

Instituto Federal de Educação, Ciência e Tecnologia do Sul de Minas
Campus - Poços de Caldas / MG

2º Trabalho Prático - Pré-Processamento e Análise de Dados

Maria Clara Sanchez de Mira - 201911020026

1 Identificação do atributo alvo (saída);

O atributo de saída da minha base de dados é o **Class**, pois ela retorna se o paciente tem ou não recorrência de eventos (no-recurrence-events, recurrence-events), o que leva a a conclusão de que essa é a classe de atributo alvo.

```
1 library(readr)
2
3 base <- read.csv("breast-cancer.csv", stringsAsFactors = FALSE)
4
5 atributo_alvo <- base$class
6 table(atributo_alvo)
```

Listing 1: Código fonte em R

Quando se utiliza o comando da linha 6 do código acima é possível visualizar a quantidade de instâncias do atributo alvo, e da um total de 70.3%: no-recurrence-events e 29.7%: recurrence-events. O que representa uma base de dados desbalanceada. E uma base de dados desbalanceada significa que o número de valores varia mais em determinados atributos.

Porém dados desbalanceados, podem ocasionar vários problemas se não forem tratados de forma correta, pode acontecer de ter **Overfitting**, ou seja quando se tem vários dados de um valor apenas, com isso a hipótese vai apenas memorizar os dados e se especializar nos dados de treinamento, ou também ter **Underfitting**, ou seja essa situação pode ocorrer, quando os exemplos de treinamento disponíveis são pouco representativos ou o modelo usado é muito simples e não captura os padrões existentes nos dados.

2 Identificação dos tipos de dados dos atributos de entrada (quantitativo, qualitativo);

Todos os dados dessa base de dados são classificados como qualitativos/simbólicos pois as instâncias de cada atributo, assumem valores que são associados a categoria. Pois mesmo tendo valores que são números na base, por exemplo o atributo *age*, que segue um padrão ou seja

uma ordem de idades que varia de *20-29*, *30-39*, *40-49*, *50-59*, *60-69*, *70-79*. É possível ver os atributos e suas respectivas classificações na tabela abaixo:

Atributos	Classificação
class	Qualitativo
age	Qualitativo
menopause	Qualitativo
tumor.size	Qualitativo
inv.nodes	Qualitativo
node.caps	Qualitativo
deg.malig	Qualitativo
breast	Qualitativo
breast.quad	Qualitativo
irradiat	Qualitativo

```

1 #Ver um por um
2
3 unique(base$class)
4 unique(base$age)
5 unique(base$menopause)
6 unique(base$tumor.size)
7 unique(base$inv.nodes)
8 unique(base$node.caps)
9 unique(base$deg.malig)
10 unique(base$breast)
11 unique(base$breast.quad)
12 unique(base$irradiat)

```

Listing 2: Código fonte em R

O código acima retorna os valores que se tem dentro de cada atributo, por exemplo no atributo *irradiat*, ele retorna *yes* ou *no* e com isso é possível montar a tabela acima também com qual tipo de cada atributo.

3 Identificação da escala de dados dos atributos de entrada (nominal, ordinal, intervalar, racional);

Na escala, como todos os dados são do tipo qualitativos, a escala vai poder assumir apenas dois valores ou **Nominal** ou **Ordinal**, pois essas duas escalas estão relacionadas ao tipo Qualitativo.

Na escala nominal, os valores assumem apenas nomes que são diferentes, ou seja trazem quase nada de informações e não existe ordem entre seus valores. Por exemplo o atributo **class**

que assume apenas dois valores ou é *no-recurrence-events* ou *recurrence-events*, ou seja ou o valor é uma coisa ou ele é outra, não tem-se relação de ordem, que um deva vir primeiro que o outro.

Na escala ordinal, os valores seguem uma certa ordem, pois mesmo os valores sendo simbólicos é possível definir uma certa ordem entre eles, por exemplo o atributo **deg.malig** que é o grau de malignidade do tumor que pode assumir 1, 2 ou 3 que é a tendência do câncer se tornar progressivamente piores e potencialmente causar a morte de uma pessoa.

Na tabela abaixo é possível ver melhor cada atributo com sua determinada classificação.

Atributos	Classificação
class	Nominal
age	Ordinal
menopause	Ordinal
tumor.size	Ordinal
inv.nodes	Ordinal
node.caps	Nominal
deg.malig	Ordinal
breast	Nominal
breast.quad	Nominal
irradiat	Nominal

Com o código abaixo, também é possível a visualização das categorias que se tem em cada atributo e com isso foi possível a determinação do tipo Nominal ou Ordinal.

```

1 library('dplyr')
2
3 #Ver como um todo o tipo da variavel
4 glimpse(base)
5
6 #Utilizando a funcao table, para encontrar as cateogiras presente
7 table(base$class)
8 table(base$age)
9 table(base$menopause)
10 table(base$tumor.size)
11 table(base$inv.nodes)
12 table(base$node.caps)
13 table(base$deg.malig)
14 table(base$breast)
15 table(base$breast.quad)
16 table(base$irradiat)

```

Listing 3: Código fonte em R

O retorno da linha 4 do código acima quando usando a função `glimpse(base)` é:

```

> glimpse(base)
Rows: 286
Columns: 10
$ class      <chr> "no-recurrence-events", "no-recurrence-events", "no-recurrence-events"...
$ age        <chr> "30-39", "40-49", "40-49", "60-69", "40-49", "60-69", "50-59", "60-69"...
$ menopause  <chr> "premeno", "premeno", "premeno", "ge40", "premeno", "ge40", "premeno", ...
$ tumor.size <chr> "30-34", "20-24", "20-24", "15-19", "0-4", "15-19", "25-29", "20-24", ...
$ inv.nodes  <chr> "0-2", "0-2", "0-2", "0-2", "0-2", "0-2", "0-2", "0-2", "0-2", "0-2", ...
$ node.caps  <chr> "no", "no", "no", "no", "no", "no", "no", "no", "no", "no", "no", "no"...
$ deg.malig  <int> 3, 2, 2, 2, 2, 2, 1, 2, 2, 3, 2, 1, 3, 3, 1, 2, 3, 1, 2, 2, 2, 2...
$ breast     <chr> "left", "right", "left", "right", "right", "left", "left", "left", "le...
$ breast.qua <chr> "left_low", "right_up", "left_low", "left_up", "right_low", "left_low"...
$ irradiat   <chr> "no", "no", "no", "no", "no", "no", "no", "no", "no", "no", "no", "no"...

```

Figura 1: Tipo de cada atributo

É possível reparar que o atributo `deg.malig` é classificado como tipo inteiro, mas mesmo ele sendo do tipo inteiro, como o câncer segue um padrão de grau 1,2 e 3, como já explicado, esse tipo acaba se tornando um categórico ordinal, e por isso também vai ocorrer a conversão desse dado.

4 Exploração dos dados através de medidas de localidade;

As medidas apontam os pontos de direção nos dados e sofrem uma variação de acordo com o tipo de dado. Por conta da base ter todos os atributos do tipo qualitativo, a única medida de localidade que se aplica é a moda.

E na tabela abaixo é possível ver cada atributo com sua determinada moda, ou seja o valor que mais se repete em cada atributo.

Atributos	Moda
class	no-recurrence-events
age	50-59
menopause	premeno
tumor.size	30-34
inv.nodes	0-2
node.caps	no
deg.malig	2
breast	left
breast.quad	left-low
irradiat	no

```

1 # Funcao que encontra a moda
2 find_mode <- function(x) {
3   u <- unique(x)
4   tab <- tabulate(match(x, u))
5   u[tab == max(tab)]
6 }

```

```
7 #Encontrando a moda de cada atributo
8 find_mode(base$class)
9 find_mode(base$age)
10 find_mode(base$menopause)
11 find_mode(base$tumor.size)
12 find_mode(base$inv.nodes)
13 find_mode(base$node.caps)
14 find_mode(base$deg.malig)
15 find_mode(base$breast)
16 find_mode(base$breast.quad)
17 find_mode(base$irradiat)
```

Listing 4: Código fonte em R

Porém depois de realizar o passo 13 letra A (*Conversão de tipos (simbólico para numérico, ordinal para numérico, nominal para binário, numérico para ordinal)*) é possível a criação de gráficos para melhor visualização da distribuição dos dados e nesse caso é possível ver essa distribuição através do boxplot.

Vale salientar que para a visualização dos dados do boxplot foi utilizada toda a base de dados, e a sua divisão em conjunto de teste e treinamento vem subsequentemente.

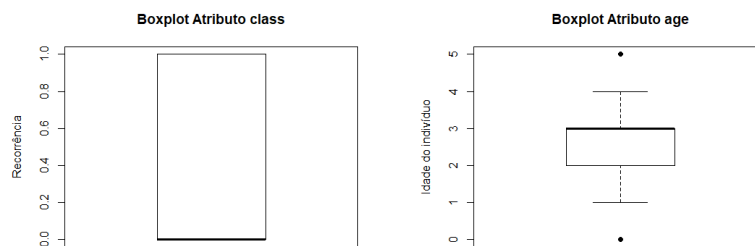


Figura 2: Boxplot atributo class e age

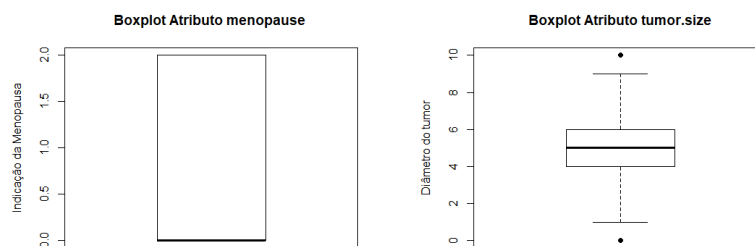


Figura 3: Boxplot atributo menopause e tumor.size

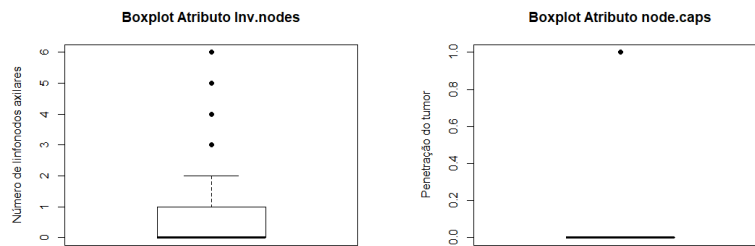


Figura 4: Boxplot atributo inv.nodes e node.caps

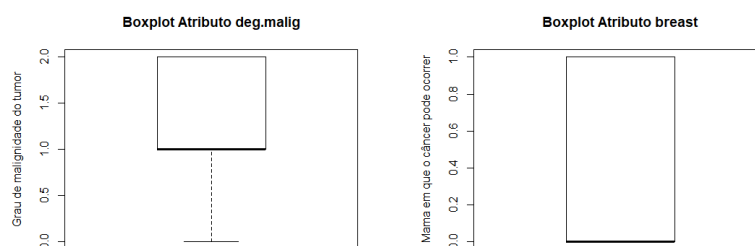


Figura 5: Boxplot atributo deg.malig e breast

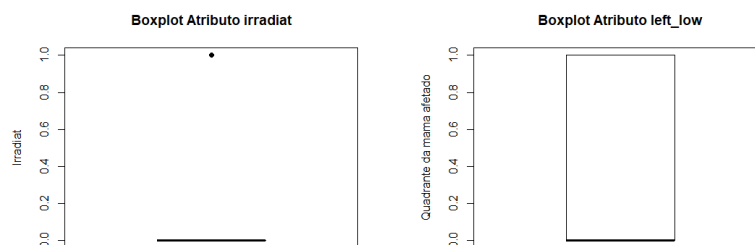


Figura 6: Boxplot atributo irradiat e left-low

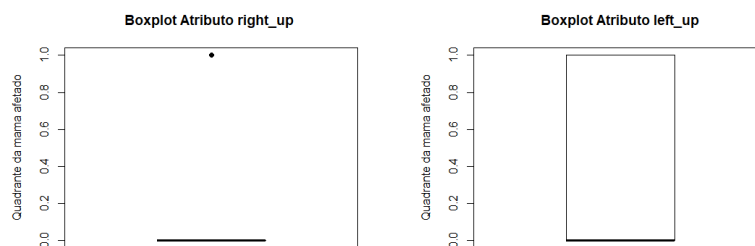


Figura 7: Boxplot atributo right-up e left-up

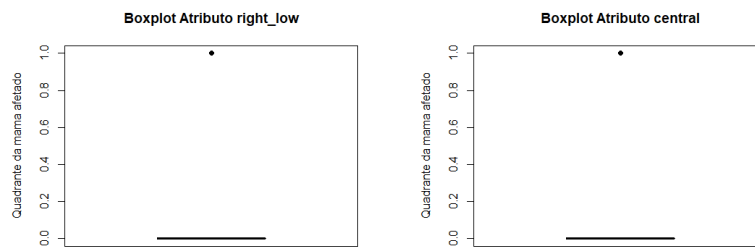


Figura 8: Boxplot atributo right-low e central

Após a visualização dos boxplots é possível notar a presença de alguns outliers em certos atributos, tais como *age*, *tumor.size* e *inv.nodes*, isso ocorre por conta de alguns valores saírem fora da normalidade dos dados e por conta disso acontece as anomalias no gráfico.

Os gráficos *node.caps*, *irradiat*, *right-up*, *right-low* e *center* apesar de apresentarem bolinhas não são considerados outliers, pois os atributos assumem apenas 2 valores, 0 ou 1, se por acaso um sair a mais que o outro, vai parecer que os dados tem uma anomalia.

Com o código abaixo, foi possível visualizar a criação de cada boxplot que está representado bem como nas figuras acima.

```

1 #Criacao dos boxplots
2 boxplot(base$class ,
3         main = "Boxplot Atributo class",
4         ylab = "Recorrencia",
5         pch = 16)
6
7 boxplot(base$age ,
8         main = "Boxplot Atributo age",
9         ylab = "Idade do individuo",
10        pch = 16)
11
12 boxplot(base$menopause ,
13         main = "Boxplot Atributo menopause ",
14         ylab = "Indicacao da Menopausa",
15         pch = 16)
16
17 boxplot(base$tumor.size ,
18         main = "Boxplot Atributo tumor.size",
19         ylab = "Diametro do tumor",
20         pch = 16)
21
22 boxplot(base$inv.nodes ,
23         main = "Boxplot Atributo Inv.nodes",
24         ylab = "Numero de linfonodos axilares",
25         pch = 16)

```

```
26 boxplot(base$node.caps ,
27         main = "Boxplot Atributo node.caps",
28         ylab = "Penetracao do tumor",
29         pch = 16)
30
31 boxplot(base$deg.malig ,
32         main = "Boxplot Atributo deg.malig",
33         ylab = "Grau de malignidade do tumor",
34         pch = 16)
35
36 boxplot(base$breast ,
37         main = "Boxplot Atributo breast",
38         ylab = "Mama em que o cancer pode ocorrer",
39         pch = 16)
40
41 boxplot(base$irradiat ,
42         main = "Boxplot Atributo irradiat",
43         ylab = "Irradiat",
44         pch = 16)
45
46 boxplot(base$left_low ,
47         main = "Boxplot Atributo left-low",
48         ylab = "Quadrante da mama afetado",
49         pch = 16)
50
51 boxplot(base$right_up ,
52         main = "Boxplot Atributo right-up ",
53         ylab = "Quadrante da mama afetado",
54         pch = 16)
55
56 boxplot(base$left_up ,
57         main = "Boxplot Atributo left-up",
58         ylab = "Quadrante da mama afetado",
59         pch = 16)
60
61 boxplot(base$right_low ,
62         main = "Boxplot Atributo right-low",
63         ylab = "Quadrante da mama afetado",
64         pch = 16)
65
66 boxplot(base$central ,
67         main = "Boxplot Atributo central",
68         ylab = "Quadrante da mama afetado",
69         pch = 16)
```

Listing 5: Código fonte em R

5 Exploração dos dados através de medidas de espalhamento;

As medidas de espalhamento são responsáveis pela medição da dispersão de um conjunto de valores. E com isso possibilita a visualização dos dados, ou seja se eles estão muito espalhados ou relativamente concentrados em torno da média por exemplo.

Com isso eu construí uma tabela simples, com a média de cada atributo para poder visualizar melhor como os gráficos se comportam em relação a média.

Atributos	Média
class	0.29720280
age	2.66136486
menopause	0.92657343
tumor.size	4.88111888
inv.nodes	0.51748252
node.caps	0.19580420
deg.malig	1.04895105
breast	0.46853147
left_low	0.38811189
right_up	0.11538462
left_up	0.33916084
right_low	0.08391608
central	0.07342657
irradiat	0.23776224

Figura 9: Tabela com as média dos atributos

Abaixo o histograma de cada atributo com sua respectiva frequência, ou seja, número de objetos que um certo atributo tem.

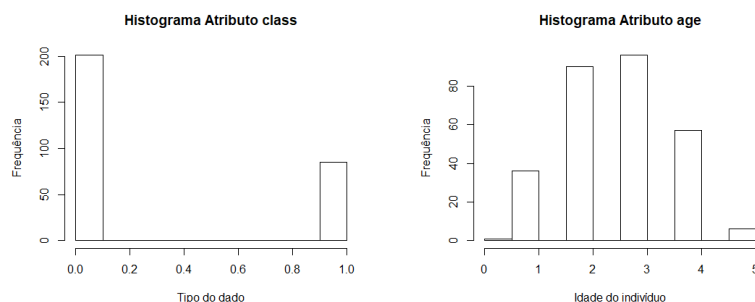


Figura 10: Histograma atributo class e age

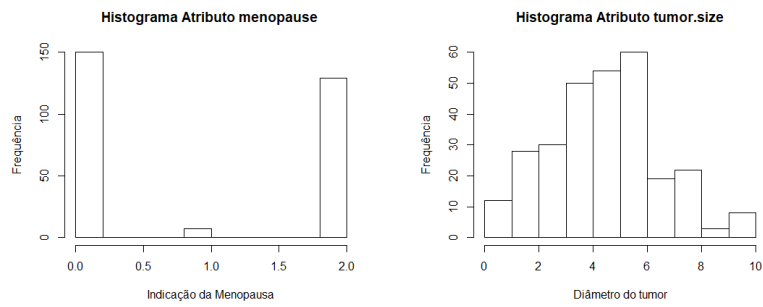


Figura 11: Histograma atributo menopause e tumor.size

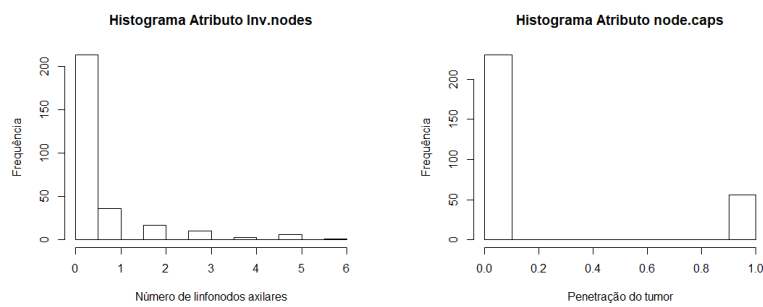


Figura 12: Histograma atributo inv.nodes e node.caps

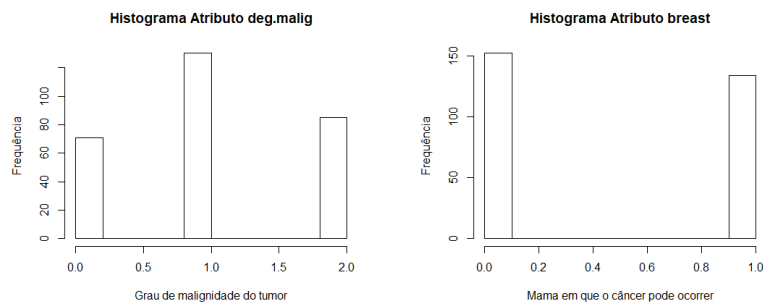


Figura 13: Histograma atributo deg.malig e breast

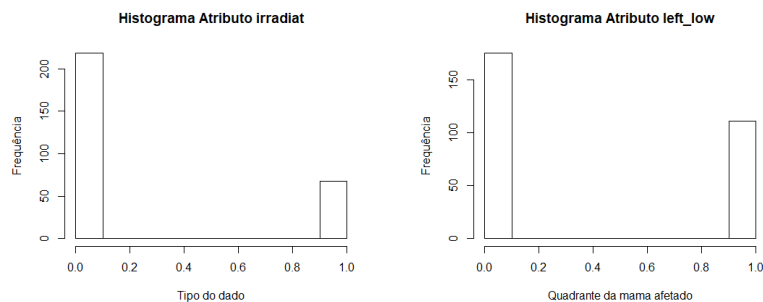


Figura 14: Histograma atributo irradiat e left-low

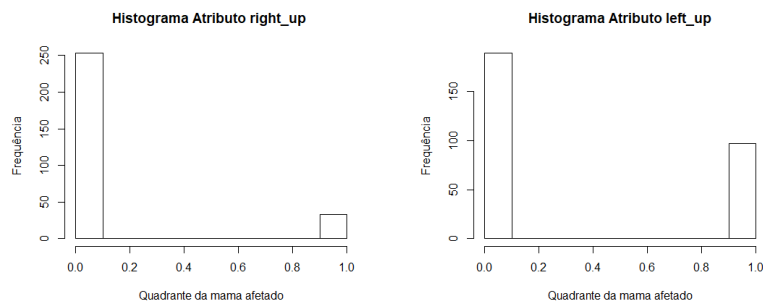


Figura 15: Histograma atributo right-up e left-up

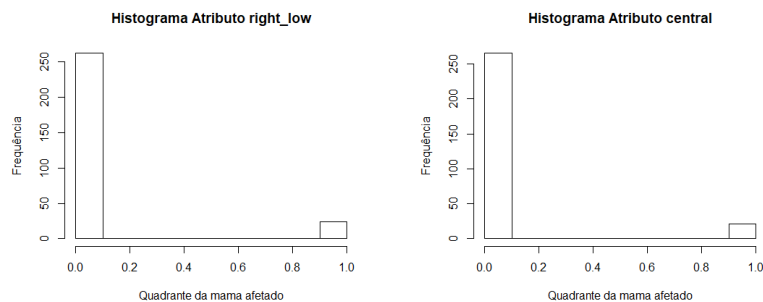


Figura 16: Histograma atributo right-low e central

Após a análise de todos os gráficos e comparando suas respectivas médias, é muito notório ver que os gráficos estão sempre concentrados em cima da média. E que os dados não seguem uma distribuição normal, então posteriormente esses dados terão que passar pela re-escala.

A média e a criação da tabela da figura 9, foi construída a partir do código fonte abaixo. Nele eu crio um novo data frame e atribuo os valores para eles, depois é importado uma biblioteca e depois eu exporto a imagem obtida.

```

1 #Tabela com a media de cada atributo
2 media <- data.frame(
3   "Atributos" = c("class", "age", "menopause", "tumor.size",
4     "inv.nodes", "node.caps", "deg.malig",
5     "breast", "left_low", "right_up",
6     "left_up", "right_low", "central", "irradiat"),
7   "Media" = c(mean(base$class),
8     mean(base$age),
9     mean(base$menopause),
10    mean(base$tumor.size),
11    mean(base$inv.nodes),
12    mean(base$node.caps),
13    mean(base$deg.malig),
14    mean(base$breast),
15    mean(base$left_low),
16    mean(base$right_up),
17    mean(base$left_up),

```

```
18         mean(base$right_low),
19         mean(base$central),
20         mean(base$irradiat)))
21
22 #biblioteca para salvar a tabela criada
23 library(gridExtra)
24
25 png("media.png", height = 50*nrow(media), width = 200*ncol(media))
26 grid.table(media)
27 dev.off()
```

Listing 6: Código fonte em R

Com o código abaixo, foi possível visualizar a criação de cada histograma que está representado as gráficos acima.

```
1 #Plotagem dos histogramas
2 hist(base$class,
3       main = "Histograma Atributo class",
4       ylab = "Frequencia",
5       xlab = "Tipo do dado")
6
7 hist(base$age,
8       main = "Histograma Atributo age",
9       ylab = "Frequencia",
10      xlab = "Idade do individuo")
11
12 hist(base$menopause,
13       main = "Histograma Atributo menopause ",
14       ylab = "Frequencia",
15       xlab = "Indicacao da Menopausa")
16
17 hist(base$tumor.size,
18       main = "Histograma Atributo tumor.size",
19       ylab = "Frequencia",
20       xlab = "Diametro do tumor")
21
22 hist(base$inv.nodes,
23       main = "Histograma Atributo Inv.nodes",
24       ylab = "Frequencia",
25       xlab = "Numero de linfonodos axilares")
26
27 hist(base$node.caps,
28       main = "Histograma Atributo node.caps",
29       ylab = "Frequencia",
30       xlab = "Penetracao do tumor")
31
32 hist(base$deg.malig,
```

```
33     main = "Histograma Atributo deg.malig",
34     ylab = "Frequencia",
35     xlab = "Grau de malignidade do tumor")
36
37 hist(base$breast,
38     main = "Histograma Atributo breast",
39     ylab = "Frequencia",
40     xlab = "Mama em que o cancer pode ocorrer")
41
42 hist(base$irradiat,
43     main = "Histograma Atributo irradiat",
44     ylab = "Frequencia",
45     xlab = "Tipo do dado")
46
47 hist(base$left_low,
48     main = "Histograma Atributo left-low",
49     ylab = "Frequencia",
50     xlab = "Quadrante da mama afetado")
51
52 hist(base$right_up,
53     main = "Histograma Atributo right-up ",
54     ylab = "Frequencia",
55     xlab = "Quadrante da mama afetado")
56
57 hist(base$left_up,
58     main = "Histograma Atributo left-up",
59     ylab = "Frequencia",
60     xlab = "Quadrante da mama afetado")
61
62 hist(base$right_low,
63     main = "Histograma Atributo right-low",
64     ylab = "Frequencia",
65     xlab = "Quadrante da mama afetado")
66
67 hist(base$central,
68     main = "Histograma Atributo central",
69     ylab = "Frequencia",
70     xlab = "Quadrante da mama afetado")
```

Listing 7: Código fonte em R

6 Exploração dos dados através de medidas de distribuição;

A medida de distribuição também estão relacionadas com a média de cada valor, nesse caso não gerei gráficos, apenas uma tabela para poder fazer a análise de acordo com os gráficos gerados no Item anterior o 5.

Atributos	Momento	Curtose	Obliquidade
class	0.29720280	-1.22488733	0.88281304
age	2.66136486	-0.97284195	-0.10329948
menopause	0.92657343	-1.96027095	0.14644434
tumor.size	4.88111888	-0.03145662	0.05513091
inv.nodes	0.51748252	6.76553585	2.59157235
node.caps	0.19580420	0.32723116	1.52513946
deg.malig	1.04895105	-1.16928931	-0.07684805
breast	0.46853147	-1.99118452	0.12546326
left_low	0.38811189	-1.79759047	0.45679117
right_up	0.11538462	3.74965237	2.39510022
left_up	0.33916084	-1.54852255	0.67591066
right_low	0.08391608	6.93840418	2.98565025
central	0.07342657	8.61662968	3.25368483
irradiat	0.23776224	0.23776224	1.22553804

Figura 17: Tabela com as Medidas de distribuição

O momento, curtose e obliquidade e a criação da tabela da figura 17, foi construída a partir do código fonte abaixo. Nele eu crio um novo data frame(assim como foi feita na média do item 6) e atribuo os valores para eles, depois exparto a imagem obtida. Nesse código criei uma variável para cada atributo para receber o momento, a curtose e a obliquidade e assim depois inserir essas variáveis na tabela.

```

1 #Momento
2 mC <- moment(base$class)
3 mA <- moment(base$age)
4 mM <- moment(base$menopause)
5 mT <- moment(base$tumor.size)
6 mI <- moment(base$inv.nodes)
7 mN <- moment(base$node.caps)
8 mD <- moment(base$deg.malig)
9 mB <- moment(base$breast)
10 mLl <- moment(base$left_low)
11 mRr <- moment(base$right_up)

```

```

12 mLu <- moment(base$left_up)
13 mRl <- moment(base$right_low)
14 mCc <- moment(base$central)
15 mIr <- moment(base$irradiat)
16
17 #Curtose
18 cC <- kurtosis(base$class)
19 cA <- kurtosis(base$age)
20 cM <- kurtosis(base$menopause)
21 cT <- kurtosis(base$tumor.size)
22 cI <- kurtosis(base$inv.nodes)
23 cN <- kurtosis(base$node.caps)
24 cD <- kurtosis(base$deg.malig)
25 cB <- kurtosis(base$breast)
26 cLl <- kurtosis(base$left_low)
27 cRr <- kurtosis(base$right_up)
28 cLu <- kurtosis(base$left_up)
29 cRl <- kurtosis(base$right_low)
30 cCc <- kurtosis(base$central)
31 cIr <- kurtosis(base$irradiat)
32
33 # obliquidade
34 oC <- skewness(base$class)
35 oA <- skewness(base$age)
36 oM <- skewness(base$menopause)
37 oT <- skewness(base$tumor.size)
38 oI <- skewness(base$inv.nodes)
39 oN <- skewness(base$node.caps)
40 oD <- skewness(base$deg.malig)
41 oB <- skewness(base$breast)
42 oLl <- skewness(base$left_low)
43 oRr <- skewness(base$right_up)
44 oLu <- skewness(base$left_up)
45 oRl <- skewness(base$right_low)
46 oCc <- skewness(base$central)
47 oIr <- skewness(base$irradiat)
48
49 #Criacao da tabela de distribuicao
50 data <- data.frame(
51   "Atributos" = c("class", "age", "menopause", "tumor.size",
52                   "inv.nodes", "node.caps", "deg.malig",
53                   "breast", "left_low", "right_up",
54                   "left_up", "right_low", "central", "irradiat"),
55   "Momento" = c(mC, mA, mM, mT, mI, mN, mD, mB, mLl, mRr, mLu, mRl, mCc, mIr),
56   "Curtose" = c(cC, cA, cM, cT, cI, cN, cD, cB, cLl, cRr, cLu, cRl, cCc, cIr),
57   "Obliquidade" = c(oC, oA, oM, oT, oI, oN, oD, oB, oLl, oRr, oLu, oRl, oCc, oIr))
58

```

```
59 png("distribuicao.png", height = 50*nrow(data), width = 200*ncol(data))
60 grid.table(data)
61 dev.off()
```

Listing 8: Código fonte em R

7 Identificação e separação do conjunto de teste, que será utilizado para testar o desempenho dos modelos – o conjunto de testes deve ser representativo e ter as características da população completa. Caso sua base de dados já tenha o conjunto de teste definido, analisar se este segue as características do conjunto de treinamento;

A divisão da base foi utilizando a estratificação dos dados, por conta da base de dados ser desbalanceada foi preciso utilizar essa abordagem afim de manter a mesma proporção do atributo de saída definido no item 1, ou seja, manter 70.3% e 29,7%, com isso destino 70.3% para a base de treinamento e 29.7% destinada para base de teste.

Com o código abaixo, foi possível realizar a divisão da base de dados, mantendo a mesma proporção.

```
1 #Biblioteca para divisao da base
2 library(caTools)
3
4 #Porcao da base de dados
5 set.seed(1)
6
7 #Quanto maior for o SplitRatio maior e a base de treinamento
8 divisao <- sample.split(base$class, SplitRatio = 0.703)
9
10 #Criando a base de teste e de treinamento
11 # 201 dados – 70.3% da base
12 base_treinamento <- subset(base, divisao == TRUE)
13
14 #Criando a variavel de teste
15 #85 dados – 29.7 % da base
16 base_teste <- subset(base, divisao == FALSE)
```

Listing 9: Código fonte em R

8 Identificação e eliminação de atributos não necessários;

No caso desta base apenas o atributo *breast.quad* será eliminada, mas isso só vai ocorrer depois de fazer o passo 13 onde ocorre a conversão dos dados simbólicos para numéricos, pois com esse atributo foi necessário fazer a divisão dos valores transformar eles usando a binarização, então após afazer isso o atributo *breast.quad* não é mais necessário para minha base.

Abaixo o código de como é feita a eliminação desse atributo:

```
1 base$breast.quad = NULL
```

Listing 10: Código fonte em R

E o restante dos atributos não são eliminados, pois como se trata de uma base relacionada ao câncer de mama e não se tem muito conhecimento do que é relevante ou não para se tirar da base, o método abordado foi a não eliminação de atributos, pois se por acaso tirar um atributo que é muito importante para detecção do câncer isso vai interferir no resultado do algoritmo.

9 Identificação e eliminação de exemplos não necessários;

O caso da eliminação de exemplos segue o mesmo padrão da eliminação de atributos, é difícil tirar uma conclusão de quais dados deve-se remover da base, sendo que todos devem ter sua certa importância, quando se trata mais de fatores biosólicas do que computacionais.

10 Análise e aplicação de técnicas de amostragem de dados (caso não seja necessário, analisar o porquê);

A amostragem de dados é muito importante pois ajuda os na criação de uma amostra precisa. Porém por conta dessa base de dados ser uma base muito pequena, pois é uma técnica aplicada mais para bases de dados grandes com mais de 1000 instâncias por exemplo e não é o caso dessa base de dados. Por isso essa base não precisou realizar a amostragem dos dados.

11 Identificação e aplicação de técnicas para minimizar problemas de desbalanceamento (caso não seja necessário, analisar o porquê);

Com o código abaixo, foi possível corrigir o problema de balanceamento, nele foi feito o seguinte, primeiro a criação de uma variável auxiliar para receber TRUE - para todos os dados que são iguais a 1, ou seja da classe minoritária e FALSE para o 0 - classe majoritária, depois

foi criada uma base auxiliar, que recebe apenas os valores que são TRUE, ou seja, da classe minoritária e por fim é feita uma junção dessa base auxiliar com a base de treinamento, e agora tem uma base de treinamento balanceada.

```
1 #variavel auxiliar , para armazenar os TRUE e FALSE
2 aux <- base_treinamento[,1]==1
3
4 #Base auxiliar para receber apenas os valores TRUE
5 base_aux <- subset(base_treinamento ,aux == TRUE)
6
7 #Base de treinamentos recebe a combinacao dela mesmo e da base auxiliar
8 base_treinamento <- rbind(base_aux,base_treinamento)
```

Listing 11: Código fonte em R

12 Limpeza de dados:

- (a) Identificação e eliminação de ruídos ou outliers;
- (b) Identificação e eliminação de dados inconsistentes;
- (c) Identificação e eliminação de dados redundantes;
- (d) Identificação e resolução de dados incompletos (ausentes) utilização de alguma técnica de preenchimento e justificar;

A) Os atributos que tem a presença de outliers são: *age*, *tumor.size*, *inv.nodes* neste caso como já se tem a divisão da base de dados, a identificação e a correção dos outliers estão a partir de agora estão relacionados a Base de **TREINAMENTO**. Antes de fazer a correção dos ruídos, utilizando o código abaixo foi possível criar a figura 18 com o boxplot de todos os atributos da base de treinamento sem a correção dos ruídos.

```
1 boxplot(base_treinamento ,
2         main = "Boxplot com os outliers",
3         ylab = "Valores",
4         xlab = "Atributos",
5         pch = 16)
```

Listing 12: Código fonte em R

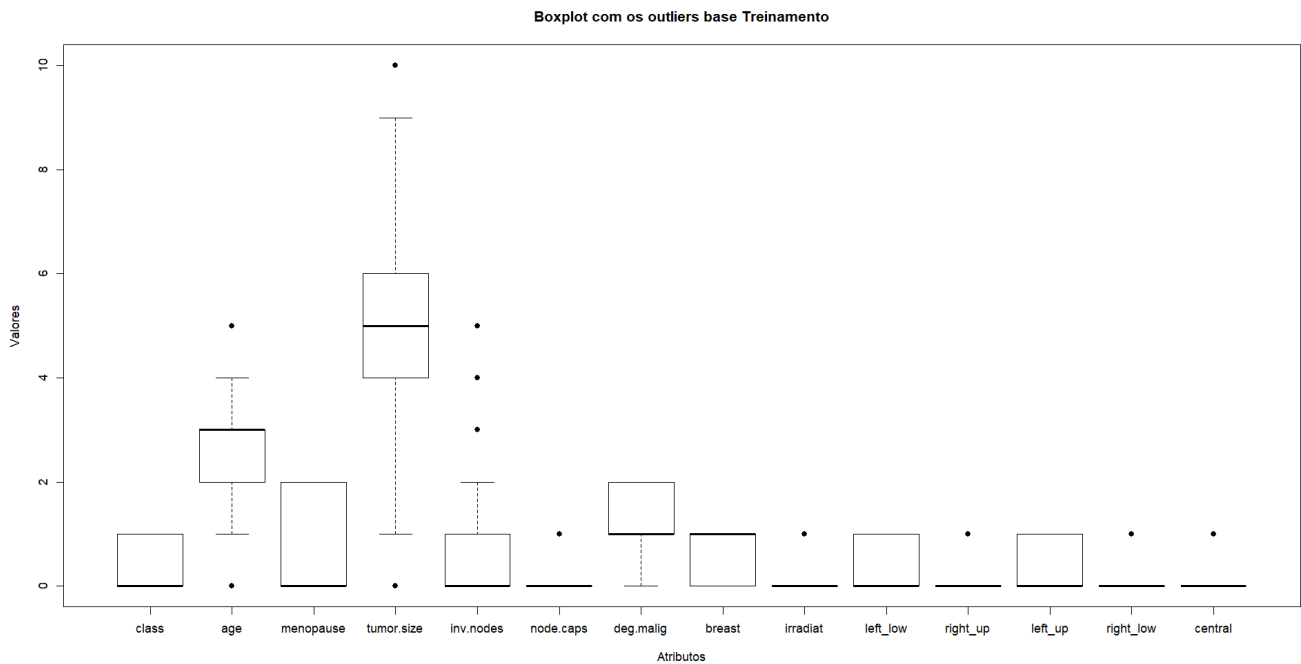


Figura 18: Boxplot com todos atributos e com todos os outliers

Agora para a correção dos ruídos foi utilizado o seguinte código:

```

1 #Atributo age
2 summary(base_treinamento$class)
3
4 #Q3 - Q1
5 IQR_age = 1
6 x_age = mean(base_treinamento$age)
7 Ls_age = x_age + 1.5*IQR_age
8 Li_age = x_age - 1.5*IQR_age
9
10 table(base_treinamento$age)
11
12 base_treinamento$age[base_treinamento$age == 0] <- Li_age
13 base_treinamento$age[base_treinamento$age == 5] <- Ls_age
14
15 #Atributo tumor.size
16 summary(base_treinamento$tumor.size)
17
18 IQR_tumor.size = 6 - 4
19 x_tumor.size = mean(base_treinamento$tumor.size)
20 Ls_tumor.size = x_tumor.size + 1.5*IQR_tumor.size
21 Li_tumor.size = x_tumor.size - 1.5*IQR_tumor.size
22
23 table(base_treinamento$tumor.size)
24
25 base_treinamento$tumor.size[base_treinamento$tumor.size == 0] <- Li_tumor.size

```

```
26 base_treinamento$tumor.size[base_treinamento$tumor.size == 10] <- Ls_tumor.size
27
28 # Atributo inv.nodes
29 summary(base_treinamento$inv.nodes)
30
31 IQR_inv.nodes = 1
32 x_inv.nodes = mean(base_treinamento$inv.nodes)
33 Ls_inv.nodes = x_inv.nodes + 1.5*IQR_inv.nodes
34
35 table(base_treinamento$inv.nodes)
36
37 base_treinamento$inv.nodes[base_treinamento$inv.nodes >= 3] <- Ls_inv.nodes
38
39 boxplot(base_treinamento,
40         main = "Boxplot sem os outliers",
41         ylab = "Valores",
42         xlab = "Atributos",
43         pch = 16)
```

Listing 13: Código fonte em R

Esse código acima faz o seguinte : Para cada atributo com outliers ele faz uma análise da amplitude interquartil utilizando sigla IQR acrescido do nome do atributo, para analisar a amplitude é feita $IQR = Q3 - Q1$, após isso analiso o limite superior chamado de Ls acrescido do nome do atributo e $Ls = media + 1.5 * IQR$, essa média é dado de acordo com cada atributo, após encontrar o limite superior, faço o mesmo para o limite inferior $Ls = media - 1.5 * IQR$, depois dou um comando `table(atributo)` e vejo no manualmente olhando para o gráfico e para os valores que tenho retornados do comando table, após isso substituo esse valor a onde eu encontro os outliers, se for no limite inferior troca o valor para o valor calculado no Li e o mesmo se for o limite superior.

Após o código acima, foi possível gerar uma imagem, que pode ser vista na figura 19, sem a presença dos ruídos. O que mostra agora que a base não tme mais dados que são completamente longe da média dos dados.

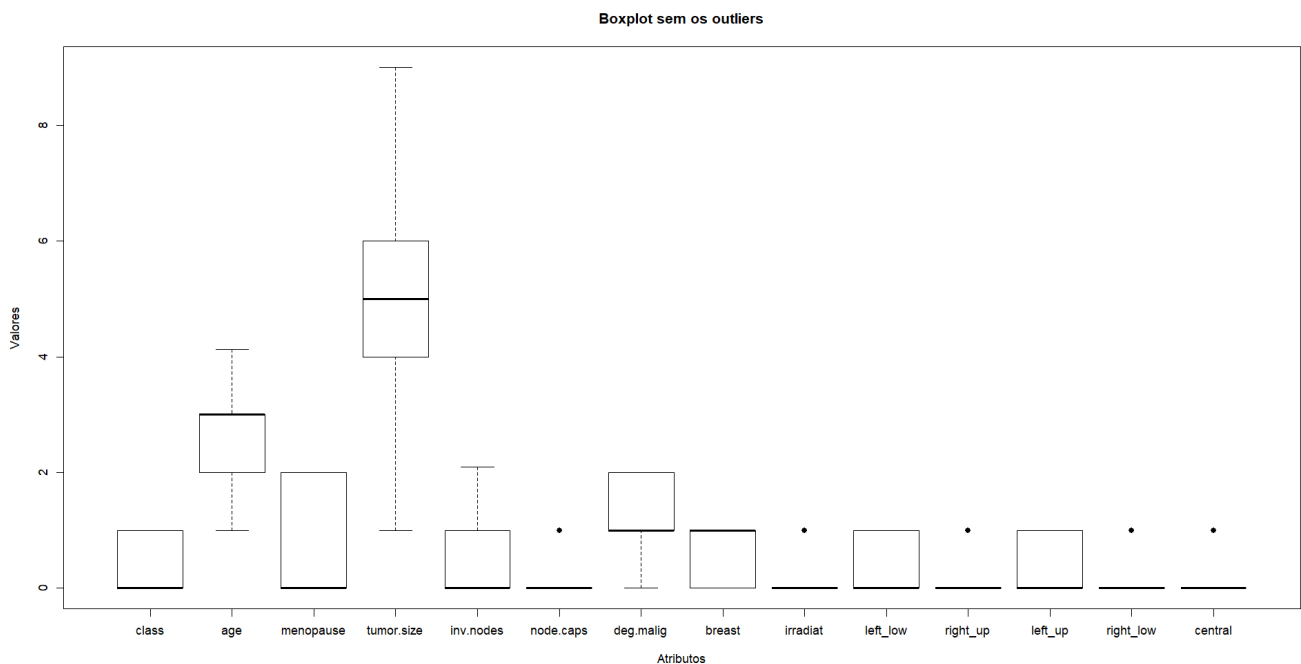


Figura 19: Boxplot com todos atributos sem os outliers

B) Os dados inconsistentes encontrados foram os dos atributos *breast.quad* e *node.caps*, como a base de dados é muito pequena optei por não eliminar esses dados e sim realizar a substituição desses dados pela moda, esse processo é realizado na letra D desse mesmo tópico.

C) A redundância de dados está relacionada a ocorrência do mesmo conjunto de dados ser armazenado em dois ou mais locais. O que não acontece nesse base de dados. Não se tem redundância de dados. Em nenhum momento os atributos se repetem e é possível ver isso no código abaixo:

```
1 #Verifica colunas duplicadas e cria uma lista
2 colunas_duplicadas <- duplicated(as.list(base))
3
4 print(colunas_duplicadas)
5
6 #Mostra o nome das colunas duplicadas
7 colnames(base[colunas_duplicadas])
```

Listing 14: Código fonte em R

O código acima faz a criação de uma lista com todas as colunas duplicadas e depois na linha 4 tem um print para ver o retorno dessa lista, ela retorna tudo *false*, o que já mostra que não se tem atributos repetidos e para confirmar na linha 7 utilizando a função *colnames* para mostrar o nome das colunas duplicadas, e ela dá um retorno de 0, ou seja confirma novamente a não existência de dados redundantes.

D) Por conta da base de dados sem composta de dados simbólicos, o método mais eficiente para resolução dos dados ausentes foi utilizando a moda, isso antes dos dados serem convertidos para numéricos e o motivo de não utilizar a média por exemplo é por conta de que se eu tenho valores que vão de 0 – 1 – 2 – 4 – 5 e a media cai em 2.5 fica difícil encaixar ela em algum lugar, sendo que os valores são definidos como inteiros, ai a média é a melhor solução nesse caso.

```
1 base$breast.quad[base$breast.quad == "?"] <- "left_low"  
2 base$node.caps[base$node.caps == "?"] <- "no"
```

Listing 15: Código fonte em R

13 Identificação e conversão dos tipos de dados (caso não seja necessário, analisar o porquê). Os tipos de conversão que podem ser utilizados são:

- (a) Conversão de tipos (simbólico para numérico, ordinal para numérico, nominal para binário, numérico para ordinal);
- (b) Normalização dos dados (re-escala ou padronização);

A) Como todos os atributos da minha base são, simbólicos, tenho que converter eles para numérico, para seja possível realizar algumas operações. Então fiz uma tabela para explicar melhor quais foram as conversões realizadas em cada atributo.

Os atributos **menopause** e **breast.quad** eu destaquei para explicar melhor sobre esses dois, pois os outros seguem um padrão de ordem que varia de 0 a quantidade de valores que tem que dentro de cada atributo.

No atributo **menopause**, a conversão dele se deu da seguinte maneira: *premeno*: está relacionado a pré-menopausa, por isso decidi atribuir o valor de **0** para ele, o *It40*: está relacionado a idade antes dos 40 anos, então para ela foi definido o valor de **1** e por fim o valor *ge40*: está relacionado a idade após 40 anos, então foi definido o valor de **2** para ela.

No atributo **premeno**, foi utilizado a técnica de *One-Hot*, que é o processo de converter uma variável categórica com várias categorias em várias variáveis, cada uma com um valor de 1 ou 0. Foi criada várias outras instâncias com os valores que tinha no atributo, ou seja dentro do atributo **breast.quad** tinha os valores do quadrante da mama afetado *left-low*, *right-up*, *left-up*, *right-low* e *central*, agora cada valor do quadrante é um novo atributo da base.

Atributos	Conversão
class	0 - 1
age	0 - 1 - 2 - 3 - 4 - 5
menopause	premeno:0 - lt40:1 - ge40:2
tumor.size	0 - 1 - 2 - 3 - 4 - 5 - 6 - 7 - 8 - 9 - 10
inv.nodes	0 - 1 - 2 - 3 - 4 - 5 - 6
node.caps	0 - 1
deg.malig	0 - 1 - 2
breast	0 - 1
breast.quad	Divide em instâncias
irradiat	0 - 1

Com o código abaixo é possível ver como foi realizada a conversão dos dados, vou destacar alguns pontos do código, da linha 28 até a 36 faço novamente uma conversão para numérico, acontece que quando eu uso o factor para conversão ele acaba convertendo os valores em números, porém o estado dele ainda continua sendo de uma variável categórica, para isso utilizo o método `as.numeric(as.character(atributo))`, para que os dados agora sejam de fato convertidos.

```

1 #A) Conversao dos atributos qualitativos para quantitativos:
2
3 #Com apenas 2 tipos
4 base$class = factor(base$class,
5     levels = c('no-recurrence-events', 'recurrence-events'), labels = c(0,1))
6 base$irradiat = factor(base$irradiat, levels = c('no', 'yes'), labels = c(0,1))
7 base$breast = factor(base$breast, levels = c('left', 'right'), labels = c(0,1))
8 base$node.caps = factor(base$node.caps, levels = c('no', 'yes'), labels = c(0,1))
9
10 #Com 3 tipos
11 base$deg.malig = factor(base$deg.malig, levels = c('1', '2', '3'), labels = c
    (0,1,2))
12
13 #Com mais de 3
14 base$age = factor(base$age,
15     levels = c('20-29', '30-39', '40-49', '50-59', '60-69', '70-79'),
16     labels = c(0,1,2,3,4,5))
17
18 base$menopause = factor(base$menopause,
19     levels = c('premeno', 'lt40', 'ge40'), labels = c(0,1,2))
20 base$tumor.size = factor(base$tumor.size,
21     levels = c('0-4', '5-9', '10-14', '15-19', '20-24', '25-29', '30-34', '35-39', '
    40-44', '45-49', '50-54'),
22     labels = c(0,1,2,3,4,5,6,7,8,9,10))
23 base$inv.nodes = factor(base$inv.nodes,

```

```

24   levels = c('0-2', '3-5', '6-8', '9-11', '12-14', '15-17', '24-26'),
25   labels = c(0,1,2,3,4,5,6))
26
27 #Converter para numerico
28 base$class <- as.numeric(as.character(base$class))
29 base$age <- as.numeric(as.character(base$age))
30 base$menopause <- as.numeric(as.character(base$menopause))
31 base$tumor.size <- as.numeric(as.character(base$tumor.size))
32 base$inv.nodes <- as.numeric(as.character(base$inv.nodes))
33 base$node.caps <- as.numeric(as.character(base$node.caps))
34 base$deg.malig <- as.numeric(as.character(base$deg.malig))
35 base$breast <- as.numeric(as.character(base$breast))
36 base$irradiat <- as.numeric(as.character(base$irradiat))
37
38 #Criando novas colunas
39 base$left_low <- 0
40 base$left_low[base$breast.quad == "left_low"] <- 1
41
42 base$right_up <- 0
43 base$right_up[base$breast.quad == "right_up"] <- 1
44
45 base$left_up <- 0
46 base$left_up[base$breast.quad == "left_up"] <- 1
47
48 base$right_low <- 0
49 base$right_low[base$breast.quad == "right_low"] <- 1
50
51 base$central <- 0
52 base$central[base$breast.quad == "central"] <- 1
53
54 #verificar se deu certo
55 base[is.na(base$inv.nodes)]

```

Listing 16: Código fonte em R

B) A técnica utilizada foi a reescala, ou seja deixei todos os dados da base de teste e treinamento na mesma escala, optei por deixar a base original com os dados normais(sem reescala), optei por isso para não ter que refazer os gráficos executados nos itens anteriores.

Com o código abaixo, foi possível fazer a reescala de todos atributos exceto o do atributo de saída. O que esse código faz é basicamente pegar a posição da base a partir do 2, pois a contagem aqui começa em 1 e ir até o último atributo da classe que termina em 14.

```

1 base_treinamento[,2:14] = scale(base_treinamento[,2:14])
2 base_teste[,2:14] = scale(base_teste[,2:14])

```

Listing 17: Código fonte em R

14 Análise e aplicação de alguma técnica para redução de dimensionalidade – pesquisar alguma técnica utilizada na literatura e aplicar;

A técnica de dimensionalidade está relacionada como o próprio nome já diz a uma redução dos dados, com isso é possível aumentar a qualidade, mas ao mesmo tempo minimizando a perda de informação. E essa técnica faz isso criando novas variáveis, que não são correlacionadas e que possam maximizar sucessivamente a medida de dispersão.

Para encontrar essas novas variáveis, os componentes principais, são definidas pelo conjunto de dados em mãos, o que torna o PCA uma técnica de análise de dados adaptativa, ou seja, o mais bacana é que sempre o usuário pode variar o número de variáveis que deseja ter.

O código abaixo representa toda essa técnica de PCA:

```
1 # 0 – no-recurrence-events
2 # 1 – recurrence-events
3
4 #Aplicar reducao
5 library('caret')
6
7 pca = preProcess(x = base_treinamento[-1], method = 'pca', pcaComp = 4)
8 base_treinamento_pca = predict(pca, base_treinamento)
9 base_teste_pca = predict(pca, base_teste)
10
11 #Testes – Extra
12
13 #Teste Com toda a base
14 classificador = naiveBayes(x = base_treinamento[-1], y = base_treinamento$class)
15
16 previsoes = predict(classificador, newdata = base_teste[-1])
17 # 64,7% de acertos
18 matriz_confusao = table(base_teste[,1], previsoes)
19 print(matriz_confusao)
20 #analisa a matriz de confusao e tras varias metricas
21 confusionMatrix(matriz_confusao)
22 #Accuracy : 0.6471
23
24
25 #Teste com apenas 4 atributos
26 classificador_pca = naiveBayes(x = base_treinamento_pca[-1], y = base_treinamento_pca$class)
27
28 previsoes_pca = predict(classificador_pca, newdata = base_teste_pca[-1])
29 # 70,6% de acertos
30 matriz_confusao_pca = table(base_teste_pca[,1], previsoes_pca)
```

```
31 print(matriz_confusao_pca)
32 #analisa a matriz de confusao e tras varias metricas
33
34 confusionMatrix(matriz_confusao_pca)
35 #Accuracy : 0.7059
```

Listing 18: Código fonte em R

Esse código acima basicamente faz o seguinte:

- Cria o pca usando a função preProcess

Essa parte é feita em cima da base de treinamentos[-1] pois é necessário retirar o atributo de saída e atribui 4 atributos aleatórios.

- Atribuição para as bases

Após criar o pca é preciso fazer as bases receberem seus novos valores e atributos. No caso foi criada uma nova variável mas com o intuito de fazer comparações ao mesmo tempo, sem ter que ficar trocando os valores, comparações no sentido de classificação e previsão e matriz de confusão, é possível até ver que esses nomes vão se repetir, o motivo foi mais para comparação dos acertos e da acurácia.

- Testes

Após isso, foi realizado teste de previsões e de matriz de confusão para toda a base de dados de teste para testar se o algoritmo está bom.

15 Conclusão

Neste trabalho foram abordadas várias técnicas de pré-processamento de dados, e com tudo que foi realizado, foi possível perceber o quão importante é a pré análise dos dados, pois se fizemos alguma análise errada isso pode repercutir em erros no trabalho.

No passo 14 foi possível ver como realmente está a porcentagem de acertos do algoritmo, com a técnica de dimensionalidade e sem ela, ainda não se pode concluir de fato nada, mas aparentemente o algoritmo está sendo acertando bem, porém como a base de dados é pequena pode ter ocorrido dele ter aprendendo mais de um determinado dado, ou seja pode ter ocorrido um Overfitting, pois ele pode ter tido um bom desempenho nos dados de treinamento porém ele pode ter generalizado mal os outros dados. A expectativa é que no próximo trabalho consiga tirar uma conclusão bem melhor sobre essa base de dados.