

Práctica 11

# Prueba de normalidad analítica

---

Luis Eduardo Galindo Amaya (1274895)

Asignatura	Estadística Avanzada
Docente	Olivia Mendoza Duarte
Fecha	23-11-2022

# Prueba de normalidad analítica

Luis Eduardo Galindo Amaya (1274895)

23-11-2022

## Información del dataset<sup>1</sup>

This is one of the best known datasets in statistics and machine learning. Fisher's paper is a classic in the field and is frequently used for tutorial and teaching purposes. The data set contains 3 classes of 50 instances each, where each class refers to a type of iris plant. One class is linearly separable from the other 2; the latter are not linearly separable from each other.

Predicted attribute: class of iris plant.

1. sepal length
2. sepal width
3. petal length
4. petal width
5. class

## Desarrollo de la práctica

Distribucion normal (sepal width)

### Observacion previa

Se puede notar como el ancho del sépalos se distribuye de manera normal, el tamaño del sépalos tiende a ser de un tamaño específico y no tanto del tamaño de los pétalos

### Resultados

Efectivamente los valores coinciden en una distribucion normal  $RQ = 0,9925113$ .

---

<sup>1</sup><https://archive-beta.ics.uci.edu/ml/datasets/iris>

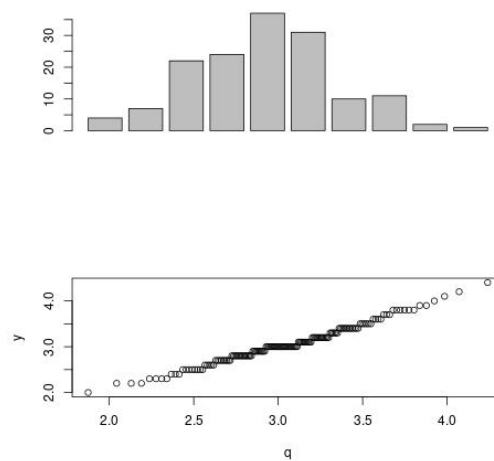


Figura 1: columna "sepal width"

Distribucion no normal (sepal length)

### Observacion previa

El largo del sépalo se distribuye de manera serrada sobre los datos, no parece haber un patrón en la gráfica

### Resultados

Aunque de a primera vista no es muy claro los datos de esta columna tambien cumplen la prueba de normalidad  $rQ = 0,9891878$ .

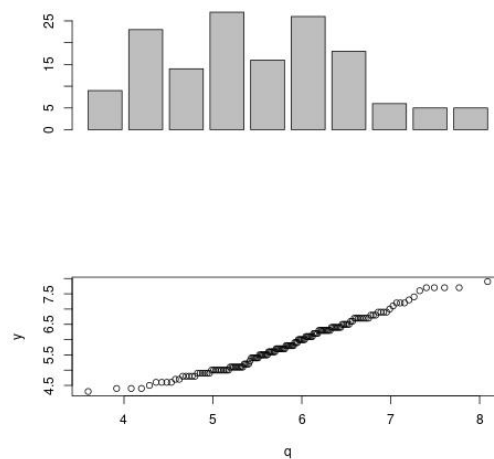


Figura 2: columna "sepal length" del dataset bezdekIris

Distribucion no normal (petal length)

### Observacion previa

Por otro lado el largo de los pétalos es muy interesante, en la parte derecha parece haber una distribución normal pero tiene unos sectores que sobresalen a la izquierda.

### Resultados

Esta no es una distribución normal, la gráfica sale completamente dividida  $rQ = 0,9378633$

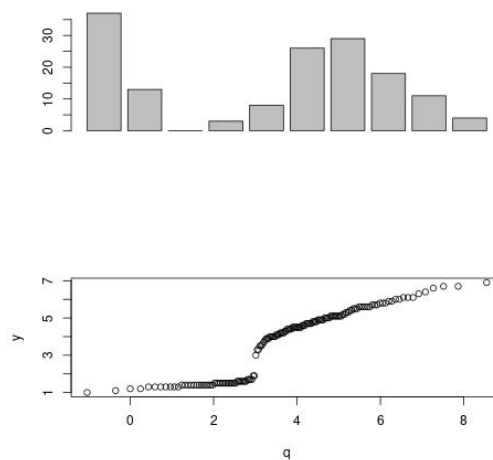


Figura 3: columna "petal length" del dataset bezdekIris

### Conclusiones

Usar métodos que solo dependen de la observación es bastante peligroso, ya que no podemos determinar con total precisión si nuestras gráficas coinciden o no y que un simple cambio en el numero de clases puede hacer que la gráfica sea completamente diferente, el método analítico nos permite tener mucha precisión para determinar si nuestros casos coinciden

### Código

```

1  ## AUTHOR: Luis Eduardo Galindo Amaya
2  ##   DATE: 24-11-2022
3  ##   DESC: Prueba de normalidad grafica y analitica
4
5  ## ENTRADAS
6  n_clases <- 10
7  nombre_del_archivo <- "../bezdekIris.csv"
8
9  columna <- 3
10 filas <- 1:150
11
12 archivo_salida <- "../practica 11/img/i3.jpeg"
13

```

```
14 ## OPERACIONES GRAFICA ####
15 archivo <- read.csv(nombre_del_archivo)
16 data <- archivo[filas, columna]
17
18 valor_minimo <- min(data)
19 valor_maximo <- max(data)
20 amplitud_de_clase <- (valor_maximo - valor_minimo) / n_clases
21
22 ## numero de elementos dentro de cada clase
23 frecuencias <- array(0, dim = (n_clases))
24
25 for (i in 1:n_clases) {
26   rango_min <- valor_minimo + amplitud_de_clase * (i - 1)
27   rango_max <- valor_minimo + amplitud_de_clase * i
28   frecuencias[i] <- sum(data >= rango_min & data < rango_max)
29 }
30
31
32 ## OPERACIONES ANALÍTICA ####
33 n <- length(data) # tamaño del arreglo
34
35 ## Se calcula las fracciones las fracciones correspondientes a
36 ## los acuartiles teoricos
37
38 fracciones <- c()
39 xdata <- c()
40
41 for (i in 1:n) {
42   fraccion <- (i - 0.5) / n
43   fracciones <- c(fracciones, fraccion)
44   xdata <- c(xdata, data[i])
45 }
46
47 # se ordena el vector x y se asiga ordenado a la variable y
48 y <- sort(xdata)
49
50 # se calcula la media y la desviación estándar de y
51 media_y <- mean(y)
52 std_y <- sd(y)
53
54 ## se calculan los cuartiles teóricos usando la distribuci
55 ## normal inversa qnorm
56
57 q <- c()
58 for (i in 1:n) {
59   qi <- media_y + std_y * qnorm(fracciones[i])
60   q <- c(q, qi)
61 }
62
63 datosQQ=data.frame(q,y)
64 rQ=cor(q,y)
65
66
67
68
69 ## grafica
70 jpeg(file = archivo_salida)
71 par(mfrow = c(2, 1))
72 barplot(frecuencias)
73 plot(datosQQ)
74 dev.off()
75
76 ## SALIDAS GRAFICAS ####
77 frecuencias
78 n_clases
79 valor_minimo
80 valor_maximo
81 amplitud_de_clase
82 length(data)
83 sum(frecuencias)
84 ## SALIDAS ANALÍTICAS ####
85 rQ
```