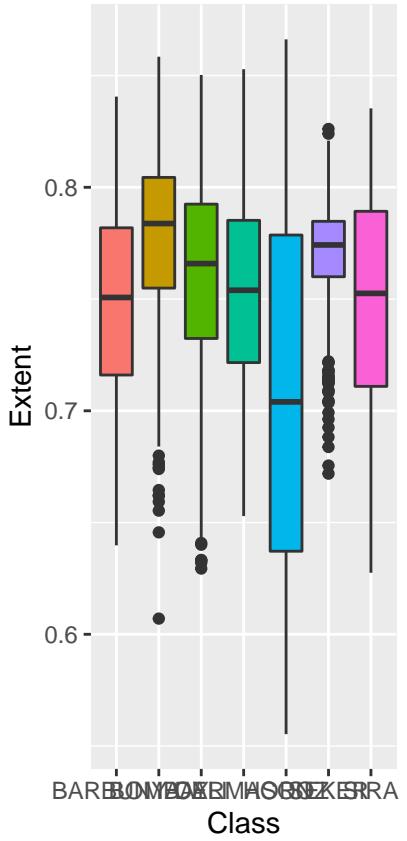
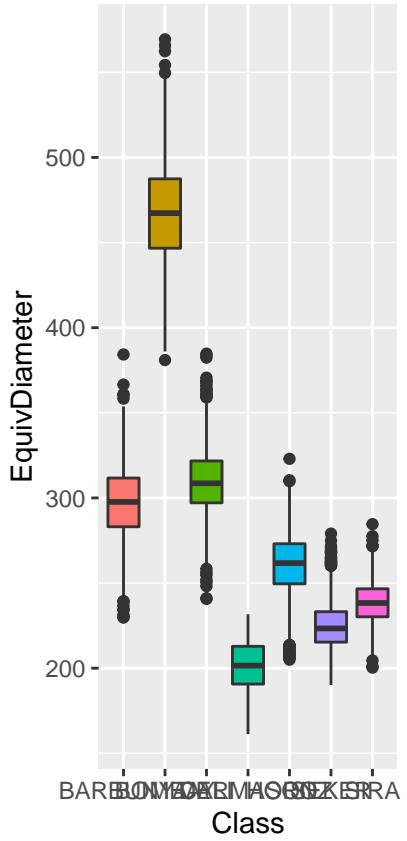
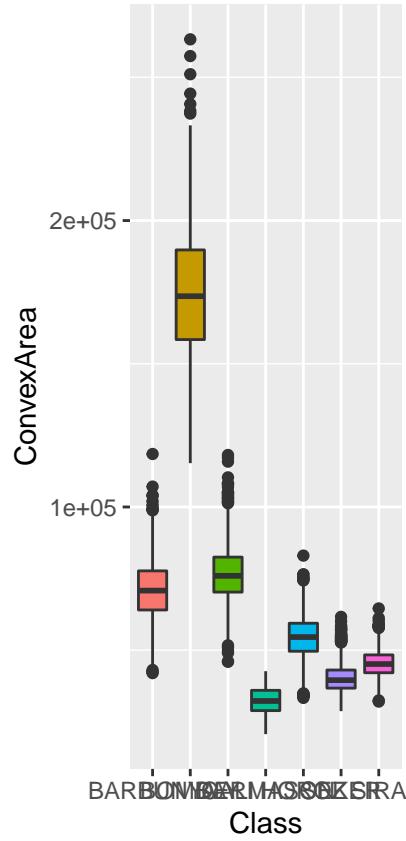
A Classe “Dermason” não apresenta outliers nas variáveis Area e EquivDiameter. Na Variável Extent as classes Barbunya, Dermason, Horoz e Sira não apresentam outliers. A classe Bombay pode ser facilmente discriminada em relação às demais classes, nas variáveis “Area”, “Perimeter”, “MinorAxisLength”, “MajorAxisLength”, “ConvexArea” e “EquivDiameter”.

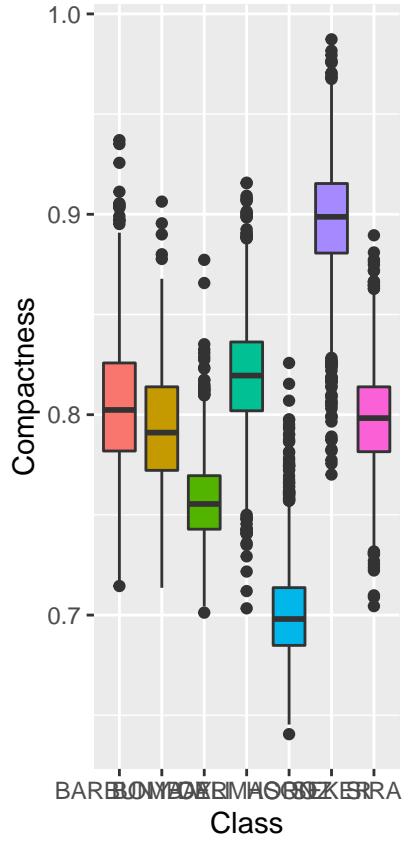
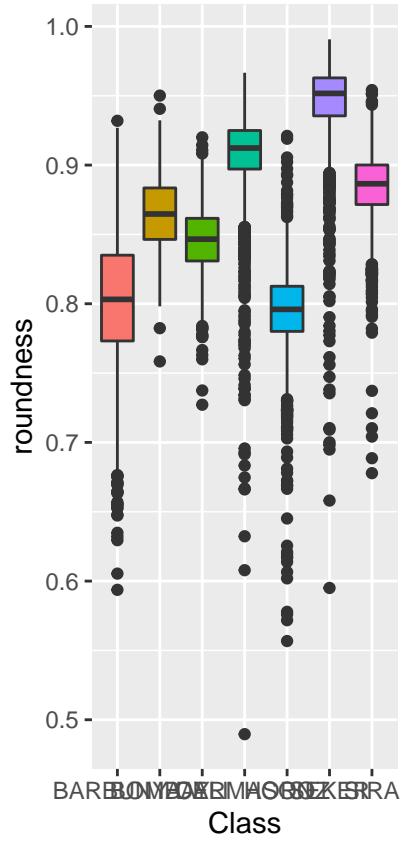
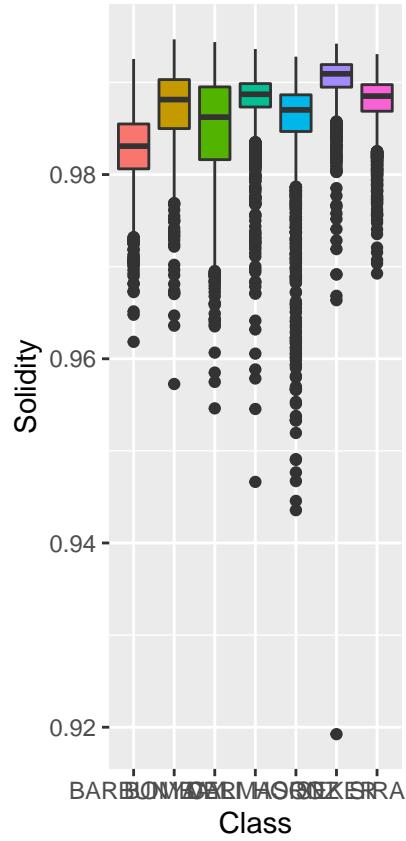
**5. Utilizando um gráfico boxplot por variável x classe (organize em 3 colunas), diga qual é a variável que teria maior poder de discriminação? Justifique a sua escolha.**

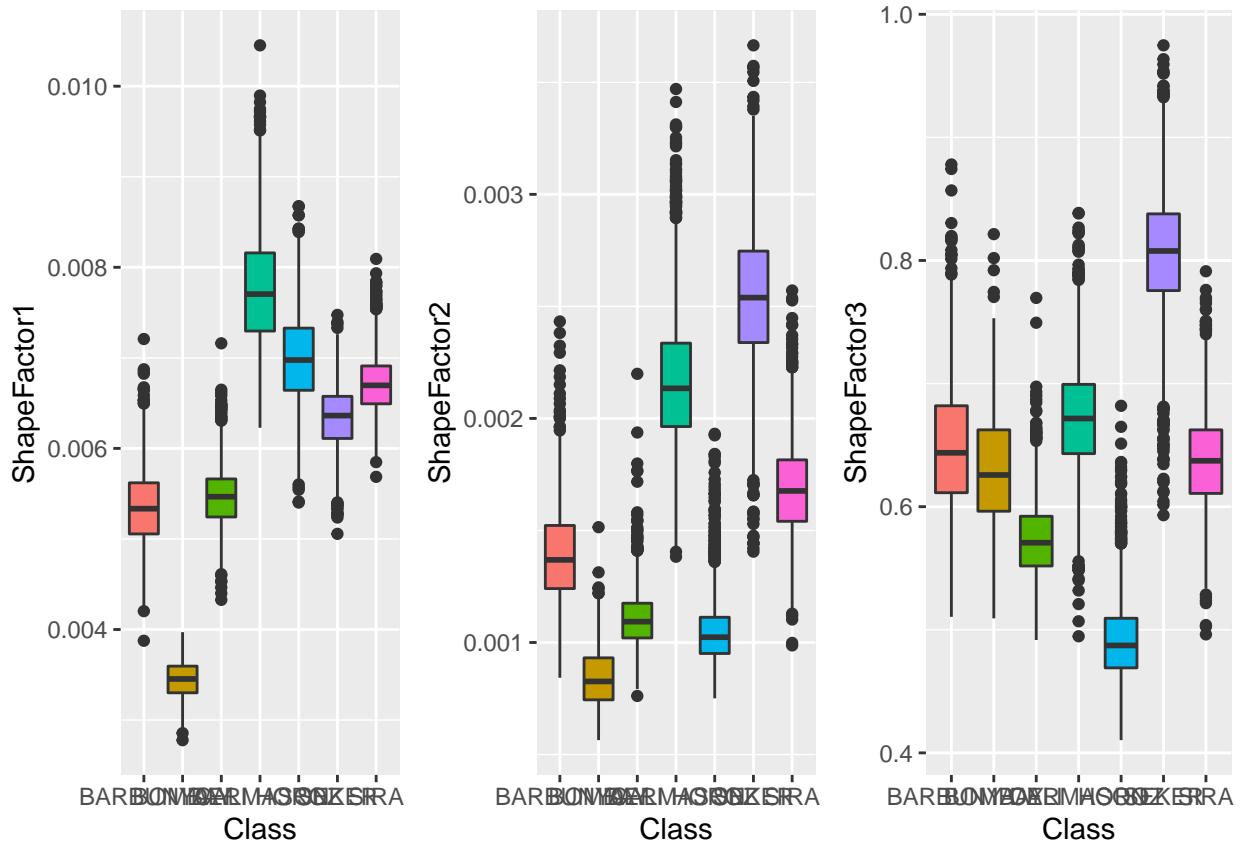
```
p = list()
for(i in 1:16){
  p[[i]] = ggplot(dados, aes_string(x="Class", y=names(dados)[i], fill="Class")) + geom_boxplot() + theme_minimal()
  if((i==3) || (i==6) || (i==9) || (i==12) || (i==15)){
    do.call(grid.arrange,c(p[(i-2):i], ncol=3))
  }
}
```







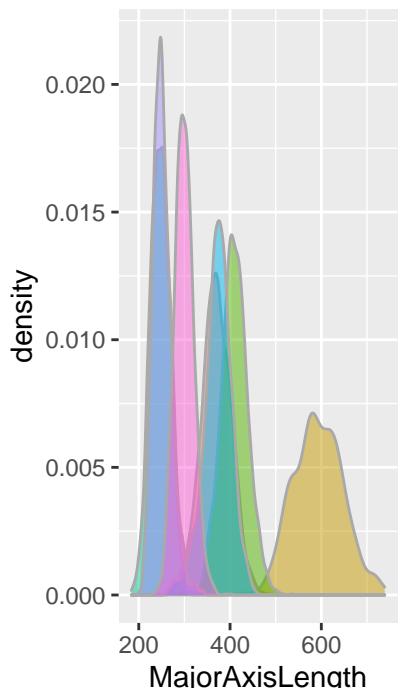
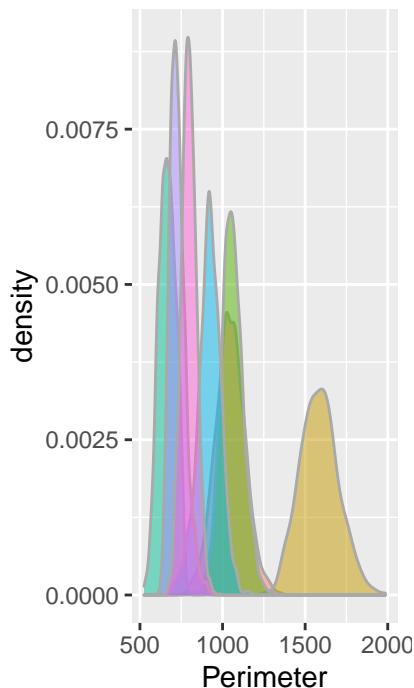
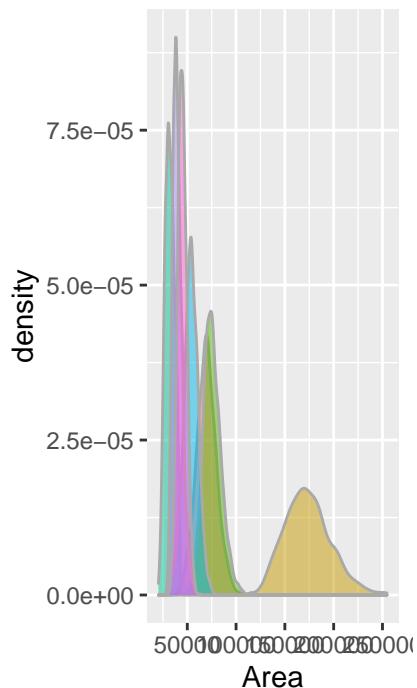


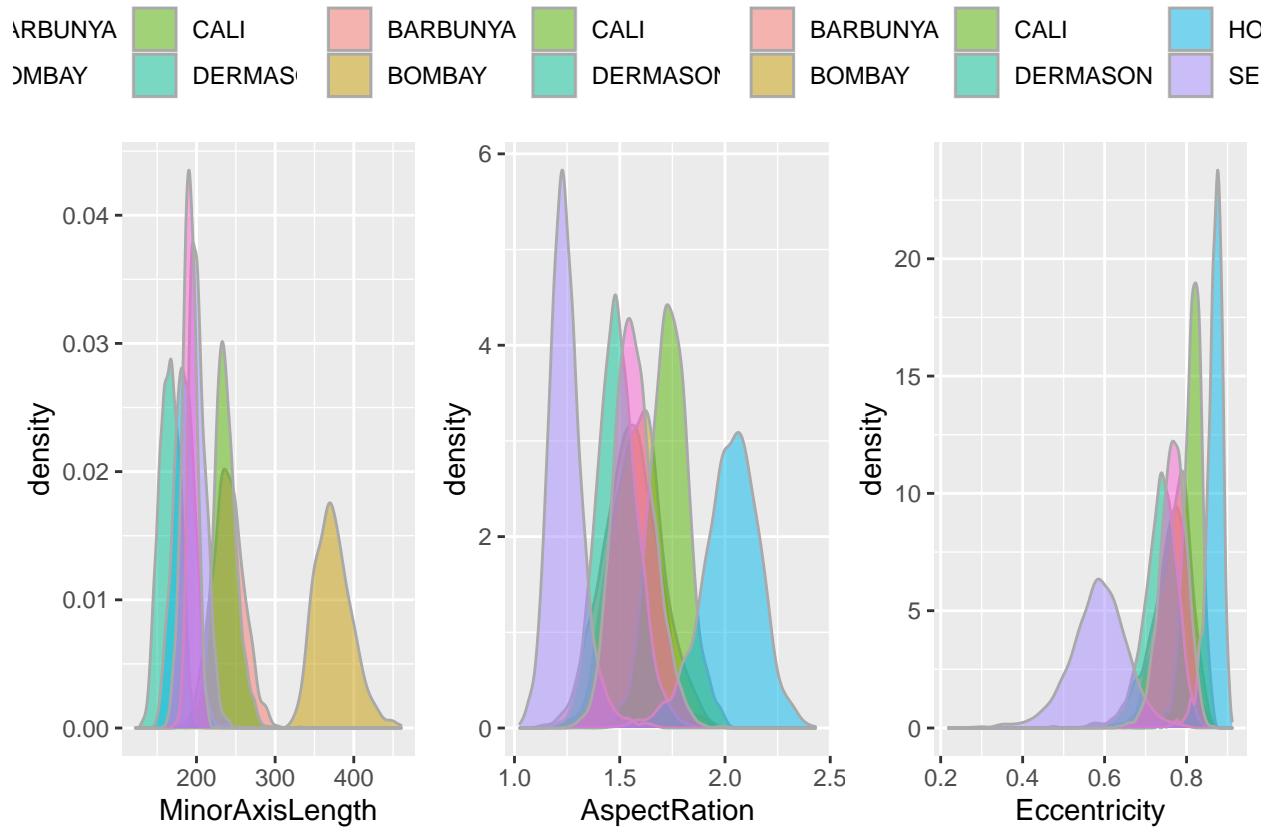


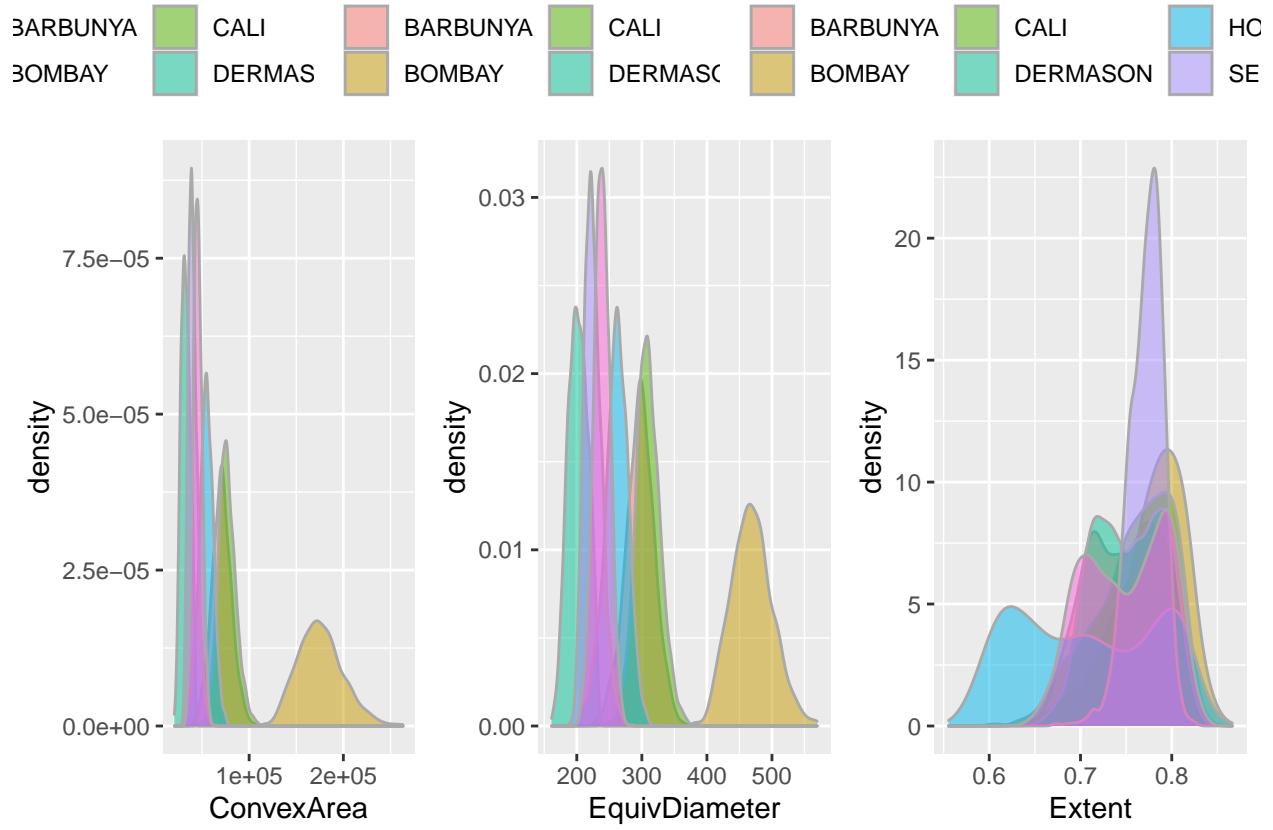
Algumas das variáveis que poderia ter mais poder de discriminação seria: “Eccentricity”, “ConvexArea”, “ShapedFactor1” e “ShapedFactor2”. Os principais pontos pelo qual estas variáveis se destacam são seus pontos mínimos e máximos, serem próximos e com poucos, ou nenhum, outlier, como também os pontos são diferentes entre as classes.

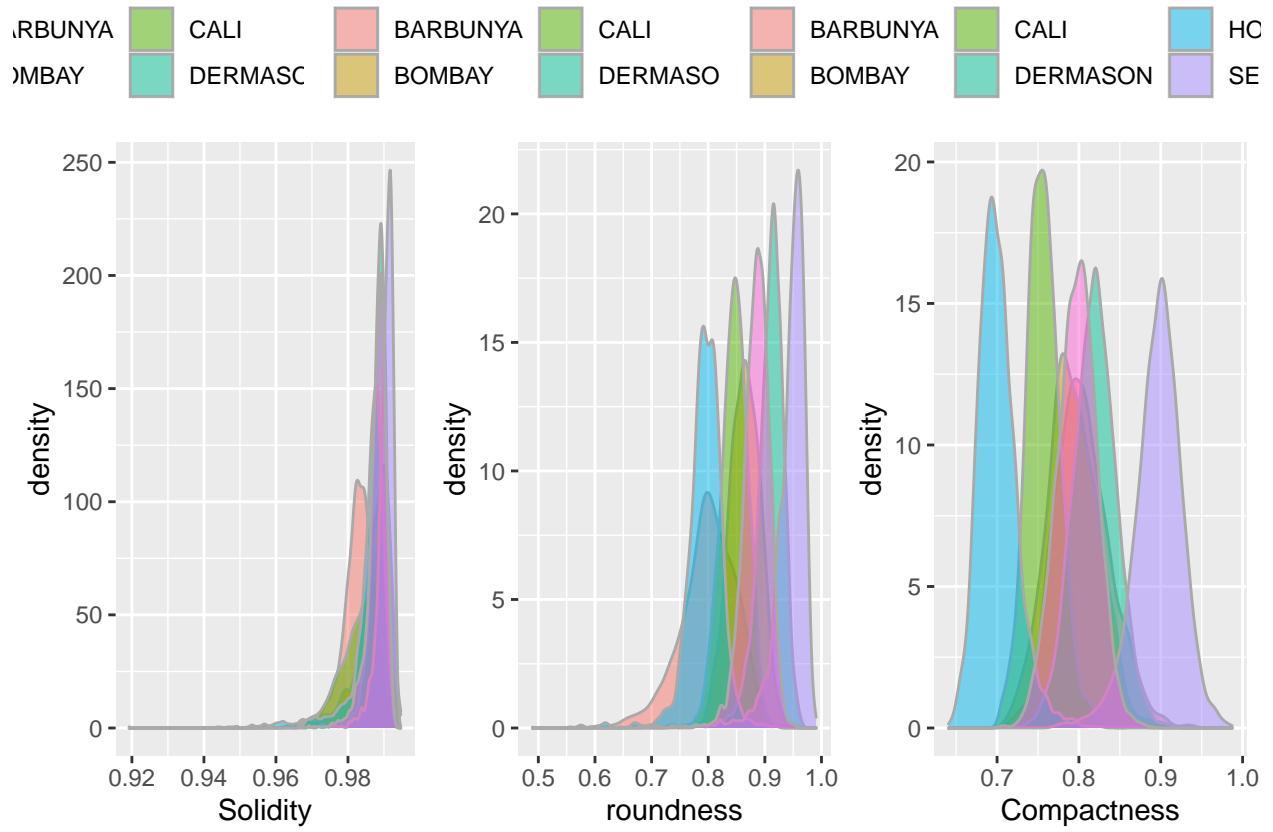
6.Utilizando gráficos de densidade por variável (organize em 3 colunas), é possível fazer alguma afirmação sobre a discriminabilidade de alguma classe? Pode utilizar os boxplots gerados na etapa anterior para auxiliar nas conclusões.

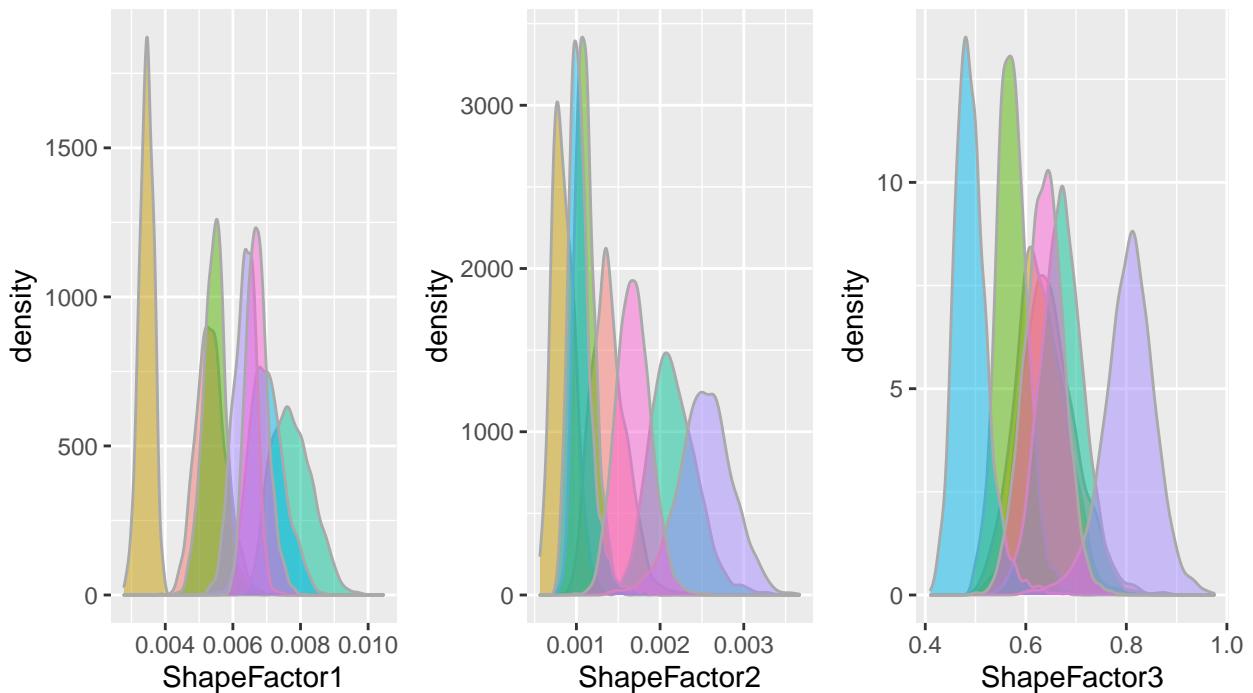
```
p = list()
for(i in 1:16){
  p[[i]] = ggplot(dados, aes_string(x=names(dados)[i], fill="Class")) +
    geom_density(alpha=0.5, color="darkgray") +
    theme(legend.position="top", legend.title = element_blank())
  if((i==3) || (i==6) || (i==9) || (i==12) || (i==15)){
    do.call(grid.arrange, c(p[(i-2):i], ncol=3))
  }
}
```







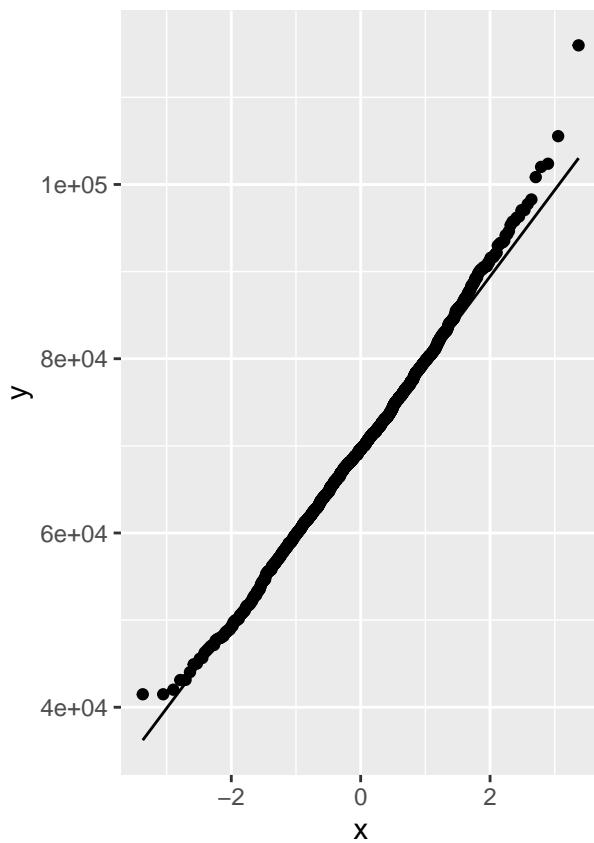
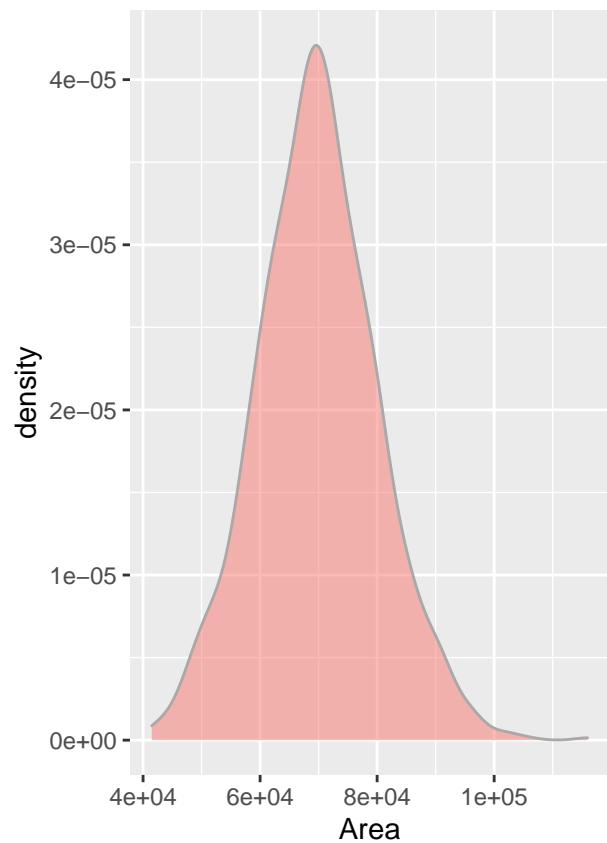


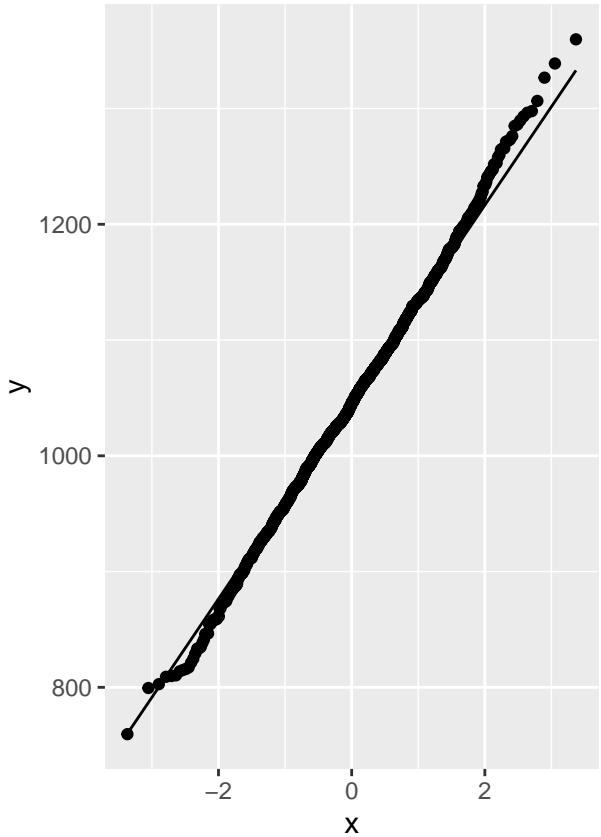
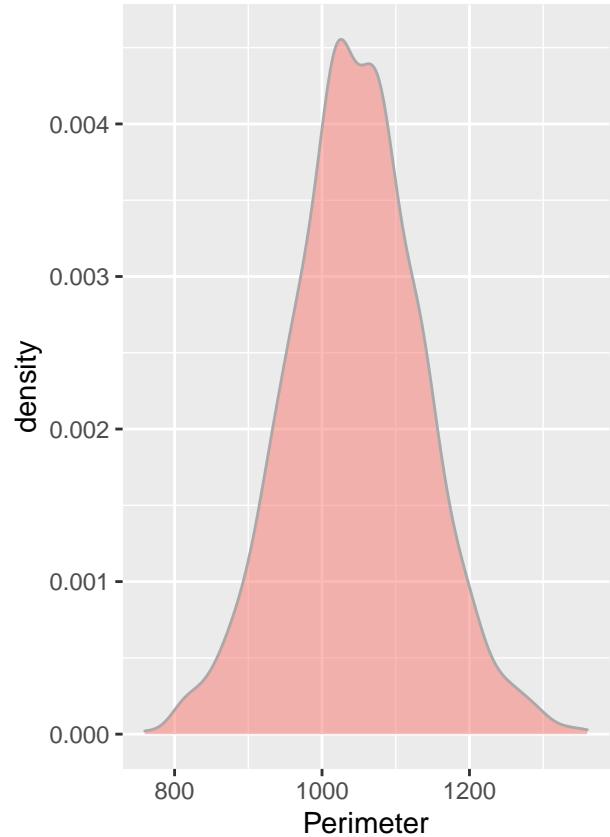


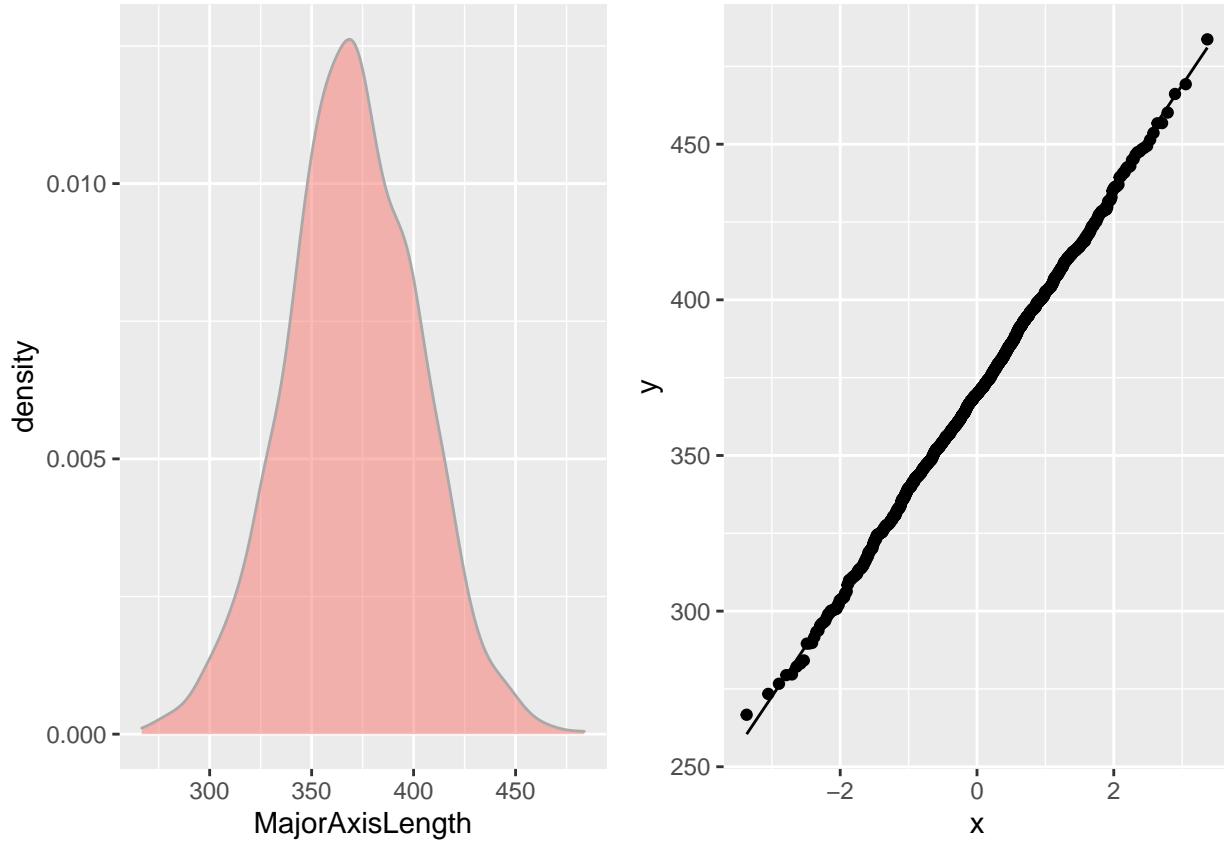
As Classes Bombay e Seker possuem grande discrepancia em algumas Variáveis.

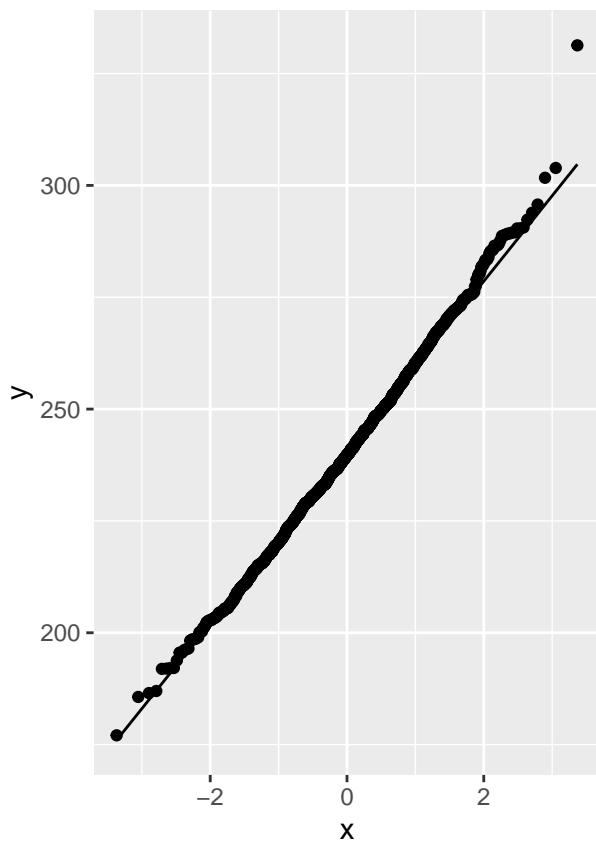
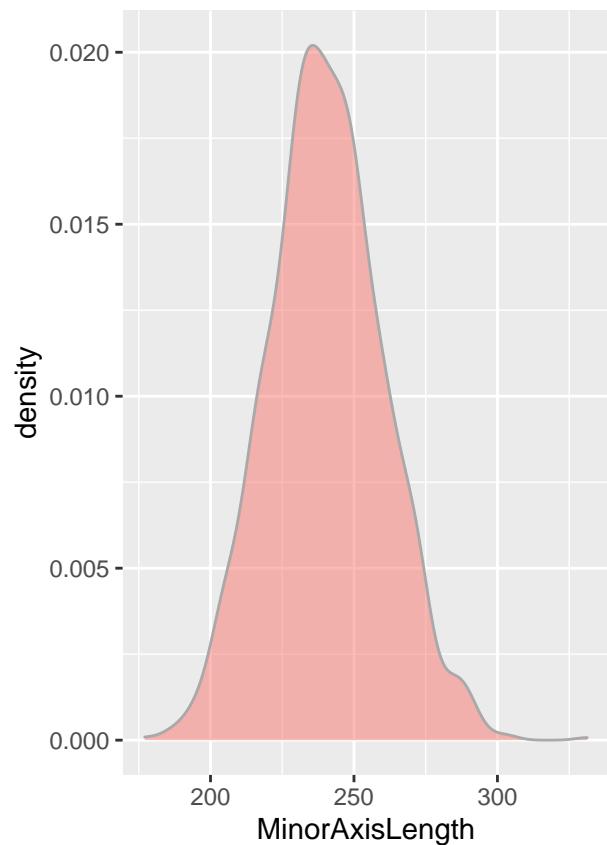
**7. Algumas variável (por classe) possui uma distribuição normal (curva do sino)?**  
É possível verificar numericamente se é verdade?

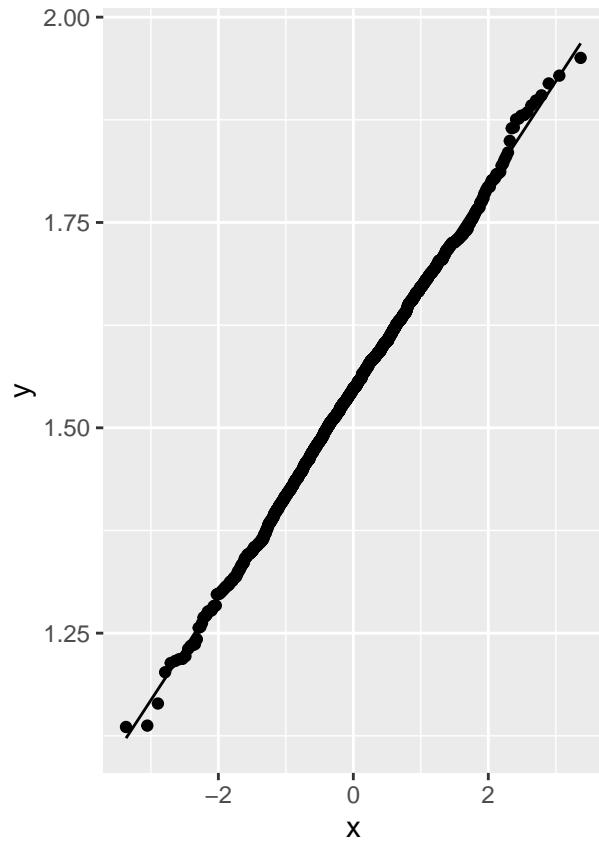
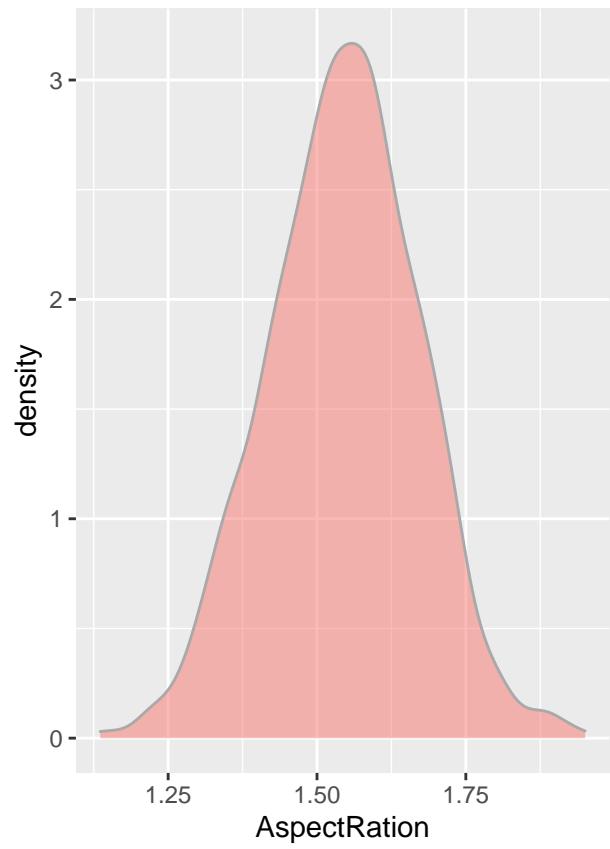
```
leves = list("BARBUNYA", "BOMBAY", "CALI", "DERMASON", "HOROZ", "SEKER", "SIRA")
p = list()
for(j in 1:7){
  print(leves[j])
  for(i in 1:16){
    p[[((i-1)*2+1)]] = ggplot(dados[dados$Class==leves[j],], aes_string(x=names(dados)[i],fill="Class"))
    p[[i*2]] = ggplot(dados[dados$Class==leves[j],], aes_string(sample=names(dados)[i]),color="Class")
    if((i%%2)==0){
      do.call(grid.arrange,c(p[[(i-1):i],ncol=2]))
    }
  }
}
## [[1]]
## [1] "BARBUNYA"
```

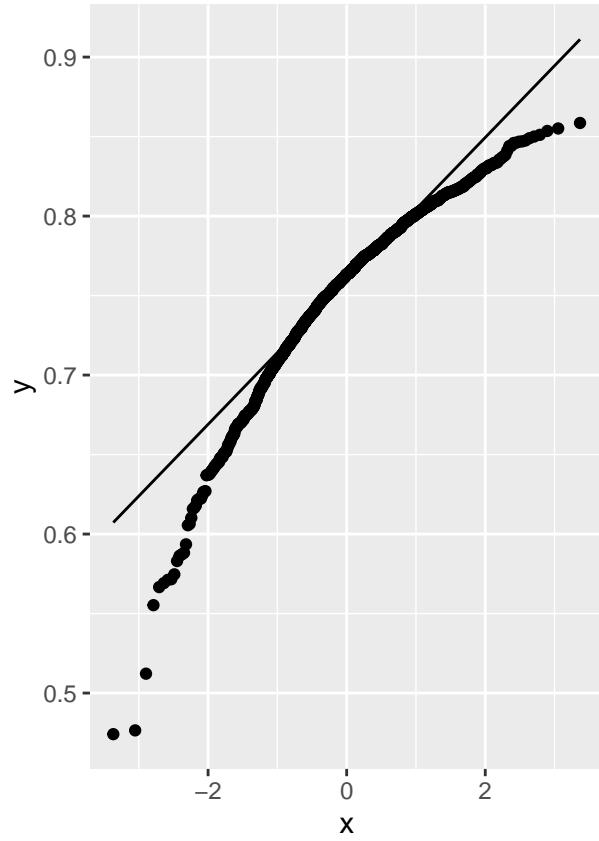
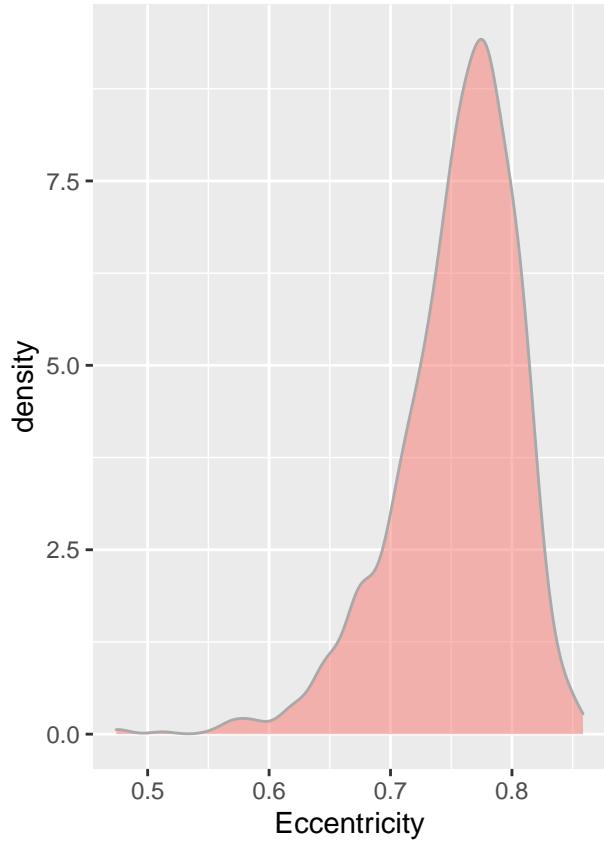


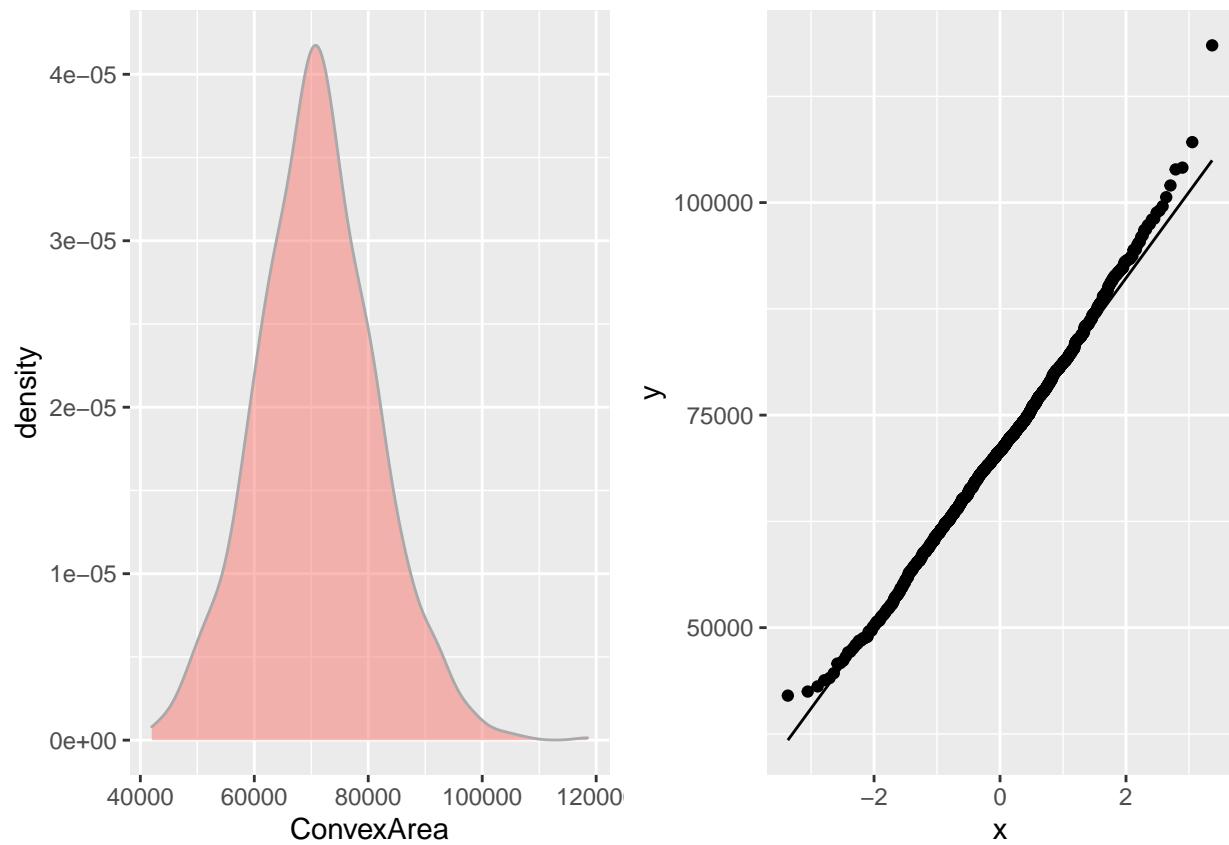


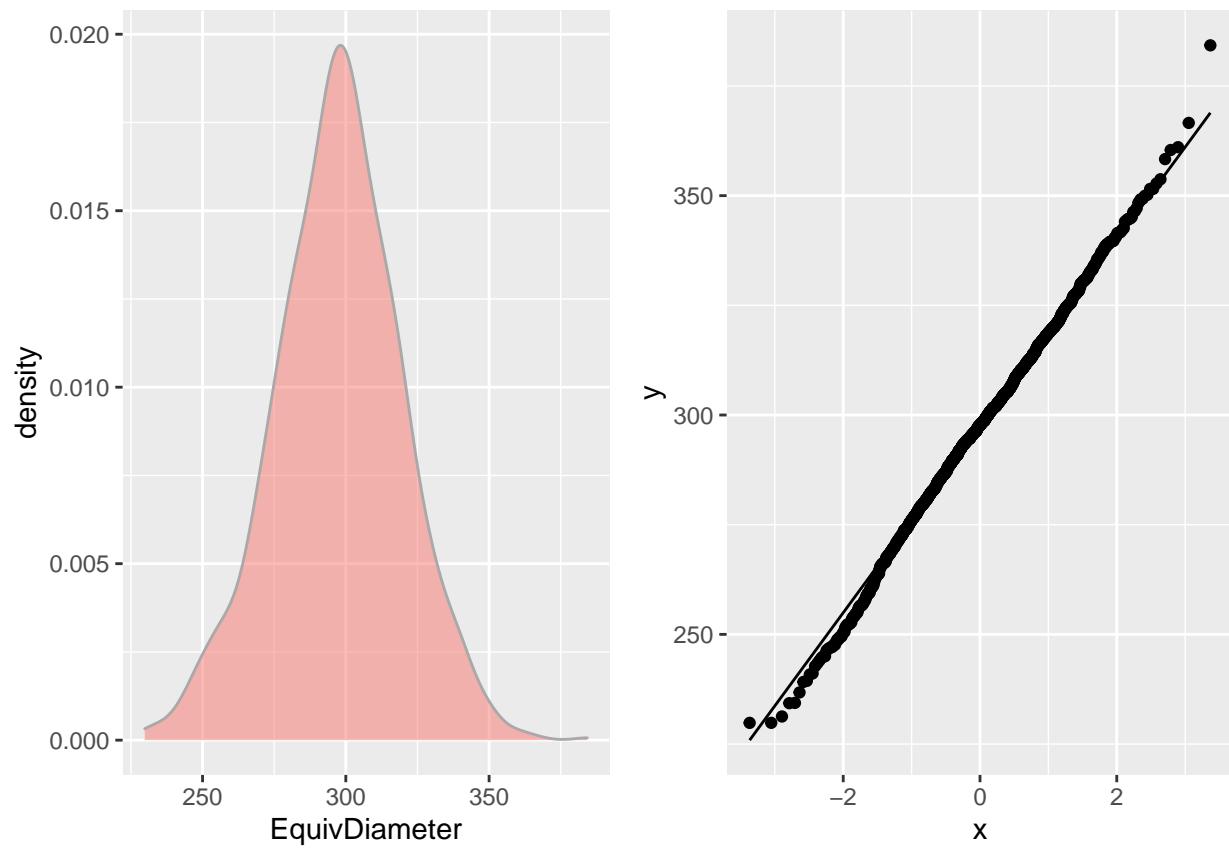




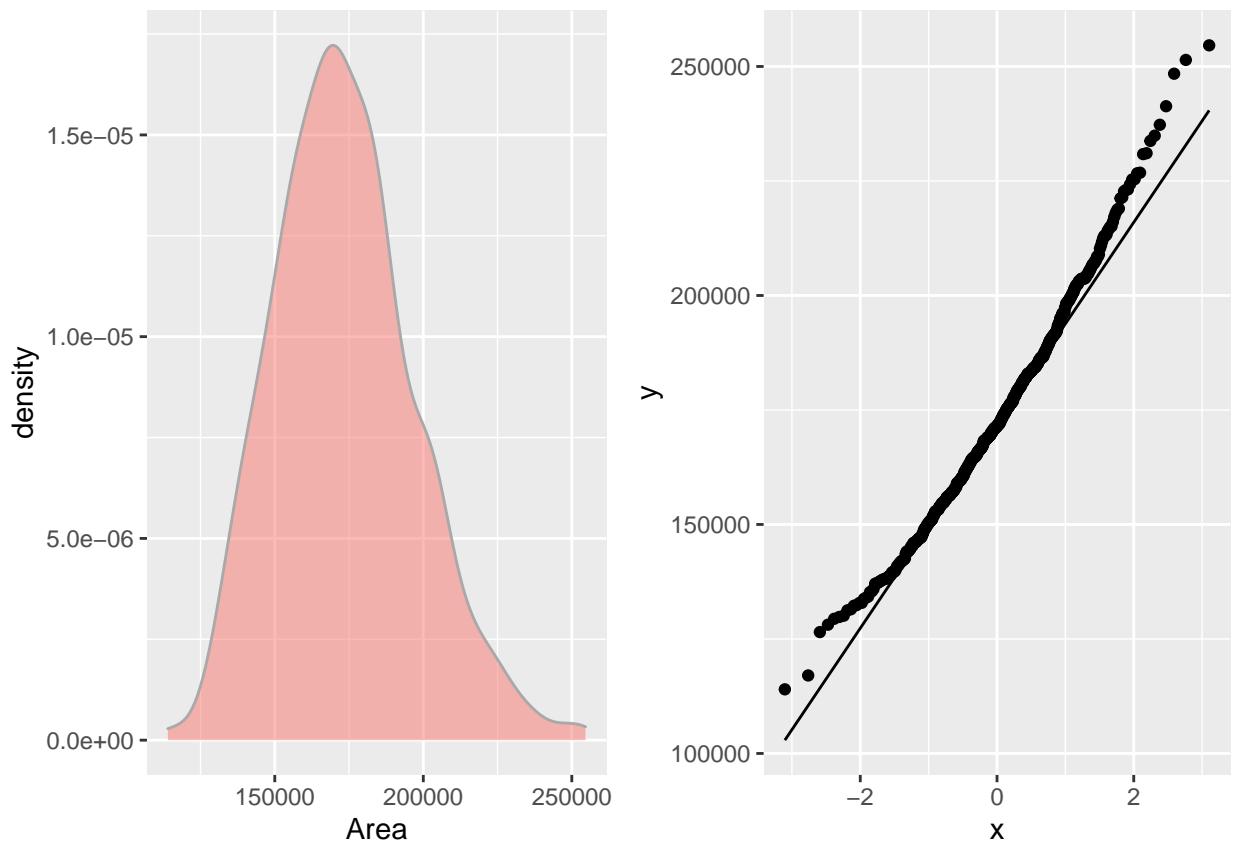


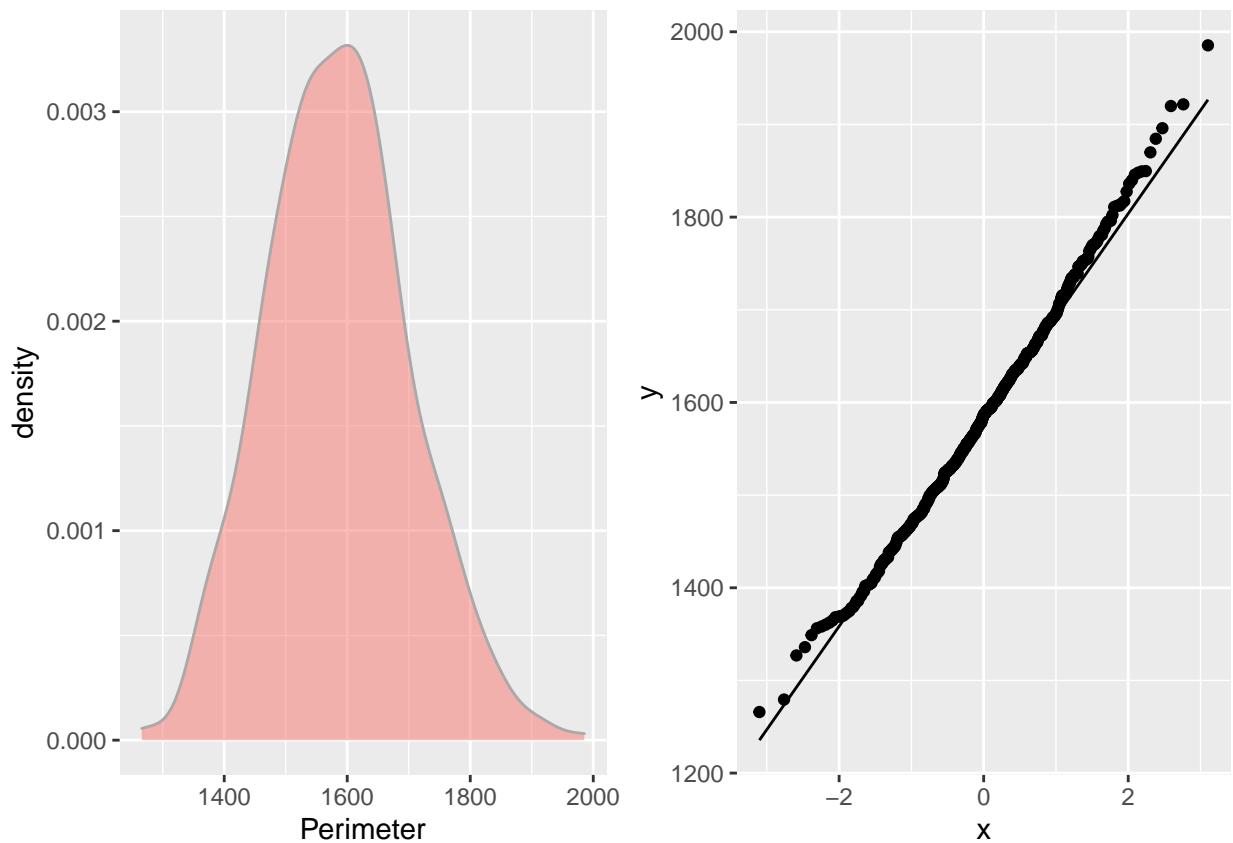


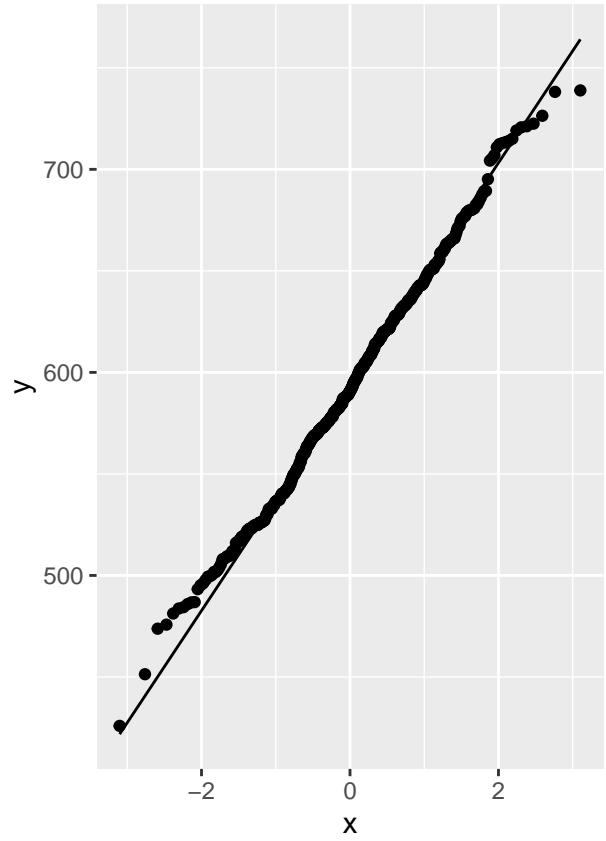
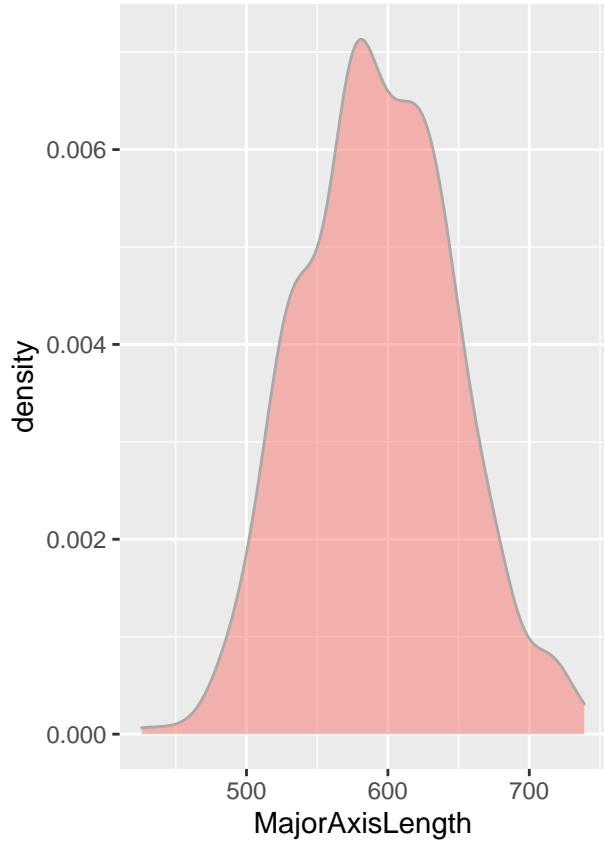


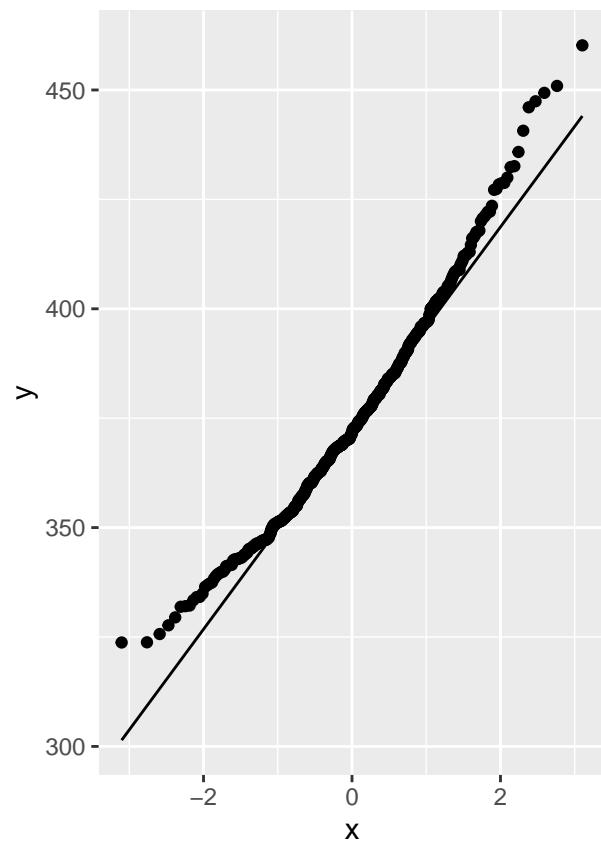
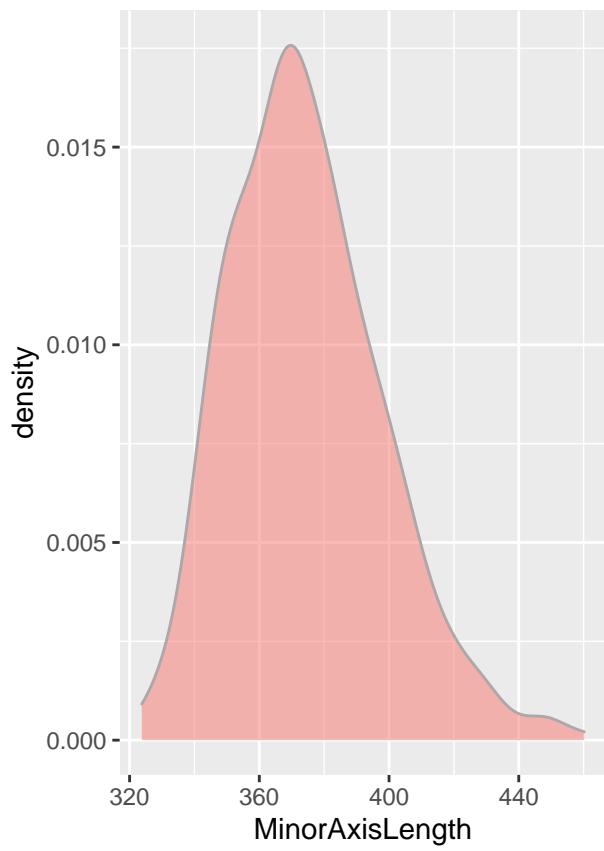


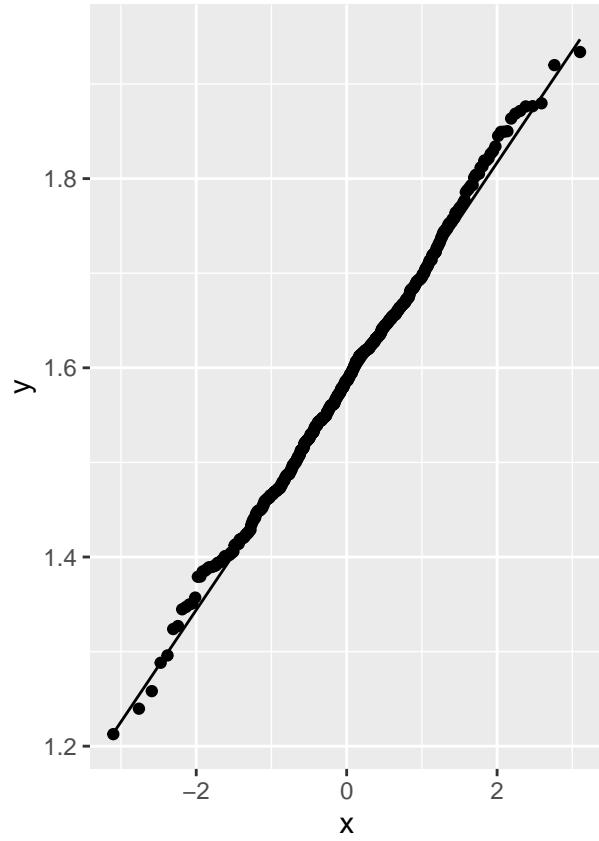
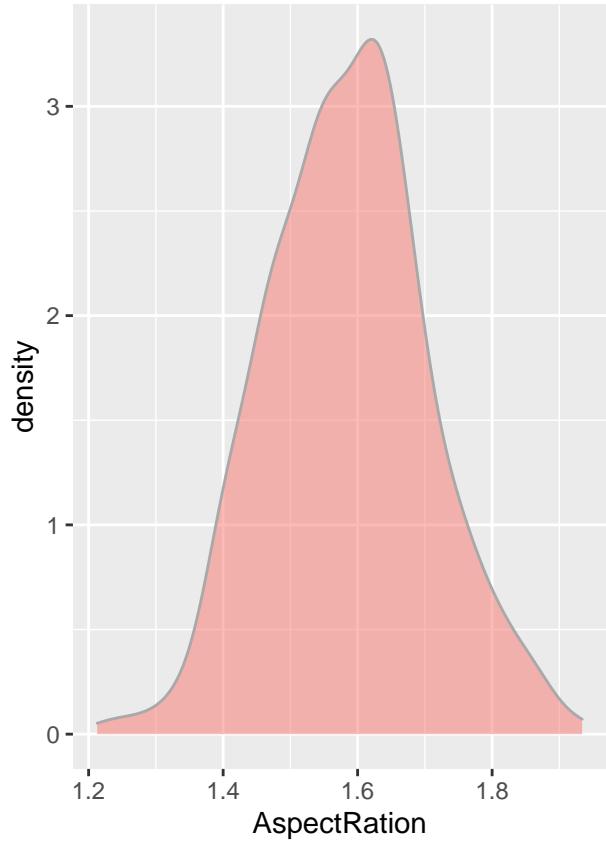
```
## [1]
## [1] "BOMBAY"
```

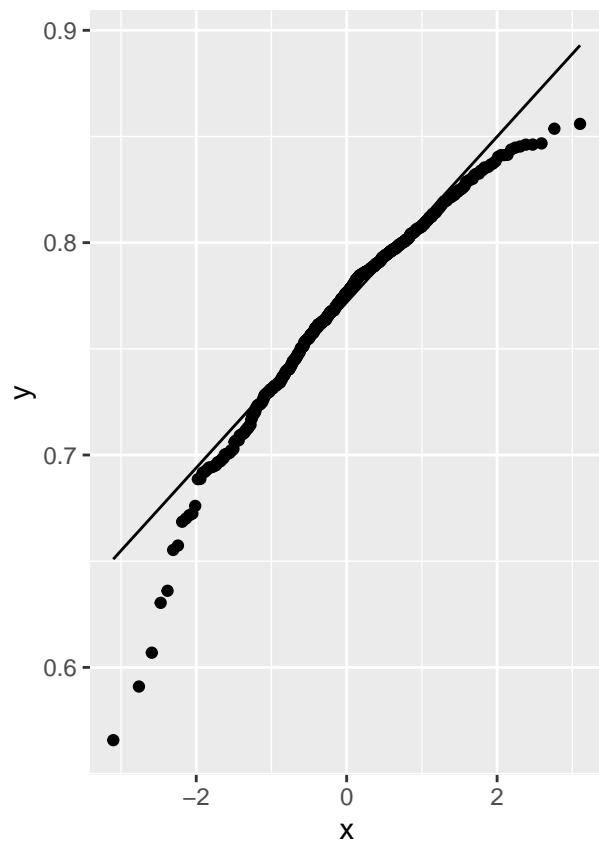
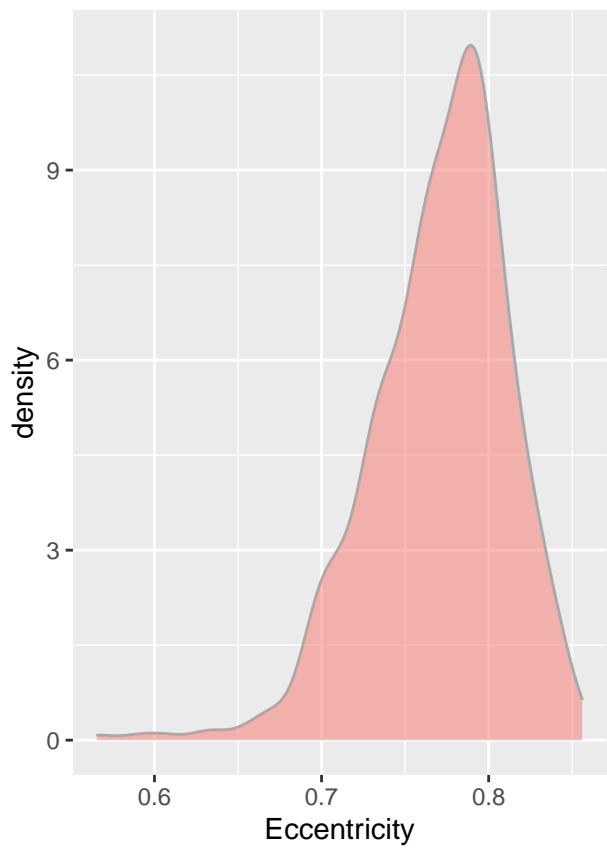


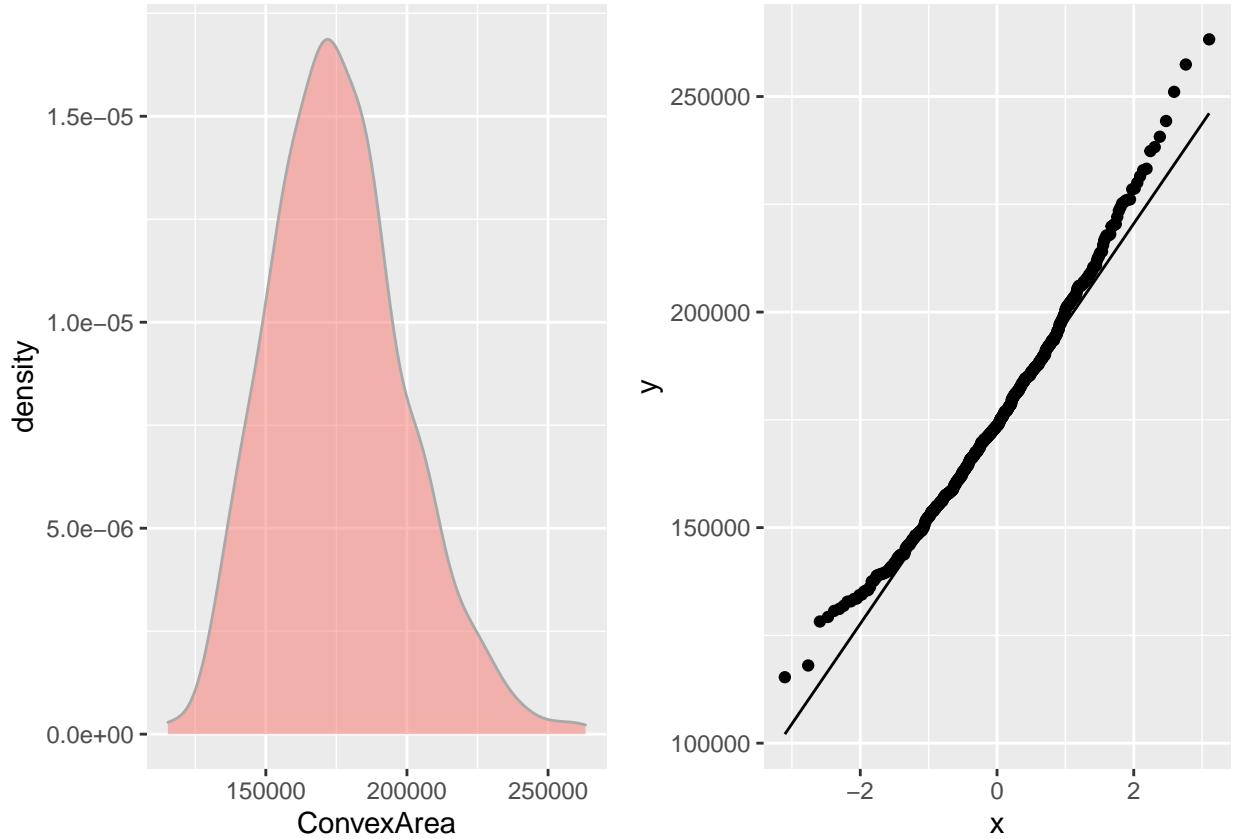


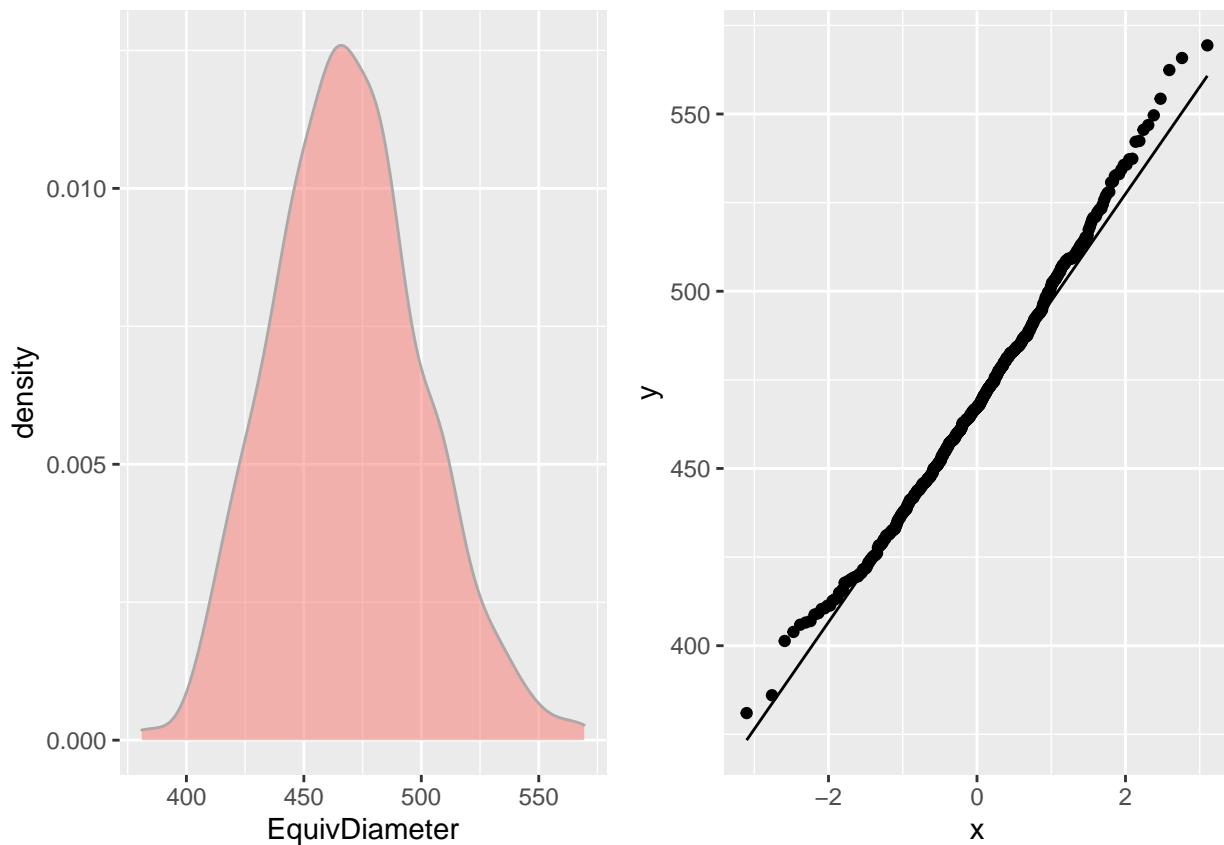




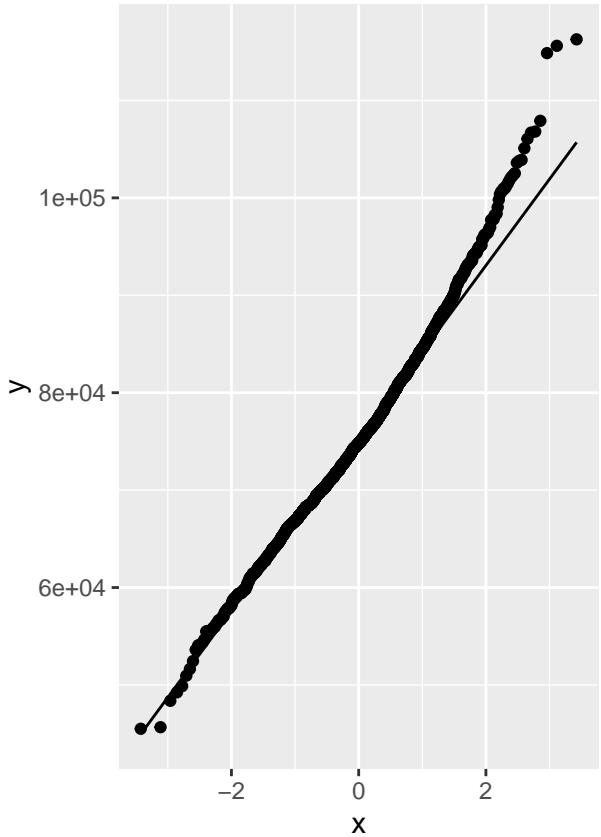
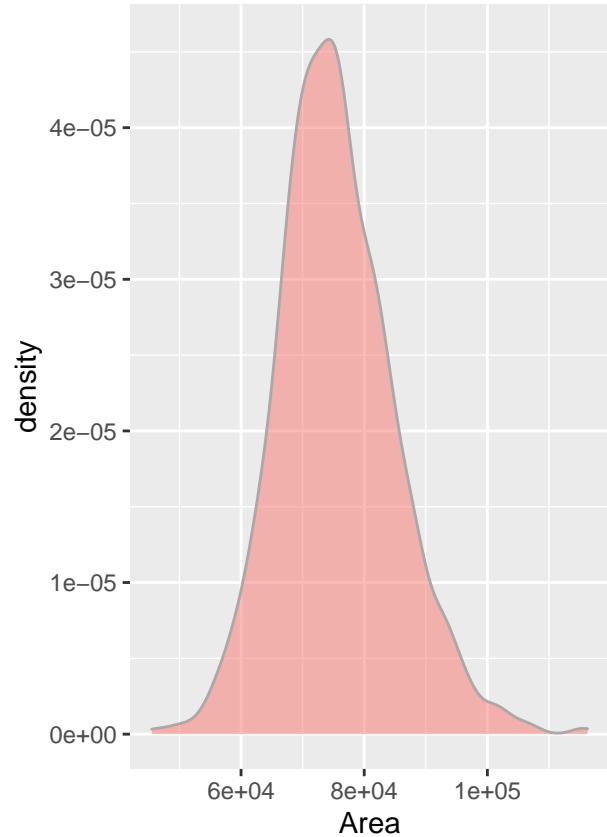


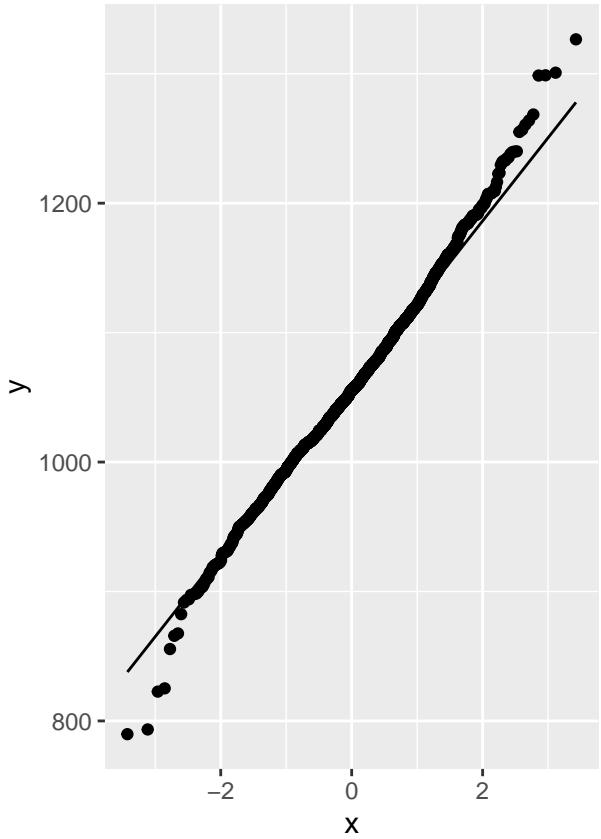
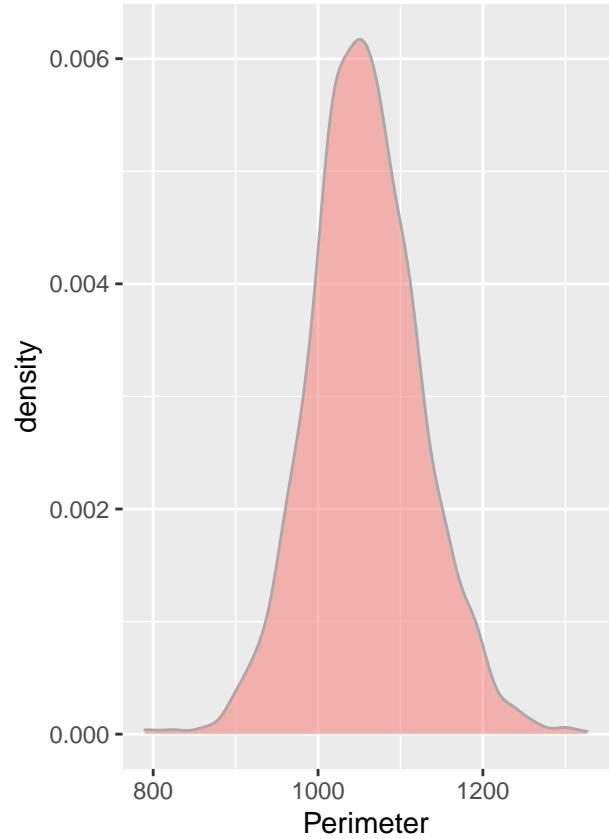


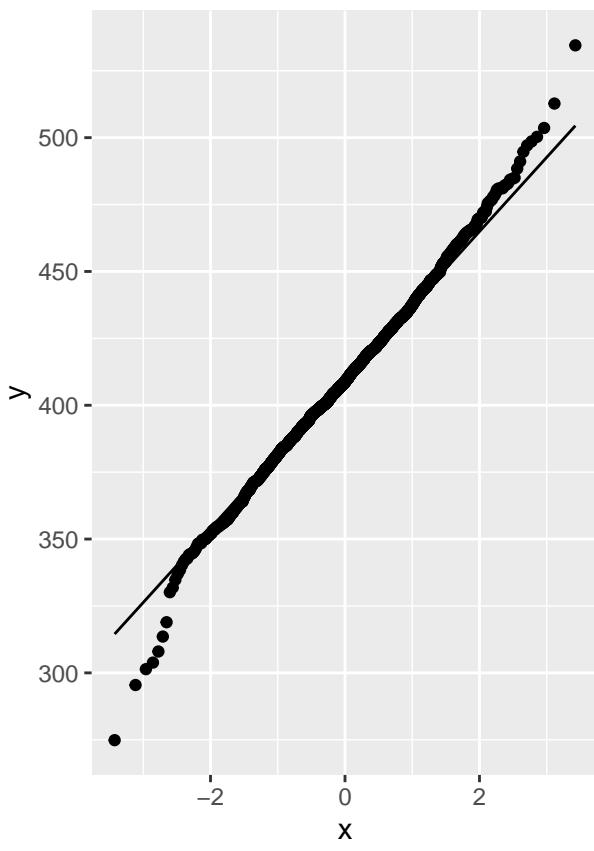
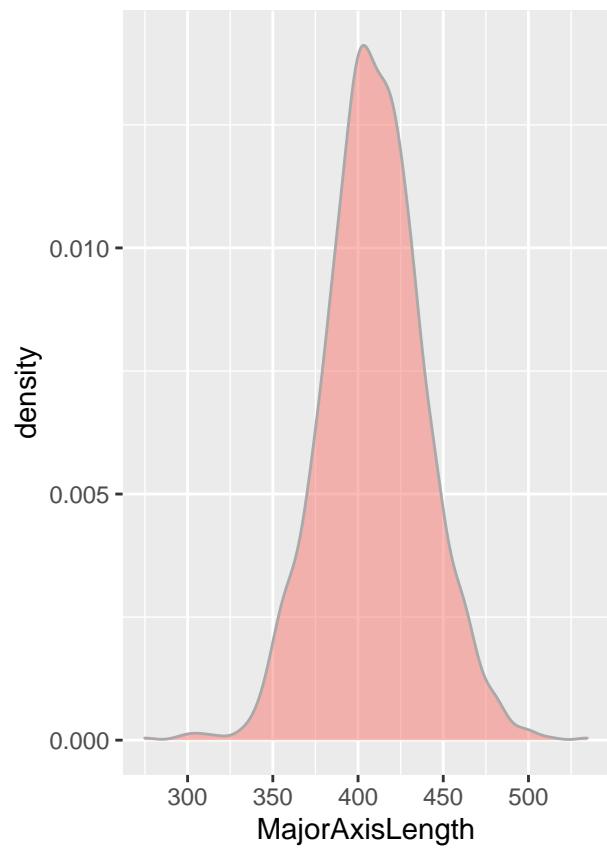


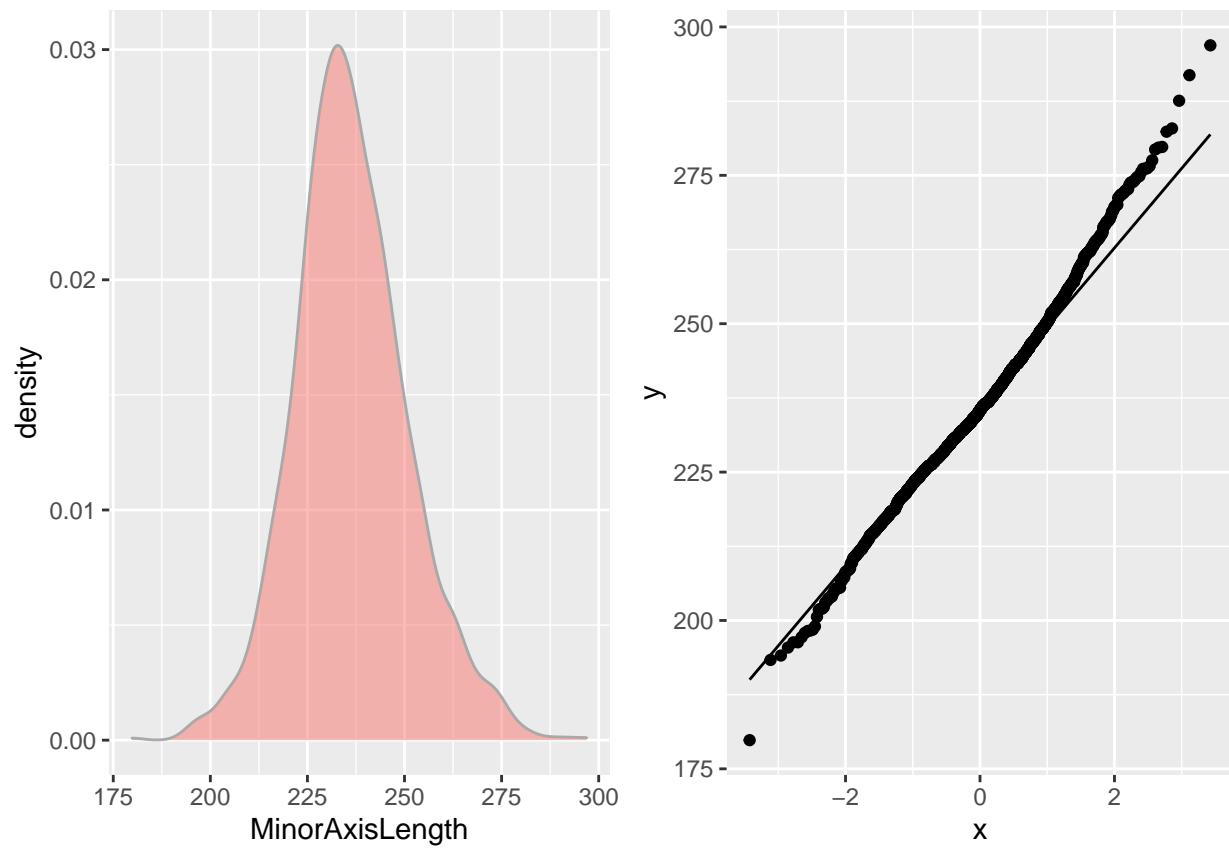


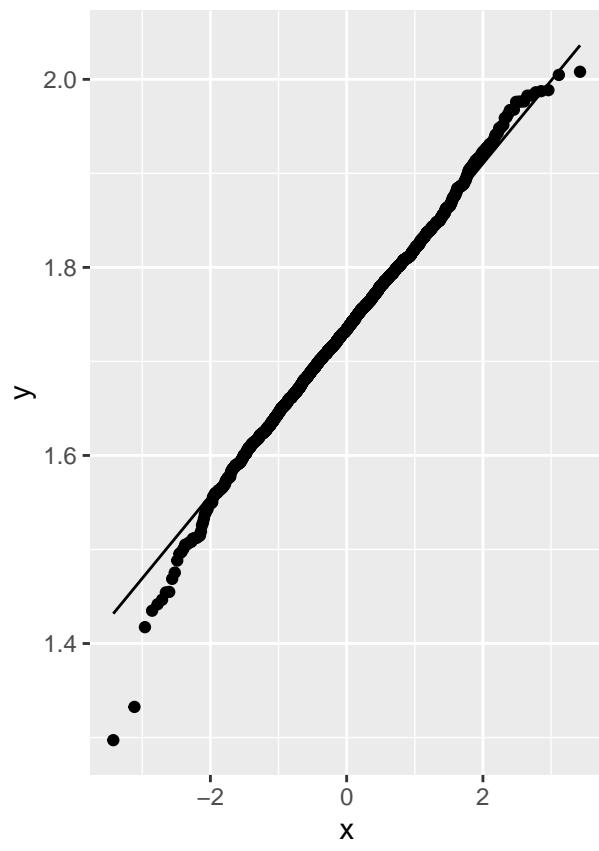
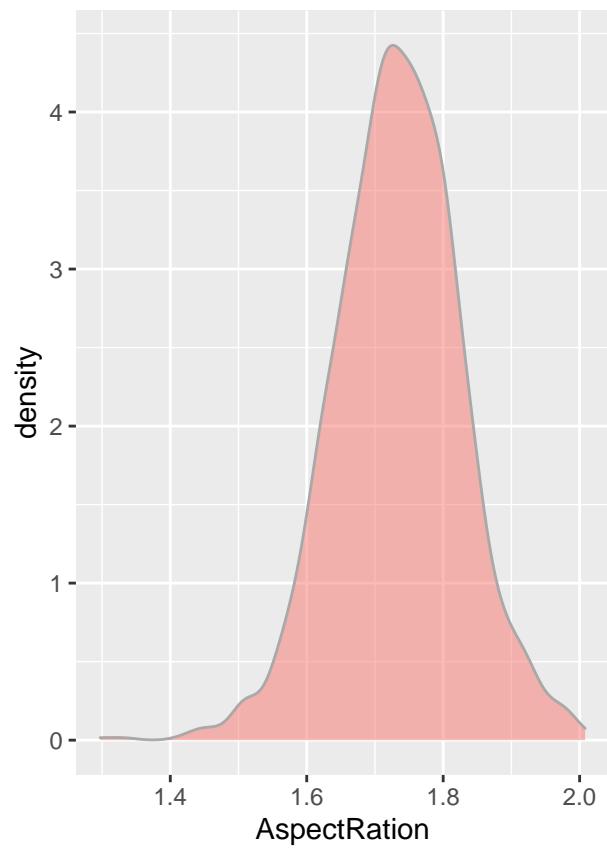
```
## [1]
## [1] "CALI"
```

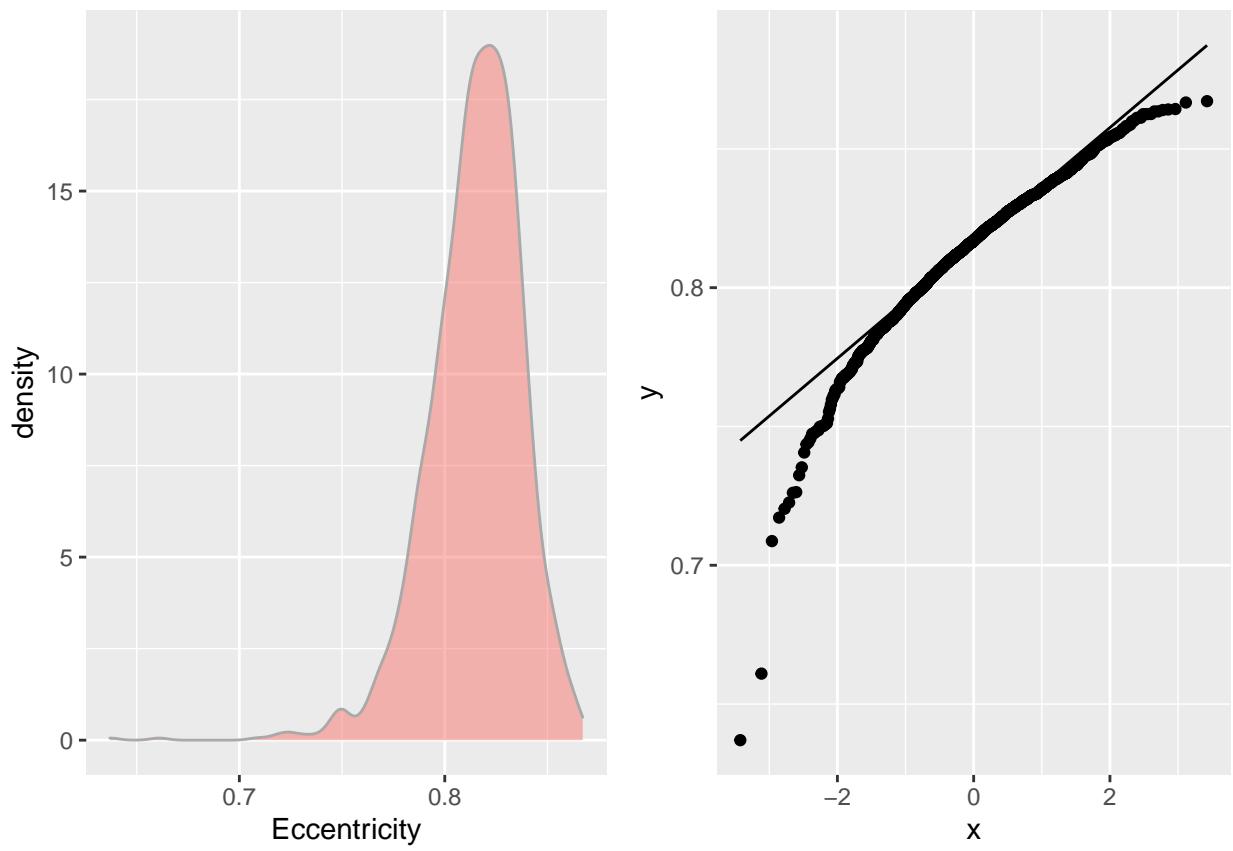


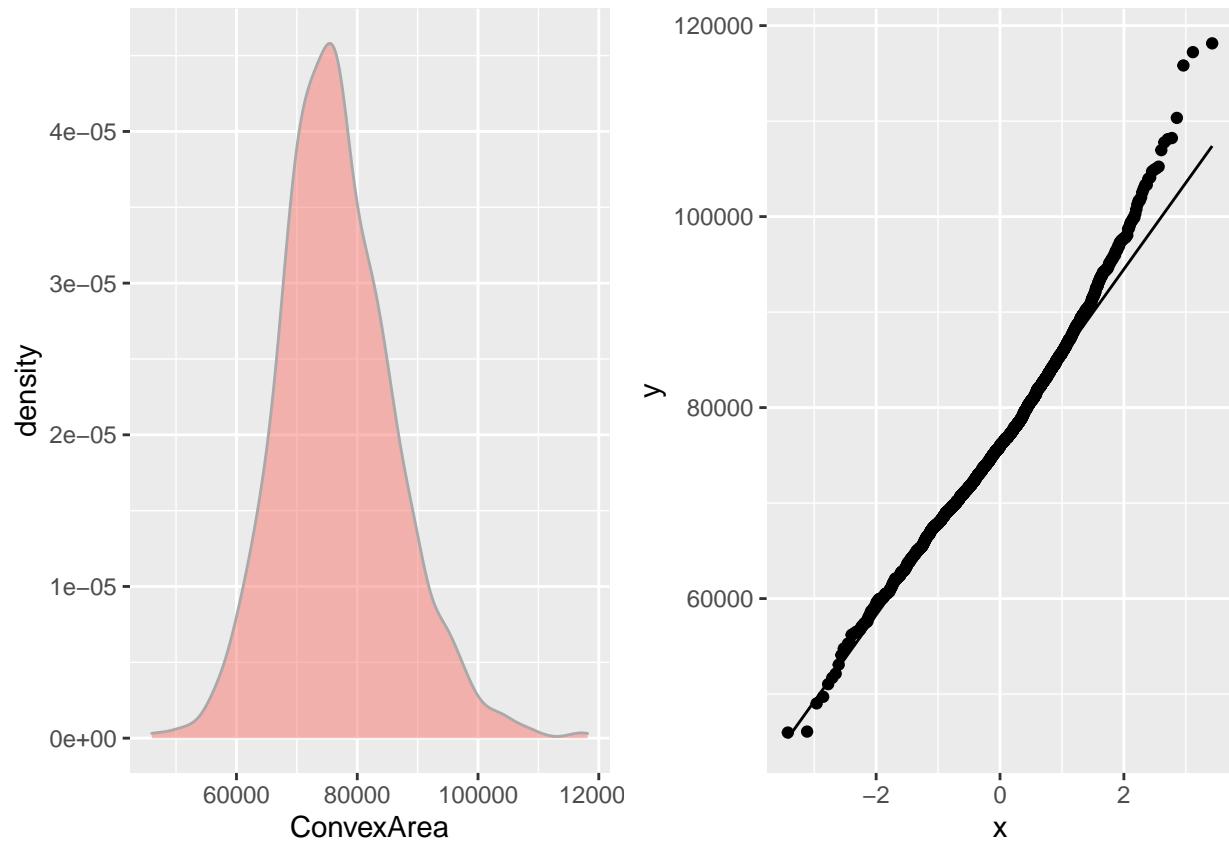


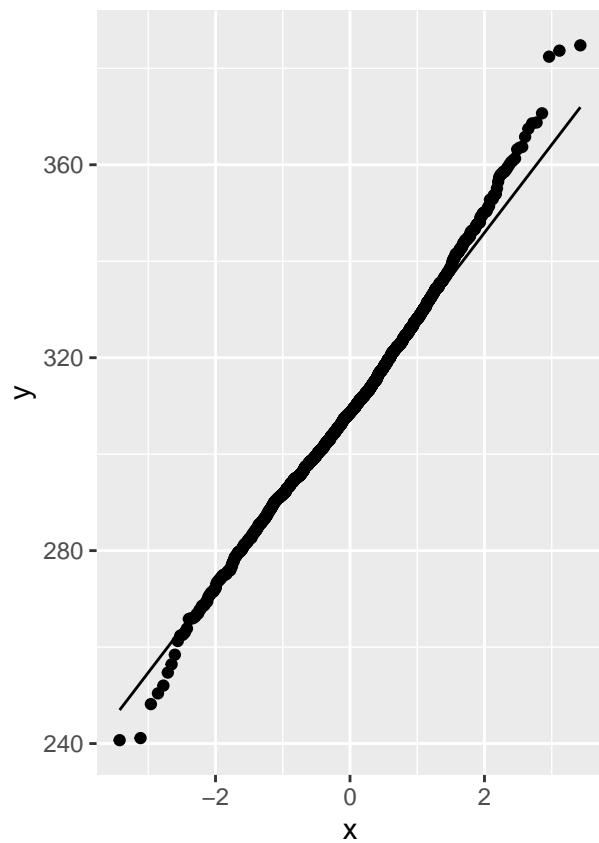
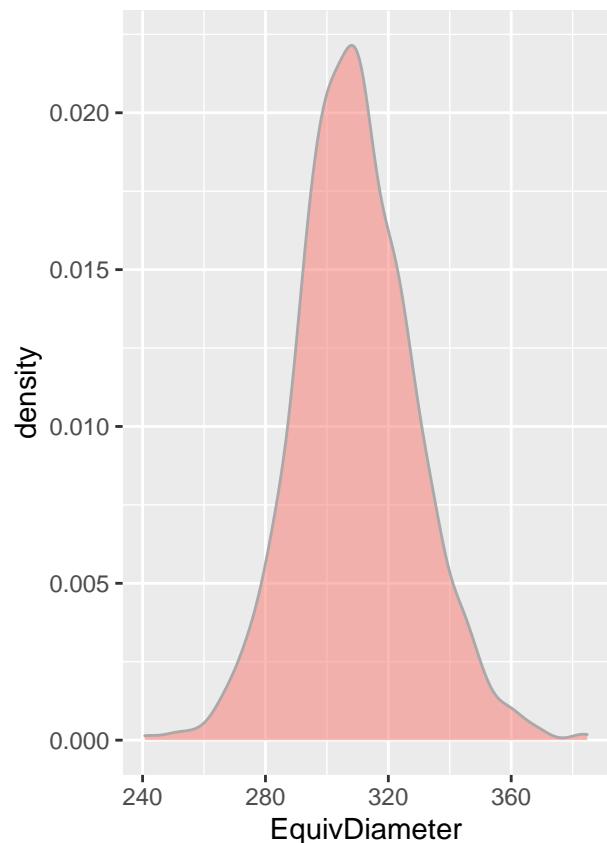




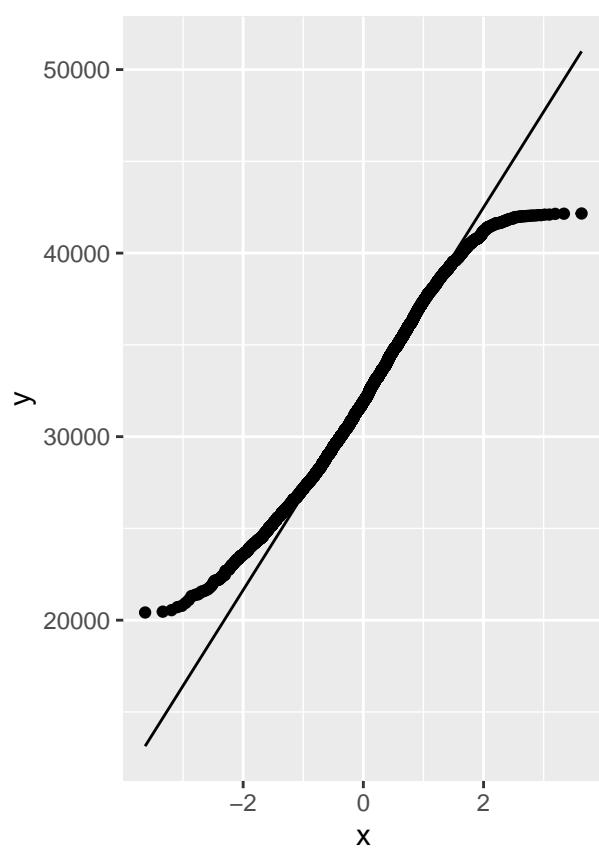
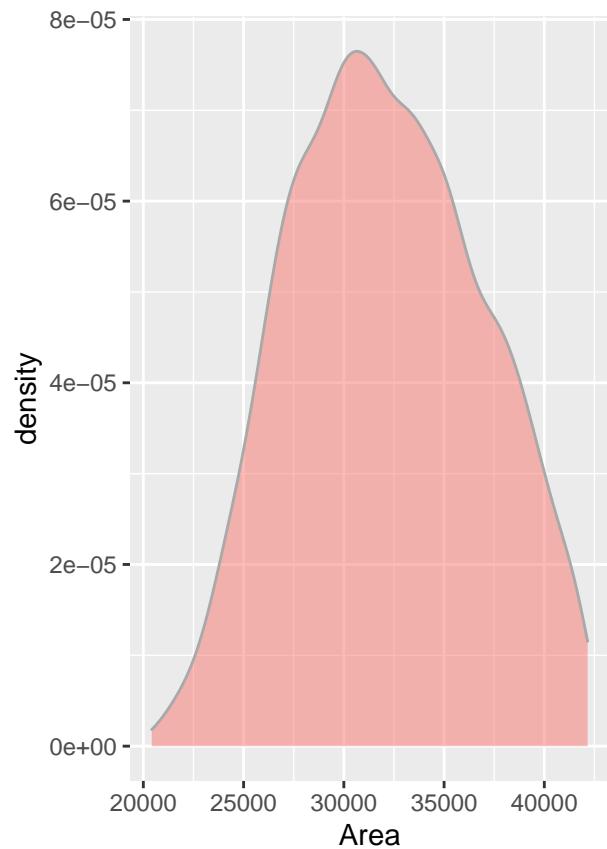


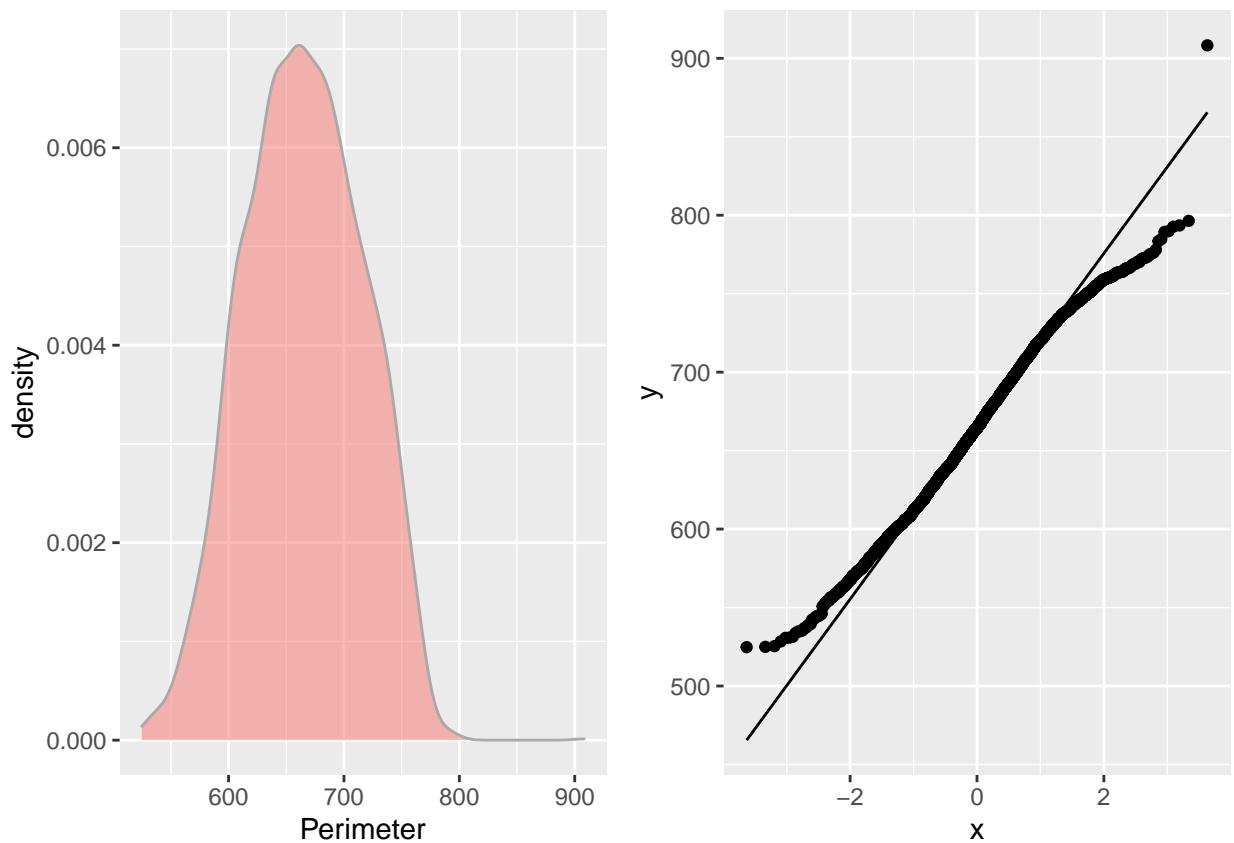


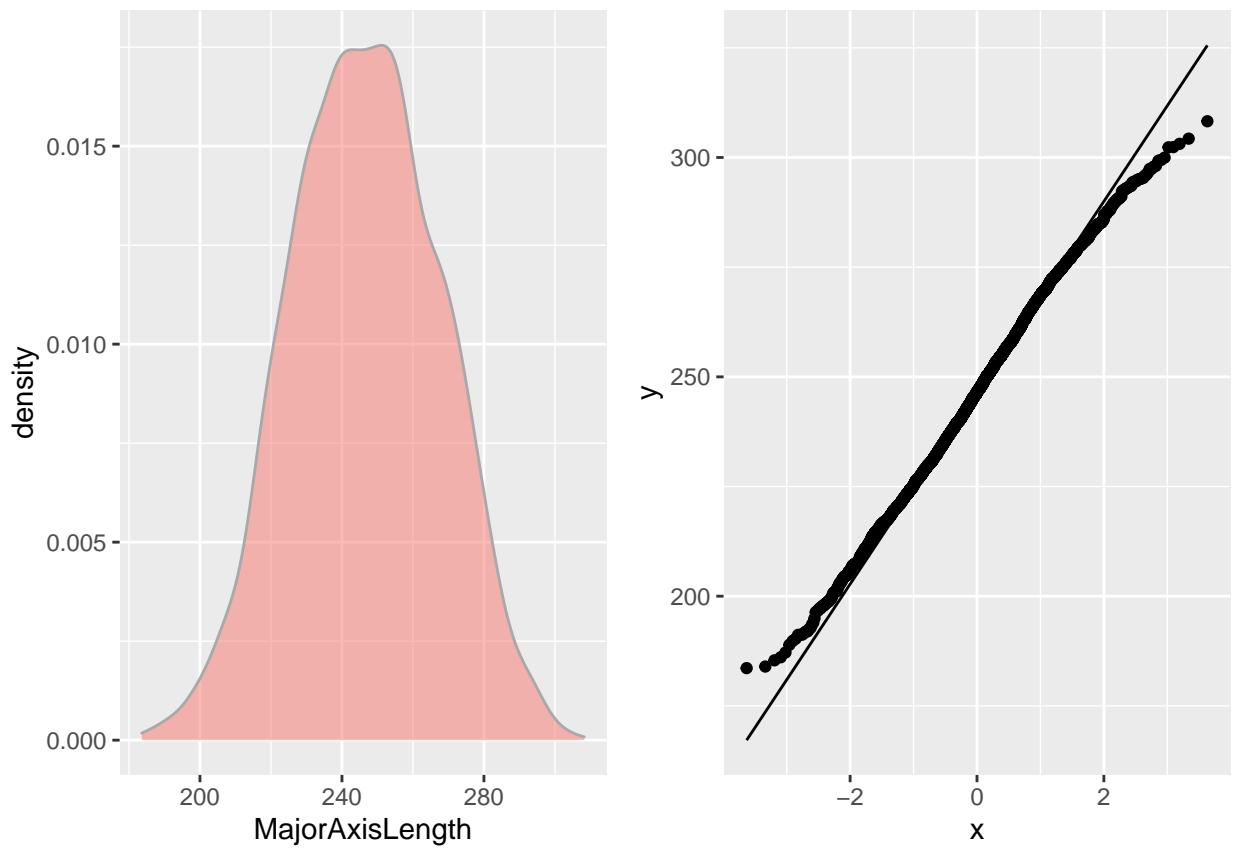


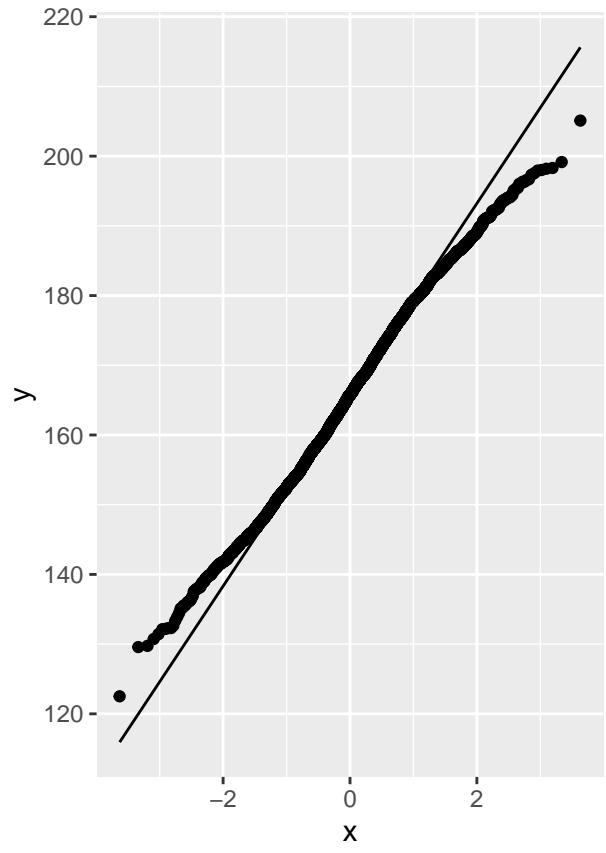
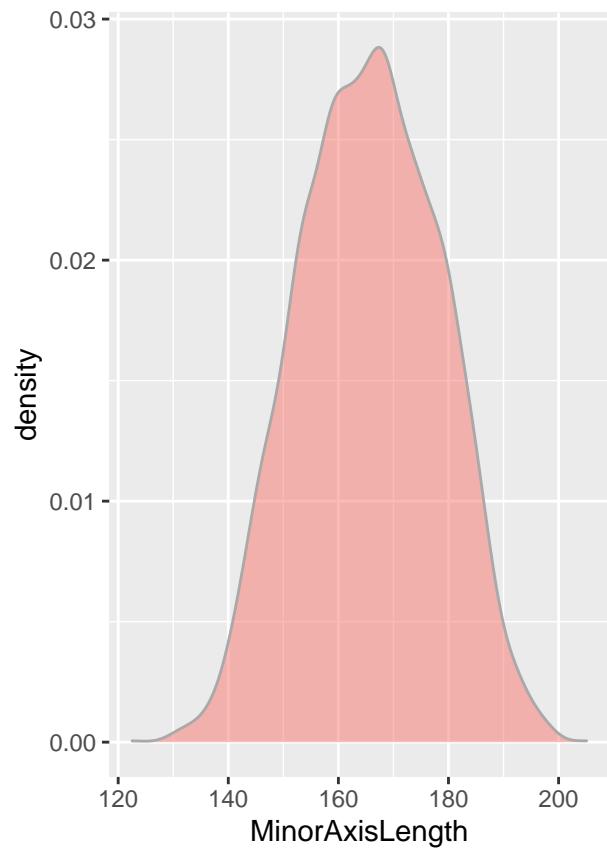


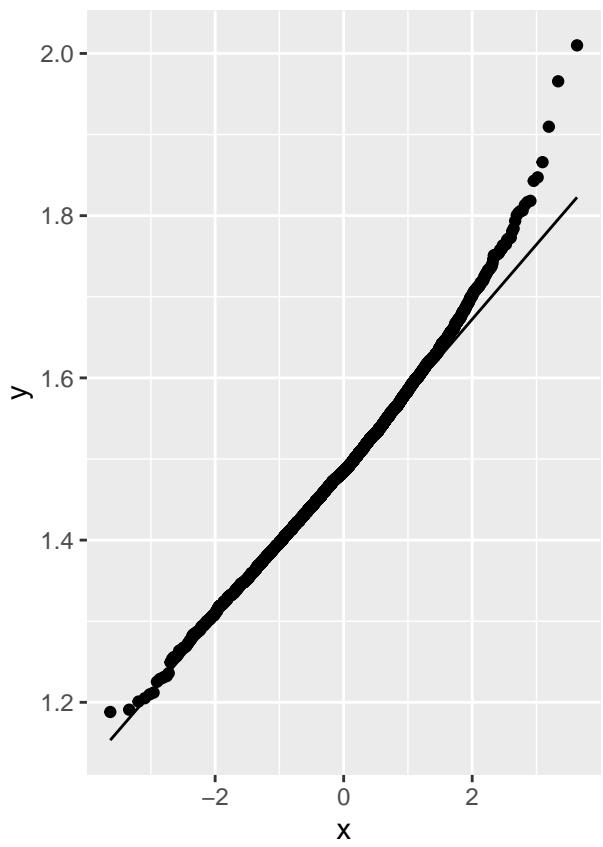
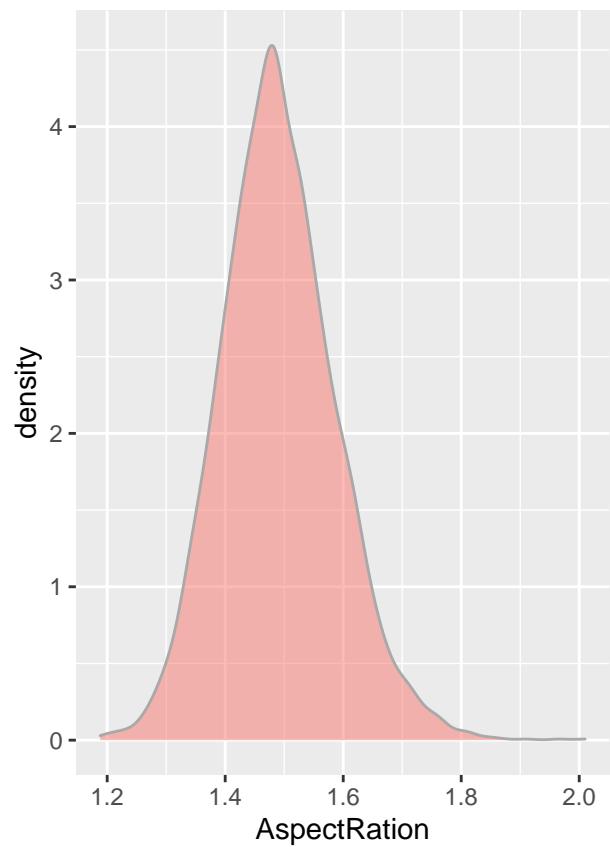
```
## [1] 
## [1] "DERMASON"
```

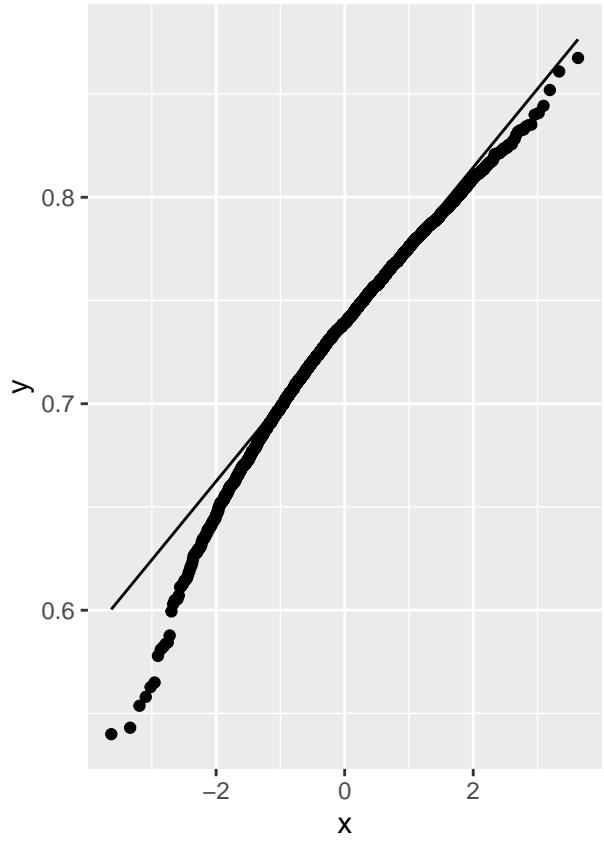
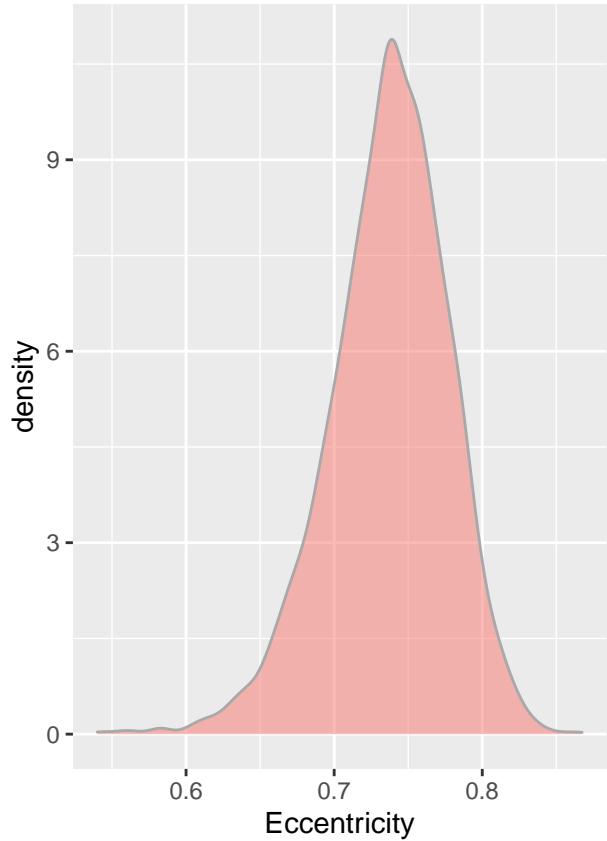


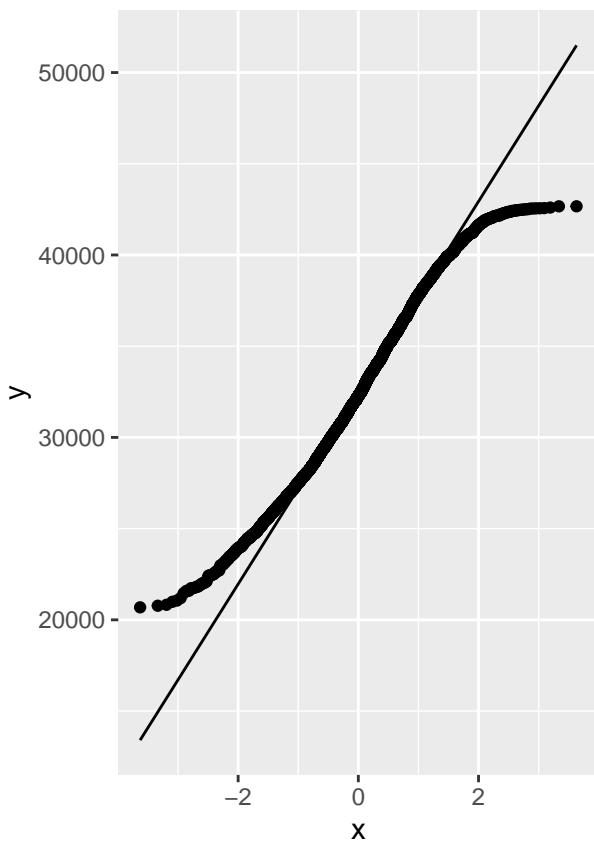
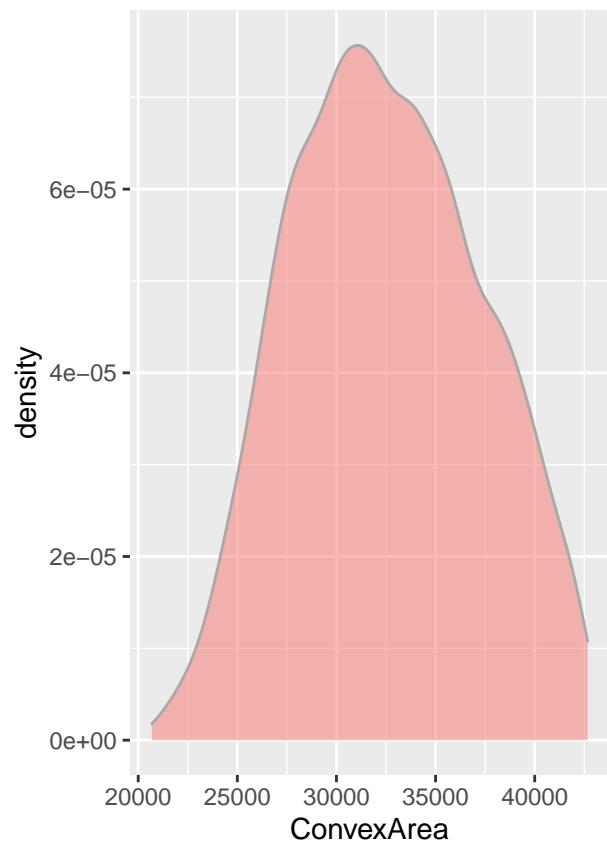


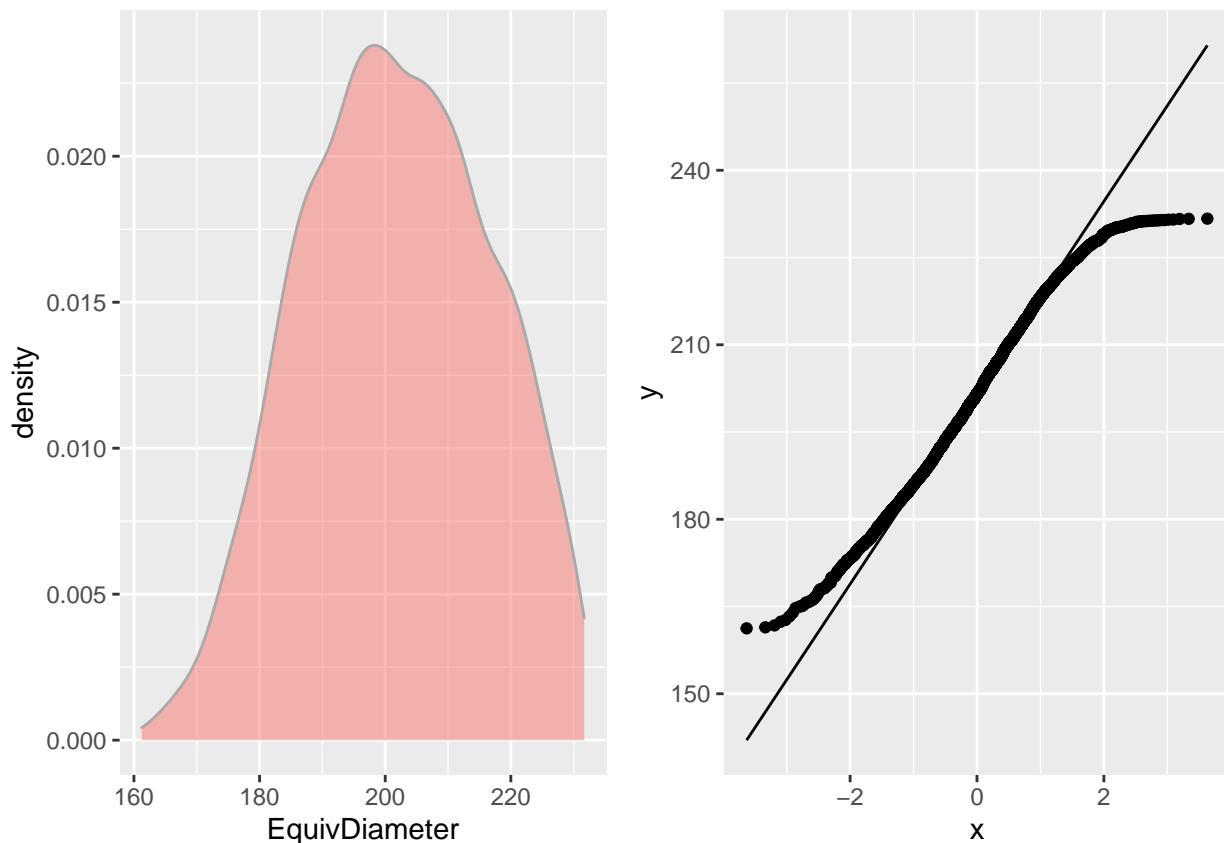




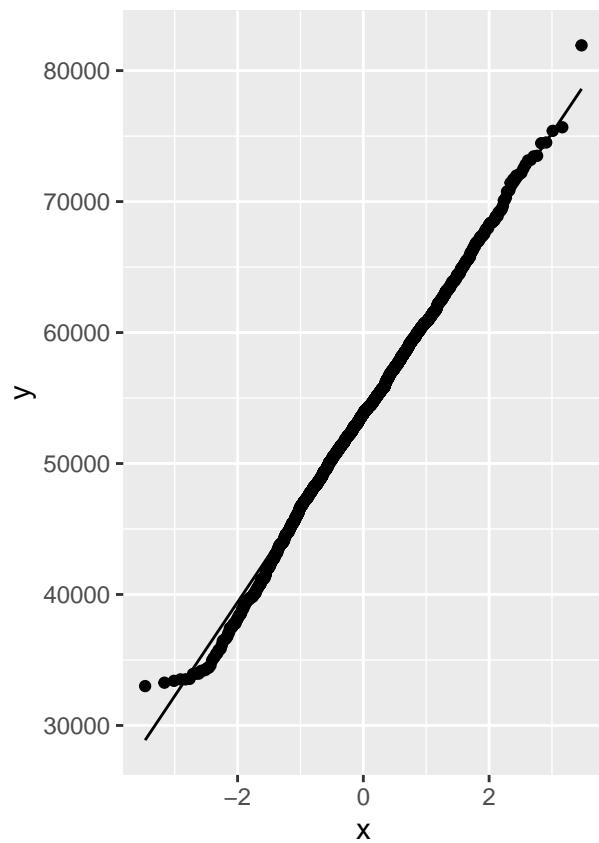
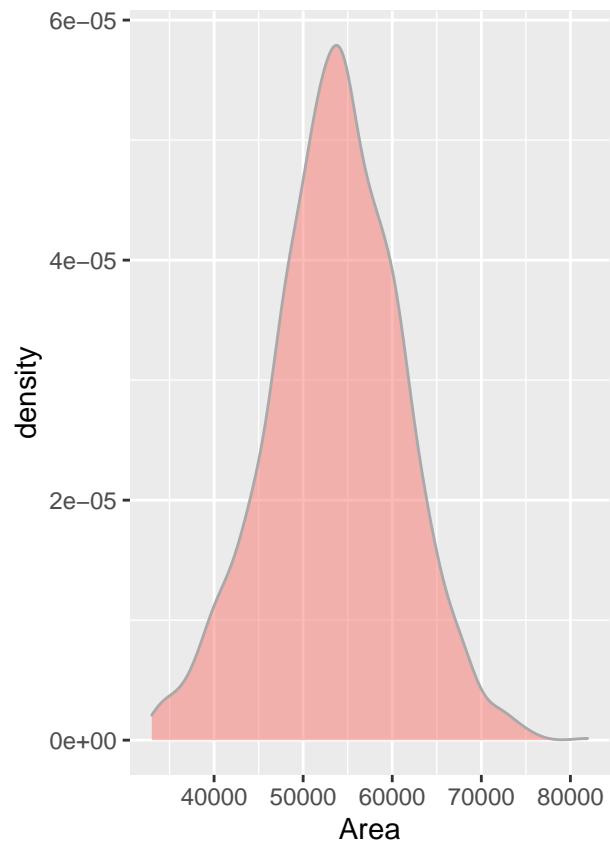


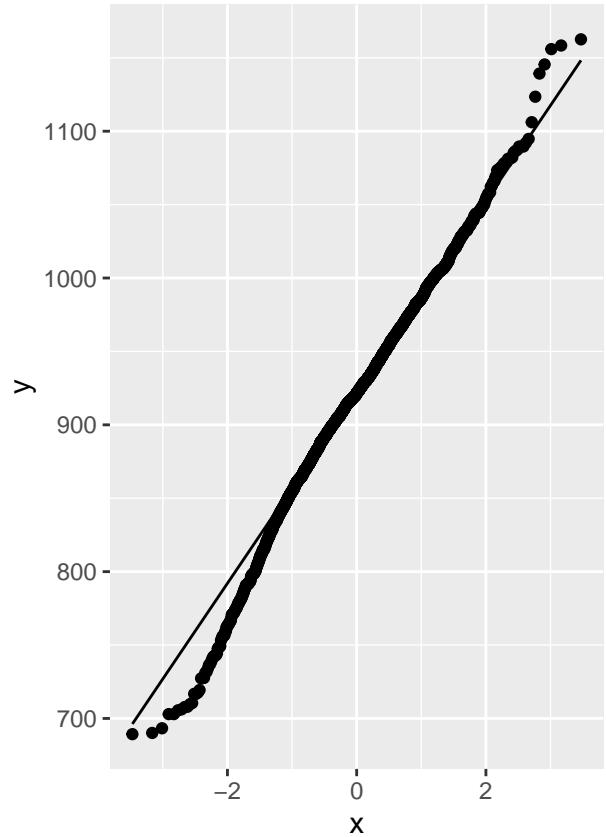
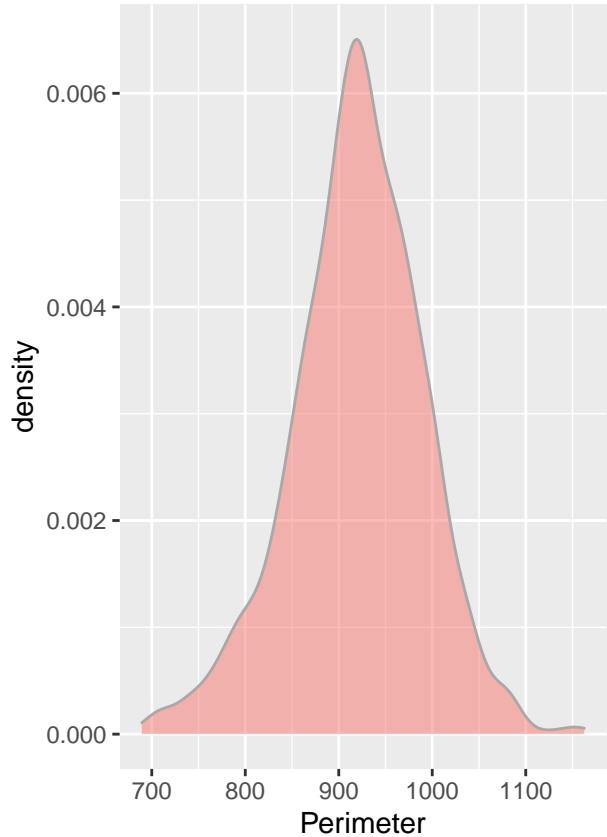


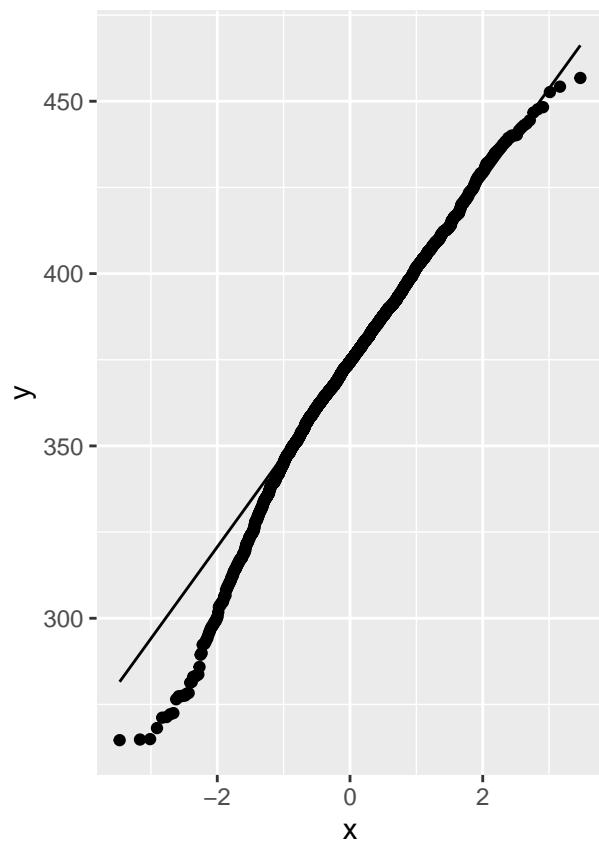
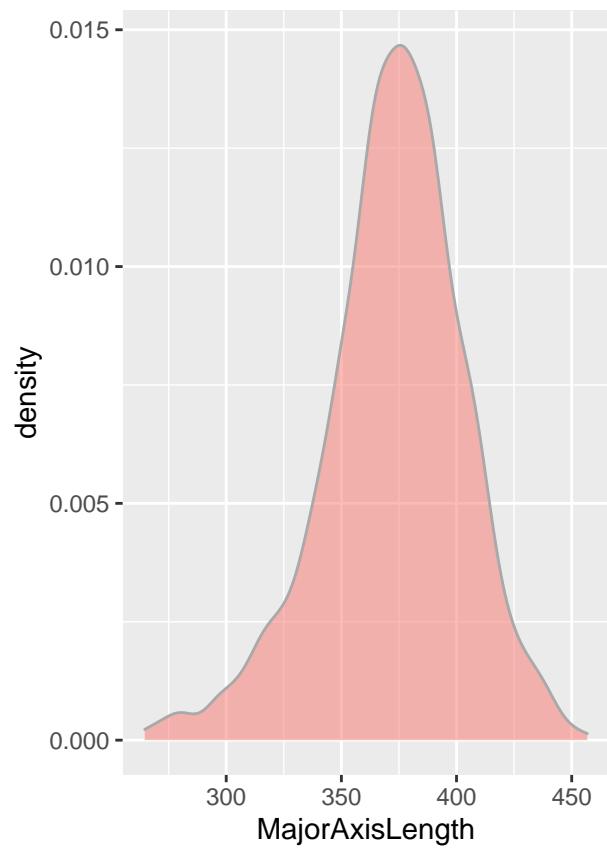


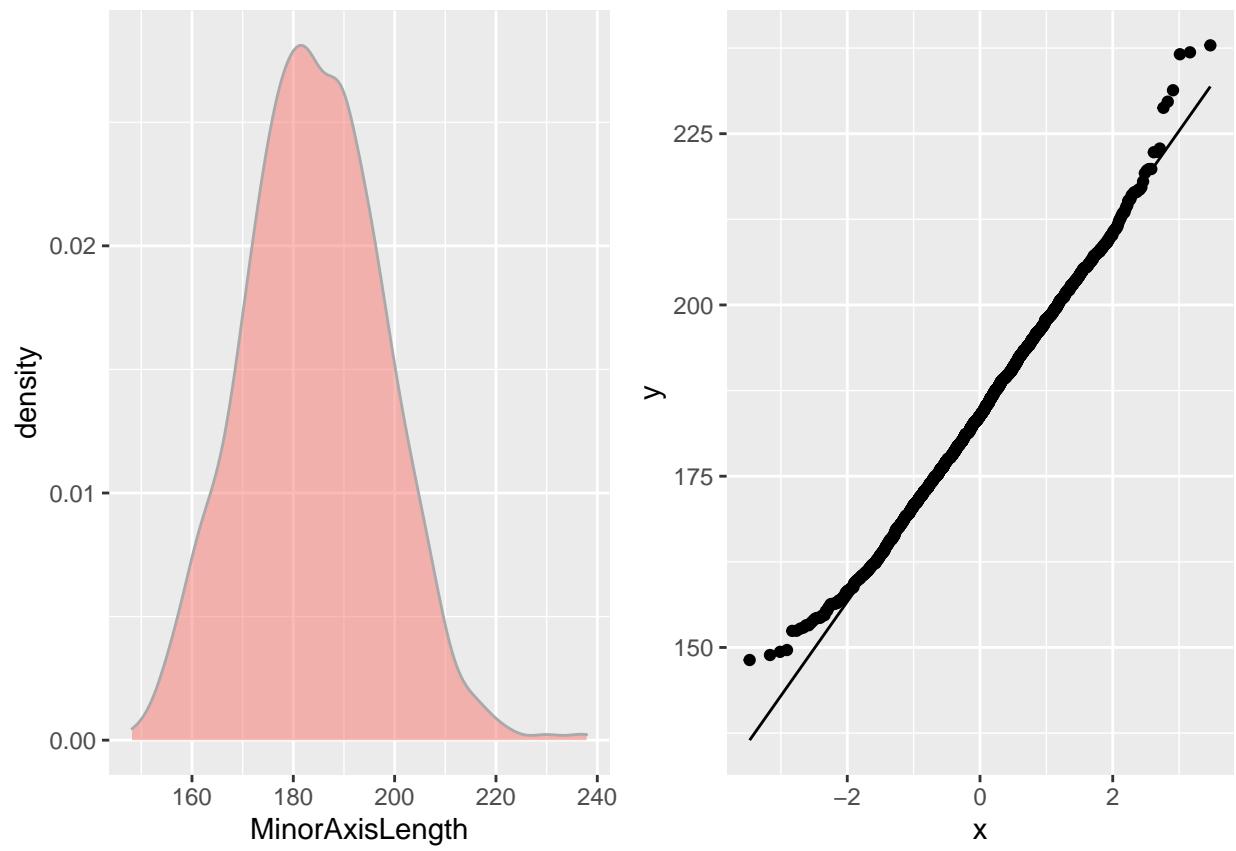


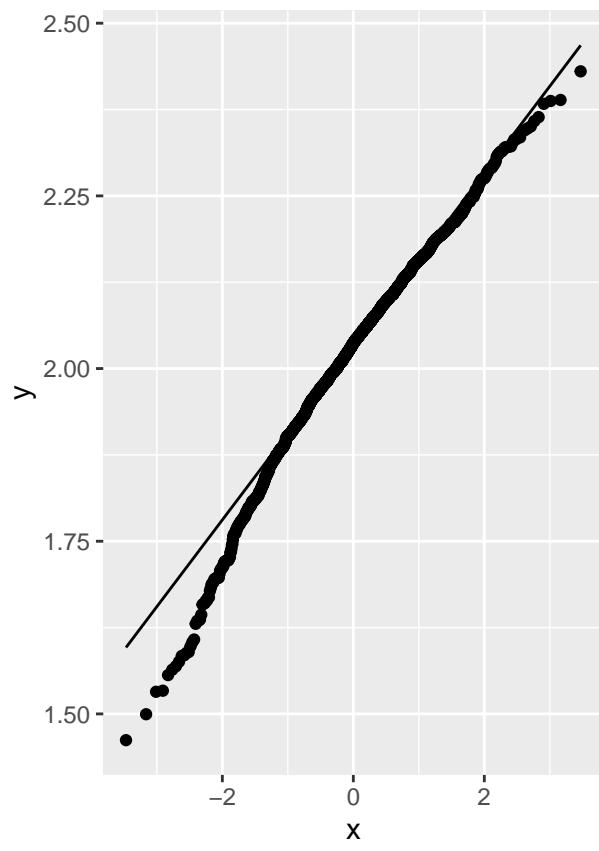
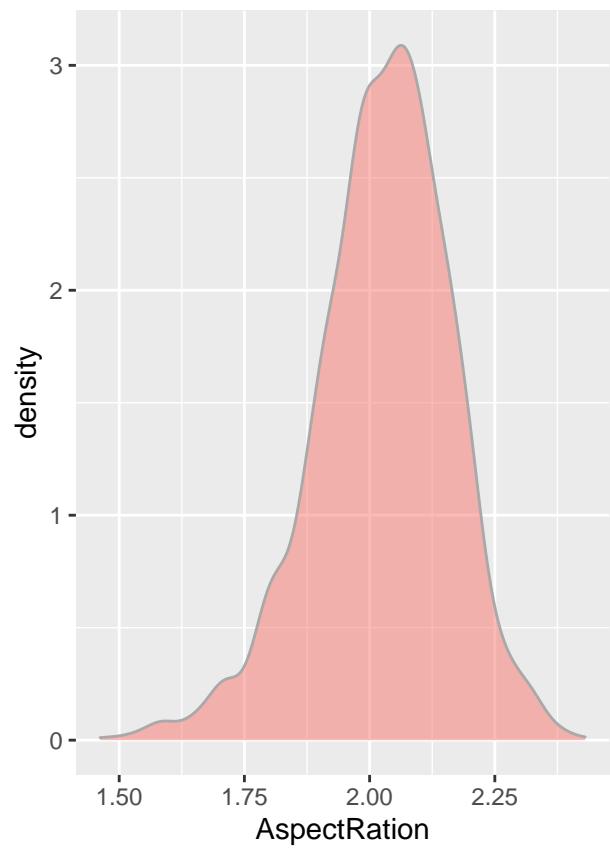
```
## [1]
## [1] "HOROZ"
```

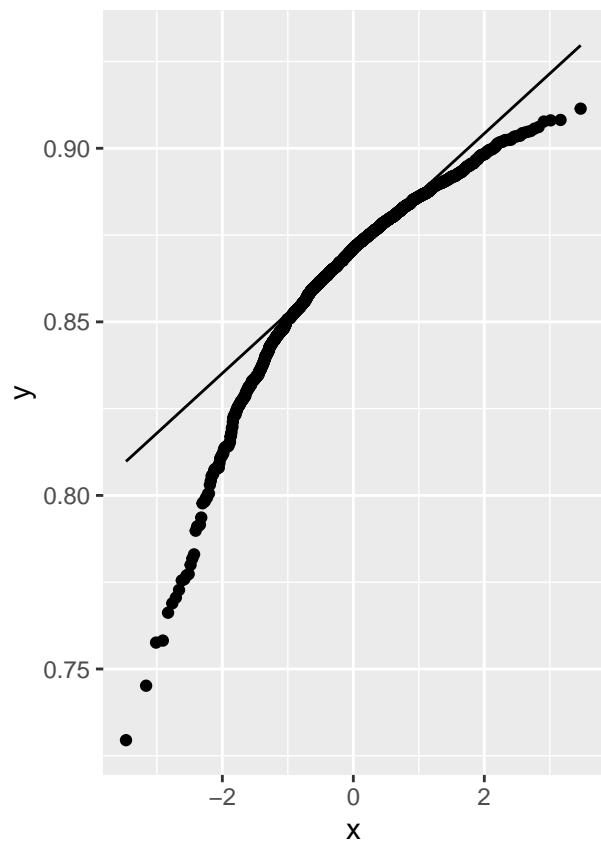
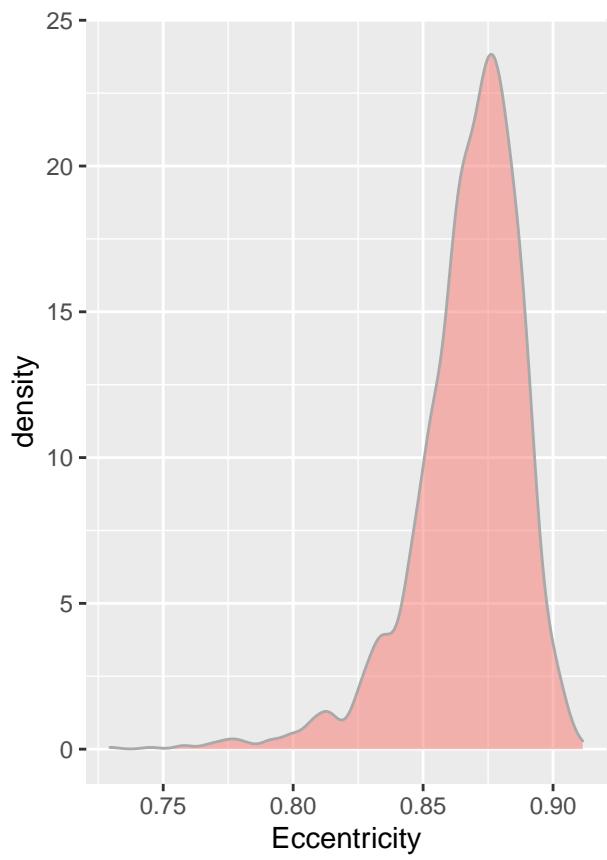


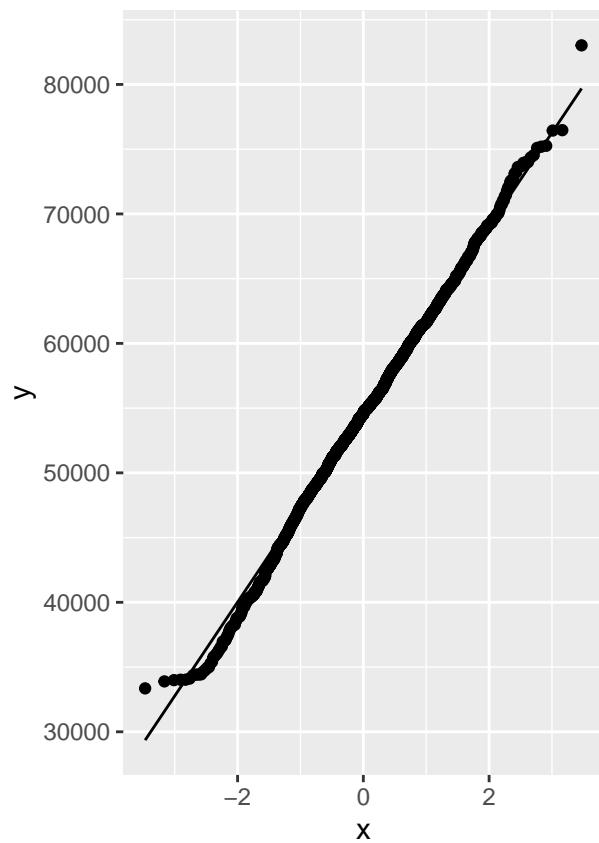
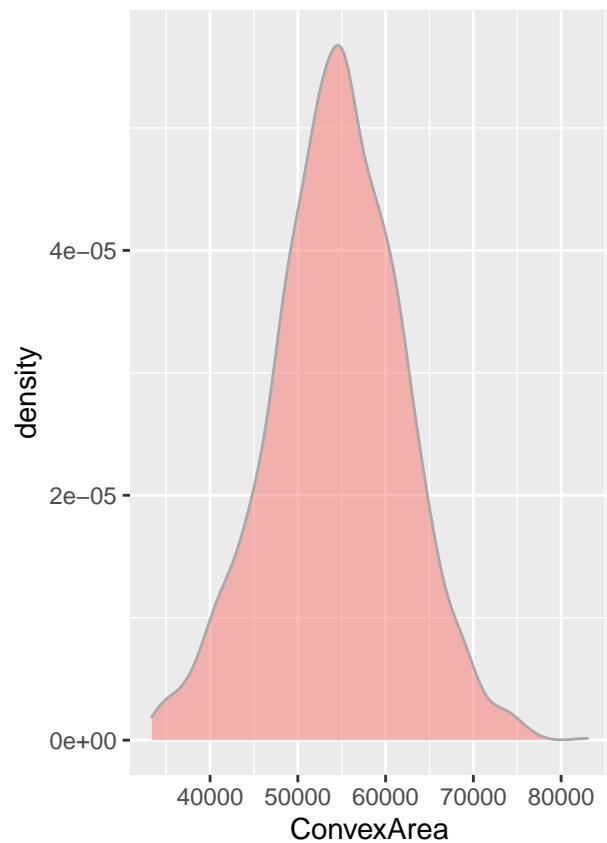


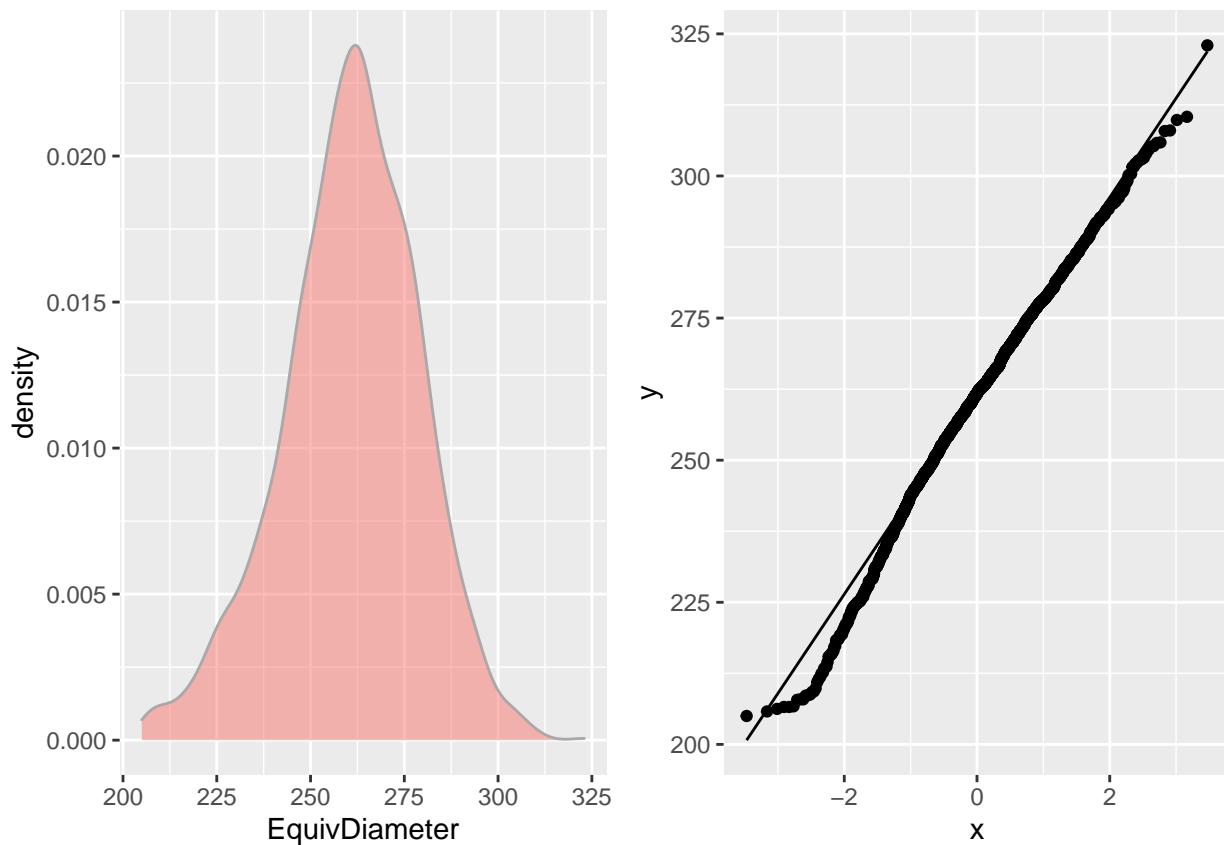




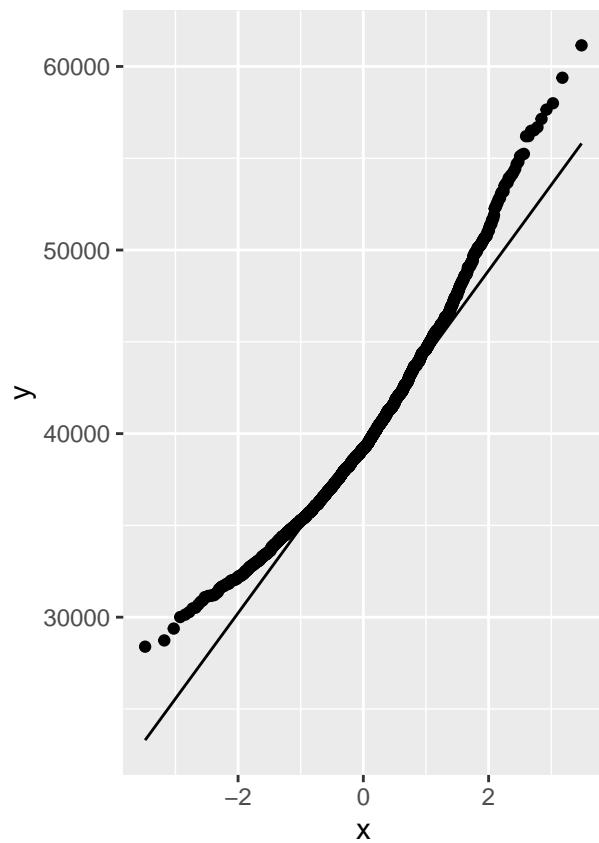
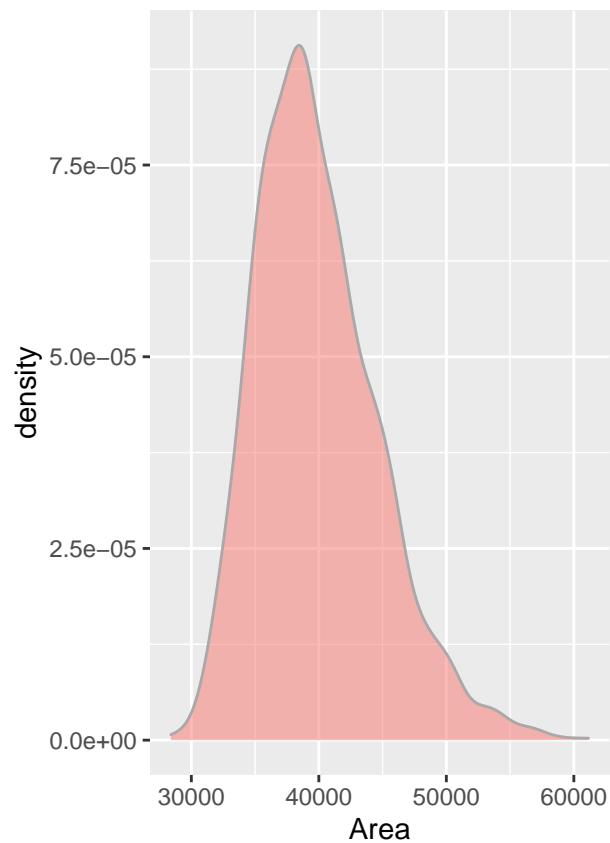


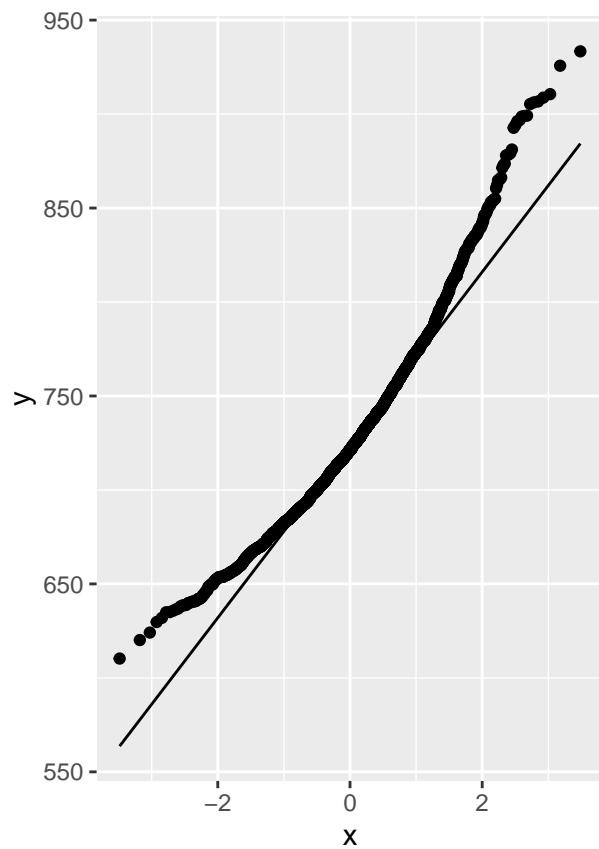
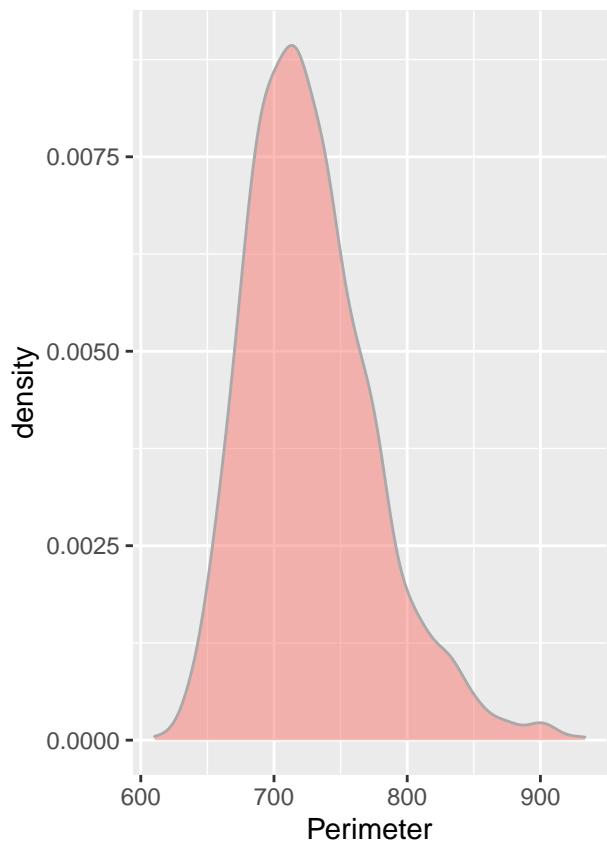


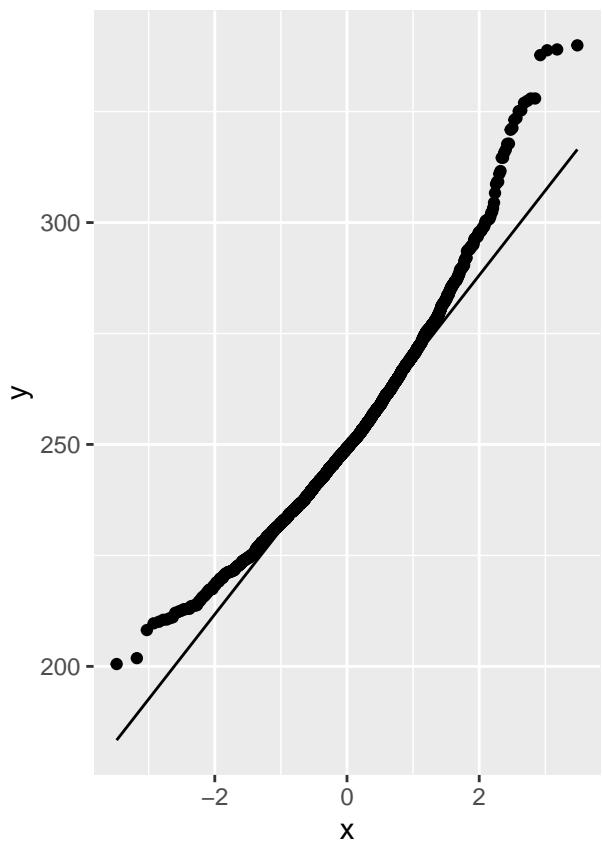
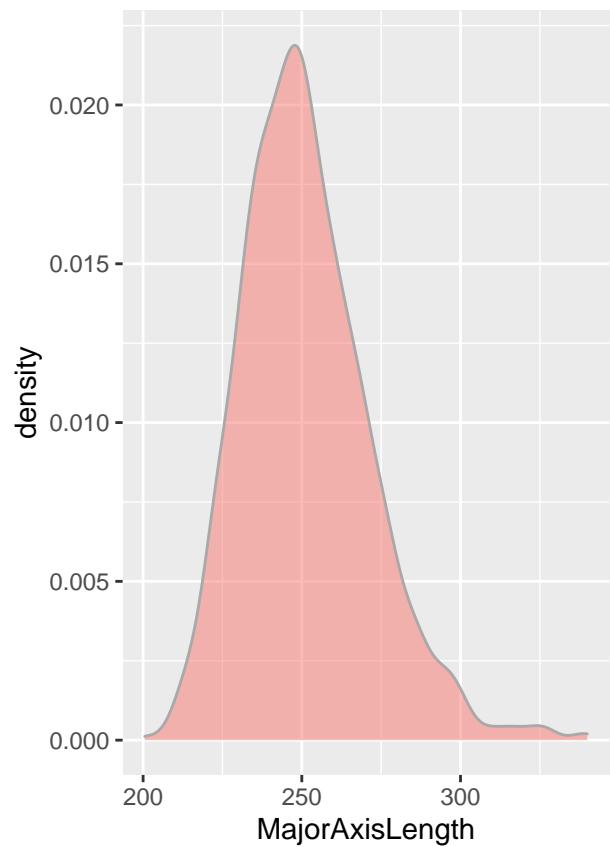


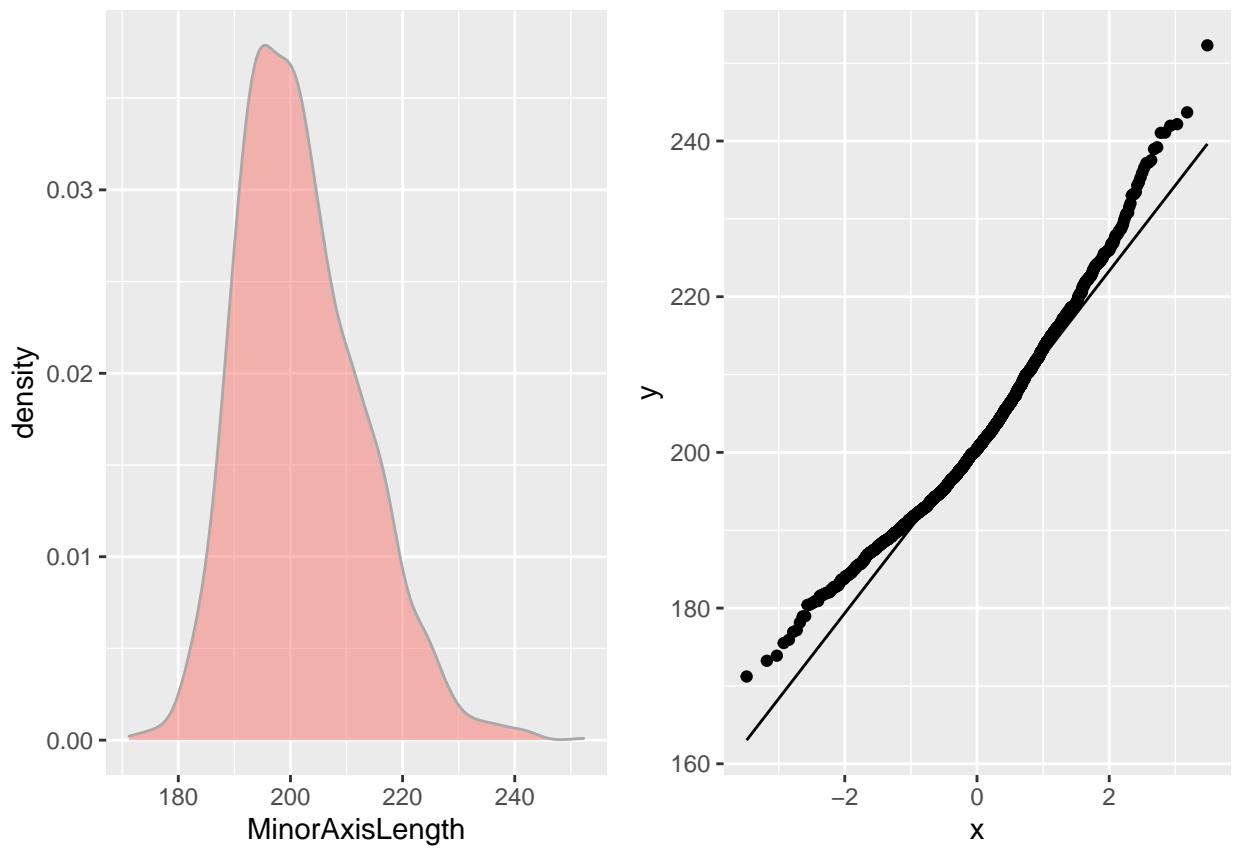


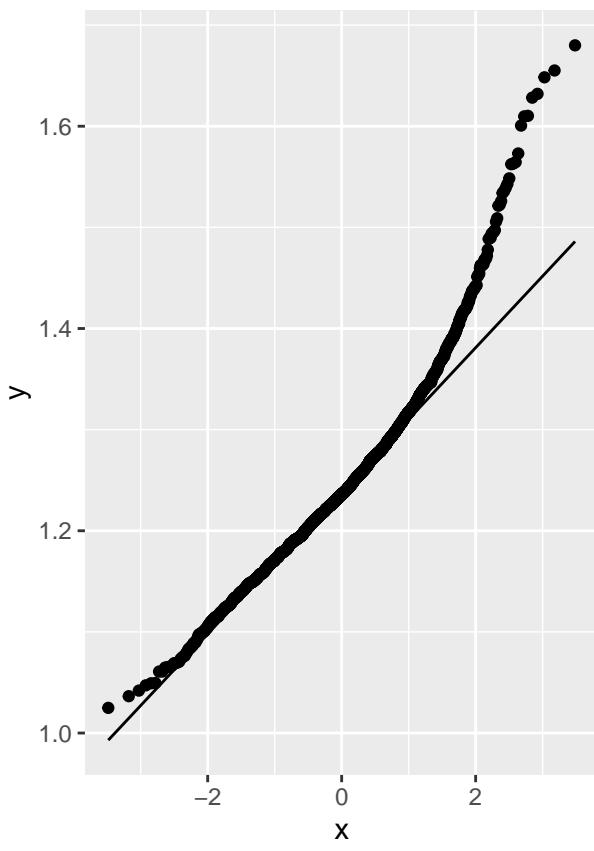
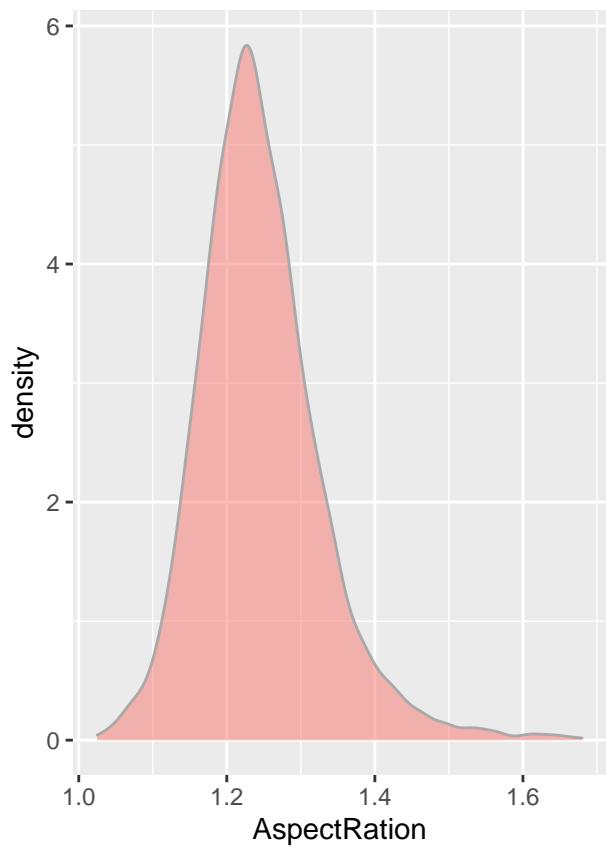
```
## [[1]]  
## [1] "SEKER"
```

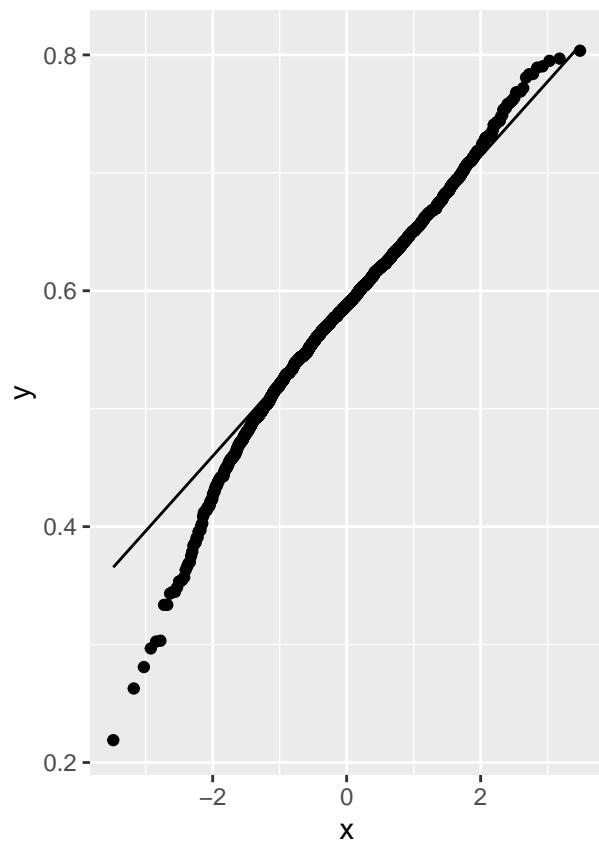
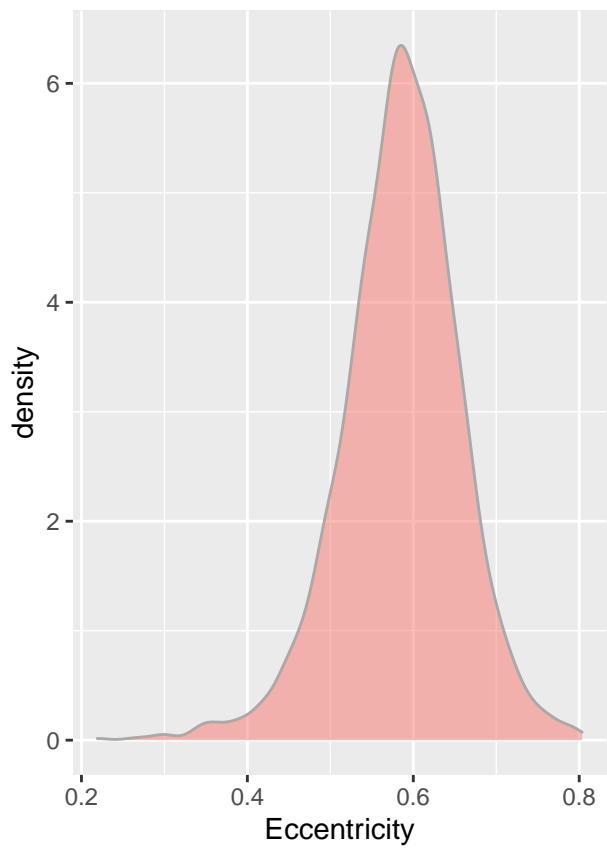


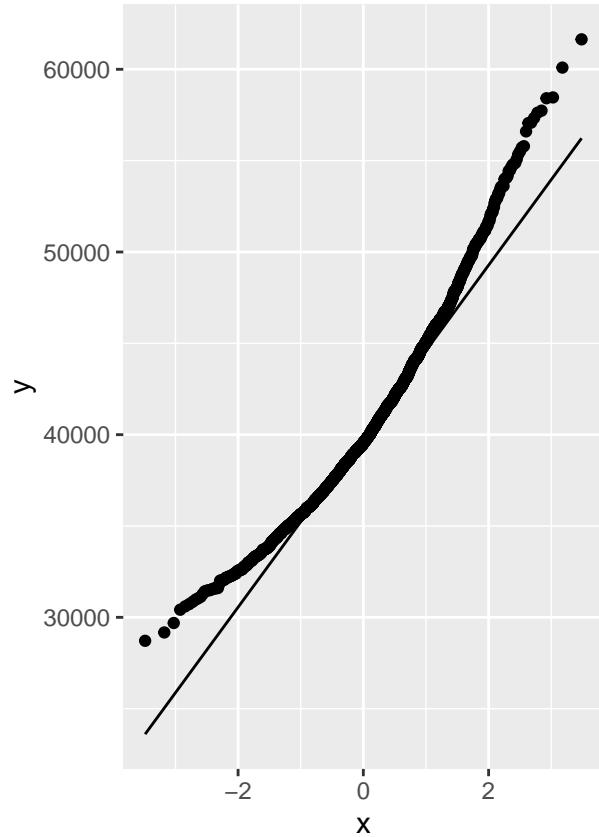
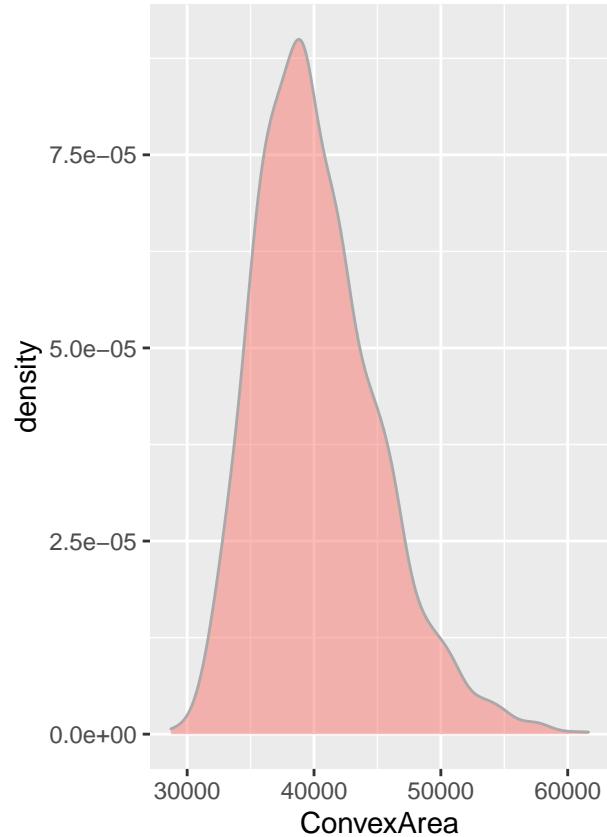


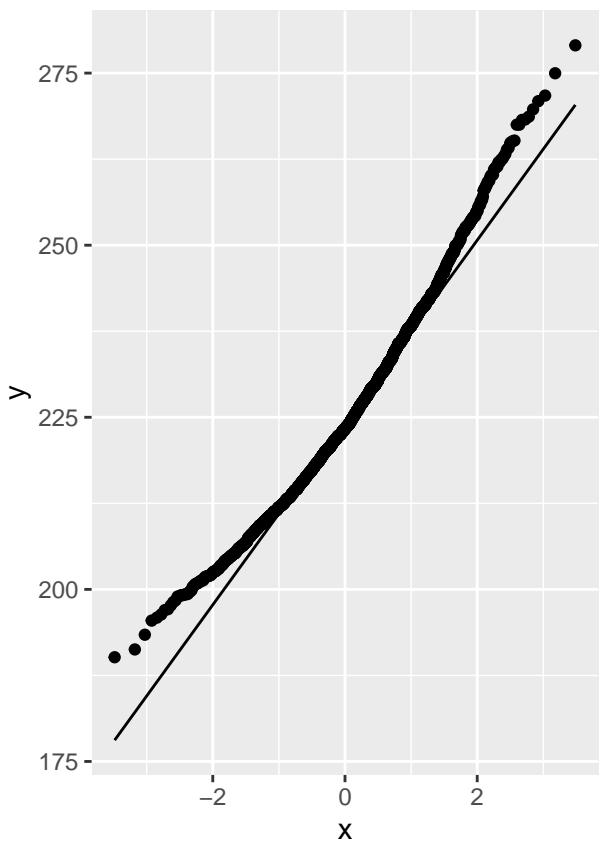
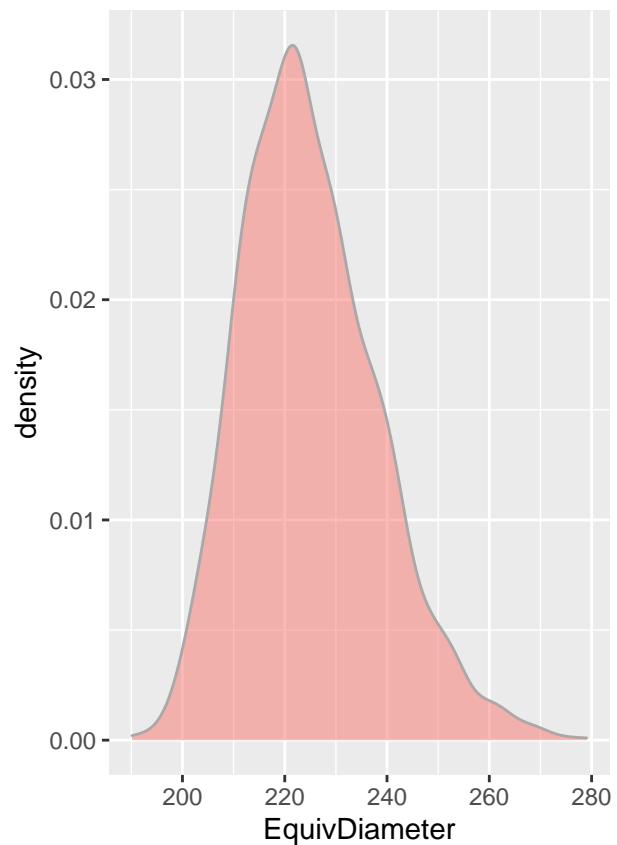




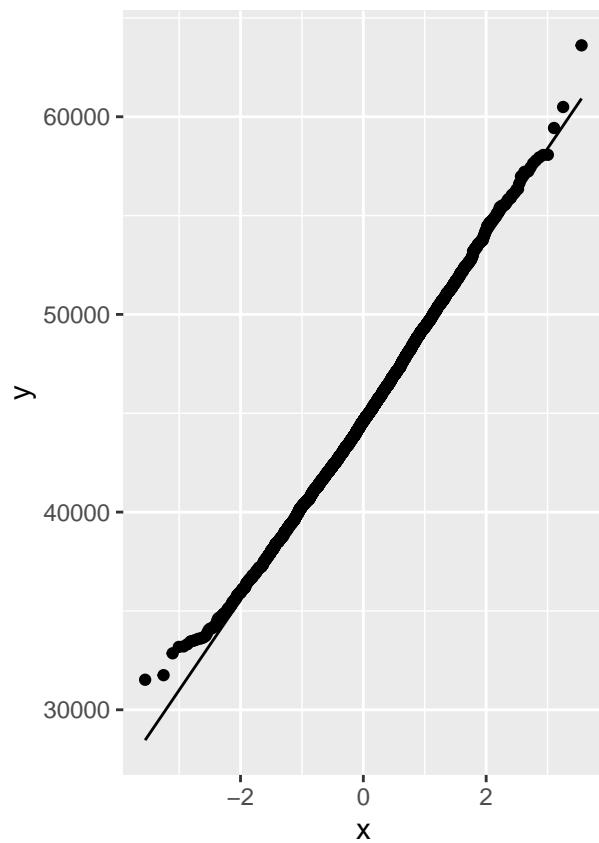
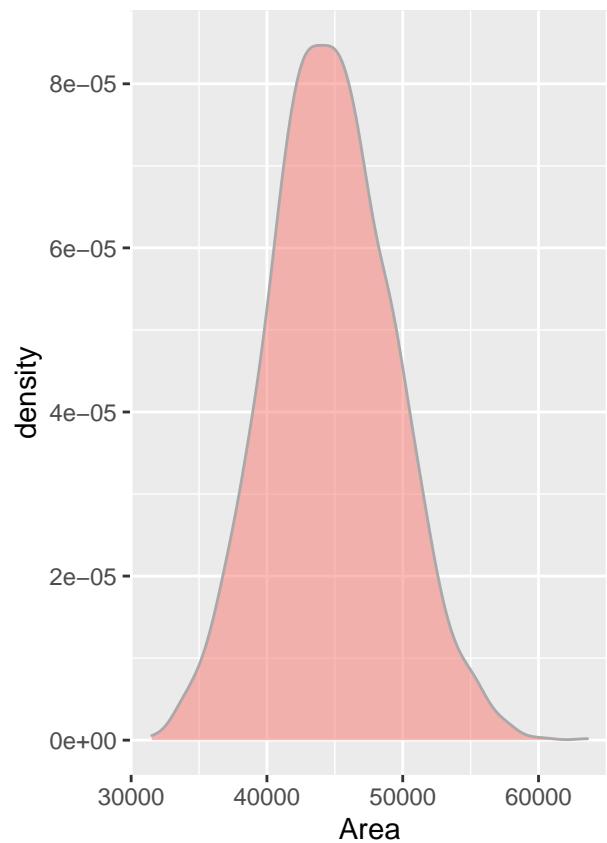


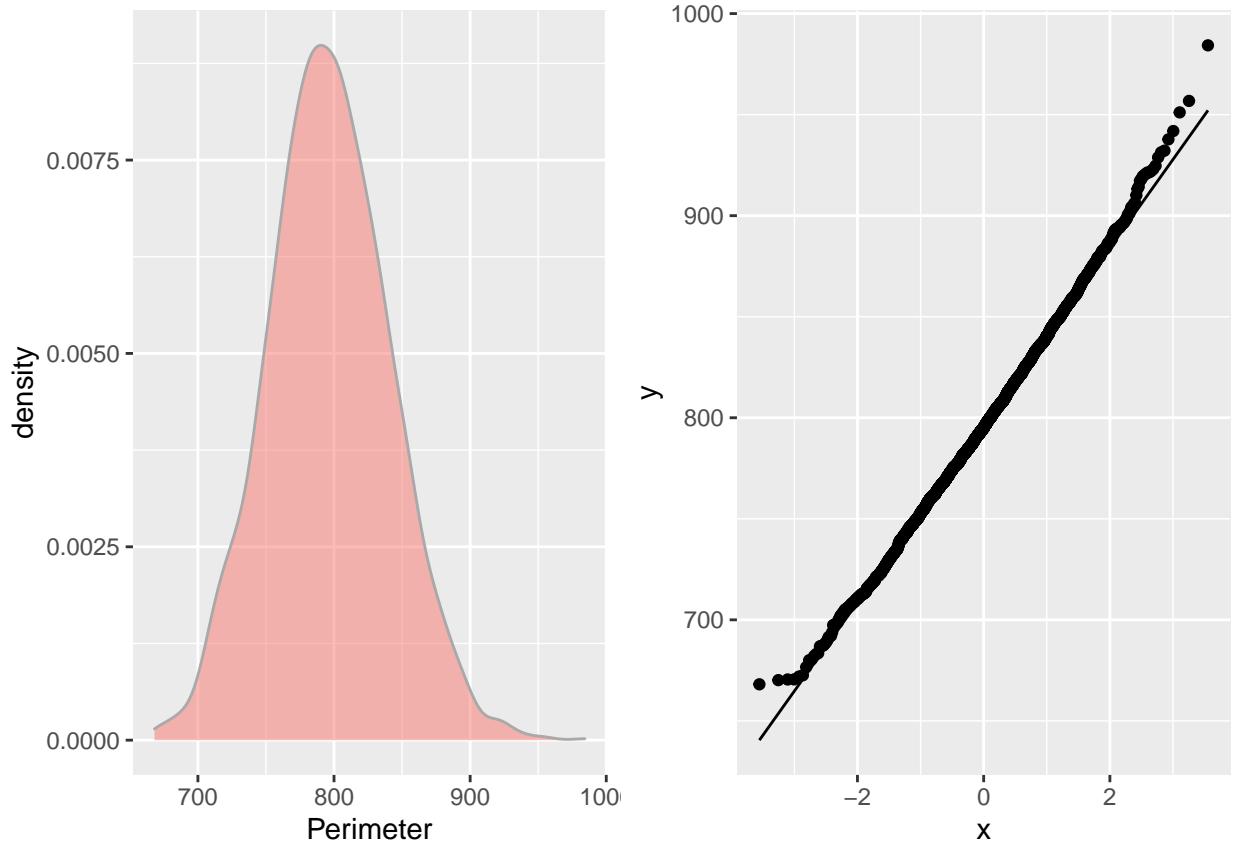


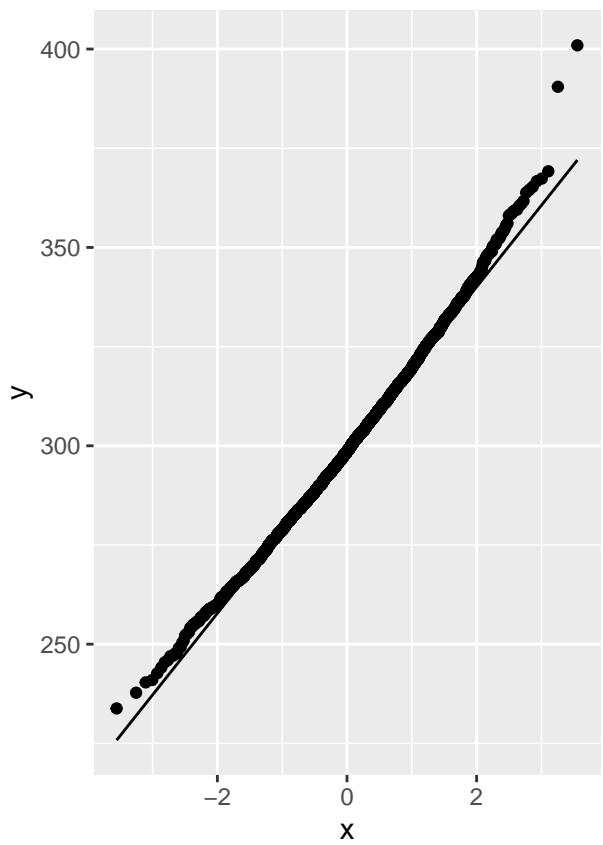
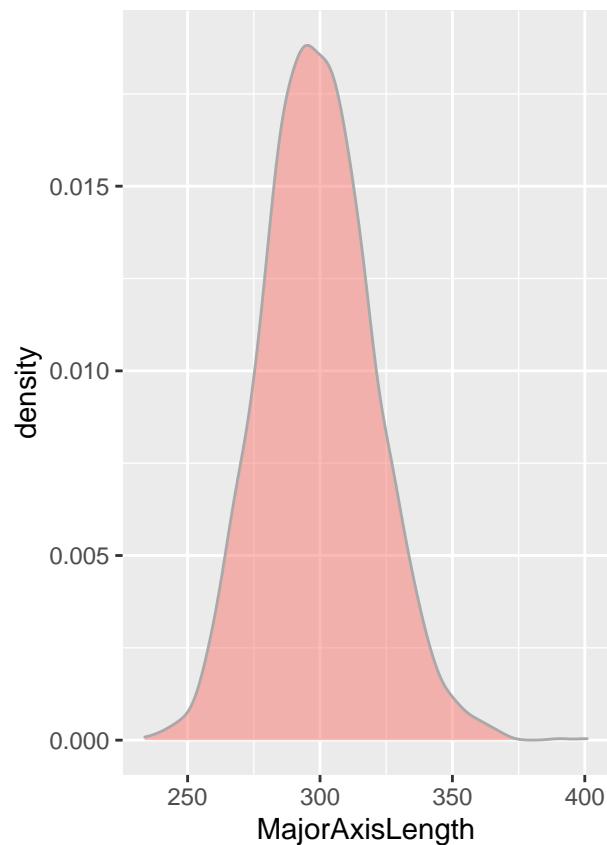


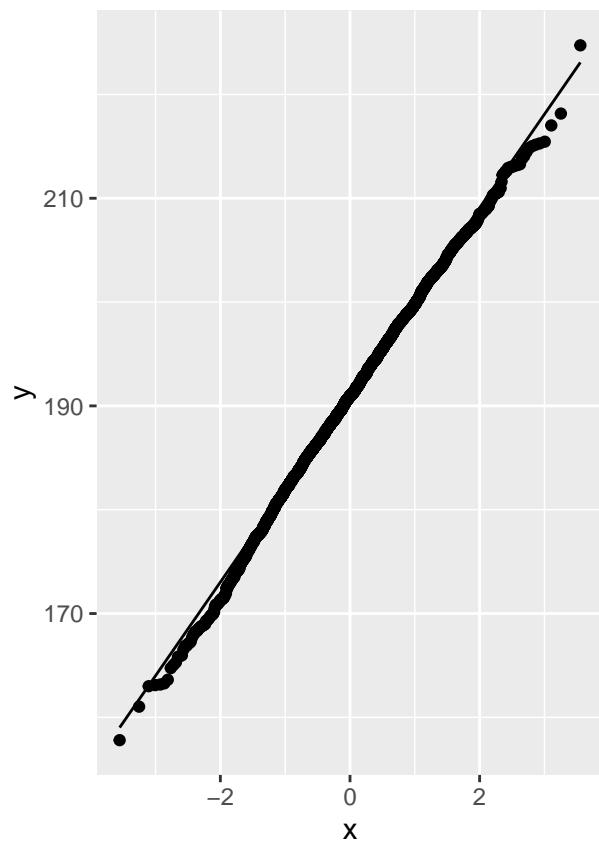
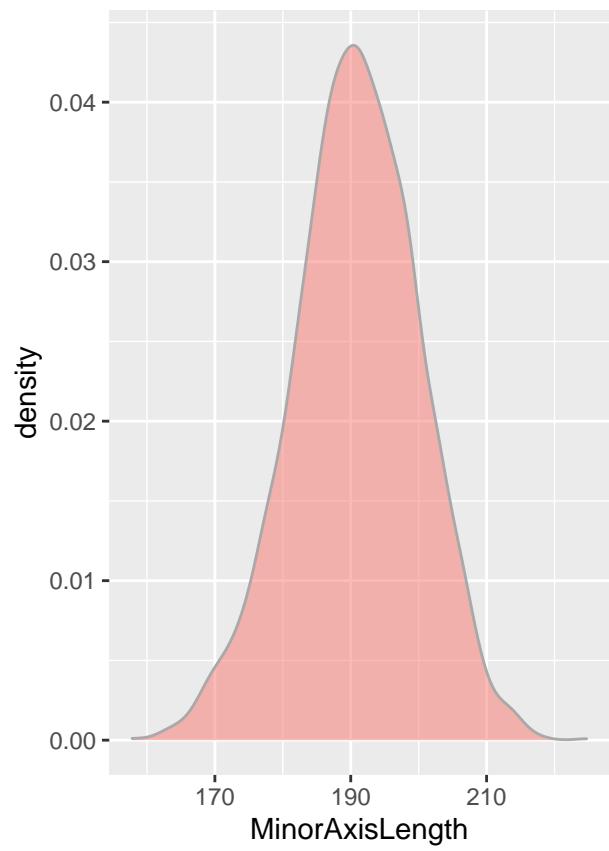


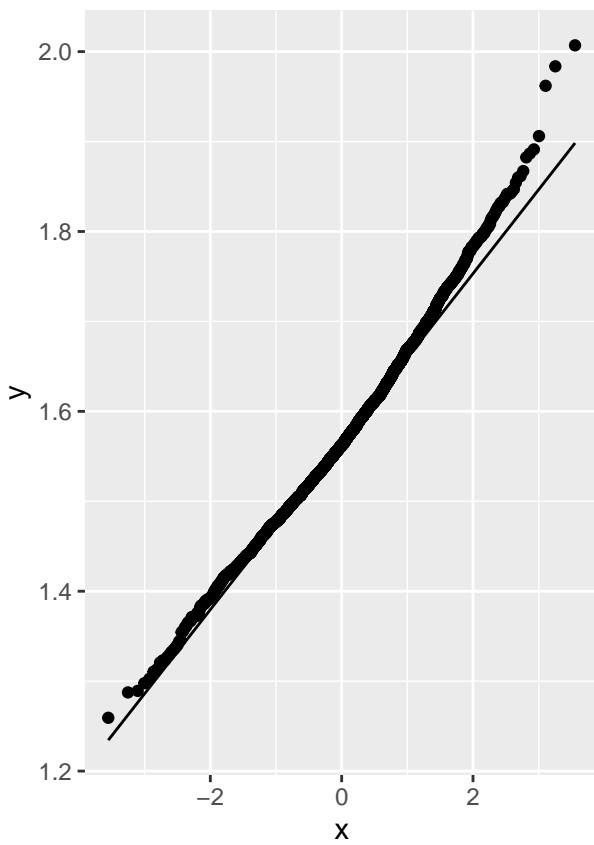
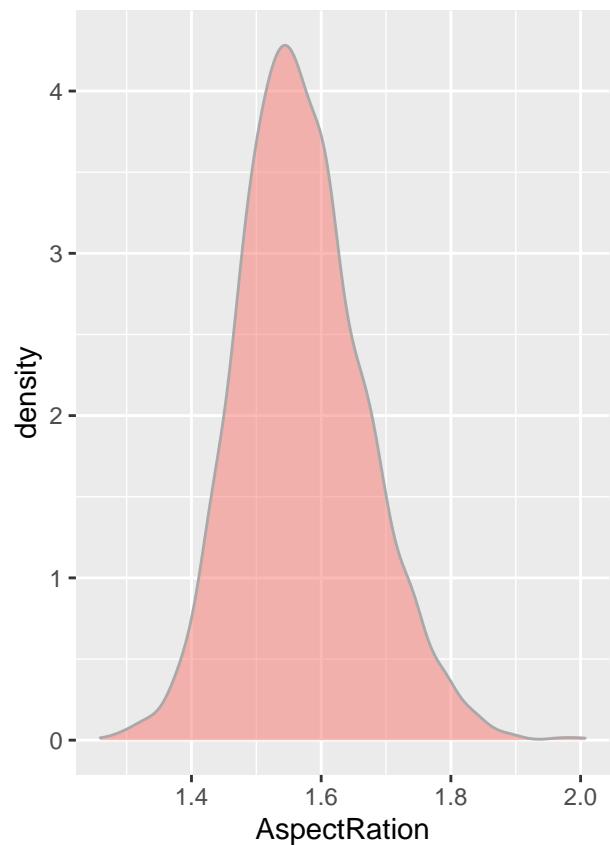
```
## [1] "SIRA"
```

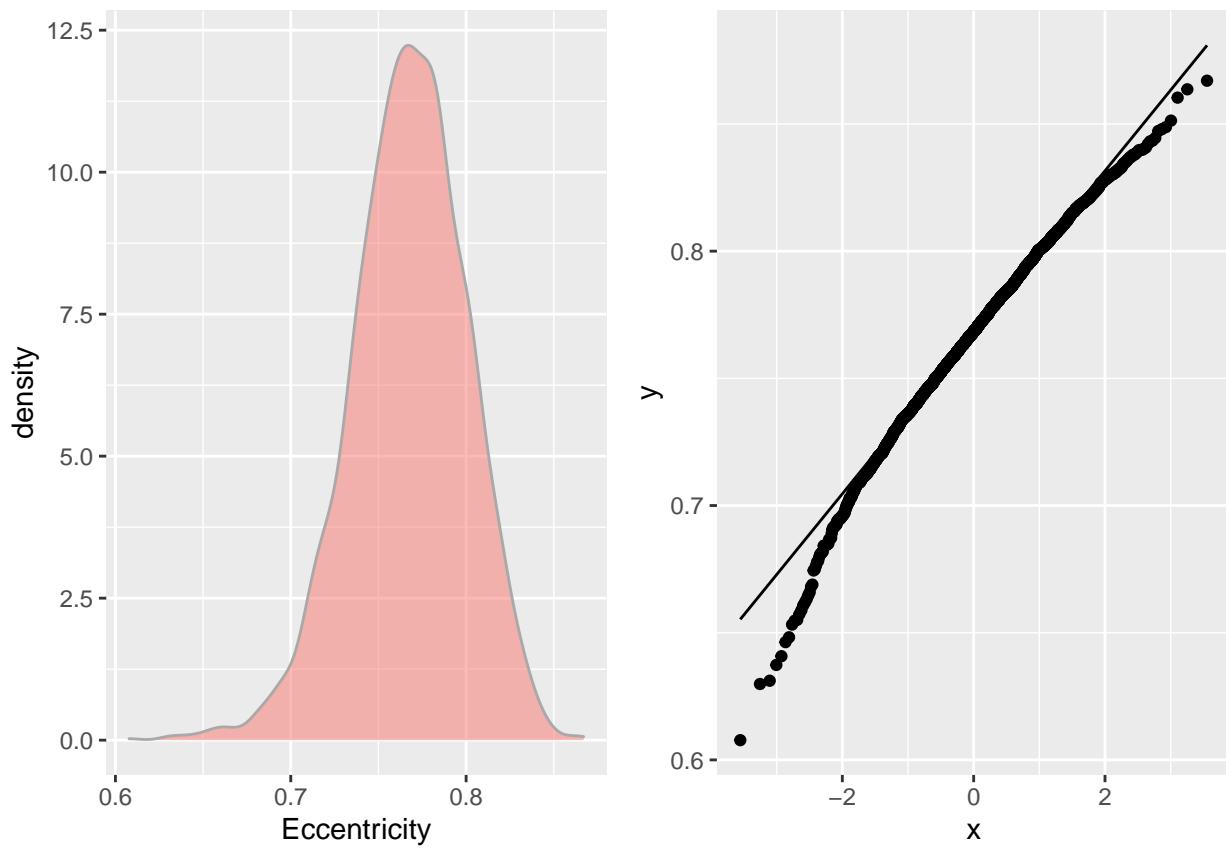


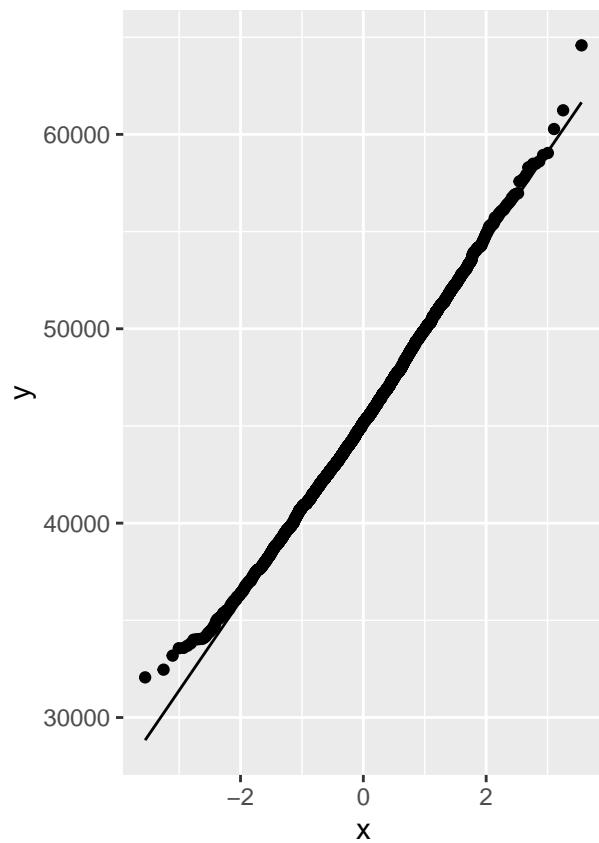
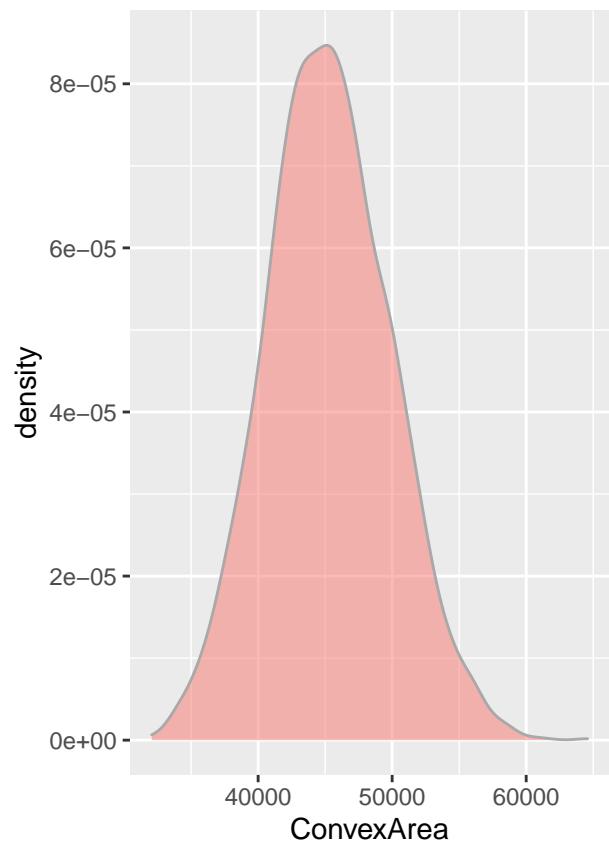


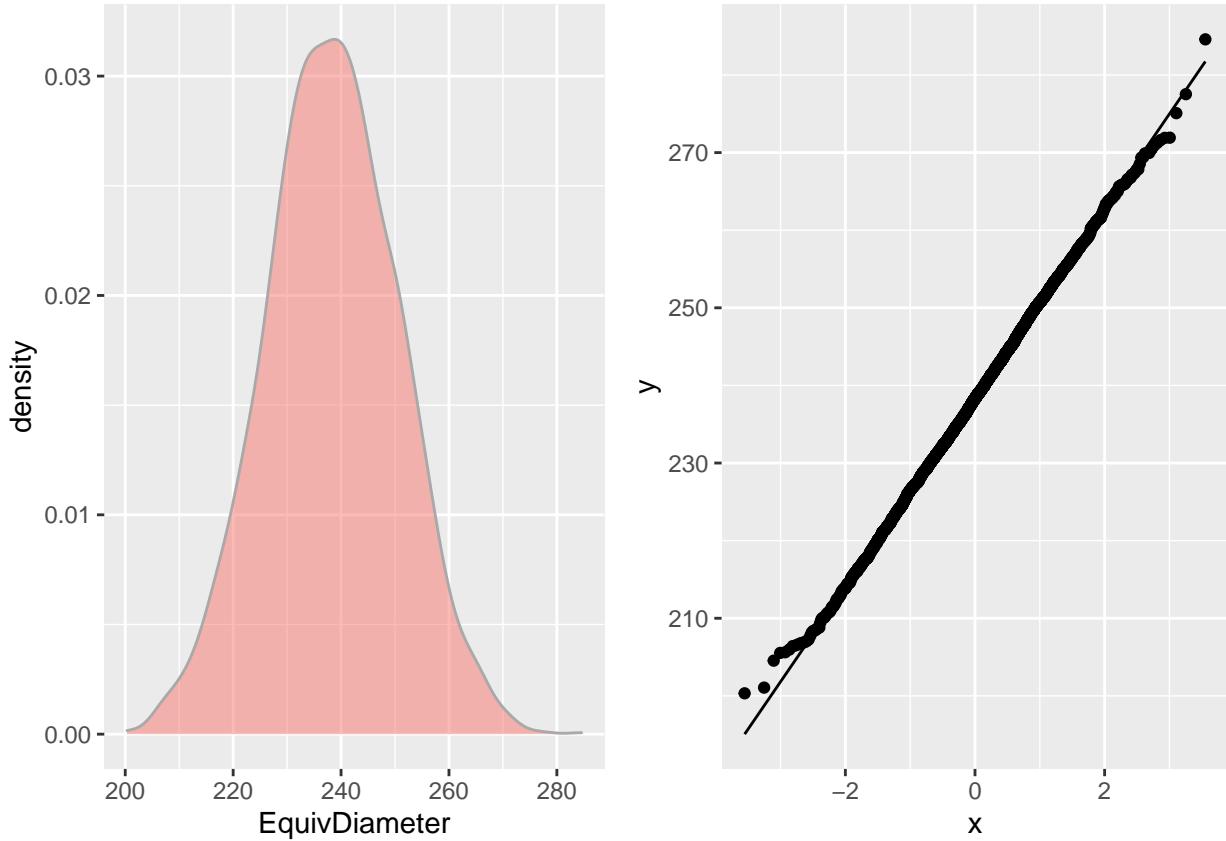








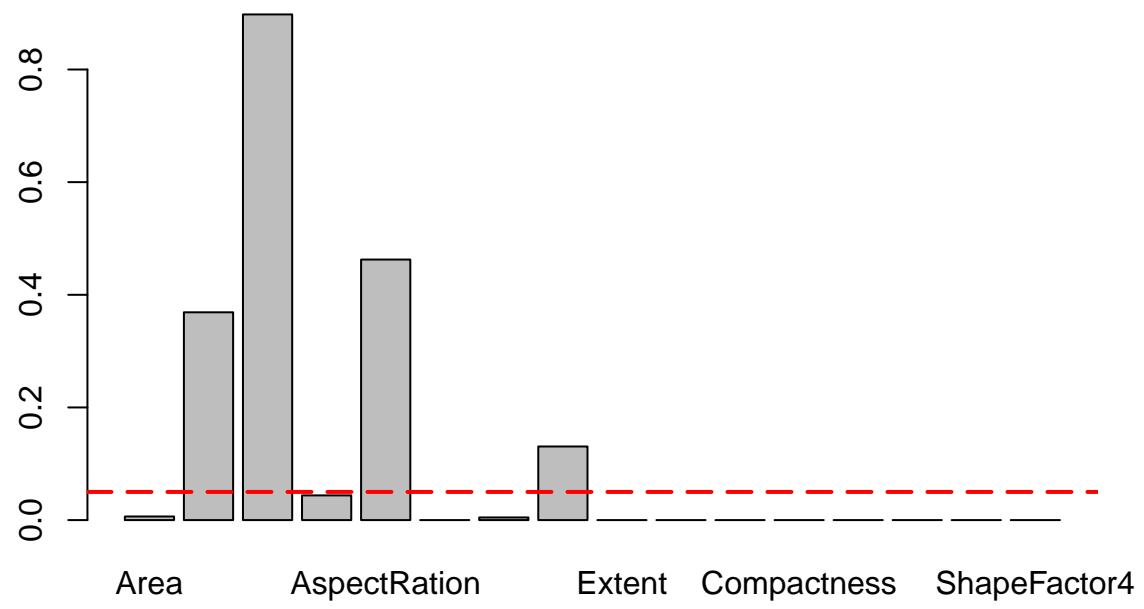


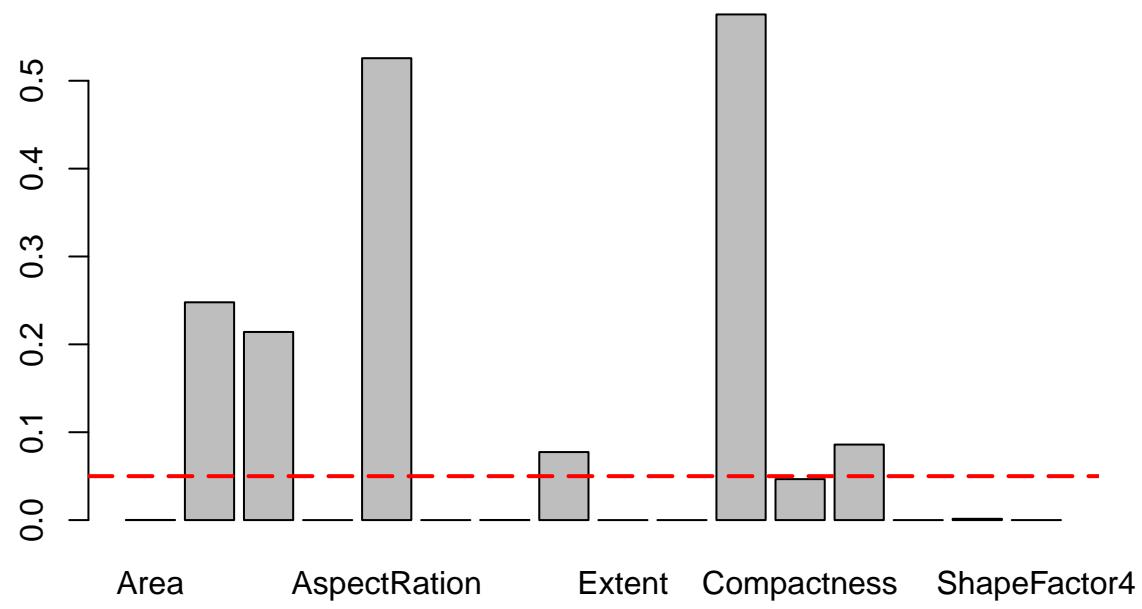


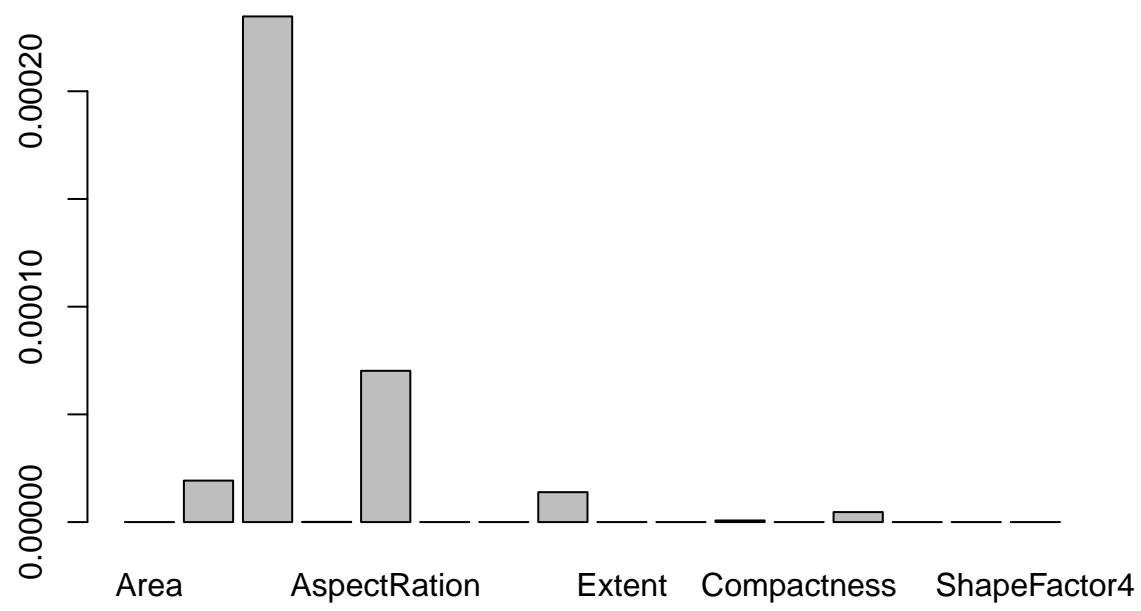
```

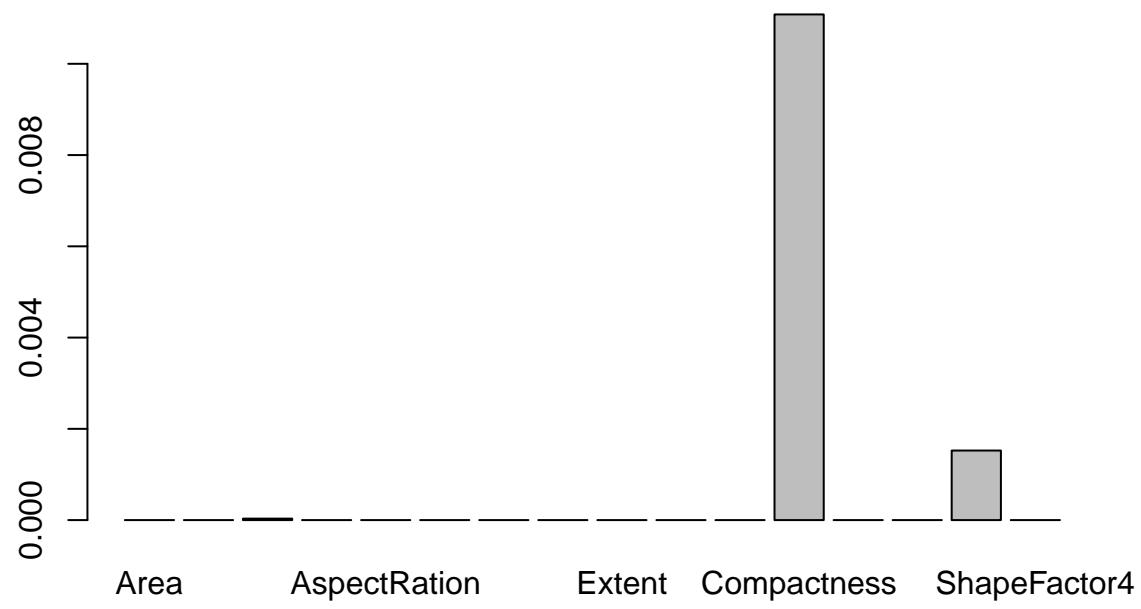
leves = list("BARBUNYA", "BOMBAY", "CALI", "DERMASON", "HOROZ", "SEKER", "SIRA")
pvalue = c()
for(j in 1:7){
  for(i in 1:16){
    pvalue[i] = shapiro.test(dados[dados$Class==leves[j],i])$p.value
  }
  barplot(pvalue,names.arg = names(dados)[1:16])
  abline(h=0.05, col = "Red", lty = 5, lwd = 2)
}

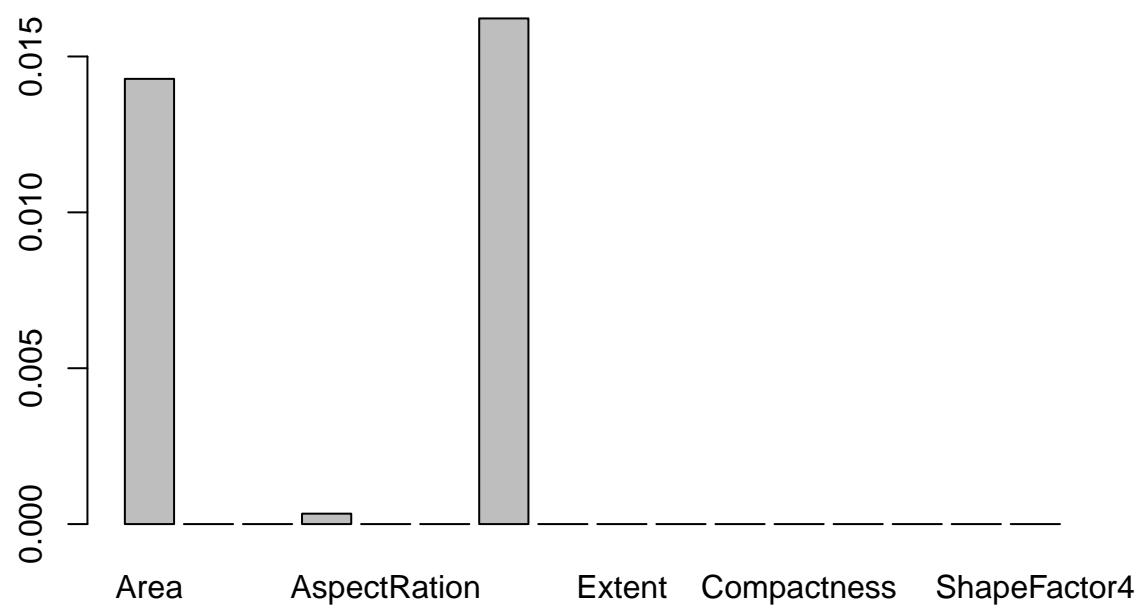
```

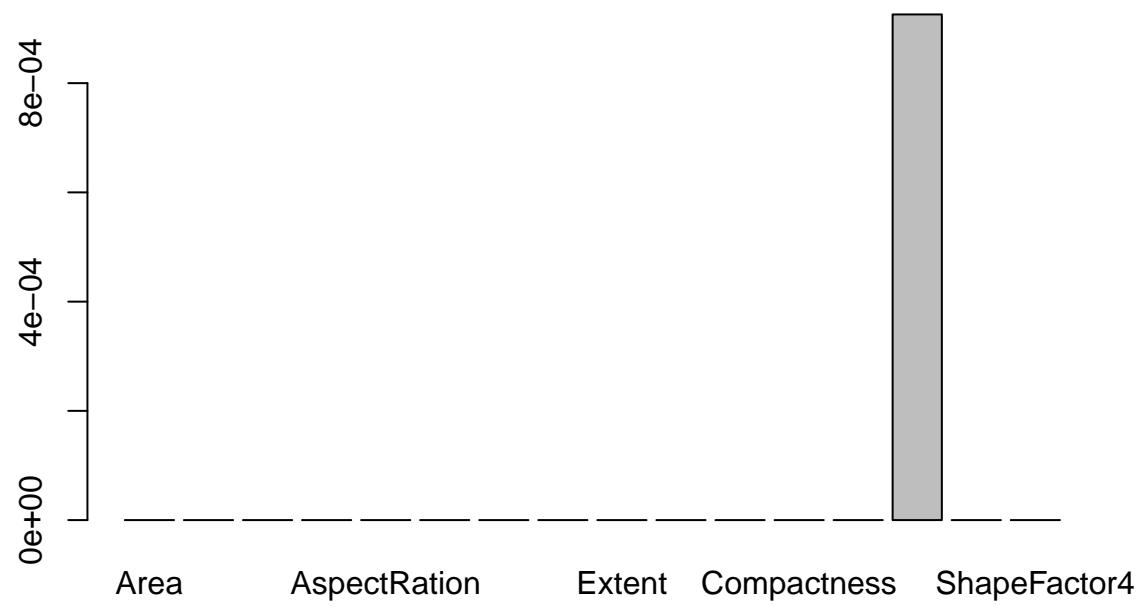


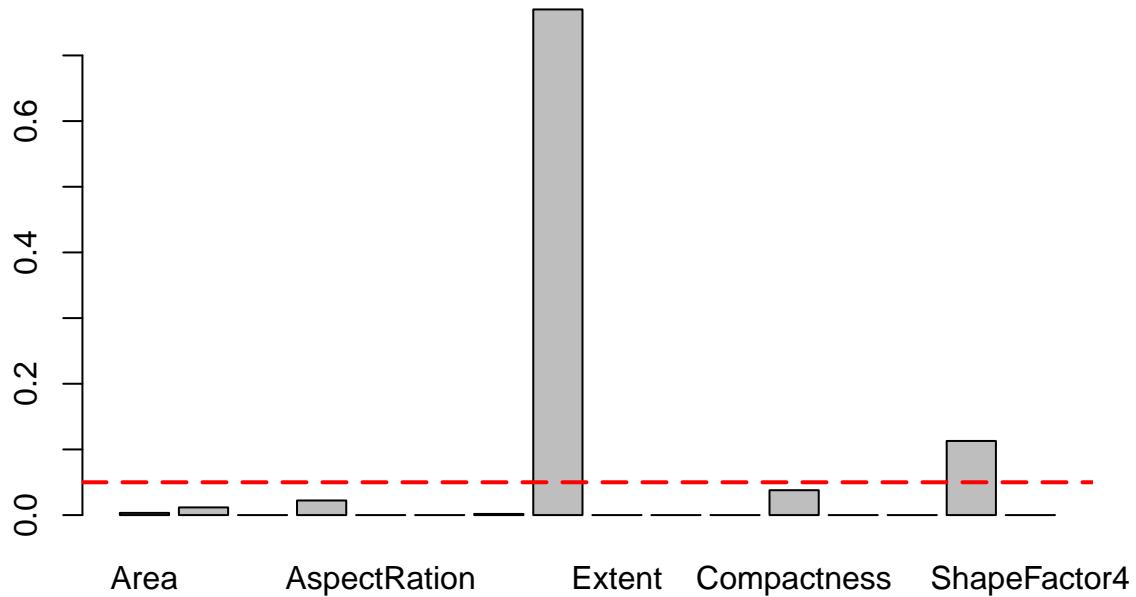








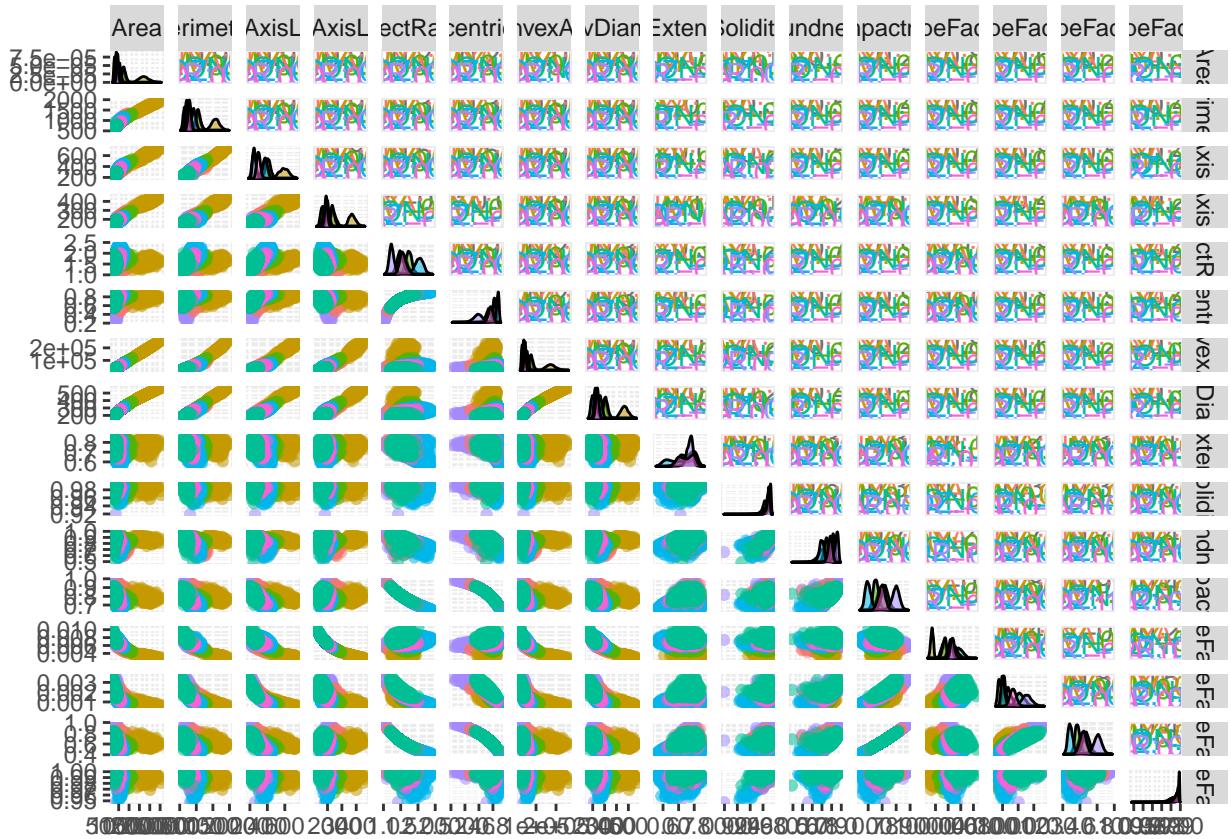




Em nenhuma das classes possui uma distribuição normal, em todas as classes há uma variável que possui uma discrepância dos dados.

**8.Utilizando um gráfico de dispersão entre pares de variáveis, diga se existe alguma associação entre variáveis que permite uma maior discriminação?**

```
ggpairs(dados[, 1:16], aes(colour = dados$Class, alpha = 0.4))
```



Há discriminação de algumas classes em determinadas variáveis com maiores diferenças que em outras mesmo desconsiderando os outliers.

## 9. Após ter analisado estas informações, quais considerações você faz sobre este conjunto de dados (ou tarefa)?

Com relação à classificação das sementes de feijão, as características de dimensão e forma das variedades de feijão não possuem características discriminatórias externas, o que torna esse processo de classificação complexo.