

Primeira Prova de Modelos de Regressão

Carlos Eduardo Alves de Souza
04/03/2022

```
# Importando livrerias
library(robustbase)

## Warning: package 'robustbase' was built under R version 4.1.2

rm(list = ls()) # removendo variáveis já salvas...

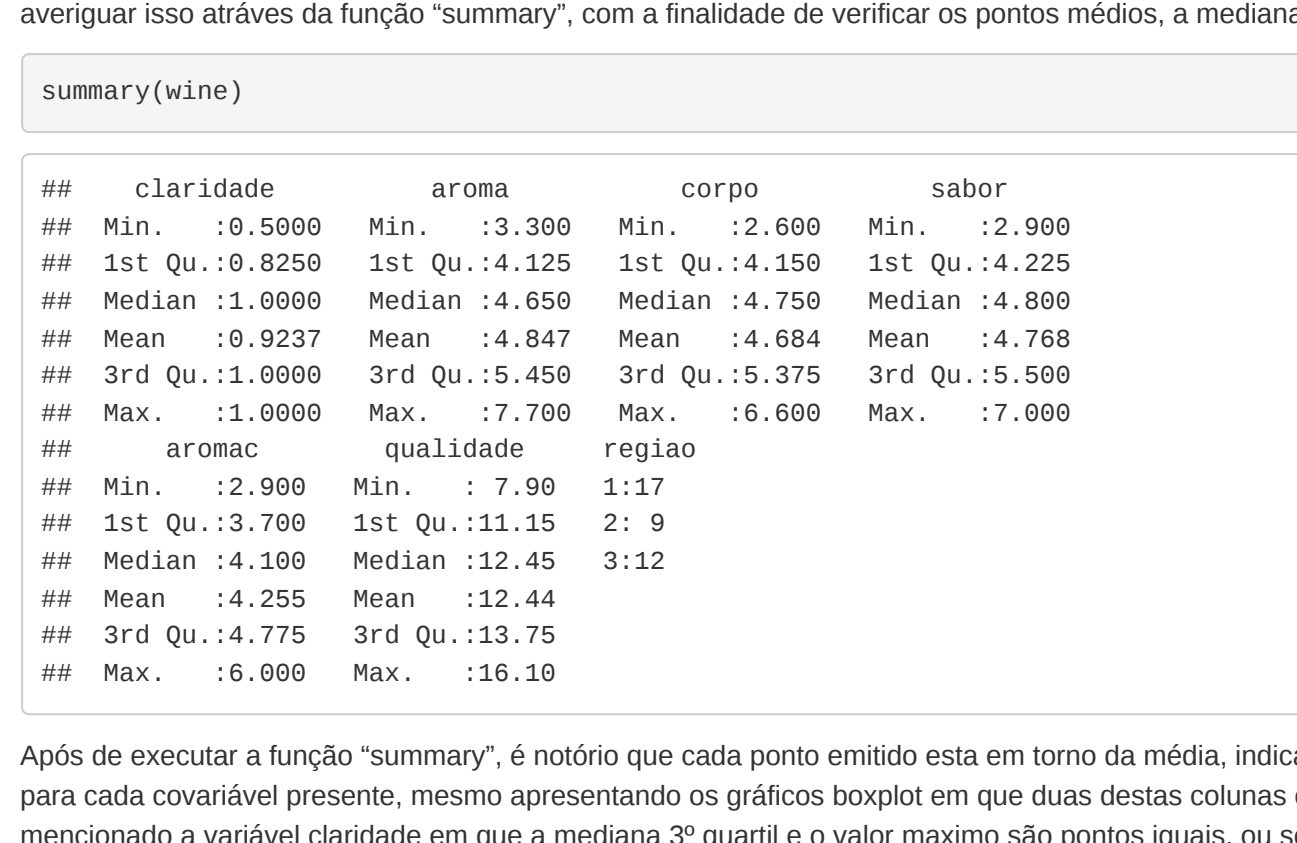
# Importando a base de dados wine
wine = read.table("C:/Users/Carlo/Downloads/wine.txt", header = T)

# Transformação da covariável "região"
wine$regiao <- factor(wine$regiao), wine

##
## clariidade aroma corpo aromac qualidade regiao
## 1 1.0 3.3 2.8 3.1 4.1 9.8 1
## 2 1.0 4.4 4.9 3.5 3.9 12.0 1
## 3 1.0 3.9 5.3 4.8 4.7 11.9 1
## 4 1.0 3.9 2.6 3.1 3.6 11.1 1
## 5 1.0 5.0 5.1 5.5 5.1 13.3 1
## 6 1.0 4.6 4.7 5.0 4.1 12.8 1
## 7 1.0 4.0 4.8 4.0 3.3 12.8 1
## 8 1.0 5.3 4.5 4.3 5.2 12.0 1
## 9 1.0 4.3 4.3 3.9 2.9 13.6 3
## 10 1.0 4.3 3.9 4.7 3.9 13.9 1
## 11 1.0 5.1 4.3 4.5 3.6 14.4 3
## 12 0.5 3.3 5.4 4.3 3.6 12.3 2
## 13 0.8 5.9 6.7 7.0 4.1 16.1 3
## 14 0.7 7.7 6.0 4.7 3.7 16.1 3
## 15 1.0 7.1 4.4 5.8 4.1 15.5 3
## 16 0.9 5.5 5.6 5.6 4.4 15.5 3
## 17 1.0 6.3 5.4 4.8 4.6 13.8 3
## 18 1.0 5.0 5.5 5.5 4.1 13.8 3
## 19 1.0 4.6 4.1 4.3 3.1 13.3 1
## 20 0.9 3.4 5.0 3.4 4.4 7.9 2
## 21 0.9 6.4 5.4 6.6 4.8 15.1 3
## 22 1.0 5.5 5.3 5.3 3.8 13.5 3
## 23 0.7 4.7 4.1 5.0 3.7 10.0 2
## 24 0.7 4.1 4.0 4.1 4.0 9.5 2
## 25 1.0 6.0 5.4 5.7 4.7 12.7 3
## 26 1.0 4.3 4.6 4.7 4.9 11.0 1
## 27 1.0 3.9 4.0 5.1 5.1 11.7 1
## 28 1.0 5.1 4.9 5.0 5.1 11.9 2
## 29 1.0 3.9 4.4 5.0 4.4 10.0 2
## 30 1.0 4.5 3.7 2.9 3.9 8.5 2
## 31 1.0 5.2 4.3 5.0 6.0 10.7 2
## 32 0.8 4.2 3.8 3.0 4.7 9.1 1
## 33 1.0 3.3 3.5 4.3 4.5 12.1 1
## 34 1.0 6.8 5.0 6.0 5.2 14.9 3
## 35 0.8 5.0 5.7 5.5 4.8 13.5 1
## 36 0.8 3.5 4.7 4.2 3.3 12.2 1
## 37 0.8 4.3 5.5 3.5 5.8 13.8 1
## 38 0.8 5.2 4.8 5.7 3.5 13.2 1
```

Análise dos dados

```
# Análise descritiva, por meio do gráfico boxplot
par(mfrow = c(1,3))
adbox(wine$clariidade, main = "Clariidade")
adbox(wine$aroma, main = "Aroma")
adbox(wine$corpo, main = "Corpo")
adbox(wine$sabor, main = "Sabor")
adbox(wine$aromac, main = "Aroma do Tonel de Carvalho")
adbox(wine$qualidade, main = "Qualidade")
```



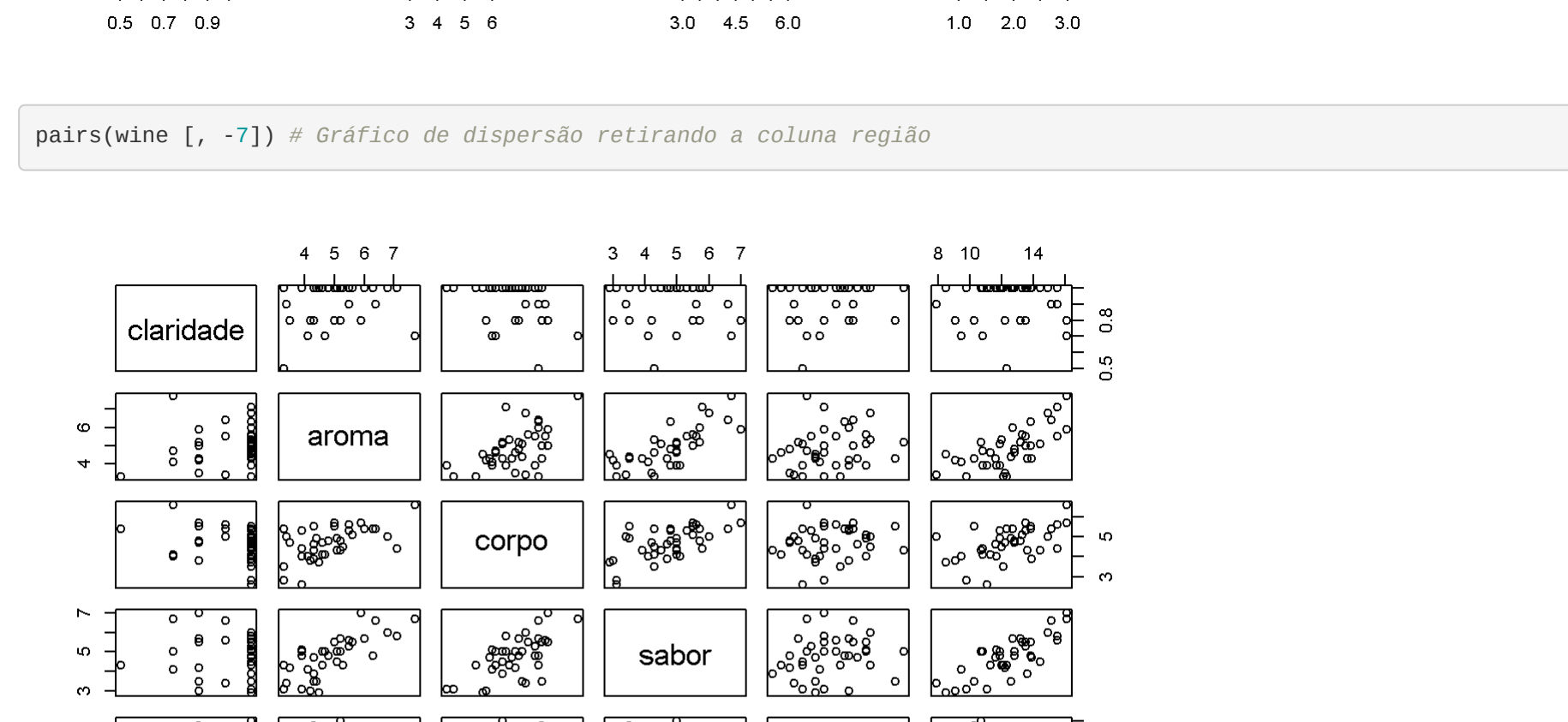
Percebemos a covariável "sabor" e "aromac" possuem outliers, sendo possível verificar se esse ponto considerado como outliers são influente ao modelo. Já as covariáveis "sabor", "qualidade" e "corpo" tem um comportamento simétrico, logo "aroma" e "aromac", exibem uma ligeira assimetria. Por fim, a covariável "clariidade" tem-se um comportamento diferente das demais, não bem representado no gráfico de boxplot. Podemos averiguar isso através da função "summary", com a finalidade de verificar os pontos médios, a mediana e/ou quantis.

```
summary(wine)

## clariidade aroma corpo aromac qualidade regiao
## Min. -0.5980 Min. -13.2000 Min. -12.6000 Min. -2.9000
## 1st Qu.:0.8250 1st Qu.:4.125 1st Qu.:4.150 1st Qu.:4.225
## Median :1.0000 Median :4.650 Median :4.750 Median :4.800
## Mean :0.9227 Mean :4.6047 Mean :4.684 Mean :4.750
## 3rd Qu.:1.0000 3rd Qu.:5.450 3rd Qu.:5.375 3rd Qu.:5.500
## Max. :1.0000 Max. :17.700 Max. :6.600 Max. :7.000
##
## aromac qualidade regiao
## Min. -12.900 Min. : 7.90 1:17
## 1st Qu.:3.700 1st Qu.:11.15 2: 9
## Median :4.100 Median :12.45 3:12
## Mean :4.255 Mean :12.44
## 3rd Qu.:4.775 3rd Qu.:13.75
## Max. :6.000 Max. :16.10
```

Após de executar a função "summary", é notório que cada ponto emitido está em torno da média, indicando que uma certa simetria nos dados, para cada covariável presente, mesmo apresentando os gráficos boxplot em que duas destas colunas exibiram valores como outliers. Como mencionado a variável "clariidade" em que a mediana 3º quartil e o valor máximo são pontos iguais, ou seja, realmente esta variável pode ser considerada não simétrica.

```
# Gráfico de dispersão
pairs(wine) # gráfico de dispersão para o banco de dados inteiro
```



```
pairs(wine[, 1:7]) # gráfico de dispersão retirando a coluna regiao
```



Apresentando os dois gráficos, podemos observar que a covariável regiao não tem tanto influencia nas demais colunas mencionadas devido ela ser uma variável categorica com 3 posições, a coluna Clariidade também mostra o mesmo comportamento da variável regiao. Após remoção da regiao, podemos ver que as variáveis Sabor e Qualidade exibem uma relação crescente com a qualidade do vinho, outrora as covariáveis não exibiam uma variabilidade maior que não permite facilmente identificar alguma dependência.

```
# Diagrama de dispersão para cada covariável
par(mfrow = c(2,3))
```

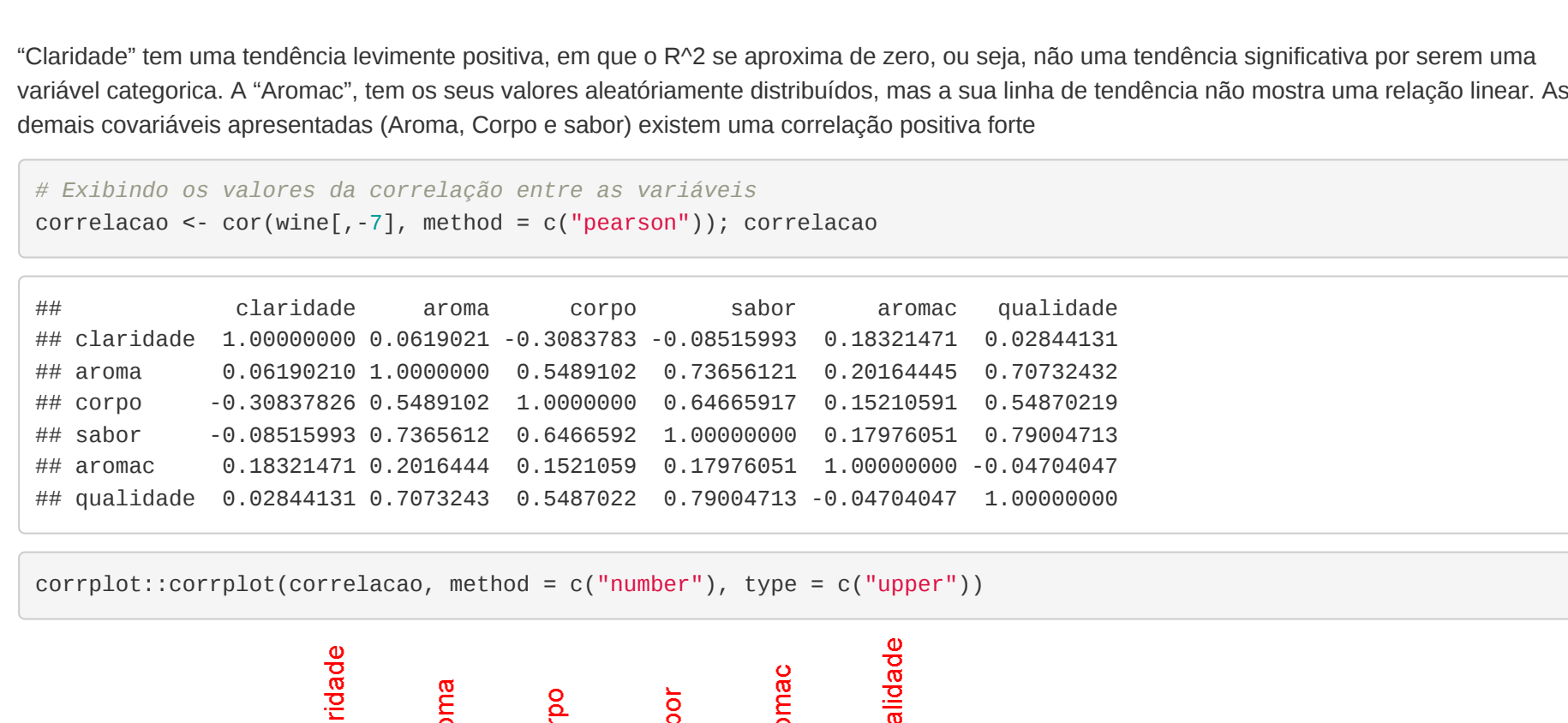
```
plot(qualidade ~ clariidade, xlab = "Clariidade", ylab = "Qualidade", data = wine)
abline(lm(qualidade ~ clariidade, data = wine), col = 2, lwd = 2)

plot(qualidade ~ aroma, xlab = "Aroma", ylab = "Qualidade", data = wine)
abline(lm(qualidade ~ aroma, data = wine), col = 2, lwd = 2)

plot(qualidade ~ corpo, xlab = "Corpo", ylab = "Qualidade", data = wine)
abline(lm(qualidade ~ corpo, data = wine), col = 2, lwd = 2)

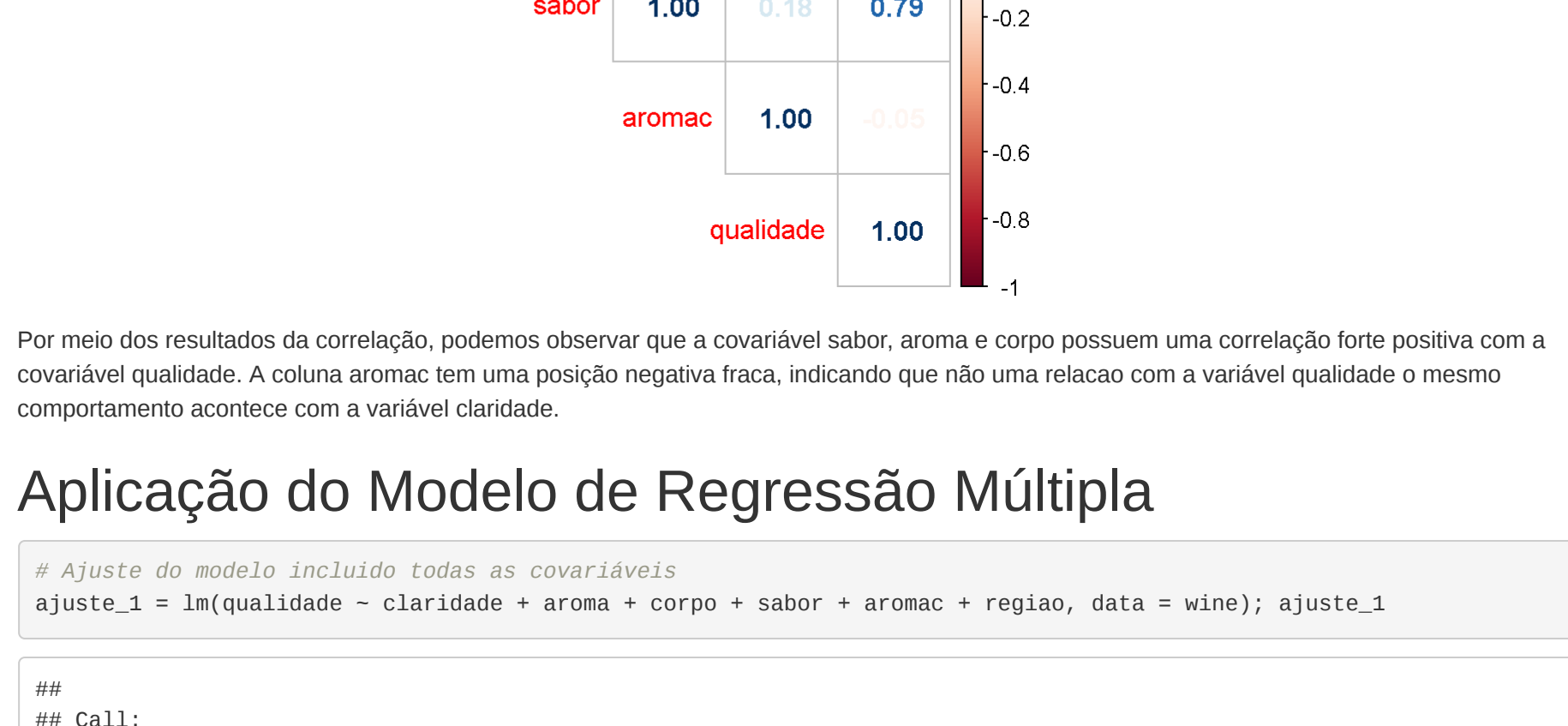
plot(qualidade ~ sabor, xlab = "Sabor", ylab = "Qualidade", data = wine)
abline(lm(qualidade ~ sabor, data = wine), col = 2, lwd = 2)

plot(qualidade ~ aromac, xlab = "Aroma do Tonel de Carvalho", ylab = "Qualidade", data = wine)
abline(lm(qualidade ~ aromac, data = wine), col = 2, lwd = 2)
```



"Clariidade" tem uma tendência levemente positiva, em que o R² se aproxima de zero, ou seja, não uma tendência significativa por serrem uma variável categorica. A "Aromac", tem os seus valores aleatoriamente distribuídos, mas a sua linha de tendência não mostra uma relação linear. As demais covariáveis apresentadas (Aroma, Corpo e Sabor) exibem uma correlação positiva forte.

```
# Exibindo os valores da correlação entre as variáveis
correlacao <- cor(wine[, 1:7], method = c("pearson")); correlacao
```



Por meio dos resultados da correlação, podemos observar que a covariável sabor, aroma e corpo possuem uma correlação forte positiva com a qualidade. A coluna aromac tem uma posição negativa fraca, indicando que não uma relação com a variável qualidade o mesmo comportamento acontece com a variável clariidade.

Aplicação do Modelo de Regressão Múltipla

```
# Ajuste do modelo incluindo todas as covariáveis
ajuste_1 <- lm(qualidade ~ clariidade + aroma + corpo + sabor + aromac + regiao, data = wine); ajuste_1
```

```
##
## Call:
## lm(formula = qualidade ~ clariidade + aroma + corpo + sabor +
## aromac + regiao, data = wine)
##
## Coefficients:
## (Intercept) clariidade aroma corpo sabor aromac
## 7.81437 0.91705 0.88901 0.07967 1.11723 -0.34644
## regiao2
## -1.51285 0.97259

## Residuals:
## Min 1Q Median 3Q Max
## -1.88824 -0.58413 -0.02861 0.48627 1.79909

## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 7.8137 1.9684 3.968 0.00041 ***
## clariidade 0.91705 1.40027 0.654 0.50979
## aroma 0.88901 0.25250 3.523 0.00080 ***
## corpo 0.07967 0.26772 0.298 0.76882
## sabor 1.11723 0.24026 4.650 6.25e-05 ***
## aromac -0.34644 0.23301 -1.487 0.14703
## regiao2 -1.51285 0.39227 -3.857 0.00050 ***
## regiao3 0.97259 0.51617 1.906 0.06618

## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 0.9154 on 38 degrees of freedom
## Multiple R-squared: 0.8276, Adjusted R-squared: 0.7997
## F-statistic: 22.1 on 7 and 30 DF, p-value: 3.295e-10
```

Segundo os valores obtidos do "ajuste_1", somente as covariáveis "sabor" e "regiao" foram significativas, em que o p_valor<0.10. As demais covariáveis não tiveram uma representatividade ao modelo.

```
# Aplicando um novo modelo de ajuste sob a covariável regiao
ajuste_2 <- lm(qualidade ~ clariidade + aroma + corpo + sabor + aromac, data = wine); ajuste_2
```

```
##
## Call:
## lm(formula = qualidade ~ clariidade + aroma + corpo + sabor +
## aromac, data = wine)
##
## Coefficients:
## (Intercept) clariidade aroma corpo sabor aromac
## 3.9960 2.3395 0.4826 0.2732 1.1683 -0.6840

## Residuals:
## Min 1Q Median 3Q Max
## -2.05552 -0.57448 -0.07892 0.67275 1.69893

## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 3.9969 2.2318 1.791 0.08275
## clariidade 2.3395 1.7348 1.349 0.18698
## aroma 0.4826 0.2724 1.772 0.08698
## corpo 0.2732 0.3326 0.821 0.41753
## sabor 1.1683 0.3845 3.037 0.00352 ***
## aromac -0.6840 0.2712 -2.537 0.01633

## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 1.163 on 32 degrees of freedom
## Multiple R-squared: 0.7206, Adjusted R-squared: 0.6769
## F-statistic: 16.51 on 5 and 32 DF, p-value: 4.703e-08
```

Obs.: Se levarmos em consideração a exclusão da covariável "regiao" como indicado na parte em análise descritiva, segundo os resultados obtidos, podemos ver que a não inclusão desta (regiao) no ajuste_2, muda totalmente os valores das estimativas (coeficientes - betas) e muda drasticamente os valores do p_valor, indicando que somente as covariáveis: corpo e clariidade não é significativo ao modelo onde o p_valor desce é > 0.10, conforme a referência na prova. É notório observar que o R-squared e o R ajustado para o modelo ajuste_2 é menor o que se encontra no modelo ajuste_1. Por fim, através do resultado do R-quadrado entre os dois modelos apresentados podemos incluir a covariável regiao.

```
# Aplicando um novo modelo de ajuste segundo os resultados obtidos do ajuste_1, excluindo as covariáveis não significativas
ajuste_3 <- lm(qualidade ~ sabor + regiao, data = wine); ajuste_3
```

```
##
## Call:
## lm(formula = qualidade ~ sabor + regiao, data = wine)
##
## Coefficients:
## (Intercept) sabor regiao2 regiao3
## 7.084 1.116 1.533 1.223

## Residuals:
## Min 1Q Median 3Q Max
## -1.97630 -0.58844 0.02184 0.51572 1.94322

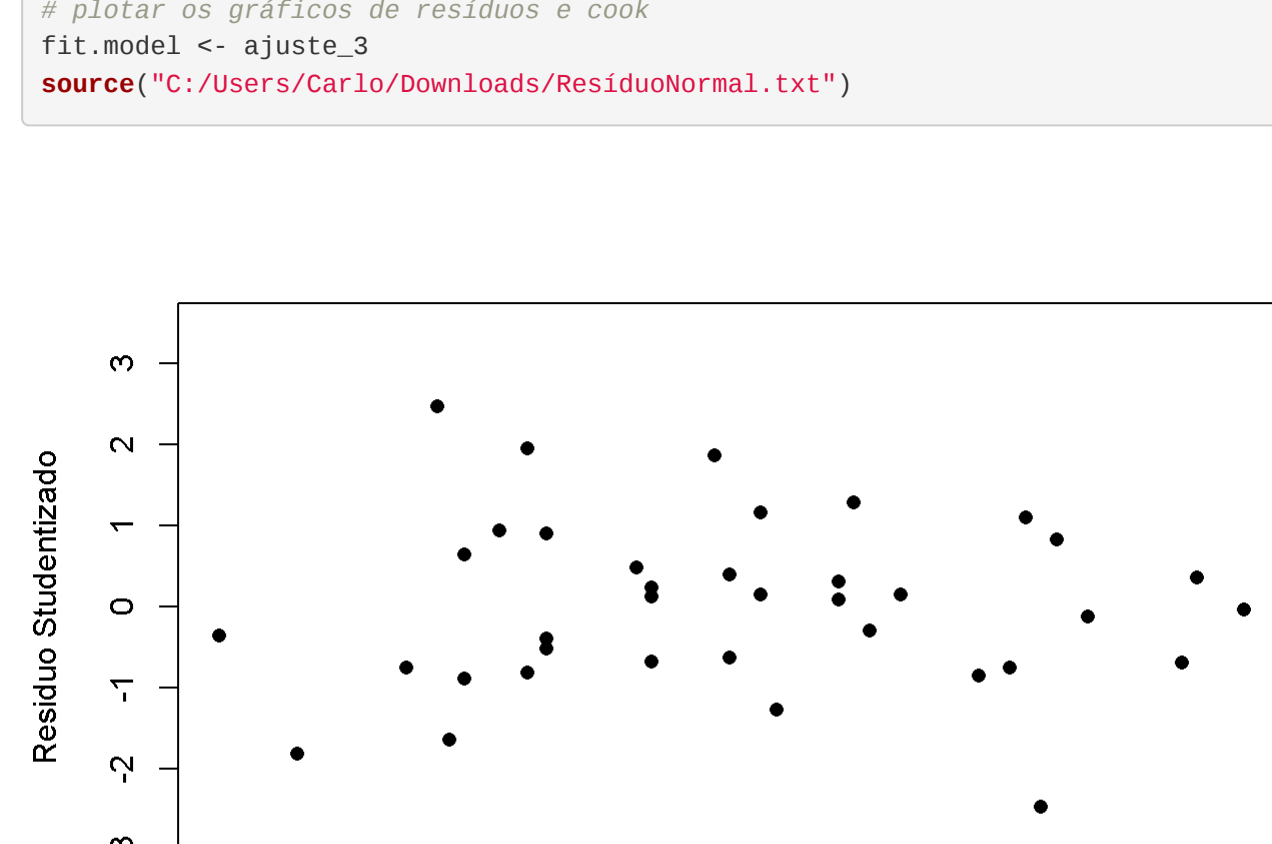
## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 7.0943 0.7912 8.967 1.70e-10 ***
## sabor 1.1155 0.1738 6.417 2.40e-07 ***
## regiao2 -1.5335 0.3688 -4.158 0.00025 ***
## regiao3 1.2234 0.4003 3.056 0.00356 ***

## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 0.8946 on 34 degrees of freedom
## Multiple R-squared: 0.8242, Adjusted R-squared: 0.8087
## F-statistic: 63.13 on 3 and 34 DF, p-value: 6.350e-13
```

Segundo os resultados do modelo ajuste_3, todas as covariáveis foram significativas onde o resultado do respectivo p_valor é < 0.10

```
# plotar os gráficos de resíduos e cook
fit.model <- ajuste_3
source("C:/Users/Carlo/Downloads/Residuonormal.txt")
```



```
source("C:/Users/Carlo/Downloads/CookNormal.txt")
```



Após executar análise de resíduos e a distância de cook, observamos que os resíduos são normalmente distribuídos no intervalo [-3, 3] e a distância de cook identifica 3 pontos fora da margem.

```
# Identificando os pontos discrepantes
distance.cook <- which(DI0>3)
wine[distance.cook, ]
```

```
## clariidade aroma corpo sabor aromac qualidade regiao
## 12 0.5 3.3 5.4 4.3 3.6 12.3 2
## 20 0.9 3.4 5.0 3.4 3.4 7.9 2
## 25 1.0 6.0 5.4 5.7 4.7 12.7 3
```

```
# Executando um novo modelo, excluindo os pontos discrepantes
ajuste_4 <- lm(qualidade ~ sabor + regiao, data = wine, subset = -c(12, 20, 25))
```

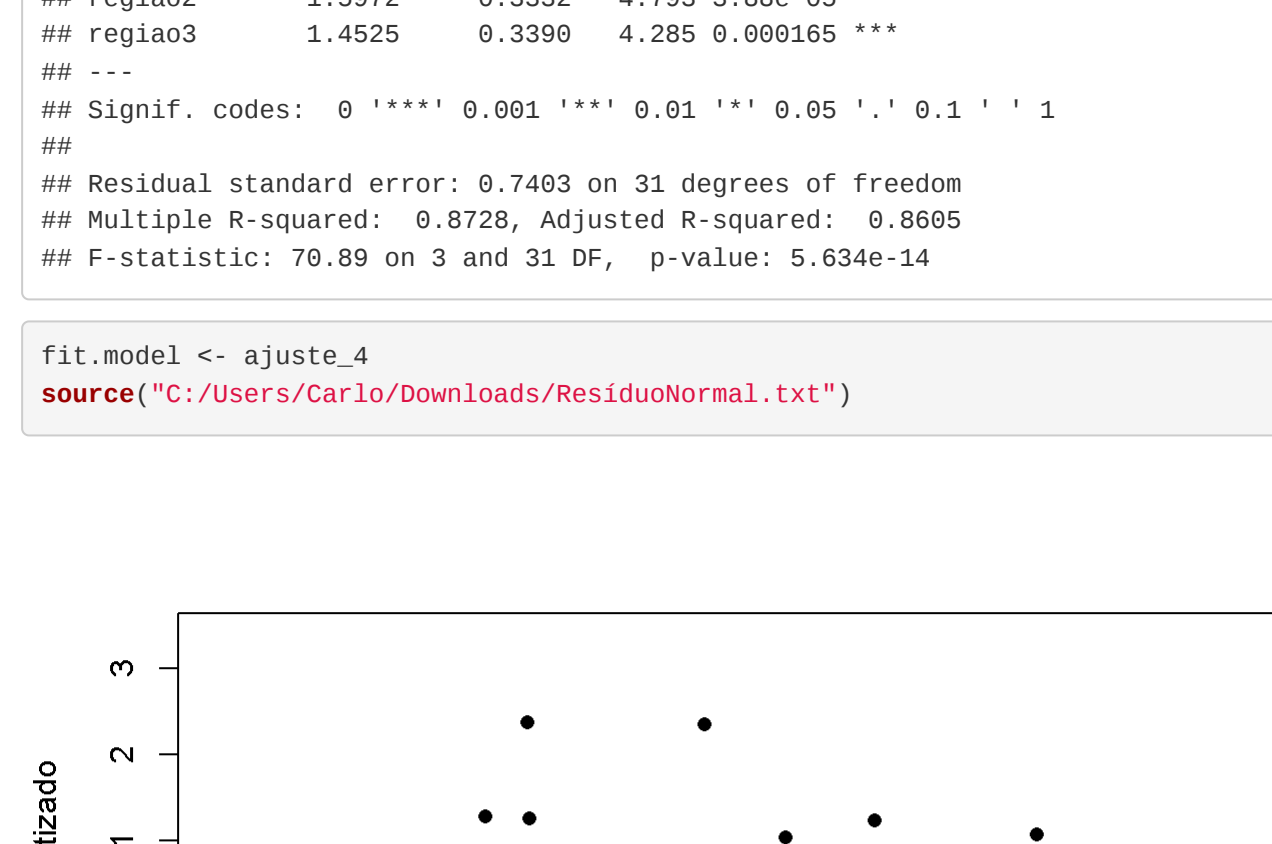
```
##
## Call:
## lm(formula = qualidade ~ sabor + regiao, data = wine, subset =
## -c(12,
## 20, 25))
##
## Residuals:
## Min 1Q Median 3Q Max
## -1.39316 -0.55059 -0.08442 0.45408 1.97559

## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 7.2699 0.6681 10.882 4.10e-12 ***
## sabor 1.0754 0.1470 7.314 1.22e-08 ***
## regiao2 -1.5072 0.3332 -4.503 0.00046 ***
## regiao3 1.4525 0.3390 4.285 0.00035 ***

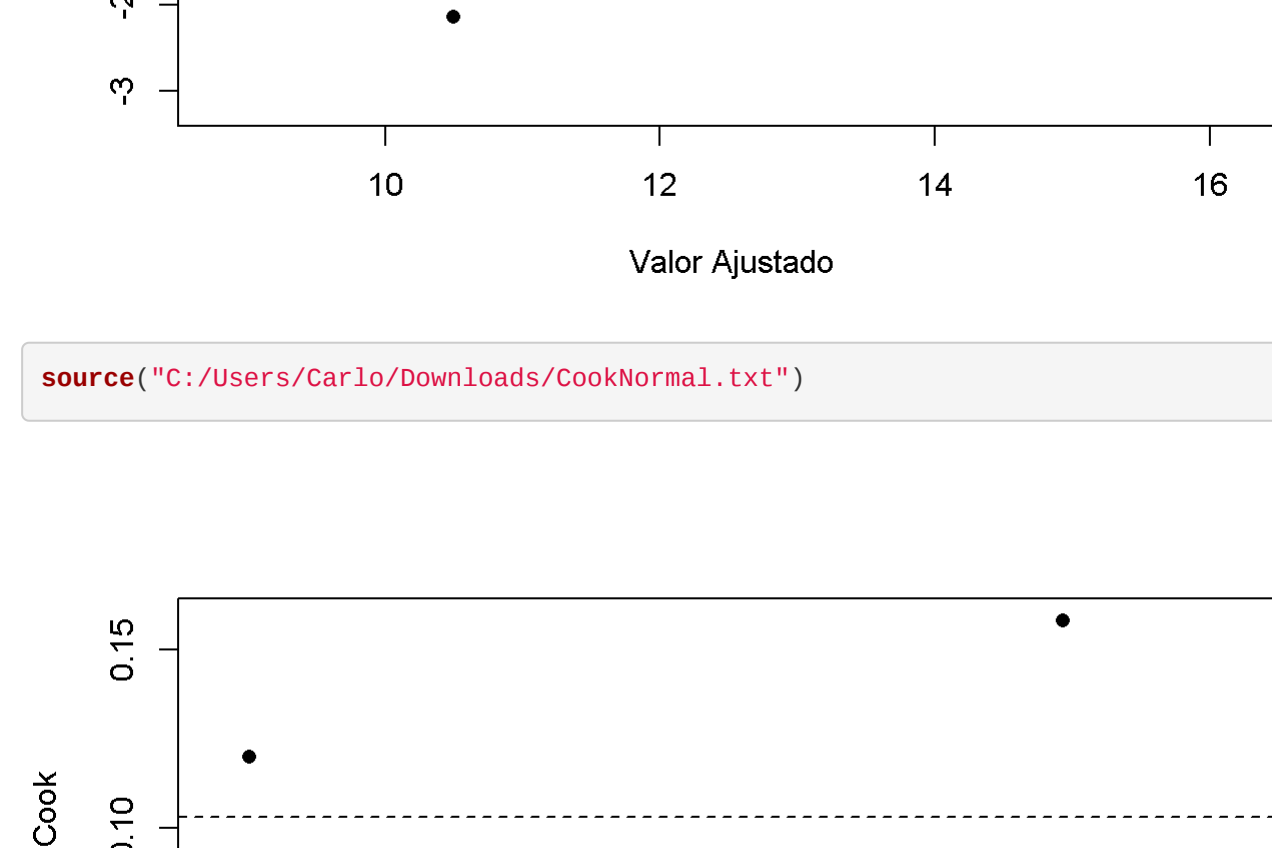
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 0.8704 on 32 degrees of freedom
## Multiple R-squared: 0.8405, Adjusted R-squared: 0.8256
## F-statistic: 96.23 on 3 and 32 DF, p-value: 7.448e-13
```

```
fit.model <- ajuste_4
source("C:/Users/Carlo/Downloads/Residuonormal.txt")
```



```
source("C:/Users/Carlo/Downloads/CookNormal.txt")
```



Aplicando um novo modelo e removendo os pontos discrepantes identificados no ajuste_5 <- lm(qualidade ~ sabor + regiao, data = wine, subset = -c(12, 20, 25))

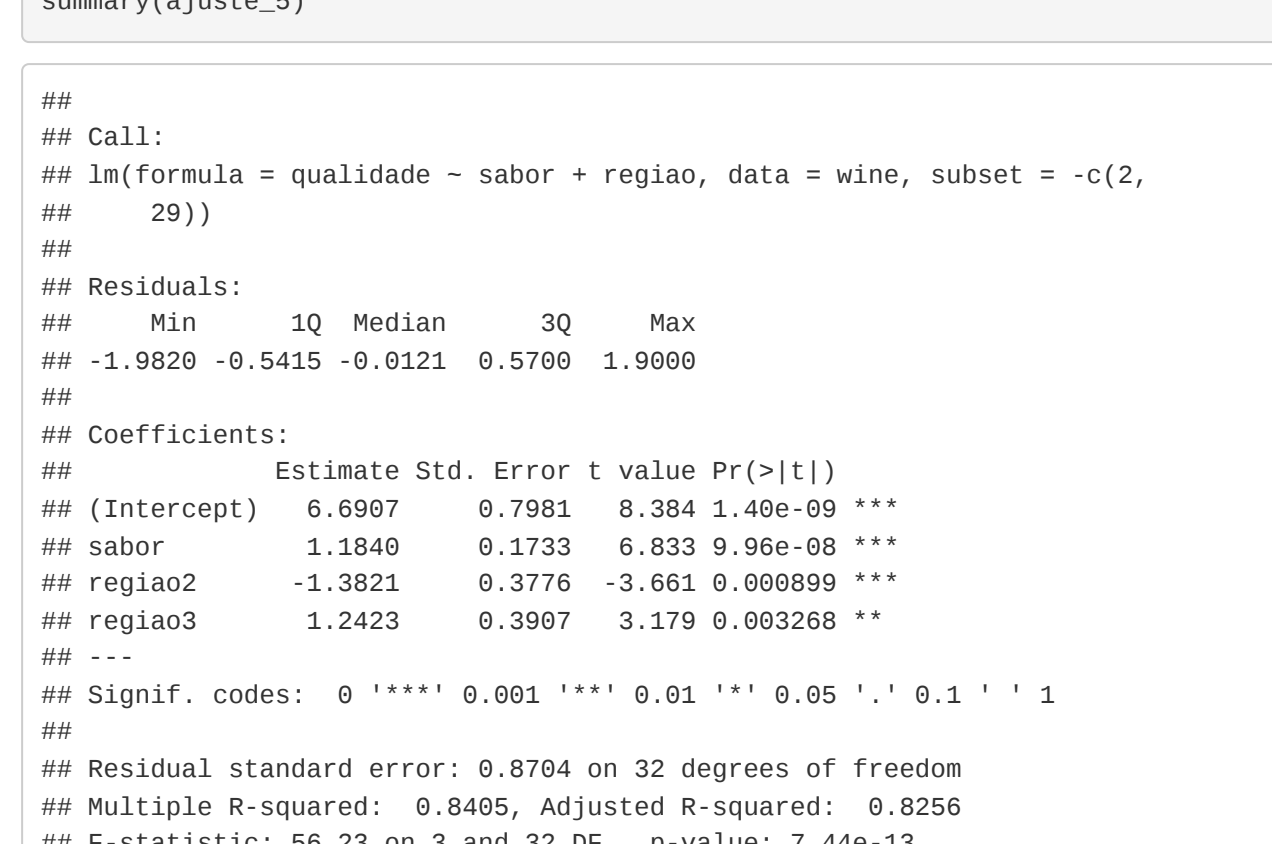
```
##
## Call:
## lm(formula = qualidade ~ sabor + regiao, data = wine, subset =
## -c(12,
## 20, 25))
##
## Residuals:
## Min 1Q Median 3Q Max
## -1.0820 -0.5415 -0.0121 0.5700 1.9000

## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) 6.9007 0.7081 9.744 1.40e-09 ***
## sabor 1.1848 0.1733 6.833 0.00000 ***
## regiao2 -1.3821 0.3777 -3.661 0.00089 ***
## regiao3 1.2423 0.3997 3.119 0.00386 ***

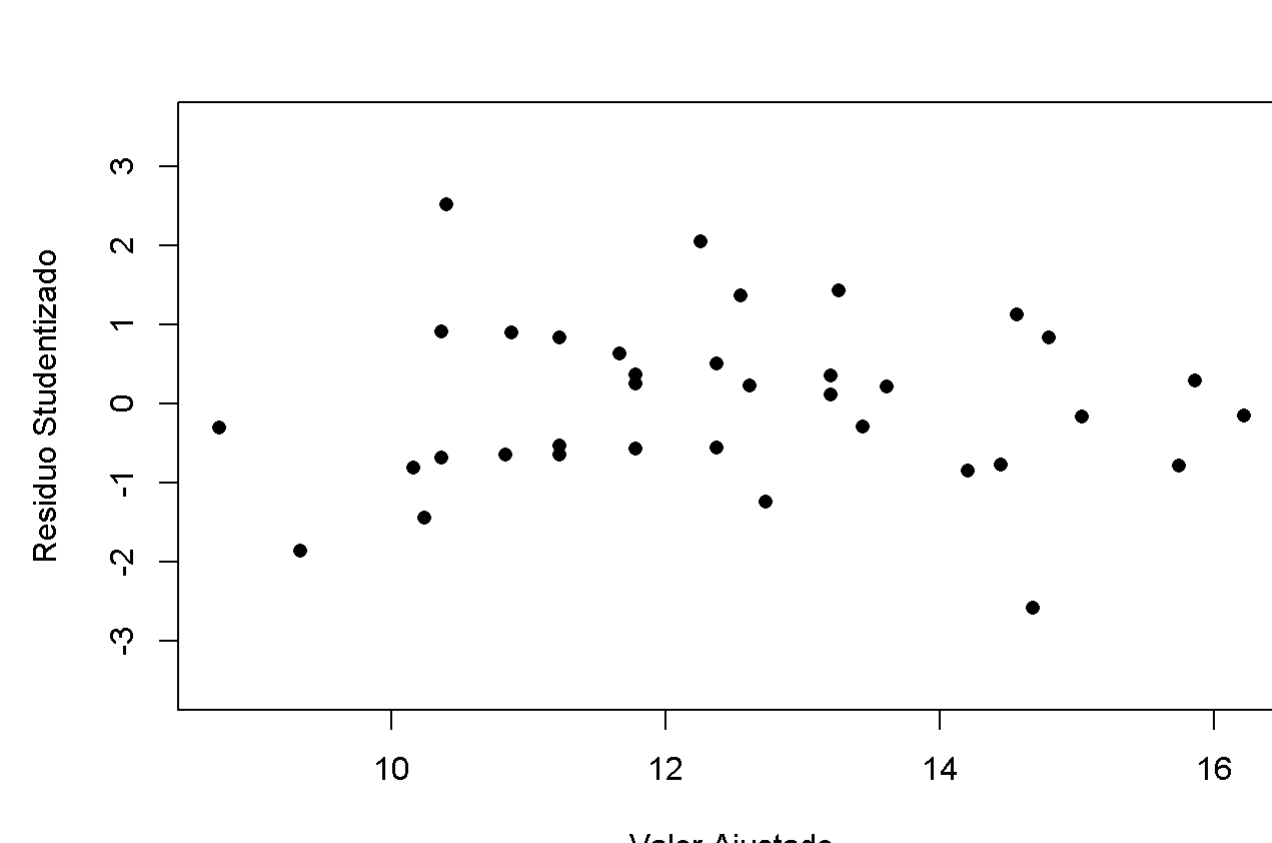
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## Residual standard error: 0.8704 on 32 degrees of freedom
## Multiple R-squared: 0.8405, Adjusted R-squared: 0.8256
## F-statistic: 96.23 on 3 and 32 DF, p-value: 7.448e-13
```

```
fit.model <- ajuste_5
source("C:/Users/Carlo/Downloads/Residuonormal.txt")
```



```
source("C:/Users/Carlo/Downloads/CookNormal.txt")
```



Executando os modelos: Ajuste_3, Ajuste_4 e Ajuste_5, foi identificando pontos discrepantes, ao identifica-los no ajuste_3 os mesmos foram removidos, então executamos o modelo ajuste_4, obtivemos o mesmo resultado em que, na análise de resíduos eles não foram distribuídos no intervalo [-3, 3] e ao plotar a distância de cook ainda mostraram pontos discrepantes no modelo. Por fim, identificamos esses pontos discrepantes realizamos um novo ajuste (ajuste_5) removendo os referidos pontos(= 29), ao executar novamente a distância de cook e resíduos identificamos mais pontos discrepantes. Por tanto, não há diferença entre os modelos executados se removemos tais pontos discrepantes encontrados. Portanto, realizando o teste de R-squared do ajuste_3 e ajuste_4, foi apontado um aumento após a remoção destes pontos (12, 20, 25), ou seja, ajuste_3 < ajuste_4. Ao executar um novo modelo ajuste_5, resultou em uma diminuição dos resultados do R-squared entre o ajuste_4 e ajuste_5, após a remoção destes pontos (2 e 29), logo ajuste_4 > ajuste_5 e o ajuste_3 < ajuste_5. Contudo, entre o modelo ajuste_3 e 4, a diferença entre eles, considerando o modelo ajuste_4 como modelo final definido.

```
# Tabela anova para ajuste_4
anova_4 <- anova(ajuste_4); anova_4
```

```
## Analysis of Variance Table
##
## Response: qualidade
## Df Sum Sq Mean Sq F value Pr(>F)
## sabor 1 83.248 83.248 151.878 1.751e-13 ***
## regiao 2 33.328 16.664 30.402 4.937e-08 ***
## Residuals 31 16.932 0.546
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

O modelo final (ajuste_4) teve 35 observações presentes, onde o seu p-valor é significativo.