

O R para Ajuste de Distribuixxes

26 de novembro de 2007

A Estatística Descritiva possui ferramentas para que o analista quantitativo elabore modelos probabilísticos para as variáveis analisadas. O R possui um conjunto de funções para o ajuste de distribuições de probabilidade aos dados.

Abaixo, apresentamos uma lista de funções do R utilizadas para auxiliar o ajuste de distribuições de probabilidade.

```

hist(): xhistogramaxdexfrequncia
xstem(): xoxramo-e-folhas
xdensity(): xajustaxdensidadexdexformaxno-paramxtrica
xcut(): xconstruxxoxdexfaixasxoxvaloresxparaxumaxvarixvel
xecdf(): xfunxxoxdexdistribuixxoxacumuladaxempxrica
xfitdistr(): xajustexdexdistribuixxoxunivariadaxporxmximaxverossimilhanxa
xchisq.test(): xttestexqui-quadrado
xshapiro.test(): xttestexparaxaxnormalidade
xjarque.bera.test(): testexparaxaxnormalidade
xks.test(): xttestexdexKolmogorov-Smirnov
xlillie.test(): xttestexparaxaxnormalidade
xad.test(): xttestexparaxaxnormalidade
xqqnorm(): xoxnormalxqqplot

```

Os grxficos exploratxrios e as medidas-resumo podem funcionar como primeira abordagem para identificar uma distribuixxo de probabilidades aderente aos dados.

Histograma

O histograma x uma das formas mais simples de inferir sobre a funxxo densidade de probabilidade de uma varixvel. Embora seja utili com freqxxncia, representa um ajuste pouco suave.

Considere o seguinte exemplo que gera dados da distribuixxo normal e explora o histograma como ferramenta de ajuste de distribuixxo.


```
xx>x#Iniciandoxaxsemente
xx>xset.seed(51007)
xx>x#Definindoxtamanhoxdaxamostra
xx>xn<-100
xx>x#Gerandoxdadosxdaxdistribuixxnormal
xx>xx.norm<-rnorm(n)
xx>x#Construindoxxhistogramaxdexfrecuencias
xx>xh<-hist(x.norm,prob=T)
xx>x#Visualizandoxxdispersioxdosxpontosxnosxintervalos.
xx>xrug(x.norm)
xx>x#Ajustandoxxcurvaxnormalxaosxdados
xx>xcurve(dnorm(x),-3,3,add=T)
```

pict

As informaxes armazenadas no objeto h são:

```
xx>xnames(h)
```

```
xx[1] x"breaks"xxxxxx"counts"xxxxxx"intensities"x"density"xxxxx"mids"  
xx[6] x"xname"xxxxxxx"equidist"
```

```

breaks: xlimitesxsuperioresxdosintervalosdexclasse
xcounts: xfreqxnciabsolutaxnosintervalosdexclasse
xintensities: xdensidadexfreqxncia
xdensity: xdensidadexfreqxncia
xmids: xpontosxmdiosdosintervalos
xxname: xnomexdaxvarixvel
xequidist: xvarixvelxlxgicaxindicandoquexosintervalosxpossuemxmesmoxcomprimento

```

A partir das informaxes geradas pelo histograma, vamos remontar a tabela de freqxncias em classes.

O comando apresentado na seqxncia transforma a varixvel `x.norm`, originalmente quantitativa, em uma varixvel qualitativa cujos `nx`veis representam as faixas de valores de uma tabela de freqxncia em classes.

```

xx>x#xCriandoxosintervalosdexclasse
xx>xclasses<-cut(x.norm,breaks=h$breaks)
xx>x#xTabelaxfreqxnciaxcomxasxclasses
xx>xtable(classes)

```

```

xxclasses
xx(-3,-2.5]x(-2.5,-2]x(-2,-1.5]x(-1.5,-1]x(-1,-0.5]xx(-0.5,0]xxx(0,0.5]xxx(0.5,1]
xxxxxxxxxx1xxxxxxxxxx2xxxxxxxxxx10xxxxxxxxxx9xxxxxxxxxx16xxxxxxxxxx7xxxxxxxxxx24xxxxxxxxxx15
xxxx(1,1.5]xxx(1.5,2]xxx(2,2.5]
xxxxxxxxxx9xxxxxxxxxx5xxxxxxxxxx2

```


Classes	Freq. Abs.
$[-3, -2.5]$	1
$[-2.5, -2]$	2
$[-2, -1.5]$	10
$[-1.5, -1]$	9
$[-1, -0.5]$	16
$[-0.5, 0]$	7
$[0, 0.5]$	24
$[0.5, 1]$	15
$[1, 1.5]$	9
$[1.5, 2]$	5
$[2, 2.5]$	2

Tabelax1: Tabela de Freqencia em Faixas de Valores

Na Tabela 1 x possível visualizar as frequências absolutas nos intervalos gerados ao construir, de modo automático, o histograma. Informações nas tabelas podem ser mais detalhadas como acontece na Tabela 2.

Classes	Ponto	Mx	di	Freq. Abs.	Freq. Rel.
11(-3,-2.5]	-2.75	11	1	11	0.01
11(-2.5,-2]	-2.25	11	2	11	0.02
11(-2,-1.5]	-1.75	11	10	11	0.1
11(-1.5,-1]	-1.25	11	9	11	0.09
11(-1,-0.5]	-0.75	11	16	11	0.16
11(-0.5,0]	-0.25	11	7	11	0.07
11(0,0.5]	0.25	11	24	11	0.24
11(0.5,1]	0.75	11	15	11	0.15
11(1,1.5]	1.25	11	9	11	0.09
11(1.5,2]	1.75	11	5	11	0.05
11(2,2.5]	2.25	11	2	11	0.02

Tabelax2: Tabela de Freqencia em Faixas de Valores

Vamos calcular algumas estatísticas com base na tabela de frequência produzida.

```
xx>x#FrequenciasAbsolutasnasclasses
xx>xni<-h$count
xx>xni
```

```
xxx[1] 1 2 10 9 16 7 24 15 9 5 2
```

```
xx>x#FrequenciasRelativasnasclasses
xx>xfi<-ni/n
xx>xfi
```

```
xxx[1] 0.01 0.02 0.10 0.09 0.16 0.07 0.24 0.15 0.09 0.05 0.02
```

```
xx>x#Mdiacombaseextabeladefrequencia
xx>xbarra<-h$mids*%fi#produtoxdevetores
xx>xbarra
```

```
xxxxxxx[,1]
xx[1,] 1.08
```

```
xx>x#Desviosdospontosmdiosxrelaxxxxxmdia
xx>xdesv<-h$mids-xbarra
xx>xdesv
```

```
xxx[1] 1.67 1.67 1.17 0.67 0.17 0.33 0.83 1.33 1.83 2.33
```

```
xx>x#Desviosquadráticosdospontosmdiosxrelaxxxxxmdia
xx>xdesv2<-desv^2
xx>xdesv2
```

```
xxx[1] 2.7889 4.7089 2.7889 1.3689 0.4489 0.0289 1.089 0.6889 1.7689 3.3489
xx[11] 5.4289
```

```
xx>x#Variancia
xx>xvarianc<-desv2*%fi
xx>xvarianc
```

```
xxxxxxx[,1]
xx[1,] 1.2061
```

Conforme a Tabela 3 x possível constatar a diferença entre estatísticas calculadas em dados brutos e dados em classes.

	dados em classes	dados brutos
mx	-0.08	-0.10
variancia	1.21	1.16

Tabela 3: Estatísticas Descritivas da Tabela em Intervalos e Dados Brutos

Ramo-e-Folhas

O Ramo-e-Folhas é um gráfico que possui as mesmas características do histograma, com a grande vantagem de reproduzir os valores numéricos e, conseqüentemente, não há perda de informação.

```
xx>xstem(x.norm)
```

```
xxxxThe decimal point is at the |
xx
xxxx-2x | x5
xxxx-2x | x410
xxxx-1x | x9987766655
xxxx-1x | x3221111000
xxxx-0x | x9999888877766655
xxxx-0x | x44220
xxxxx0x | x000001112222333344444
xxxxx0x | x55555556677778889
xxxxx1x | x0111112233
xxxxx1x | x66689
xxxxx2x | x02
```

Repare que, na figura acima, o formato da distribuição é semelhante ao exibido pelo histograma e o número de intervalos é exatamente o mesmo.

Densidade Não-Paramétrica

Uma forma mais sofisticada de estimar a densidade de probabilidade dos dados é feita através do comando `density`.


```
xx>xhist(x.norm,prob=T)
xx>xrug(x.norm)
xx>xcurve(dnorm(x),-3,3,add=T,col='red')
xx>xlines(density(x.norm),col='blue')
```

pict

O resultado pode ser visto como uma variante do histograma, porém suavizado.

Função de Distribuição Acumulada Empírica

Através do comando `ecdf` é possível encontrar, e plotar, F_e , a Função de Distribuição Acumulada Empírica, ou seja, as probabilidades acumuladas com base nos dados.

Veja o exemplo da amostra abaixo

```
xx>xx<-c(3,5,7,9,11,13,17)
xx>xx
```

```
xx[1] xx3xx5xx7xx9x11x13x17
```

O comportamento de F_e para este conjunto de dados x obtido da seguinte forma:

```
xx>x#xTamanhoxdaxamostra
xx>xn<-length(x)
xx>x#xFunçõxodexDistribuiçõxEmpírica
xx>xFe<-seq(1:n)/n
xx>xFe
```

```
xx[1] x0.1428571x0.2857143x0.4285714x0.5714286x0.7142857x0.8571429x1.0000000
```

O gráfico desta função x obtido por :


```
xx>x#xGrxficoxdaxFunxxoxdexDistribuixxoxAcumuladaxEmpxrica
xx>xplot(ecdf(x),ylab="ProbabilidadexAcumulada")
```

pict

```
xx>x#FunxxdexDistribuixxoxAcumuladaxEmpxricaxSuavizada
xx>xFe.s<-((seq(1:n)-0.5)/n)
xx>xFe.s<-c(0,Fe.s,1)
xx>xFe.s
```

```
xx[1]x0.00000000x0.07142857x0.21428571x0.35714286x0.50000000x0.64285714x0.78571429
xx[8]x0.92857143x1.00000000
```



```
xx>xplot(ecdf(x),ylab="ProbabilidadexAcumulada")
xx>xlines(c(min(x)-0.5,x,max(x)+0.5),Fe.s,col='red')
```

pict

Com a Função de Distribuição Acumulada Suavizada, podemos estimar qualquer quantil que for de interesse.

Exercícios

Os dados no arquivo a ser carregado contém o desempenho dos alunos no Processo Seletivo Estendido no ano de 2007.

```
xx>x#xLeitura do conjunto de dados
xx>xpse<-read.csv('http://www.leg.ufpr.br/~joel/dados/pse2007.csv')
xx>x#xExibir o conteúdo das 6 primeiras linhas
xx>xhead(pse)
```

```
xxxxXxCURS0xmatxbioxquixgeoxfisporxlitxhistxLEM
xx1x1xESTxNxxx0xxx0xxx5xxx3xxx2xx10xxx3xxxx1xxx5
xx2x2xESTxNxxx3xxx3xxx4xxx7xxx4xxx8xxx2xxxx4xxx1
xx3x3xESTxNxxx5xxx2xxx2xxx3xxx3xxx7xxx4xxxx2xxx2
xx4x4xESTxNxxx2xxx2xxx3xxx4xxx6xxx7xxx1xxxx4xxx2
xx5x5xESTxNxxx3xxx4xxx4xxx9xxx3xxx9xxx3xxxx7xxx8
xx6x6xESTxNxxx2xxx3xxx1xxx3xxx5xxx8xxx1xxxx5xxx1
```

```
xx>x#xExibir o conteúdo das 6 últimas linhas
xx>xtail(pse)
```

```
xxxxxxxxXxxxxxCURS0xmatxbioxquixgeoxfisporxlitxhistxLEM
xx459x459xMatxIndxTxxx2xxx3xxx3xxx2xxx4xxx3xxx1xxxx3xxx3
xx460x460xMatxIndxTxxx1xxx1xxx2xxx2xxx4xxx5xxx1xxxx3xxx2
xx461x461xMatxIndxTxxx3xxx1xxx3xxx3xxx8xxx4xxxx3xxx5
xx462x462xMatxIndxTxxx4xxx3xxx3xxx3xxx5xxx7xxx0xxxx4xxx2
xx463x463xMatxIndxTxxx3xxx4xxx3xxx6xxx6xxx8xxx2xxxx4xxx4
xx464x464xMatxIndxTxxx4xxx3xxx3xxx5xxx7xxx8xxx1xxxx2xxx3
```

Calcule as principais medidas resumo para estes dados

```
xx>x#xResumo das informações de todas as variáveis.
xx>xsummary(pse)
```



```
xx>x#xSelecionandoxtodasxasxlinhasxexasxcolunasxdex3xatex11
xx>xdiscip<-pse[,3:11]
xx>x#xCalculandoxmatrizdexcorrelacaotentrexasxnotas
xx>xcorrelacoes<-cor(discip)
xx>xcorrelacoes
```

```
xxxxxxxxxxxxxxxxmatxxxxxxxxbioxxxxxxxxquixxxxxxxxxgeoxxxxxxxxxfisxxxxxxxxpor
xxmatxx1.00000000x0.2579604xx0.13570161x0.2789203x0.34611361x0.16180803
xxbioxx0.25796043x1.0000000xx0.15974264x0.2527182x0.23427131x0.15794839
xxquixx0.13570161x0.1597426xx1.00000000x0.1237300x0.14197504x0.09533673
xxgeoxx0.27892025x0.2527182xx0.12372997x1.00000000x0.21084688x0.40538627
xxfisxx0.34611361x0.2342713xx0.14197504x0.2108469x1.00000000x0.13704096
xxporxx0.16180803x0.1579484xx0.09533673x0.4053863x0.13704096x1.00000000
xxlitxx0.07527067x0.1160417xx0.15972211x0.1640895x0.07056438x0.16296187
xxhistx0.13068988x0.2051098xx0.14710358x0.3606197x0.13940618x0.29867822
xxLEMxx0.15263761x0.1433847x-0.03162456x0.2150677x0.11234875x0.24979889
xxxxxxxxxxxxxxxxlitxxxxxxxxhistxxxxxxxxLEM
xxmatxx0.07527067x0.1306899xx0.15263761
xxbioxx0.11604169x0.2051098xx0.14338466
xxquixx0.15972211x0.1471036x-0.03162456
xxgeoxx0.16408949x0.3606197xx0.21506770
xxfisxx0.07056438x0.1394062xx0.11234875
xxporxx0.16296187x0.2986782xx0.24979889
xxlitxx1.00000000x0.2250333xx0.17445784
xxhistx0.22503329x1.0000000xx0.22367091
xxLEMxx0.17445784x0.2236709xx1.00000000
```