

Cwiczenie 3 - Rozkłady empiryczne i teoretyczne

Michał Marosz

31 października 2015

Spis treści

| | |
|---|-----------|
| Rozkład empiryczny i dystrybuanta empiryczna | 6 |
| Estymacja parametrów rozkładów teoretycznych | 8 |
| Zmienne dyskretne - rozkłady teoretyczne | 15 |
| Rozkład Bernouli'ego | 15 |
| Rozkład Poisson'a | 15 |
| Zmienne ciągłe - rozkłady teoretyczne | 16 |
| Rozkład Gauss'a | 16 |
| Rozkład Gamma | 16 |
| Rozkład Weibull'a | 16 |
| Rozkład GEV (Generalized Extreme Value) | 16 |
| Rozkłady wykorzystywane często we wnioskowaniu statystycznym | 17 |
| Rozkład t-Studenta | 17 |
| Rozkład χ^2 | 17 |
| Rozkład Fishera-Snedecora | 17 |

Analiza właściwości zmiennych jest jednym z podstawowych zadań z jakimi przyjdzie się Wam zmierzyć, w trakcie analizy danych. Dlatego rozpoczniemy od analizy rozkładów empirycznych a następnie wprowadzimy pojęcia i praktyczne zastosowanie dostępnych w R rozkładów teoretycznych, które najczęściej znajdują zastosowanie w analizach z zakresu Klimatologii. W R dostępnych jest cała gama rozkładów teoretycznych. Wystarczy w *pomocy* Rstudio wpisać *distributions* i otrzymamy dostęp do informacji z tego zakresu.

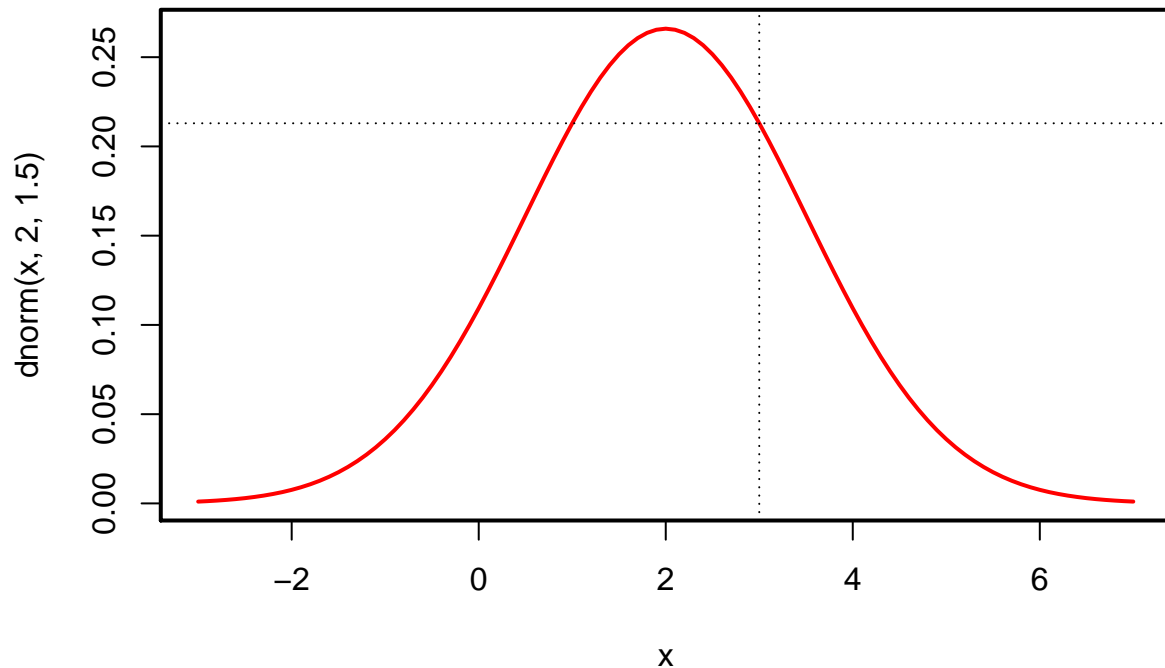
Istotnym jest aby nauczyć się “zadawać” R odpowiednie *pytania* w analizach rozkładów. I tak, jeżeli chcemy uzyskać informację o wartości gęstości prawdopodobieństwa w określonym rozkładzie i dla konkretnej wartości skrót nazwy rozkładu poprzedzimy literą *d* (od Density) np. wpisując:

```
dnorm(3, 2, 1.5)
```

```
## [1] 0.2129653
```

jako rezultat otrzymamy prawdopodobieństwo wystąpienia wartości zmiennej 3, jeżeli ma ona rozkład normalny (Gauss’a) o $\mu = 2$ i $\sigma = 1.5$. Można to graficznie przedstawić w sposób następujący:

```
curve(dnorm(x, 2, 1.5), xlim=c(-3, 7), col=2, lwd=2, add=F)
abline(v=3, h=dnorm(3, 2, 1.5), lty=3)
box(lwd=2)
```



Kolejną informację której możemy pożądać stanowi wartość prawdopodobieństwa, że przy założeniu określonego rozkładu przekroczone zostanie (lub też nie), określona wartość. Wówczas posłużymy się przedrostkiem *p* np.:

```
pnorm(3, 2, 1.5)
```

```
## [1] 0.7475075
```

udzieli nam to odpowiedzi na pytanie jakie jest prawdopodobieństwo że w rozkładzie normalnym o $\mu = 2$ i $\sigma = 1.5$ wartość będzie niższa od 3. Jeżeli natomiast interesuje nas to, czy będzie ona większa użyjemy dodatkowego argumentu *lower.tail=FALSE* np.:

```
pnorm(3, 2, 1.5, lower.tail = FALSE)
```

```
## [1] 0.2524925
```

można oczywiście wykorzystać poprzedni kod ale wynik trzeba odjąć od 1.

```
1-pnorm(3, 2, 1.5)
```

```
## [1] 0.2524925
```

Np. analizujemy rozkład wartości średniej dobowej temperatury powietrza dla jednego z miesięcy letnich. Załóżmy, że jest to rozkład normalny o $\mu = 15$ i $\sigma = 3.2$. jakie jest prawdopodobieństwo, że średnia dobowa temperatura powietrza spadnie poniżej 10°C

```
pnorm(10, 15, 3.2)
```

```
## [1] 0.05908512
```

albo że przekroczy 22°C

```
pnorm(22, 15, 3.2, lower.tail = FALSE)
```

```
## [1] 0.01435302
```

Można oczywiście przemnożyć wynik przez 100 i ładnie zaokrąglić, aby otrzymać wartość w %

```
round(100*pnorm(22, 15, 3.2, lower.tail = FALSE), digits=1)
```

```
## [1] 1.4
```

Kolejnym z pytań które można zadawać R odnosi się do wartości o konkretnym prawdopodobieństwie przekroczenia, czyli innymi słowy poszukujemy wartości kwantyla. W tym wypadku nasz predrostek to q a dla rozkładu normalnego o $\mu = 15$ i $\sigma = 3.2$ o wartość której prawdopodobieństwo przekroczenia wynosi 1% można “zapytać się” w sposób następujący:

```
qnorm(0.99, 15, 3.2)
```

```
## [1] 22.44431
```

lub z uwzględnieniem argumentu *lower.tail*

```
qnorm(0.01, 15, 3.2, lower.tail = FALSE)
```

```
## [1] 22.44431
```

Rozkład empiryczny i dystrybuanta empiryczna

Rozkład empiryczny zazwyczaj przedstawia się z wykorzystaniem histogramu natomiast dystrybuantę empiryczną z wykorzystaniem funkcji *ecdf*

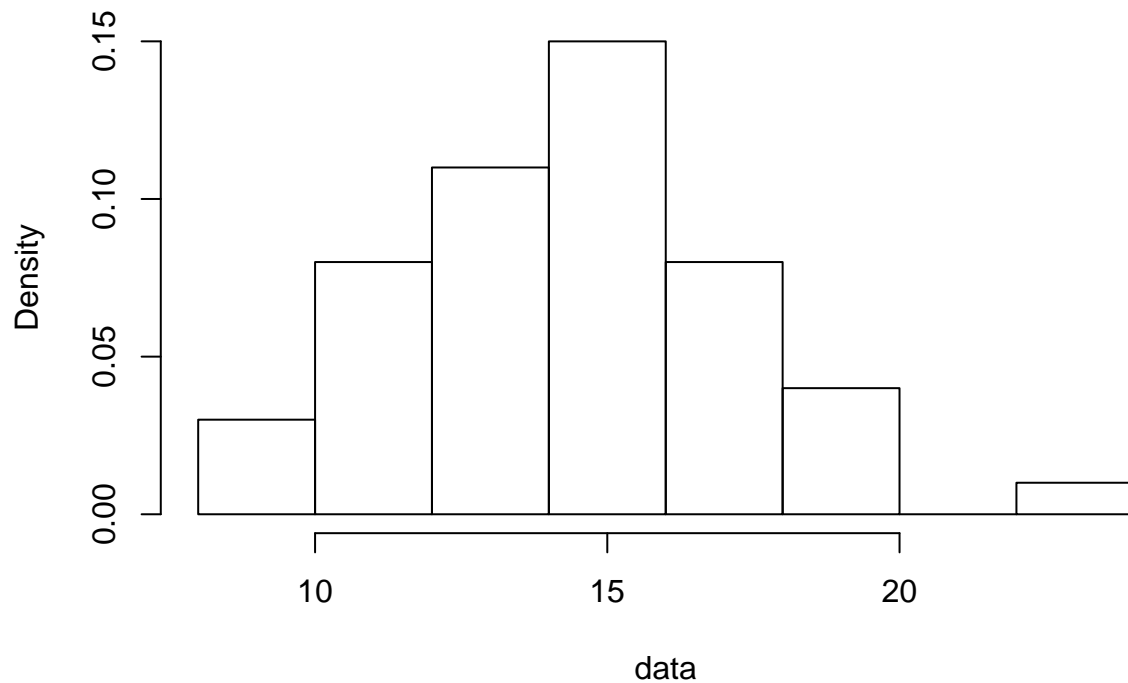
Utwórzmy wektor 100 losowych wartości o rozkładzie normalnym o $\mu = 15$ i $\sigma = 3.2$

```
set.seed(1000)
```

```
data=rnorm(50, 15, 3.2)
```

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

Histogram of data

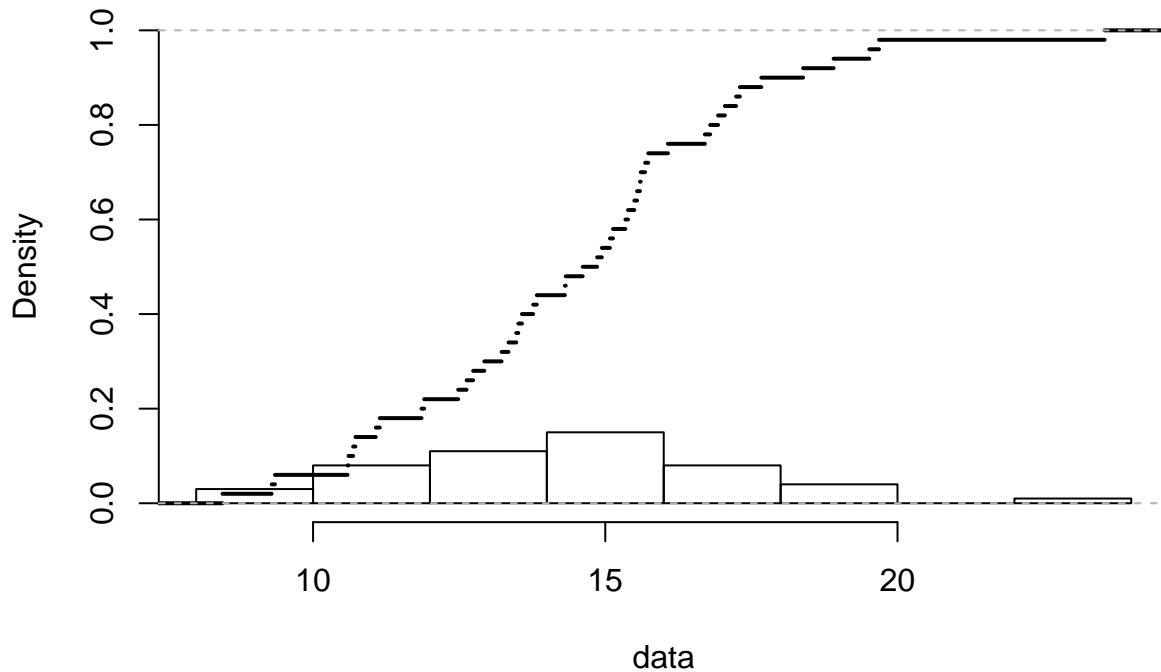


Aby do-
dać do powyższego histogramu dystrybuantę empiryczną posłużymy się funkcją *ecdf*

```
hist(data, prob=T, ylim=c(0,1))  
plot(ecdf(data), vertical=FALSE, pch="", add=T, lwd=2)
```

```
## Warning in segments(ti.l, y, ti.r, y, col = col.hor, lty = lty, lwd =  
## lwd, : 'vertical' nie jest parametrem graficznym
```

Histogram of data



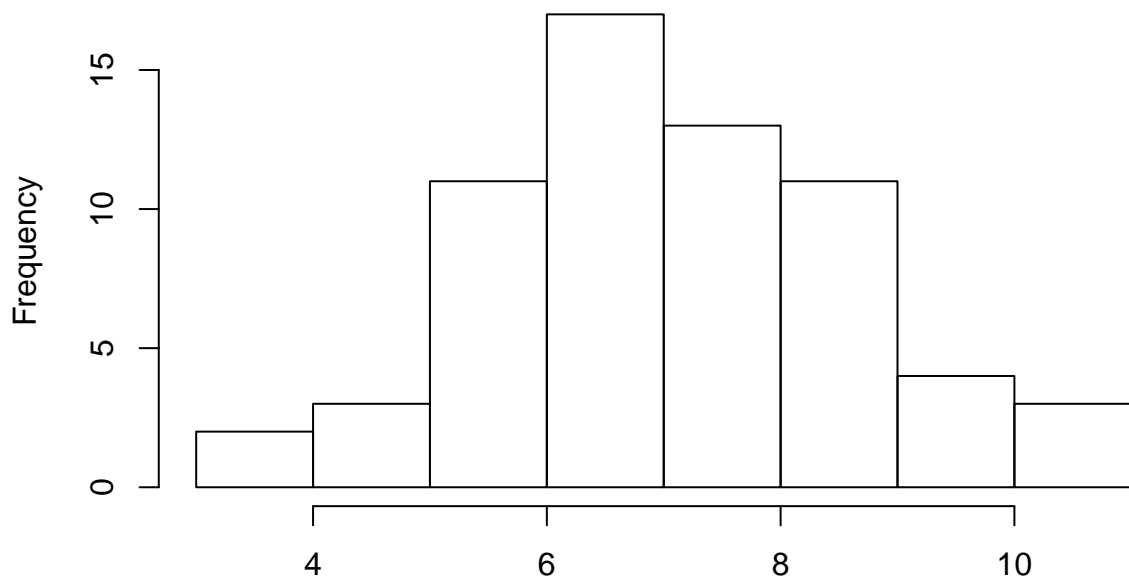
Estymacja parametrów rozkładów teoretycznych

W przypadku rozkładu normalnego (Gauss'a) estymacja parametrów nie następuje większych problemów. Średnia μ oraz odchylenie standardowe σ są wystarczająco dobrymi estymatorami parametrów rozkładu normalnego. Jednak dla pozostałych rozkładów niezbędne jest posłużenie się dodatkowymi funkcjami pozwalającymi na dopasowanie parametrów rozkładów w tym celu niezbędne jest zainstalowanie *paczki fitdistplus*

Dopasujemy parametry rozkładu Gaussa za pomocą funkcji *fitdist* do temperatur powietrza w kwietniu.

```
dane=read.table("air.txt", header=T)
attach(dane)
temp04=TEMP[which(MC==4)]
hist(temp04)
```


Histogram of temp04



temp04

policzmy

wartość średnią μ oraz odchylenie standardowe σ

```
mean(temp04)
```

```
## [1] 7.184375
```

```
sd(temp04)
```

```
## [1] 1.564408
```

Teraz sprawdzimy jakie wartości parametrów dopasuje funkcja *ftdis* z wykorzystaniem metody największej wiarygodności (*mle*), będącej standardem we współczesnych analizach.

Wpiszmy następujący kod:

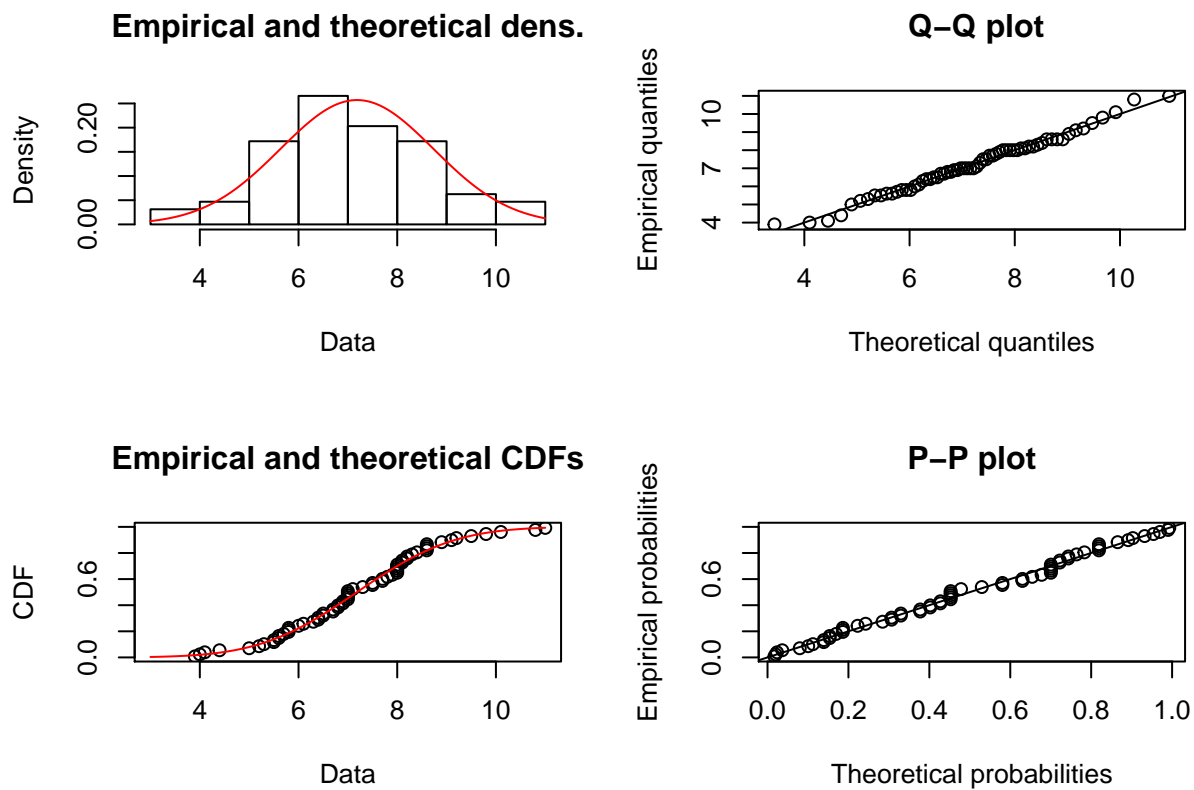
```
library("fitdistrplus", lib.loc="~/R/x86_64-pc-linux-gnu-library/3.2")
```

```
## Loading required package: MASS
```

```
normal_fit = fitdist(temp04, dnorm, method = "mle")
```

Zwróćcie uwagę wynik analiz jest obiektem, który można poddać dalszej analizie w celu oceny dopadowania np.:

```
plot(normal_fit)
```



Moż-

na również “wyciągnąć” wartości parametrów ze zmiennej *estimate* będącej jedną ze składowych obiektu.

```
normal_fit$estimate
```

```
##      mean      sd  
## 7.184375 1.552138
```

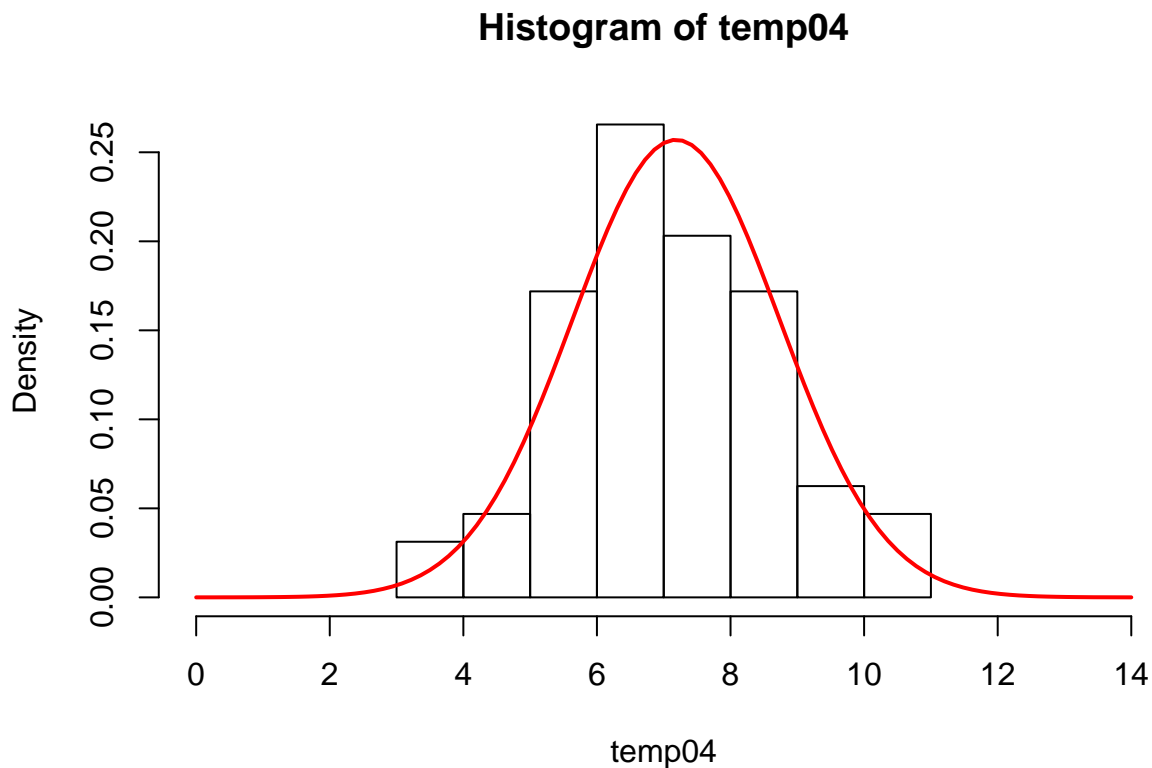
Teraz możemy wykreślić ponownie histogram i dodać do niego krzywą rozkładu normalnego wykreslona na podstawie dopasowanych parametrów. “Wypreparujmy” je najpierw z obiektu

```
mean04=as.numeric(normal_fit$estimate[1])
```

```
sd04=as.numeric(normal_fit$estimate[2])
```

```
hist(temp04, prob=T, xlim=c(0,14))
```

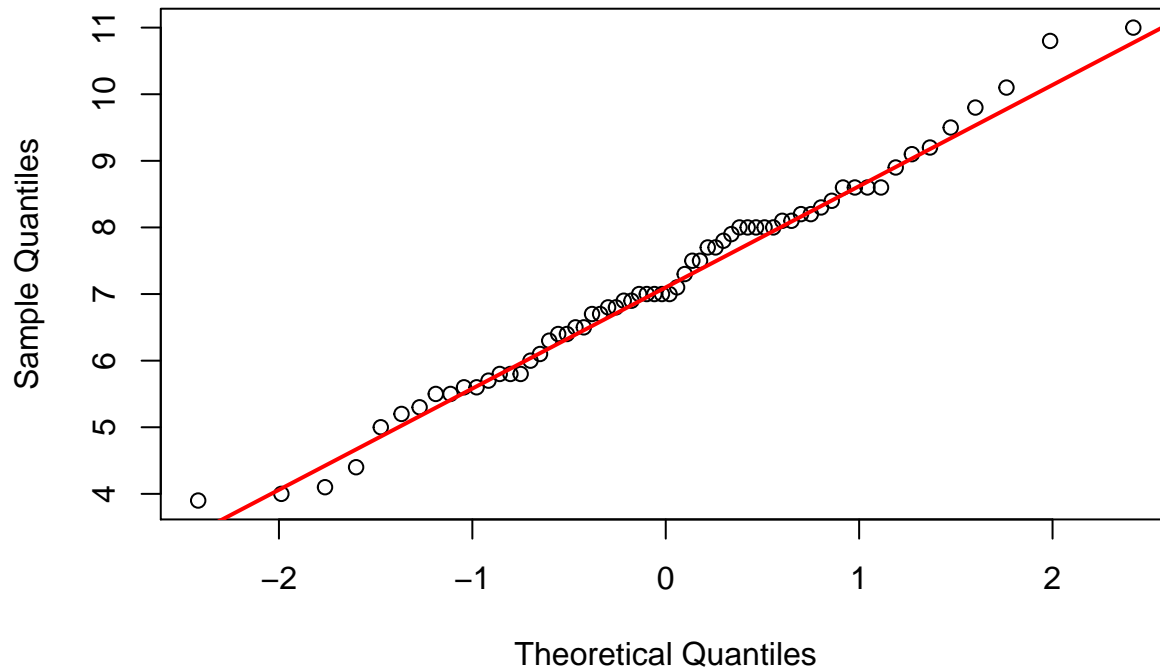
```
curve(dnorm(x, mean04, sd04), add=T, col=2, lwd=2)
```



Istnieje również klasa wykresów przeznaczona specjalnie do wizualnego porównywania rozkładów empirycznych z teoretycznym normalnym: *qqnorm*

```
qqnorm(temp04)
qqline(temp04, col=2, lwd=2)
```

Normal Q-Q Plot



Weryfikację jakości dopasowania można przeprowadzić wizualnie za pomocą uprzednio wywołanej funkcji `plot/qqnorm` lub zaprząć do tego nieco bardziej sformalizowane testy np.: Shapiro-Wilk'a `shapiro.test`

```
shapiro.test(temp04)
```

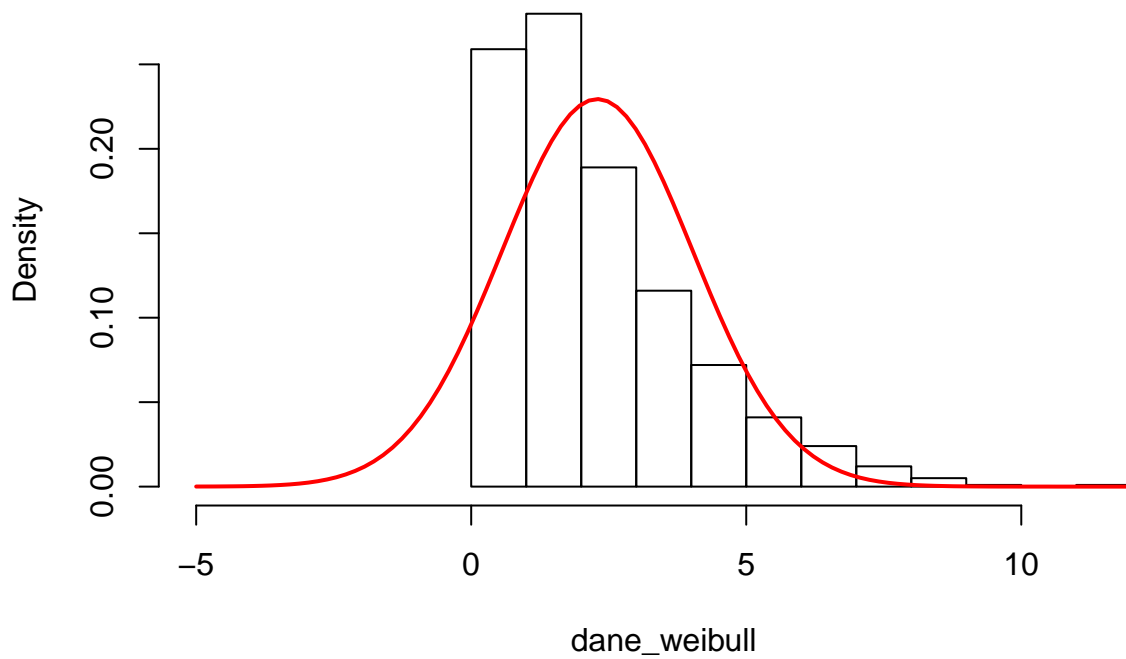
```
##
## Shapiro-Wilk normality test
##
## data: temp04
## W = 0.98799, p-value = 0.7913
```

Dla jasności wygenerujemy dane o rozkładzie innym niż normalny i sprawdzimy wykresy oraz wartości z powyższego testu.

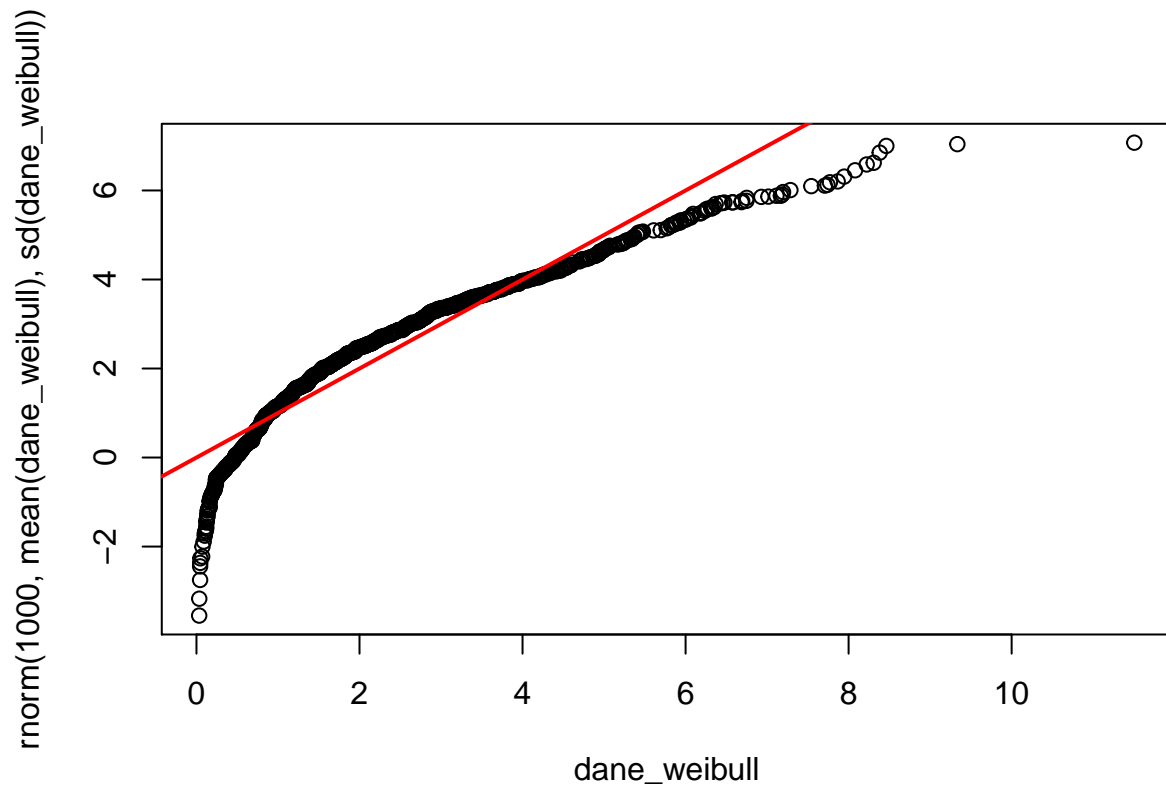
```
set.seed(1000)
dane_weibull=rweibull(1000, 1.3, 2.5)

hist(dane_weibull, prob=T, xlim=c(-5, 12))
curve(dnorm(x, mean(dane_weibull), sd(dane_weibull)), col=2, lwd=2, add=T)
```

Histogram of dane_weibull



```
qqplot(dane_weibull, rnorm(1000, mean(dane_weibull), sd(dane_weibull)))
abline(0,1, col=2, lwd=2)
```



```
shapiro.test(dane_weibull)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data:  dane_weibull  
## W = 0.90707, p-value < 2.2e-16
```

Zmienne dyskretne - rozkłady teoretyczne

Rozkład Bernouli'ego

Rozkład Poisson'a

Zmienne ciągłe - rozkłady teoretyczne

Rozkład Gauss'a

Rozkład Gamma

Rozkład Weibull'a

Rozkład GEV (Generalized Extreme Value)

Rozkłady wykorzystywane często we wnioskowaniu statystycznym

Rozkład t-Studenta

Rozkład χ^2

Rozkład Fishera-Snedecora