

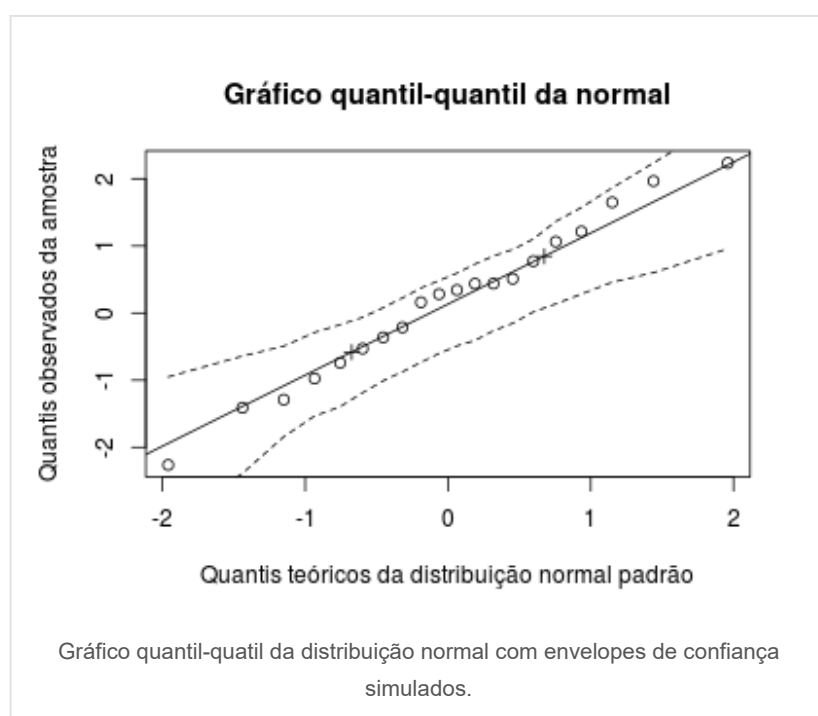
# R blogs / lang

Articles about R, in your own language

## Como fazer e interpretar o gráfico quantil-quantil

by [Walmes Zeviani](#) • November 30, 2012

This post was kindly contributed by [Ridículas](#) - go there to comment and to read [the full post](#).



O gráfico quantil-quantil (q-q) é uma ferramenta muito útil para checar adequação de distribuição de frequência dos dados à uma distribuição de probabilidades. Situações como essa ocorrem principalmente na análise de resíduos de modelos de regressão onde o gráfico q-q é usado para verificar se os resíduos apresentam distribuição normal. O gráfico q-q é melhor que o histograma e o gráfico de distribuição acumulada empírica porque nós temos mais habilidade para verificar se uma reta se ajusta aos pontos do que se uma curva de densidade se ajusta a um histograma ou uma curva de probabilidade acumulada se ajusta à acumulada empírica. Compare às três visualizações com o código a seguir.

```
1  #-----  
2  # q-q vs histograma vs ecdf  
3  
4  y <- rnorm(50)  
5  par(mfrow=c(1,3))  
6  qqnorm(y); qqline(y)  
7  plot(density(y))  
8  curve(dnorm(x, mean(y), sd(y)), add=TRUE, col=2)  
9  plot(ecdf(y))  
10 curve(pnorm(x, mean(y), sd(y)), add=TRUE, col=2)
```

```

11
12 # qual você sente mais segurança para verificar adequação?
13 # nossos olhos têm mais habilidade para comparar alinhamentos
14 #-----

```

Apesar de muito usado, poucos usuários conhecem o procedimento para fazê-lo e interpretá-lo. O procedimento é simples e pode ser estendido para outras distribuições de probabilidade, não apenas para a distribuição normal como muitos podem pensar. Além do mais, alguns padrões do gráfico q-q obedecem à certas características dados dados, como assimetria, curtose e discrecicidade. Saber identificar essas características é fundamental para indicar uma transformação aos dados. Abaixo o código para o gráfico q-q para distribuição normal, q-q para qualquer distribuição e o q-q com envelope de confiança obtido por simulação. A execução e estudo do código esclarece o procedimento. O envelope de confiança por simulação torna-se proibitivo para grandes amostras pelo tempo gasto na simulação.

```

1 #-----
2 # função para fazer o gráfico quantil-quantil da normal
3
4 qqn <- function(x, ref.line=TRUE){
5   x <- na.omit(x)           # remove NA
6   xo <- sort(x)             # ordena a amostra
7   n <- length(x)           # número de elementos
8   i <- seq_along(x)         # índices posicionais
9   pteo <- (i-0.5)/n         # probabilidades teóricas
10  qteo <- qnorm(pteo)        # quantis teóricos sob a normal padrão
11  plot(xo~qteo)              # quantis observados ~ quantis teóricos
12  if(ref.line){
13    qrto <- quantile(x, c(1,3)/4) # 1º e 3º quartis observados
14    qrtt <- qnorm(c(1,3)/4)      # 1º e 3º quartis teóricos
15    points(qrtt, qrto, pch=3)    # quartis, passa uma reta de referência
16    b <- diff(qrto)/diff(qrtt)  # coeficiente de inclinação da reta
17    a <- b*(0-qrtt[1])+qrto[1]  # intercepto da reta
18    abline(a=a, b=b)           # reta de referência
19  }
20 }
21
22 x <- rnorm(20)
23 par(mfrow=c(1,2))
24 qqn(x)
25 qqnorm(x); qqline(x)
26 layout(1)
27
28 #-----
29 # função para fazer o gráfico quantil-quantil de qualquer distribuição
30
31 qqq <- function(x, ref.line=TRUE, distr=qnorm, param=list(mean=0, sd=1)){
32   x <- na.omit(x)           # remove NA
33   xo <- sort(x)             # ordena a amostra
34   n <- length(x)           # número de elementos
35   i <- seq_along(x)         # índices posicionais
36   pteo <- (i-0.5)/n         # probabilidades teóricas
37   qteo <- do.call(distr,    # quantis teóricos sob a distribuição
38                 c(list(p=pteo), param))
39   plot(xo~qteo)              # quantis observados ~ quantis teóricos
40   if(ref.line){
41     qrto <- quantile(x, c(1,3)/4) # 1º e 3º quartis observados
42     qrtt <- do.call(distr,      # 1º e 3º quartis teóricos
43                   c(list(p=c(1,3)/4), param))
44     points(qrtt, qrto, pch=3)    # quartis, por eles passa uma reta de referênc
45     b <- diff(qrto)/diff(qrtt)  # coeficiente de inclinação da reta
46     a <- b*(0-qrtt[1])+qrto[1]  # intercepto da reta
47     abline(a=a, b=b)           # reta de referência
48   }
49 }
50
51 x <- rnorm(20)

```

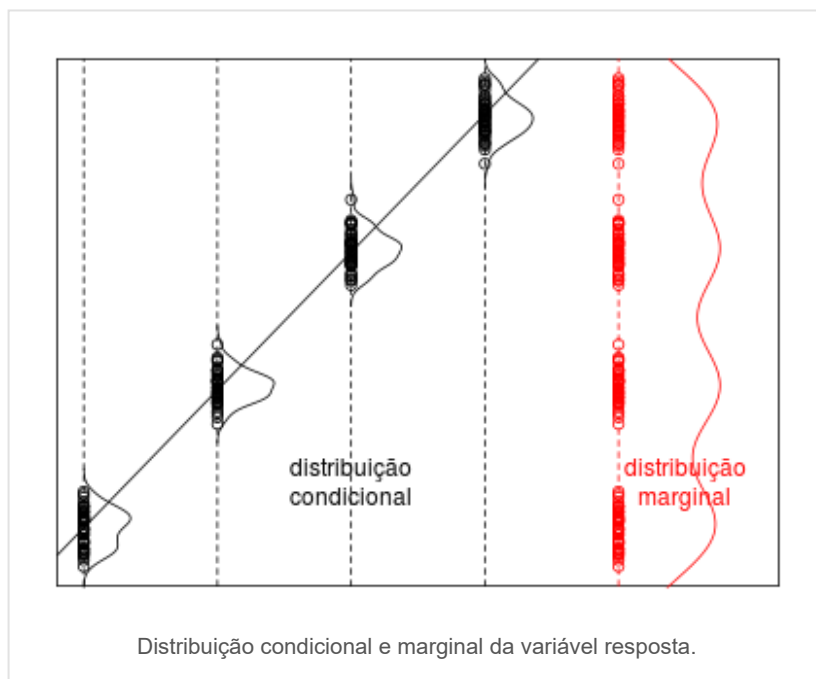
```

52 par(mfrow=c(1,2))
53 qqnorm(x)
54 qqnorm(x); qqline(x)
55 layout(1)
56
57 x <- runif(20)
58 qq(x, ref.line=TRUE, distr=qunif, param=list(min=0, max=1))
59
60 x <- rgamma(20, shape=4, rate=1/2)
61 qq(x, ref.line=TRUE, distr=qgamma, param=list(shape=4, rate=1/2))
62
63 #-----
64 # envelope para o gráfico de quantis (simulated bands)
65
66 qqsb <- function(x, ref.line=TRUE, distr=qnorm, param=list(mean=0, sd=1),
67                  sb=TRUE, nsim=500, alpha=0.95, ...){
68   x <- na.omit(x) # remove NA
69   xo <- sort(x) # ordena a amostra
70   n <- length(x) # número de elementos
71   i <- seq_along(x) # índices posicionais
72   pteo <- (i-0.5)/n # probabilidades teóricas
73   qteo <- do.call(distr, # quantis teóricos sob a distribuição
74                  c(list(p=pteo), param))
75   plot(xo~qteo, ...) # quantis observados ~ quantis teóricos
76   if(ref.line){
77     qrto <- quantile(x, c(1,3)/4) # 1º e 3º quartis observados
78     qrtt <- do.call(distr, # 1º e 3º quartis teóricos
79                    c(list(p=c(1,3)/4), param))
80     points(qrtt, qrto, pch=3) # quartis, passa uma reta de referência
81     b <- diff(qrto)/diff(qrtt) # coeficiente de inclinação da reta
82     a <- b*(0-qrtt[1])+qrto[1] # intercepto da reta
83     abline(a=a, b=b) # reta de referência
84   }
85   if(sb){
86     rdistr <- sub("q", "r", # função que gera números aleatórios
87                  deparse(substitute(distr)))
88     aa <- replicate(nsim, # amostra da distribuição de referência
89                    sort(do.call(rdistr, c(list(n=n), param))))
90     lim <- apply(aa, 1, # limites das bandas 100*alpha%
91                 quantile, probs=c((1-alpha)/2,(alpha+1)/2))
92     matlines(qteo, t(lim), # coloca as bandas do envelope simulado
93              lty=2, col=1)
94   }
95 }
96
97 x <- rnorm(20)
98
99 #png("f047.png", 400, 300)
100 qqsb(x, xlab="Quantis teóricos da distribuição normal padrão",
101      ylab="Quantis observados da amostra",
102      main="Gráfico quantil-quantil da normal")
103 #dev.off()
104
105 x <- rpois(50, lambda=20)
106 qqsb(x, distr=qpois, param=list(lambda=20))
107
108 #-----

```

A normalidade dos resíduos é um dos pressupostos da análise de regressão. **São os resíduos e não os dados que devem apresentar normalidade.** Se a distribuição dos dados, ou melhor, da sua variável resposta (Y) condicional ao efeito das suas variáveis explicativas for normal, os resíduos terão distribuição normal. Porém, se você aplica um teste de normalidade aos dados (Y) você não está considerando os efeitos das variáveis explicativas, ou seja, você está aplicando um teste na distribuição marginal de Y que não tem porque atender a normalidade. Todo teste de normalidade supõe que os dados têm uma média e

uma variância e só os resíduos atendem essa premissa porque os dados (Y) têm média dependente do efeito das variáveis explicativas.



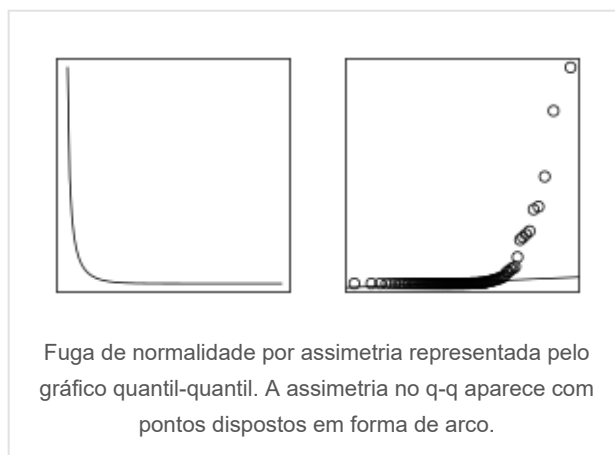
```

1  #-----
2  # distribuição condicional vs distribuição marginal
3
4  layout(1)
5  x <- rep(1:4, e=50)
6  y <- rnorm(x, mean=0+0.75*x, sd=0.1)
7  da <- data.frame(x, y)
8  plot(y~x, da)
9
10 db <- split(da, f=da$x)
11
12 #png("f046.png", 400, 300); par(mar=c(1,1,1,1))
13 plot(y~x, da, xlim=c(1,6), xaxt="n", yaxt="n", xlab="", ylab="")
14 abline(a=0, b=0.75)
15 lapply(db,
16         function(d){
17             dnst <- density(d$y)
18             lines(d$x[1]+dnst$y*0.1, dnst$x)
19             abline(v=d$x[1], lty=2)
20         })
21 points(rep(5, length(y)), da$y, col=2)
22 abline(v=5, lty=2, col=2)
23 dnst <- density(da$y)
24 lines(5+dnst$y*2, dnst$x, col=2)
25 text(3, 1, "distribuição\ndicional")
26 text(5.5, 1, "distribuição\nmarginal", col=2)
27 #dev.off()
28
29 par(mfrow=c(1,2))
30 qqnorm(da$y)
31 qqnorm(residuals(lm(y~x, da)))
32 layout(1)
33
34 #-----

```

Muitos usuários preferem aplicar um teste de normalidade do que olhar para o gráfico q-q. Isso tem duas razões: (1) costume, o usuário sempre usou aplicativos para análise de dados que não dispõem de recursos gráficos, eles conduzem toda análise sem sequer ver os dados em um gráfico, (2) consideram subjetiva a análise gráfica. Do meu ponto de vista, a subjetividade está presente ao aplicar um teste também pois o

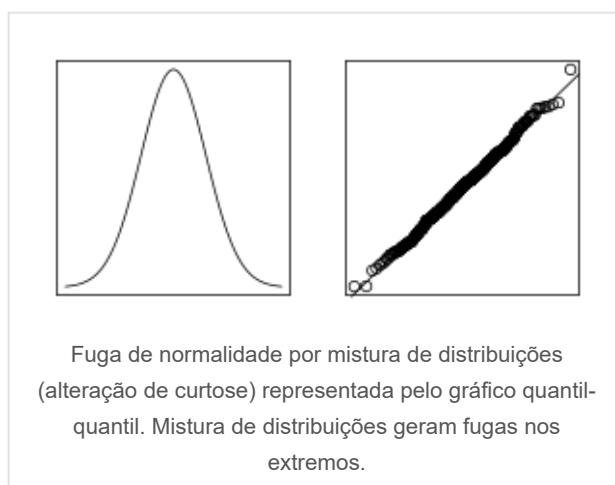
usuário é quem escolhe o teste e o nível de significância. Mas o que de fato eu defendo é que a análise gráfica é indiscutivelmente mais informativa. Veja, se o teste rejeita a normalidade é porque os dados não apresentam distribuição normal por algum motivo. Quando você visualiza o q-q é possível explicar a fuga de normalidade que pode ser sistemática: (1) devido à desvio de assimetria, (2) de curtose, (3) à mistura de distribuições, (4) à presença de um outlier e (5) ao dados serem discretos. Essas são as principais causas de afastamento. Cada uma delas sugere uma alternativa para corrigir a fuga: transformação, modelagem da variância, deleção de outlier, etc. Nesse sentido, como identificar esses padrões de fuga? É o que os gráficos animados vão mostrar. Até a próxima ridícula.



```

1  #-----
2  # assimetria
3
4  dir.create("frames")
5  setwd("frames")
6
7  png(file="assimet%03d.png", width=300, height=150)
8  par(mar=c(1,1,1,1))
9  par(mfrow=c(1,2))
10 for(i in 10*seq(0.01, pi-0.01, l=100)){
11   curve(dbeta(x, i, 10-i), 0, 1, xaxt="n", yaxt="n", xlab="", ylab="")
12   y <- rbeta(100, i, 10-i)
13   qqnorm(y, xaxt="n", yaxt="n", main=NULL, xlab="", ylab=""); qqline(y)
14 }
15 dev.off()
16
17 # converte os pngs para um gif usando ImageMagick
18 system("convert -delay 10 assimet*.png assimet.gif")
19
20 # remove os arquivos png
21 file.remove(list.files(pattern=".png"))
22
23 #-----

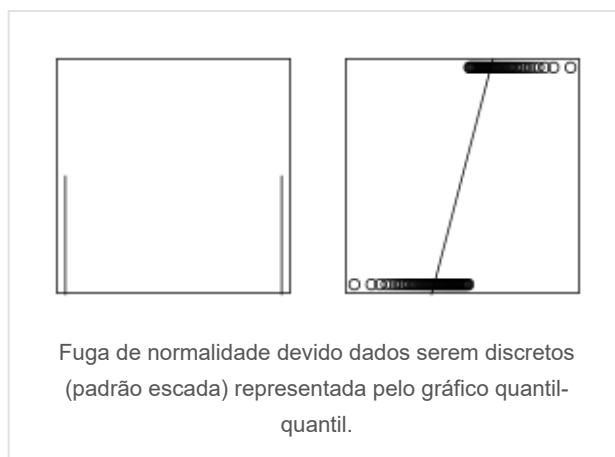
```



```

1  #-----
2  # mistura de distribuições
3
4  png(file="mistura%03d.png", width=300, height=150)
5  par(mar=c(1,1,1,1))
6  par(mfrow=c(1,2))
7  for(i in seq(0, pi, l=100)){
8    curve(i*dnorm(x,0,1)+(1-i)*dnorm(x,0,6), -20, 20,
9          xaxt="n", yaxt="n", xlab="", ylab="")
10     y <- c(rnorm(ceiling(500*i),0,1), rnorm(floor(500*(1-i)),0,6))
11     qqnorm(y, xaxt="n", yaxt="n", main=NULL, xlab="", ylab=""); qqline(y)
12   }
13   dev.off()
14
15   # converte os pngs para um gif usando ImageMagick
16   system("convert -delay 10 mistura*.png mistura.gif")
17
18   # remove os arquivos png
19   file.remove(list.files(pattern=".png"))
20
21   #-----

```



```

1  #-----
2  # discreticidade
3
4  png(file="discret%03d.png", width=300, height=150)
5  par(mar=c(1,1,1,1))
6  par(mfrow=c(1,2))
7  for(i in c(1:100, 99:1)){
8    x <- 0:i
9    px <- dbinom(x, i, 0.5)
10    plot(px~x, type="h", xaxt="n", yaxt="n", xlab="", ylab="")
11    y <- rbinom(100, i, 0.5)
12    qqnorm(y, xaxt="n", yaxt="n", main=NULL, xlab="", ylab=""); qqline(y)
13  }
14  dev.off()
15
16  # converte os pngs para um gif usando ImageMagick
17  system("convert -delay 10 discret*.png discret.gif")
18
19  # remove os arquivos png
20  file.remove(list.files(pattern=".png"))
21
22  #-----

```

- Ile punktów potrzeba by się dostać do szkoły średniej w Warszawie?

---

Copyright © 2019 R blogs / lang. All Rights Reserved.

The Magazine Basic Theme by bavotasan.com.

☺

[← 决策树之三国争霸](#)[Eine kurze Geschichte über R →](#)

## LANGUAGES

- [Chinese](#)
- [Dutch](#)
- [French](#)
- [German](#)
- [Indonesian](#)
- [Italian](#)
- [Korean](#)
- [Polish](#)
- [portuguese](#)
- [Russian](#)
- [Serbian](#)
- [Spanish](#)
- [Uncategorized](#)

## RECENT POSTS

- [Tell Me a Story: How to Generate Textual Explanations for Predictive Models](#)
- [dime: Deep Interactive Model Explanations](#)
- [Learn about XAI in R with „Predictive Models: Explore, Explain, and Debug”](#)
- [Communautés de prénoms](#)