

Ciencia de datos, práctica 5

Juan Casado Ballesteros, Samuel García Gonzalez, Iván Anaya Martín

December 3, 2019

Abstract

Pendiente

```

;;functions, echo=False;frecuenciaAbsoluta j- function(data) uniquedataj-
unique(data) uniquedataj- sort(uniquedata) newdataj- vector(mode="numeric",
length=0) for (value in uniquedata) ammountj-length(data[data==value]) newdataj-
c(newdata, ammount) setNames(newdata, uniquedata) frecuenciaRelativa j-
function(data) frecuencia j-frecuenciaAbsoluta(data) uniquedataj-unique(data)
uniquedataj- sort(uniquedata) newdata j- vector(mode="numeric", length=0)
for (value in 1:length(uniquedata)) newdata j- c(newdata, frecuencia[value]/length(data))
setNames(newdata, uniquedata) plotFrequencyData j- function(data, xlabel="")
uniquedataj-unique(data) frecuenciarelativa < -as.table(frecuenciaRelativa(data))media<-
mediaAritmetica(data)desviaciontípica<-desviacionTípica(data)tchebychevmin<
-media-2*desviaciontípicatchebychevmax<-media+2*desviaciontípicainrange=
min(tchebychevmin, min(uniquedata))maxrange=max(tchebychevmin, max(uniquedata))plot(frecuencia
"h", xlab = xlabel, xlim = c(minrange, maxrange))abline(v = c(media, tchebychevmin, tchebychevmax), col
c("blue", "gray", "gray"), lty = c(2, 2, 3, 3), lwd = c(2, 2, 2, 1))legend("topright", legend =
c("media", "2*desviación típica"), fill = c("blue", "gray"))mediaAritmetica <
-function(data)acc < -0for(valueindata)acc < -acc + valueacc/length(data)varianza <
-function(data)vmedia < -mediaAritmetica(data)acc = 0for(valueindata)acc < -acc + (value - vmedia
-function(data)varianza(data)(1/2)@

```

1 K-vecinos sobre la muestra proporcionada para obtener outliers

Aplicaremos el algoritmo k-vecinos sobre la muestra que tenemos. Este algoritmo identificará de forma supervisada en la muestra datos anómalos, para poder obtener los outliers. En primer lugar deberemos cargar esta desde un archivo .txt. `##cargar_datos >>= datos1 <- read.table("datos1.txt")datos1@`

En segundo lugar calculamos las distancias euclídeas entre todos los puntos `##distancias_i_i= distancias_i- as.matrix(dist(datos1)) distancias @`

Posteriormente, ordenamos las distancias de cada punto a todos los demás.

`##distancias_ordenar >>= for(i in 1 : length(distancias_ordenadas[, 1]))distancias[, i] <- sort(distancias[, i]) -distancias@`

Como último paso, reordenamos la matriz para organizarla en función de la distancia de cada punto a su vecino número 1,2,3...etc. Tras haber organizado la matriz, buscamos en el tercer vecino, que es el valor k que hemos usado en nuestro análisis para poder identificar los outliers. `##ordenacionfinalconoutlier_i_i= outliers_k_means = list()for(i in 1 : length(distancias_ordenadas[, 1]))if(distancias_ordenadas[4, i] > 2.5)outliers_k_means[i] = list(datos1[i,], distancias_ordenadas[i,], distancias[i, i])`

2 Deteccion de datos anómalos sobre la resistencia

`##cargar_datos >>= datos2 <- read.table("datos2.txt")datos2 <- data.frame(datos2)datos2@`

`##plot_caja_bigotes, fig = TRUE >>= boxplot(datos2Resistencia, range=1.5, horizontal = TRUE) @`

`##plot_caja_bigotes, fig = TRUE >>= boxplot(datos2Densidad, range=1.5, horizontal = TRUE) @`

3 Dispersión sobre la densidad, desviación típica

```
;;desviacion_típica >>= datos2 <- read.table("datos2.txt")@
```

```
;;desviacion_típica_densidad_plot, fig = TRUE >>=
plotFrequencyData(datos2Densidad) @
```

```
;;desviacion_típica_densidad >>= int <- -c(mediaAritmetica(datos2Densidad)
- 2*desviacionTipica(datos2Densidad), mediaAritmetica(datos2Densidad) + 2*desviacionTipica(datos2Densidad))
list())for(iin1 : length(datos2Densidad)) if ((datos2Densidad[i] < int[1]) || (datos2Densidad[i] > int[2]))
outliers_desviacion[[length(outliers_desviacion)+1]] <- -t(matrix(c(i, datos2[i, ]Densidad),
dimnames=list(c("Indice", "Densidad")))) outliers_desviacion@
```

```
;;desviacion_típica_resistencia_plot, fig = TRUE >>=
plotFrequencyData(datos2Resistencia) @
```

```
;;desviacion_típica_resistencia >>= int <- -c(mediaAritmetica(datos2Resistencia)
- 2*desviacionTipica(datos2Resistencia), mediaAritmetica(datos2Resistencia)
+ 2*desviacionTipica(datos2Resistencia))outliers_desviacion = list()for(iin1 :
length(datos2Resistencia)) if ((datos2Resistencia[i] < int[1]) || (datos2Resistencia[i] > int[2]))
outliers_desviacion[[length(outliers_desviacion)+1]] <- -t(matrix(c(i, datos2[i, ]Resistencia),
dimnames=list(c("Indice", "Resistencia")))) outliers_desviacion@
```