

Módulo VII: Estadística elemental

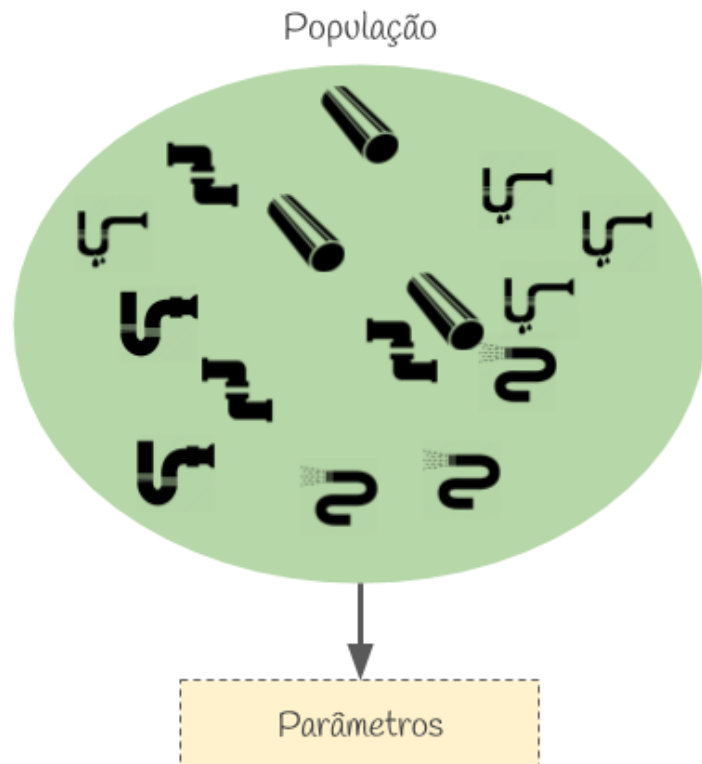
Revisão

A estatística é composta por um conjunto de métodos destinados à coleta, à apresentação, à análise e à interpretação de dados, tendo como objetivo compreender fenômenos de natureza aleatória e ajudar o usuário na tomada de decisão.

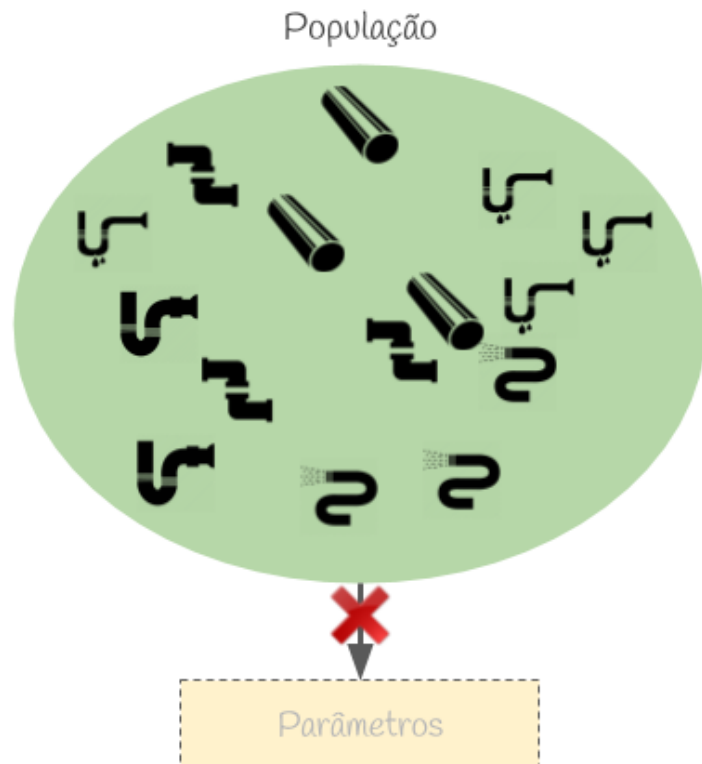
Palavras-chave

- **Variável:** O objeto/fenômeno/característica de interesse em estudo
- **População:** Conjunto de todas as observações da variável em estudo
- **Amostra:** Subconjunto de observações retiradas da população
- **Parâmetro:** Uma característica da população
- **Inferência:** Com base numa amostra inferimos os parâmetros associados à população

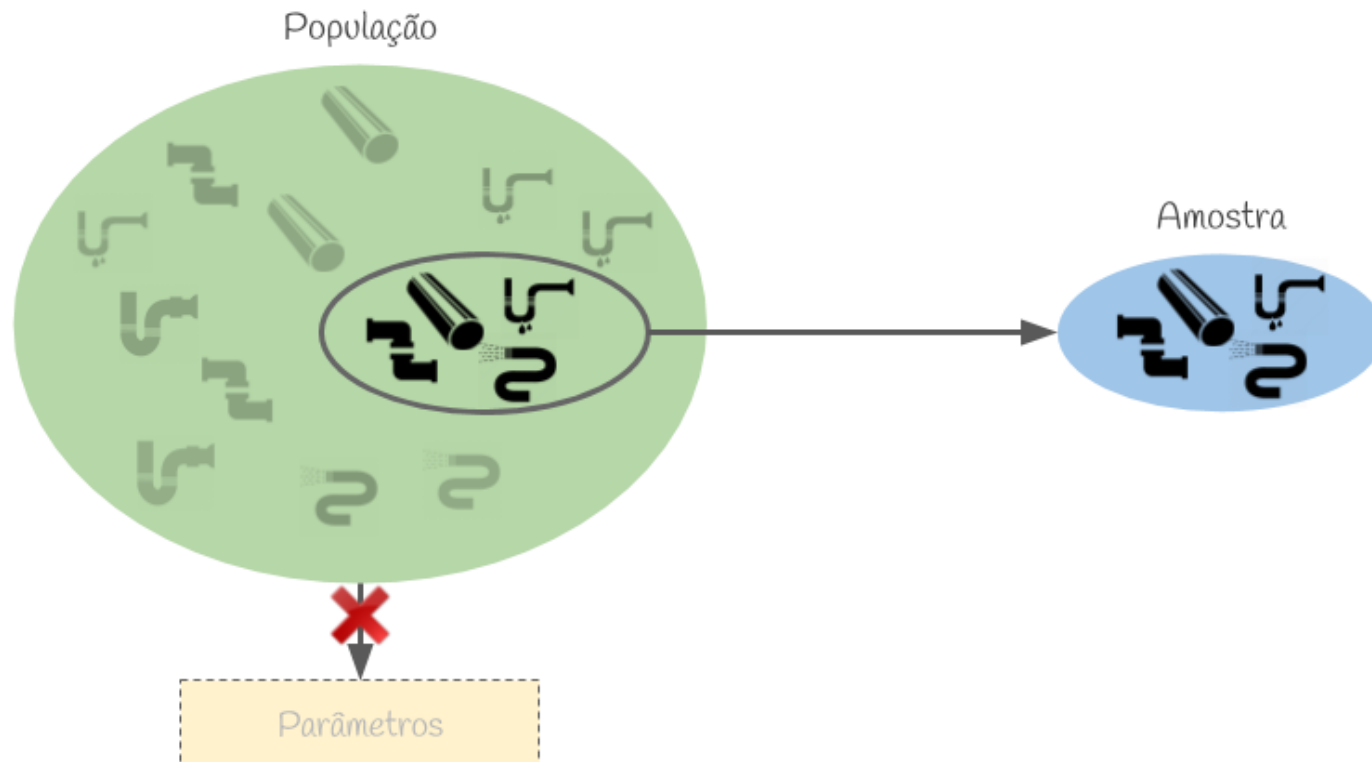
Estatística em um *slide*



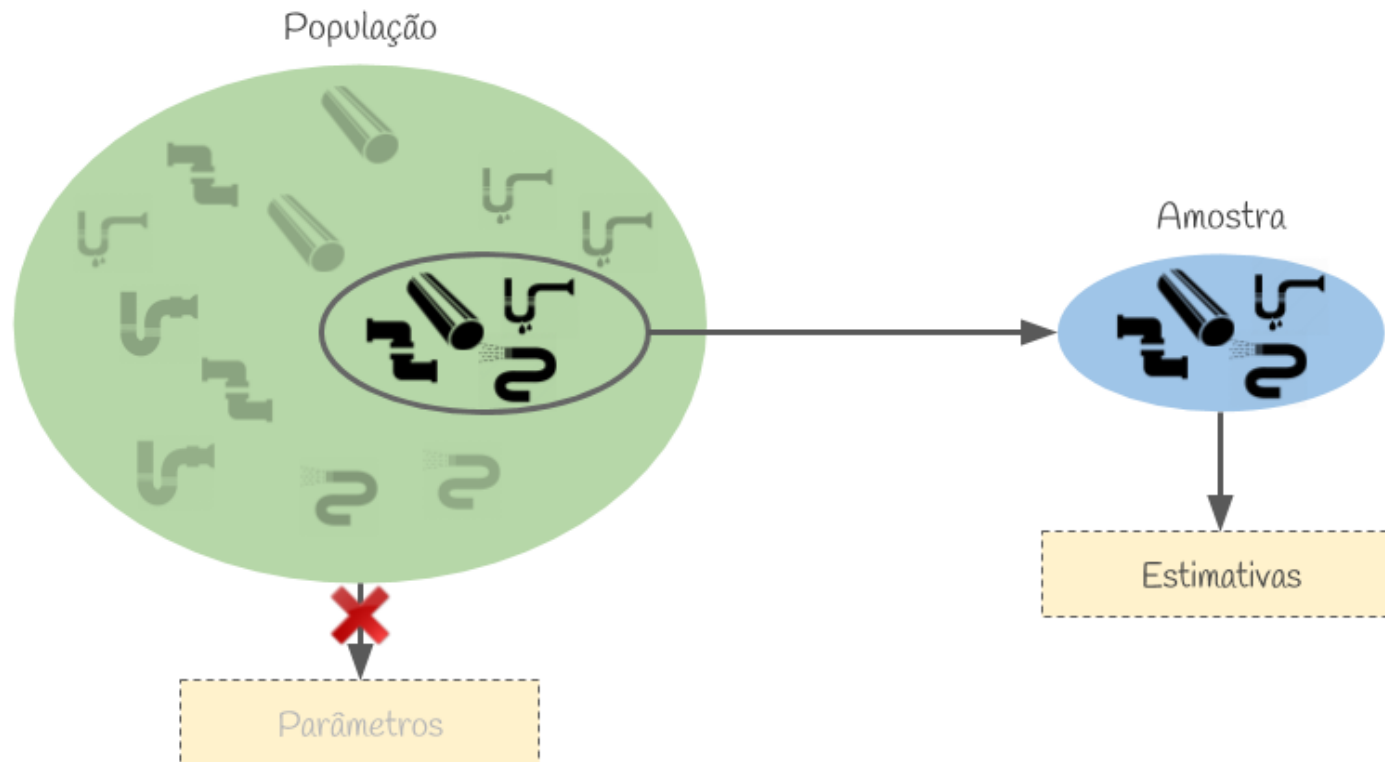
Estatística em um *slide*



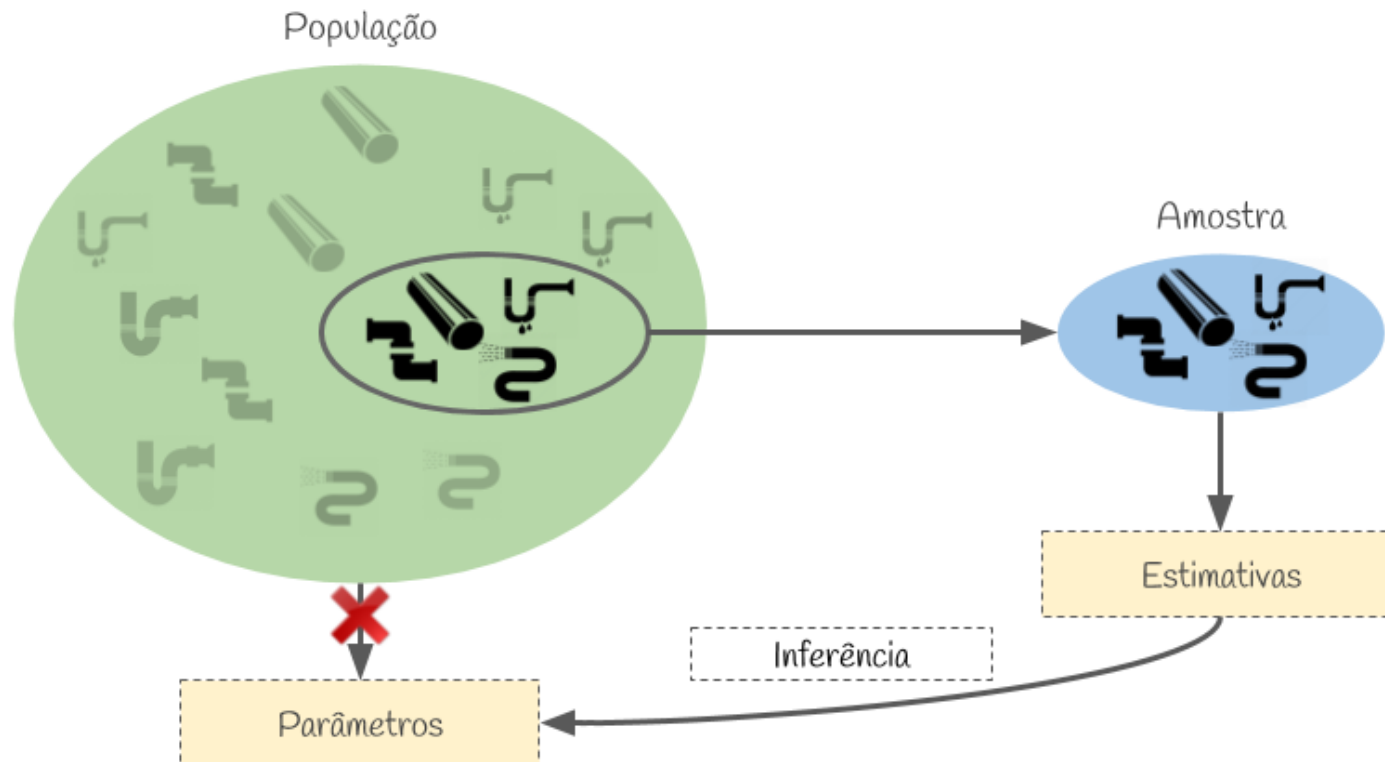
Estatística em um *slide*



Estatística em um *slide*



Estatística em um *slide*

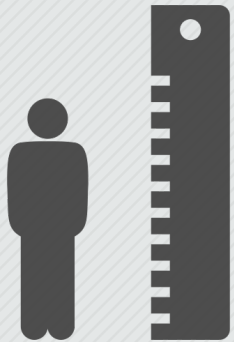


Tipos de variáveis

Variáveis quantitativas

Variáveis qualitativas

Contínuas



Discretas



Ordinais



Nominais



Para cada tipo de variável, diferentes tipos de análises (e gráficos) são indicados.

Estadística descriptiva

A estatística descritiva é uma área da estatística responsável por descrever e sumarizar um conjunto de dados.

Estatística descritiva

- **Medidas de tendência central:** Média, moda, mediana, ...
- **Medidas de variabilidade:** Variância, desvio padrão, amplitude, ...
- **Gráficos:** Boxplot, histograma, gráfico de barras, ...

Base para exemplificação

Vamos trabalhar com a base `dadosForno` e principalmente com as variáveis `PRS_GAS` e `PRS_FORNO`.

```
> load("/home/douglas/Documents/Dropbox/Cursos/Curso_R/data/Dados_VSB/dadosForno.RData")
> PRS_GAS <- dadosForno$PRS_GAS
> head(PRS_GAS)
[1] 100.1  99.5 100.1  99.9  99.5  99.9
```

```
> PRS_FORNO <- dadosForno$PRS_FORNO
> head(PRS_FORNO)
[1] 0.011 0.012 0.014 0.007 0.006 0.007
```

```
> miniBase <- cbind(PRS_FORNO,
+                   PRS_GAS)
```

Medidas de tendência central

As medidas de tendência central nos dão uma ideia de onde os dados estão concentrados.

Medidas de tendência central

- `mean`: Média de um vetor
- `median`: Mediana de um vetor
- `table`: Tabela de frequências/contingência
- `mlv`: Moda de um vetor
- `colSums`: Soma por colunas de base de dados
- `colMeans`: Média por coluna de uma base de dados

As medidas de variabilidade nos dão uma ideia da dispersão dos dados.

Medidas de variabilidade

- `min`: Menor valor de um vetor
- `max`: Maior valor de um vetor
- `which.min`: Posição em que o mínimo ocorre
- `which.max`: Posição em que o máximo ocorre
- `quantile`: Quantis de um vetor
- `var`: Variância de uma matriz/vetor
- `sd`: Desvio padrão de um vetor
- `cov`: Covariância entre duas variáveis
- `cor`: Correlação entre duas ou mais variáveis

Uma maneira simples de resumir os dados é utilizando a função `summary`

```
> summary(miniBase)
  PRS_FORNO          PRS_GAS
Min.      :-0.037000   Min.      :  0.00
1st Qu.:  0.007000   1st Qu.: 91.40
Median   :  0.010000   Median  : 95.30
Mean     :  0.009762   Mean     : 95.63
3rd Qu.:  0.012000   3rd Qu.:100.00
Max.     :  0.042000   Max.     :125.10
```

Agora é a sua vez

Análise descritiva

Utilizando a base de dados `dadosForno` encontre (quando possível):

- A média, mediana e variância
- O maior e o menor valor juntamente com a data em que entraram no forno
- Existem valores atípicos? O que ocorre com a média retirando 1
- Divida o desvio padrão pela média e multiplique por 100. O que podemos concluir?

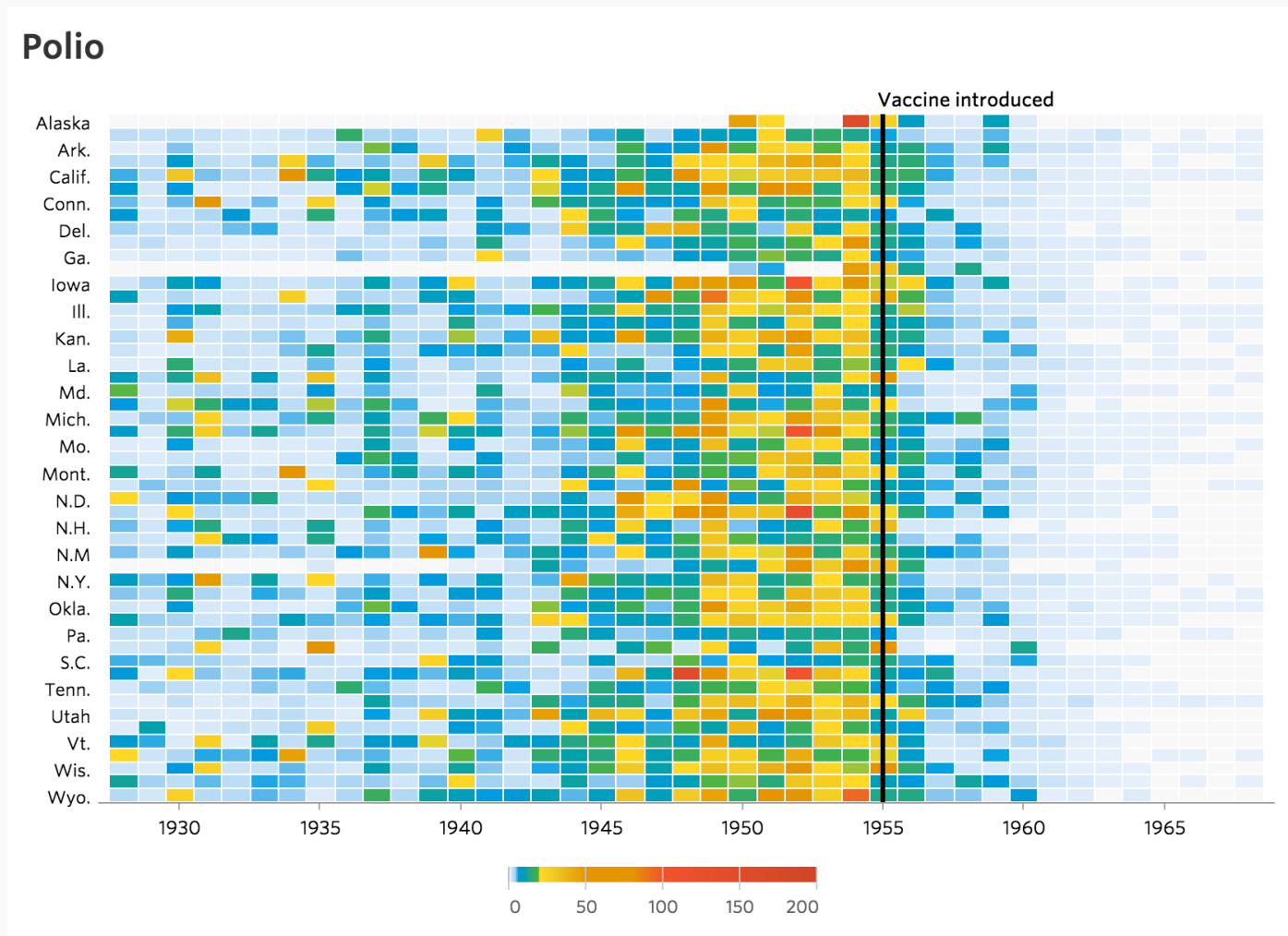
Utilize as variáveis:

- `DIS_ULT_BLOCO_ENFORNADO`,
- `CMP_BLOCO_ENFORNAMENTO` e
- `PER_02_CANAL_FUMOS`

Gráficos estadísticos

Gráficos estatísticos

Visualizar os dados ajuda no entendimento do problema e também na apresentação dos resultados. Sendo assim é parte **essencial** de uma análise.



Gráficos mostram de uma forma clara e objetiva o comportamento dos dados.

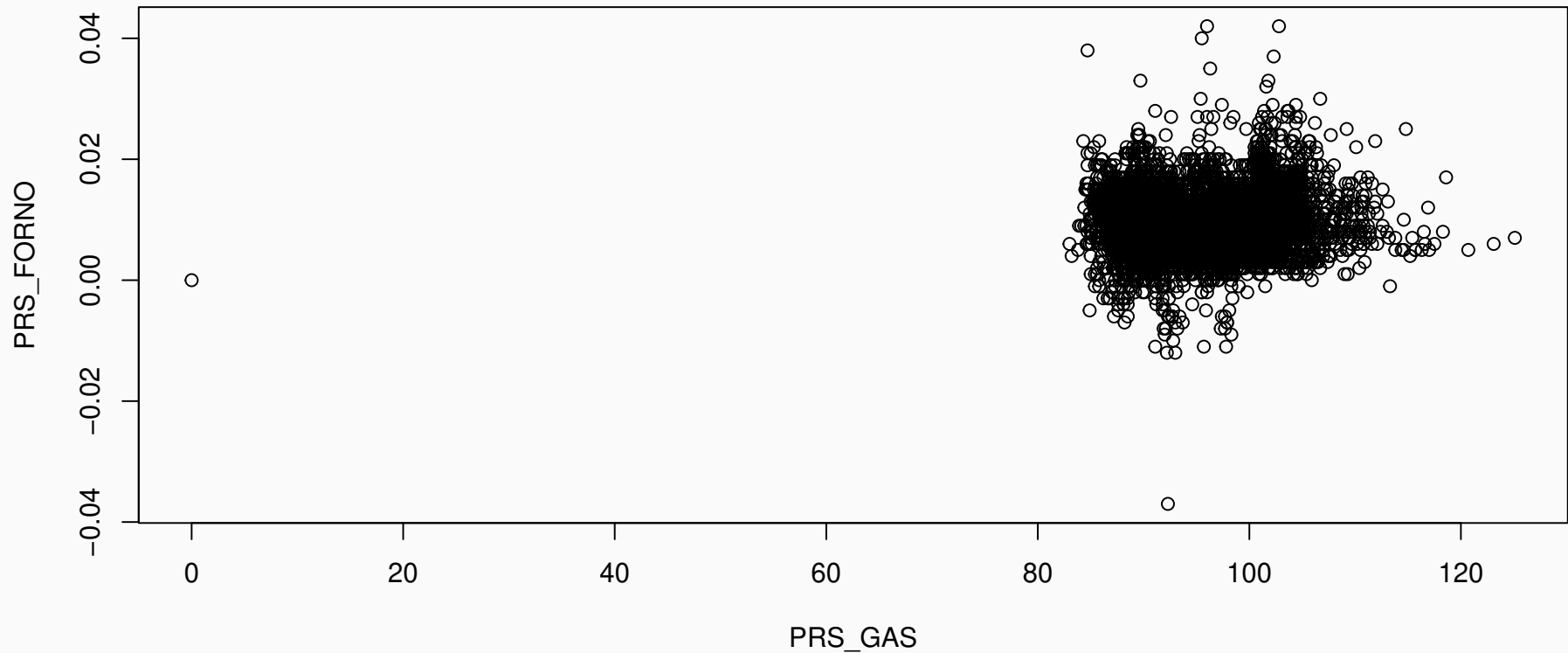
Gráficos estatísticos

- Dispersão
- Linhas
- Barras
- Setores/Pizza
- Histograma
- Boxplot

Veremos no próximo módulo como customizar estes e outros gráficos.

Gráfico de dispersão: Gráfico contrapondo duas variáveis que pode evidenciar correlação entre as variáveis.

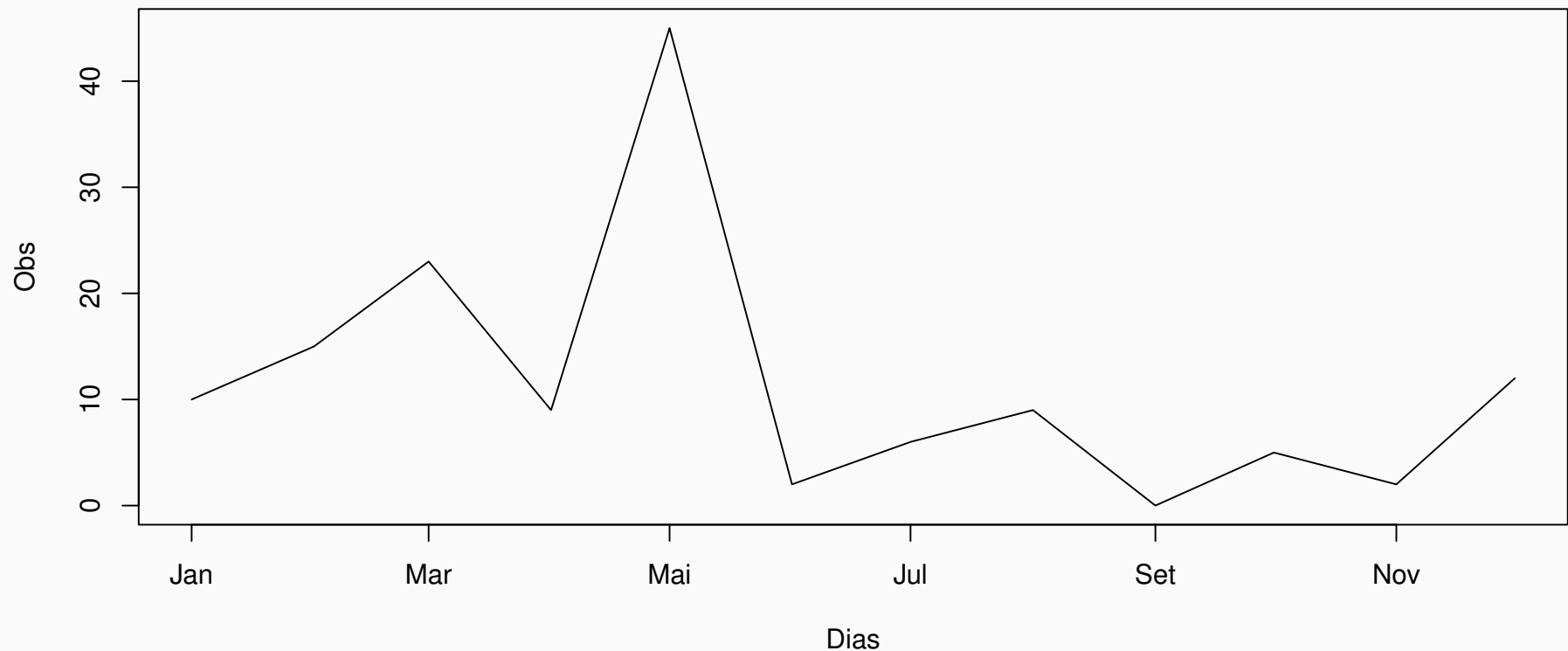
```
> plot(x = PRS_GAS, y = PRS_FORNO)
```



Gráficos estatísticos - Linhas

Gráfico de linhas: Gráfico contrapondo duas variáveis em que normalmente o eixo x apresenta uma variável temporal.

```
> Dias <- seq(as.Date("2016/01/01"), by = "month", length.out = 12)
> Obs <- c(10, 15, 23, 9, 45, 2, 6, 9, 0, 5, 2, 12)
> plot(x = Dias, y = Obs, type = "l")
```



Gráficos estatísticos - Barras

Gráfico de barras: É uma forma de resumir um conjunto de dados categóricos. Cada barra representa uma categoria particular e a altura de cada barra é proporcional à frequência de cada categoria.

```
> DIM_BLOCO_ENFORNAMENTO <- dadosForno$DIM_BLOCO_ENFORNAMENTO  
> cont <- table(DIM_BLOCO_ENFORNAMENTO)  
> barplot(height = cont)
```

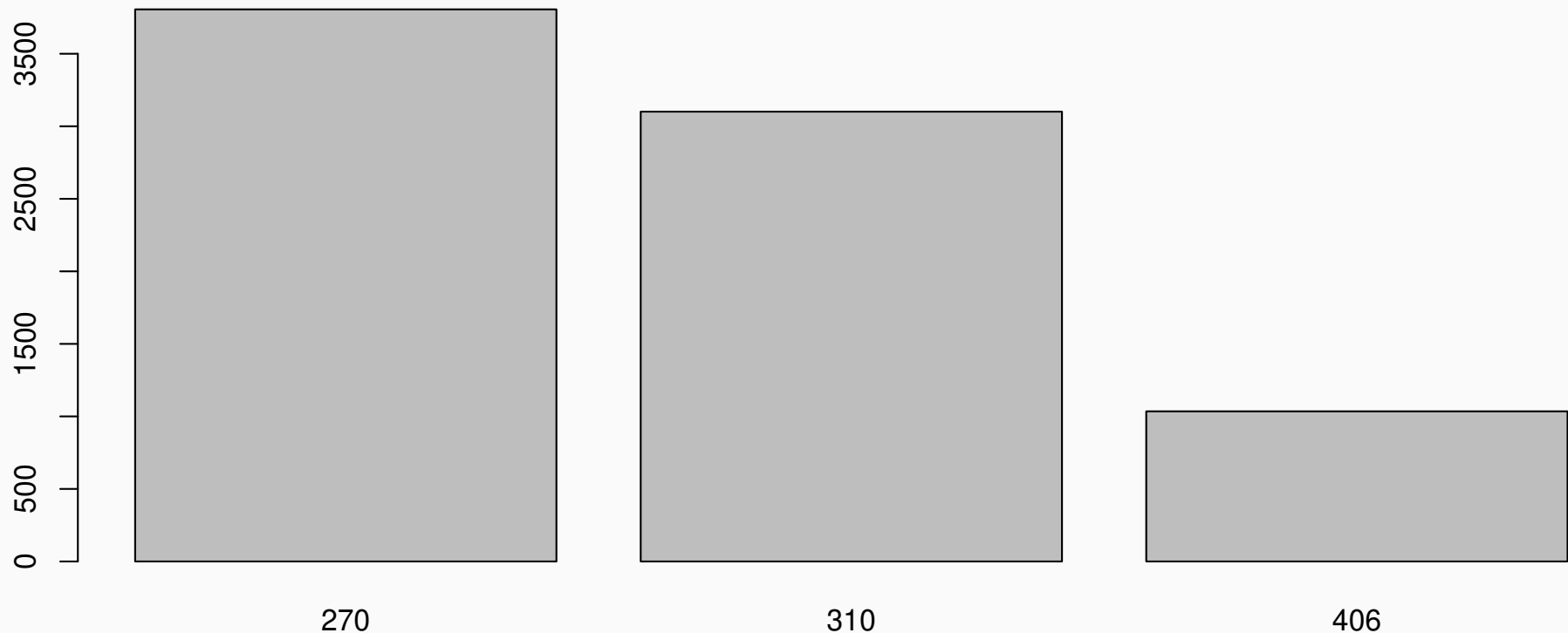
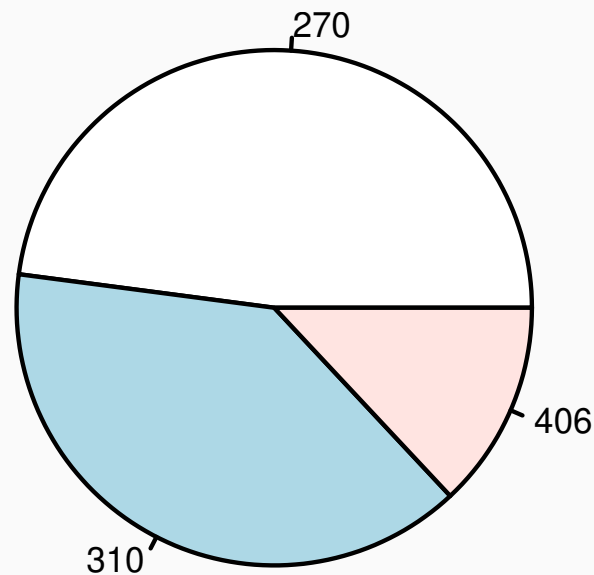


Gráfico de setores: Útil para visualizar a participação de cada categoria no total geral.

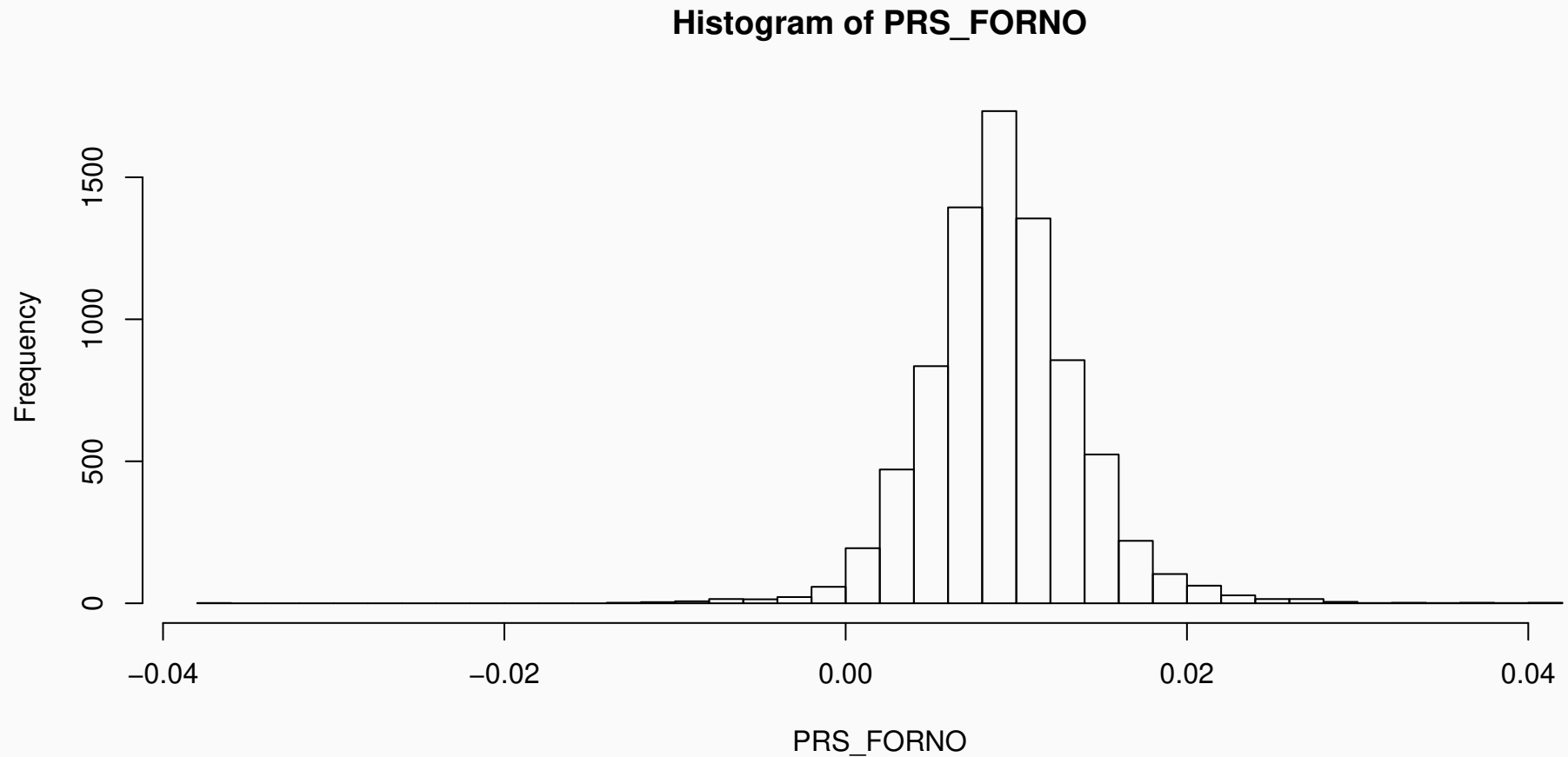
```
> pie(x = cont, cex = 0.5)
```



Gráficos estatísticos - Histograma

Histograma: Útil para caracterizar a distribuição dos dados.

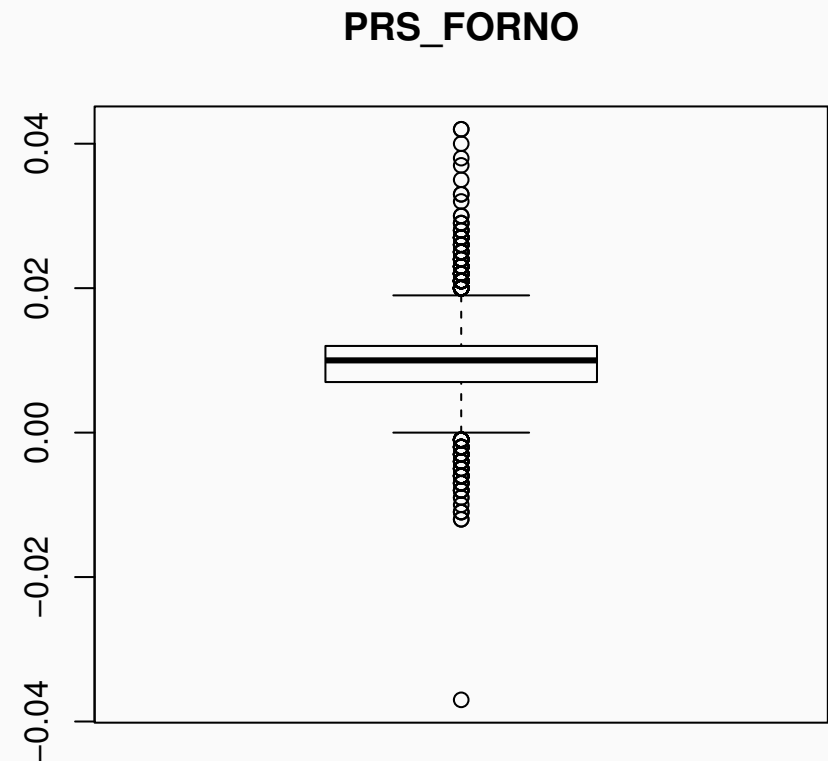
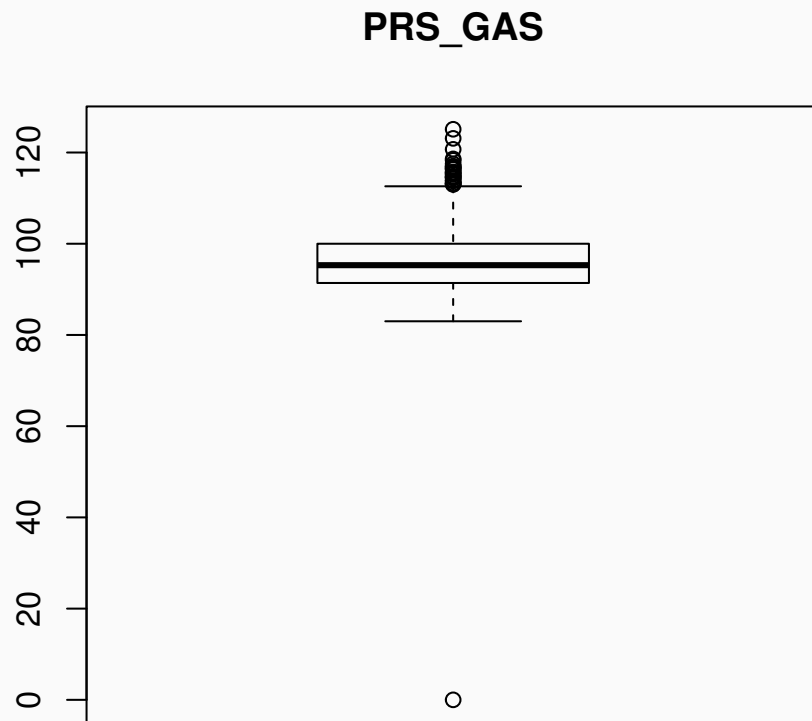
```
> hist(PRS_FORNO, breaks = 50)
```



Gráficos estatísticos - Boxplot

Boxplot: Útil para caracterizar a distribuição dos dados. Cada elemento do gráfico representa uma estatística e os pontos os *outliers*.

```
> par(mfrow=c(1,2))
> boxplot(PRS_GAS, main = "PRS_GAS")
> boxplot(PRS_FORNO, main = "PRS_FORNO")
```



Agora é a sua vez

Gráficos descritivos

Utilizando a base de dados `dadosForno` e as variáveis `DIS_ULT_BLOCO_ENFORNADO`, `CMP_BLOCO_ENFORNAMENTO` e `PER_02_CANAL_FUMOS`, faça:

- O gráfico de dispersão entre DIS e CMP
- O histograma e o boxplot das variáveis DIS e CMP
- O gráfico de barras e de setores da variável PER

Distribuições de probabilidade

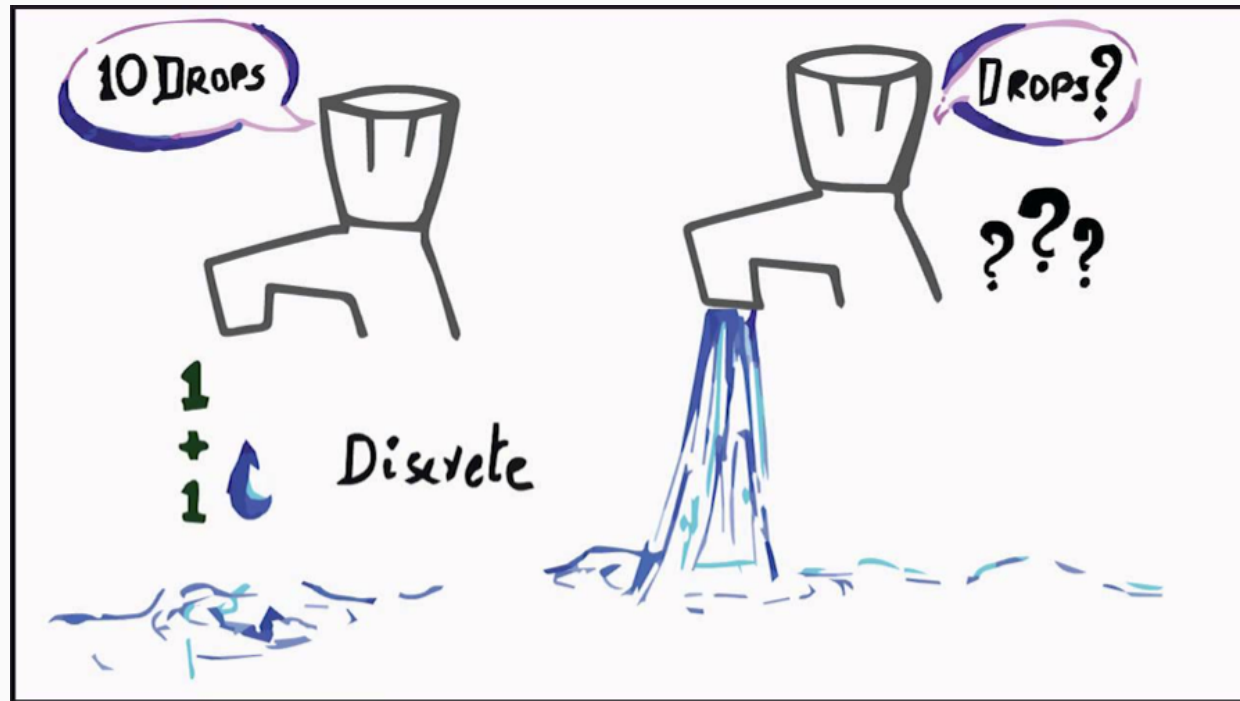
Distribuições de probabilidade

Distribuições contínuas

Assumem valores em uma escala contínua. Exemplo: Altura

Distribuições discretas

Assumem um número finito de valores. Exemplo: Número de filhos



Distribuições discretas

- `sample`: Uniforme
- `binom`: Binomial
- `multinom`: Multinomial
- `pois`: Poisson
- `geom`: Geométrica
- `nbinom`: Binomial negativa
- ...

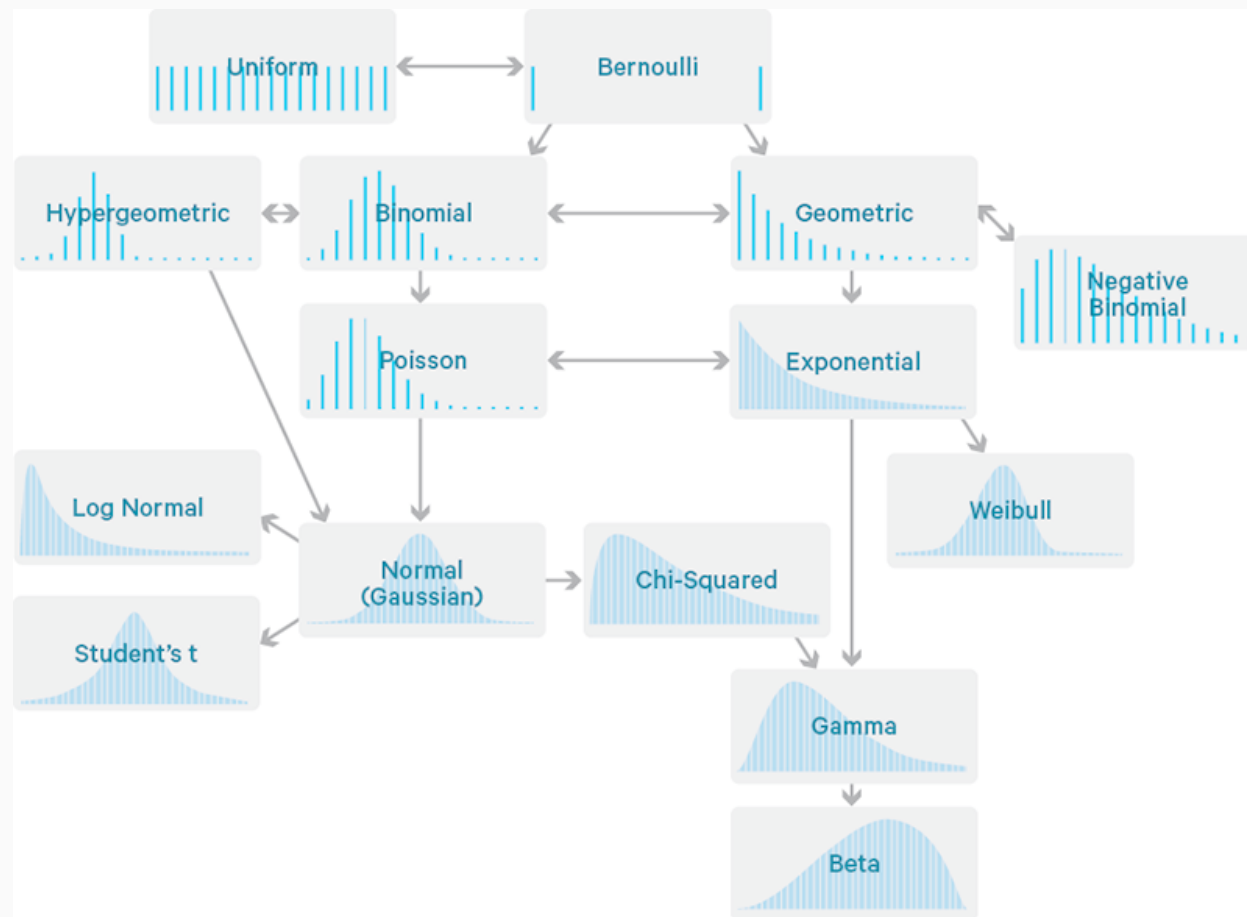
Distribuições contínuas

- `unif`: Uniforme
- `norm`: Normal
- `exp`: Exponencial
- `chisq`: Qui-quadrado
- `t`: t-student
- `f`: F de snedecor
- ...

O uso de uma distribuição de probabilidade depende da natureza da variável em análise. Para isso é interessante conhecer o suporte das distribuições de probabilidade.

Distribuições de probabilidade

Existem muitas (**Muitas mesmo**) distribuições de probabilidade. Na prática algumas delas são mais utilizadas (como as citadas no *slide* anterior).



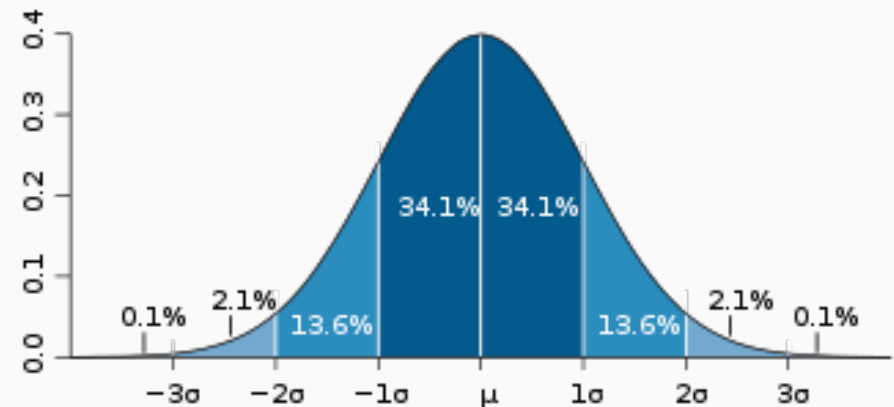
Distribuições de probabilidade

No R as distribuições são utilizadas com um dos prefixos abaixo.

Distribuições de probabilidade no R

- Prefixo **d**: Densidade
- Prefixo **q**: Quantil
- Prefixo **p**: Probabilidade
- Prefixo **r**: Números aleatórios

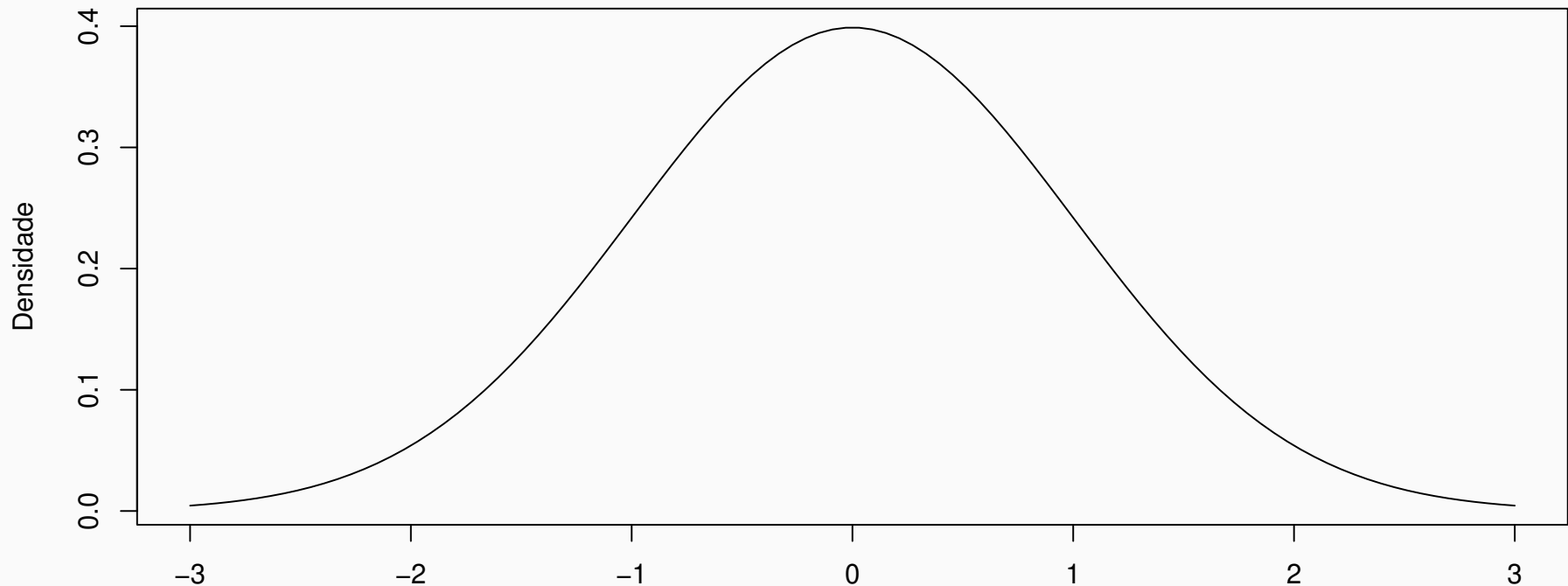
```
8 2 0 3 1 4 5 8 2 1 7 2 7 3 8 5 5 2 9 0 6 3 1 6 4
0 8 7 3 3 1 9 7 5 2 5 7 6 9 8 0 3 6 2 5 1 2 7 5 2
2 3 3 8 6 1 4 2 4 0 2 6 1 8 9 5 2 6 9 8 3 4 0 1 0
4 7 5 5 6 3 0 7 7 1 9 1 6 1 7 4 1 7 1 3 7 9 3 3 7
1 9 3 9 5 3 4 9 5 5 2 7 5 8 0 3 4 8 8 1 2 7 5 3 4
2 8 7 8 1 4 1 4 9 4 2 4 1 5 2 9 4 6 2 1 5 2 8 1 9
8 4 8 5 1 3 9 6 6 0 7 2 1 9 0 2 0 6 7 0 6 0 1 3 0
0 3 8 8 4 7 5 1 5 1 7 3 4 5 2 0 7 4 7 9 6 6 7 7 4
3 5 3 1 9 3 7 4 9 5 0 2 0 1 4 6 2 5 4 5 8 5 0 9 2
3 4 5 9 5 2 7 9 8 9 0 5 5 8 5 1 7 7 3 5 5 4 7 7 2
4 1 5 3 0 9 1 3 7 2 5 8 7 7 1 3 6 3 9 7 8 7 9 1 7
7 2 9 5 6 7 8 5 4 5 3 4 5 4 1 9 8 6 7 5 7 9 3 1 8
5 9 2 8 9 8 6 4 4 1 5 3 7 7 0 8 0 2 5 6 0 6 1 2 0
1 3 3 3 9 0 5 2 8 7 4 0 9 0 3 7 3 1 7 9 4 5 5 2 8
4 6 0 1 0 8 6 2 1 0 0 5 0 3 1 5 4 9 0 3 7 4 7 0 1
7 7 0 6 6 3 2 8 8 5 8 9 5 6 4 0 5 9 1 8 0 5 4 9 4
3 3 8 5 7 5 7 4 3 4 5 7 9 6 9 5 0 7 7 6 6 8 8 5 9
9 1 7 1 3 6 9 2 9 1 9 4 2 3 3 0 8 1 8 7 7 6 4 7 2
6 2 2 8 0 9 4 5 3 7 2 5 4 6 6 5 6 6 5 0 4 6 5 6 8
1 7 5 9 0 0 2 0 5 6 5 8 5 1 9 5 3 3 7 4 0 5 8 2 4
0 3 9 6 9 4 7 3 5 7 0 6 5 4 7 1 1 8 5 3 2 8 0 9 8
```



Distribuições de probabilidade - O prefixo d

O prefixo 'd' seguido do nome de uma distribuição corresponde à **densidade** da variável.

```
> x <- seq(-3, 3, length.out = 100)
> densidade <- dnorm(x = x, mean = 0, sd = 1)
> plot(x = x, y = densidade, type = "l",
+       xlab = "", ylab = "Densidade")
```



O prefixo 'q' seguido do nome de uma distribuição corresponde ao **quantil** da variável.

```
> quantil1 <- qnorm(p = 0.025, mean = 0, sd = 1)
> quantil1
[1] -1.959964
```

```
> quantil2 <- qnorm(p = 0.975, mean = 0, sd = 1)
> quantil2
[1] 1.959964
```

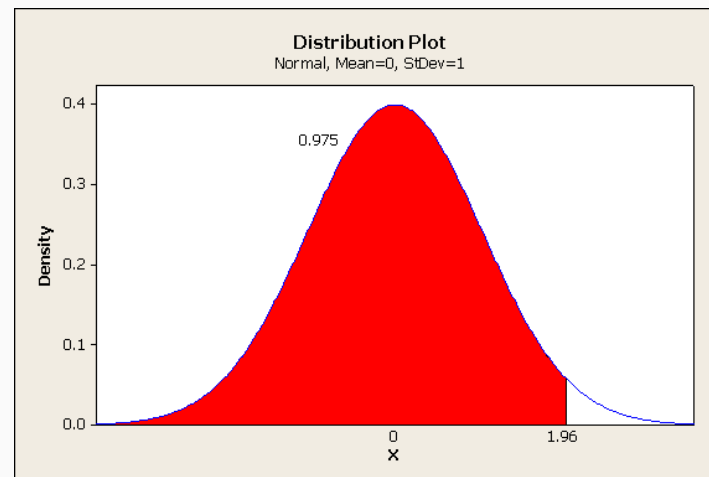
Ou seja, em uma distribuição Normal padrão, a probabilidade acumulada em $(-\infty, -1.96] \cup [1.96, \infty)$ é de 0.05.

Distribuições de probabilidade - O prefixo p

O prefixo 'p' seguido do nome de uma distribuição corresponde à **probabilidade acumulada** até determinado quantil.

```
> prob1 <- pnorm(q = quantil1, mean = 0, sd = 1)
> prob1
[1] 0.025
```

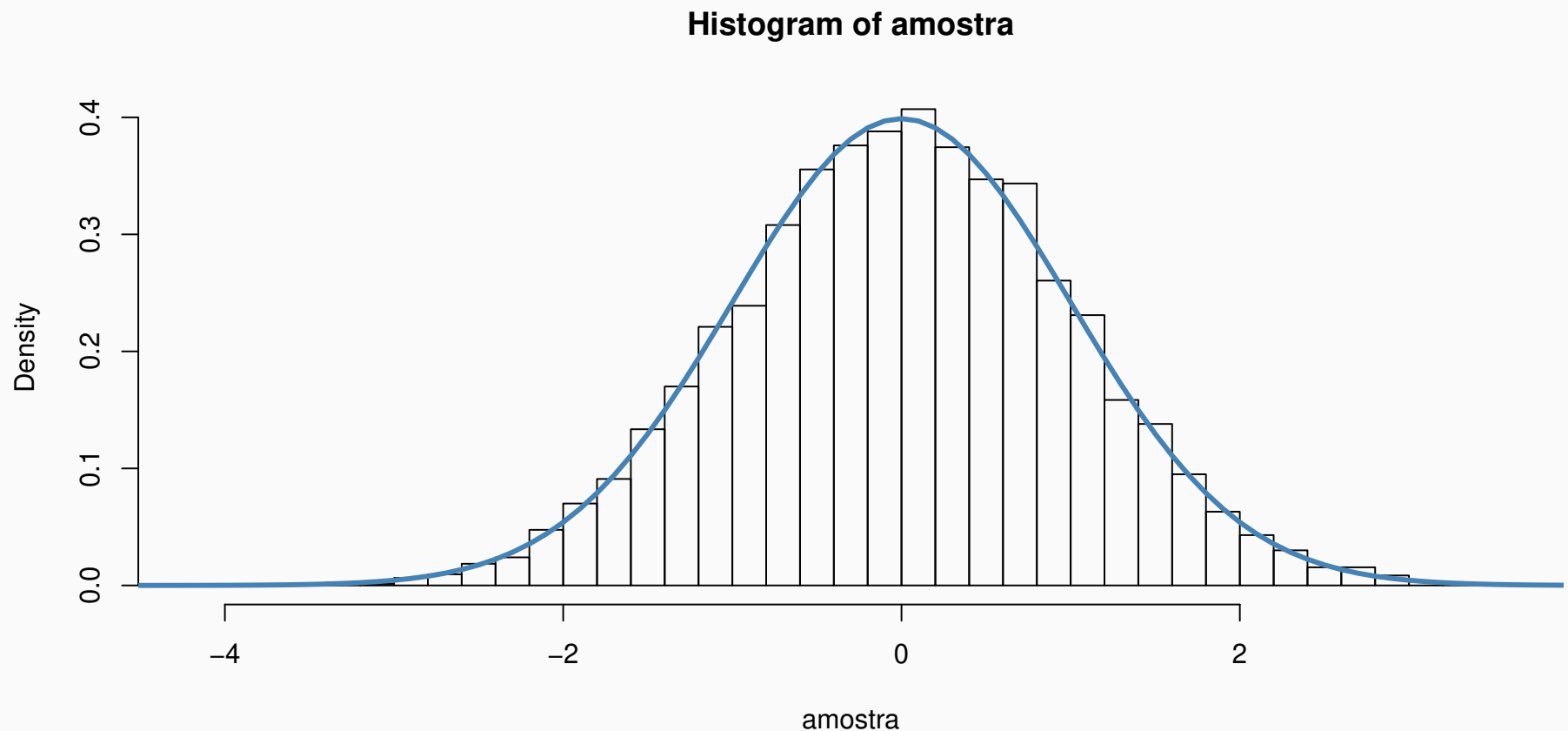
```
> prob2 <- pnorm(q = quantil2, mean = 0, sd = 1)
> prob2
[1] 0.975
```



Distribuições de probabilidade - O prefixo r

O prefixo 'r' seguido do nome de uma distribuição permite a geração de **números aleatórios**.

```
> amostra <- rnorm(n = 10000, mean = 0, sd = 1)
> hist(amostra, probability = T, breaks = 30)
> curve(dnorm, from = -5, to = 5, add = T,
+       col = "steelblue", lwd = 3)
```



Agora é a sua vez

Gráficos de uma distribuição

Utilizando os prefixos 'd', 'q', 'p' e 'r' e se valendo dos nomes das distribuições vistas no slide Distribuições de probabilidade:

- Gere 1.000 observações de duas distribuições à sua escolha
- Faça o histograma e a curva da distribuição (ao mesmo tempo)
- Encontre o quantil 0.025 e o quantil 0.975 destas distribuições
- Qual a proporção de observações maiores e menores que estes quantis (2.5
- Qual a probabilidade de observarmos valores entre -1 e 1?
- Qual a proporção encontrada na nossa amostra?