

Lista 7 [AD-UFPE-2019]

Antonio Fernandes

28 de maio de 2019

Conteúdo

Apresentação	1
Questão 1	1
a) Análise descritiva de todas as variáveis da base de dados	1
b) Modelo de regressão	7
Resultado da regressão	8
Capacidade explicativa do modelo	8
Ajuste do modelo	8
c) Modelo de regressão com mais de uma VI	11
Resultados do modelo	12
Capacidade explicativa do modelo	12
Comparando resultados dos modelos	12

Apresentação

Este documento apresenta as respostas da lista de exercícios 7 da disciplina de Análise de dados.

O link está disponível no GitHub: https://github.com/alvesat/AD_7

A lista envolve a execução de diversos modelos de regressão linear multivariada envolvendo um banco de dados específicos.

Questão 1

O primeiro passo para responder a questão é abrir o banco de dados contendo os dados necessários. É importante verificar que os dados estão no formato *.dta*, que é o formato relacionado ao *software* Stata. Devido a isso, teremos que abrir o pacote **haven** e executar o comando para abrir o banco.

```
library(haven)
```

```
## Warning: package 'haven' was built under R version 3.5.3
```

```
fair <- read_dta("~/Dados/Listas/AD_7/DATA_L7/fair.dta")
```

a) Análise descritiva de todas as variáveis da base de dados

Vamos verificar o nome das variáveis presentes no banco

```
names(fair)
```

```
## [1] "YEAR"      "VOTE"      "PARTY"     "PERSON"    "DURATION"  "WAR"
## [7] "GROWTH"    "INFLATION" "GOODNEWS"
```

Podemos observar que o banco apresenta 9 variáveis: *Year*, *Vote*, *Party*, *Person*, *Duration*, *War*, *Growth*, *Inflation* e *Good news*.

Um outro passo envolvendo o processo de descrição das variáveis é identificar a estrutura de cada variável:

```
str(fair)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':   32 obs. of  9 variables:
## $ YEAR      : num  1880 1884 1888 1892 1896 ...
##   .. attr(*, "format.stata")= chr "%8.0g"
## $ VOTE       : num  50.2 49.8 50.4 48.3 47.8 ...
##   .. attr(*, "format.stata")= chr "%9.0g"
## $ PARTY      : num  -1 -1 1 -1 1 -1 -1 -1 -1 1 ...
##   .. attr(*, "format.stata")= chr "%8.0g"
## $ PERSON     : num   0 0 1 1 0 1 0 0 1 1 ...
##   .. attr(*, "format.stata")= chr "%8.0g"
## $ DURATION   : num   1.75 2 0 0 0 0 1 1.25 1.5 0 ...
##   .. attr(*, "format.stata")= chr "%9.0g"
## $ WAR        : num   0 0 0 0 0 0 0 0 0 0 ...
##   .. attr(*, "format.stata")= chr "%8.0g"
## $ GROWTH     : num   3.88 1.59 -5.55 2.76 -10.02 ...
##   .. attr(*, "format.stata")= chr "%9.0g"
## $ INFLATION  : num   1.974 1.055 0.604 2.274 3.41 ...
##   .. attr(*, "format.stata")= chr "%9.0g"
## $ GOODNEWS   : num    9 2 3 7 6 7 5 8 8 3 ...
##   .. attr(*, "format.stata")= chr "%8.0g"
```

É possível identificar que todas as variáveis presentes no banco são numéricas e que o banco possui 32 observações. Agora vamos fazer uma análise descritiva de cada variável:

Year

A variável ano apresenta os anos relacionados as outras variáveis. Por meio do comando `fivenum`, podemos obter um sumário da variável (valores mínimos, 1º quartil, mediana, 3º quartil e valor máximo). No caso dessa variável, o importante é verificarmos o começo da análise e o máximo.

```
fivenum(fair$YEAR)
```

```
## [1] 1880 1910 1942 1974 2004
```

Em relação a variável Year, percebemos que o primeiro ano do banco é 1880 e o ano final do banco é 2004.

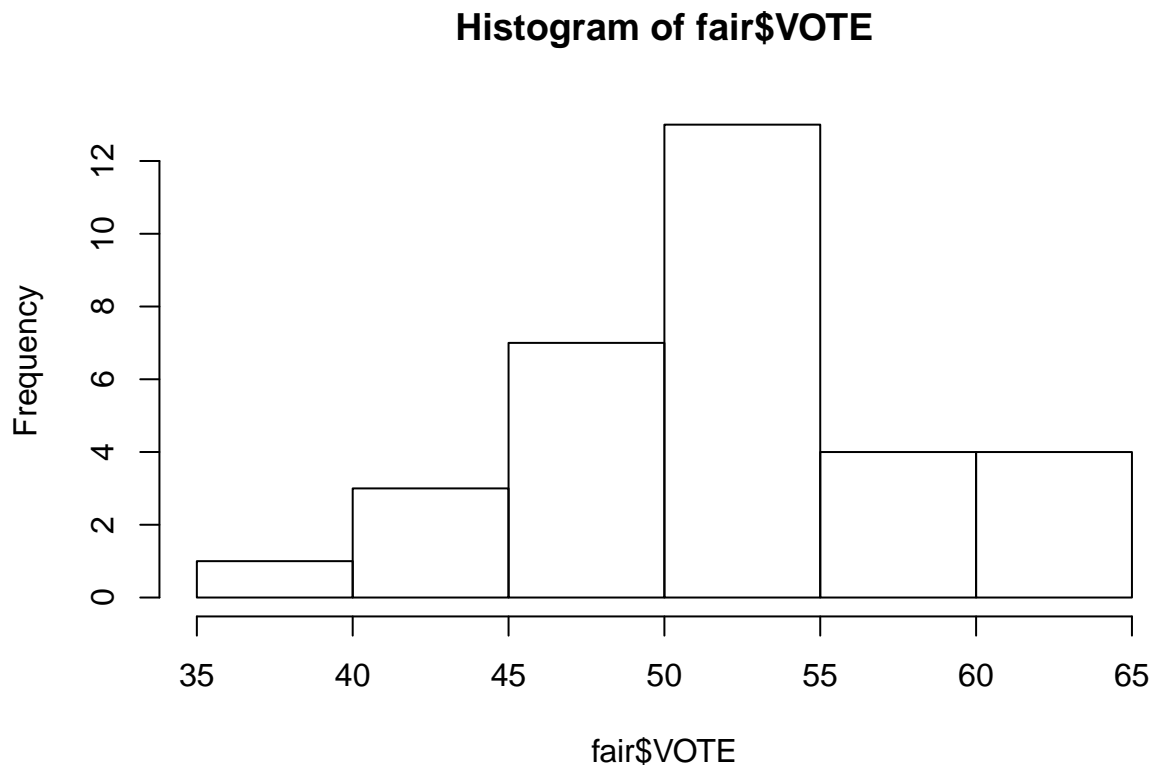
Vote

```
fivenum(fair$VOTE)
```

```
## [1] 36.1190 49.2720 52.0260 56.3815 62.4580
```

Em rela  o a vari  vel vote, a vari  vel representa a porcentagem de votos recebidas pelo partido do incumbente. Observando os valores obtidos do comando `fivenum`, o valor m  nimo da vari  vel foi de 36.12, o 1   quartil    de 49.27, a mediana    de 52.03, o 3      de 56.38 e o valor m  ximo    de 62.46. Podemos analisar a vari  vel vote por meio de um histograma:

```
hist(fair$VOTE)
```



Com o histograma, verificamos que os valores est  o mais concentrados entre 45 e 55, ou seja, no centro da distribui  o.

Party

A vari  vel party    bin  ria, apresentando como valores -1 e 1.

```
ftable(fair$PARTY)
```

```
## -1  1  
##  
## 18 14
```

Com o comando `ftable`, podemos fazer uma tabela de frequ  ncia da vari  vel party e identificar que das 32 observa  es do banco, 18 s  o do valor -1 e 14 s  o do valor 1.

Person

A variável person também é binária, apresentando como valores 1 e 0.

```
fable(fair$PERSON)
```

```
##    0    1  
##  
##   13   19
```

O resultado mostra que 13 observações apresentam o valor 0 e 19 observações apresentam o valor 1.

Duration

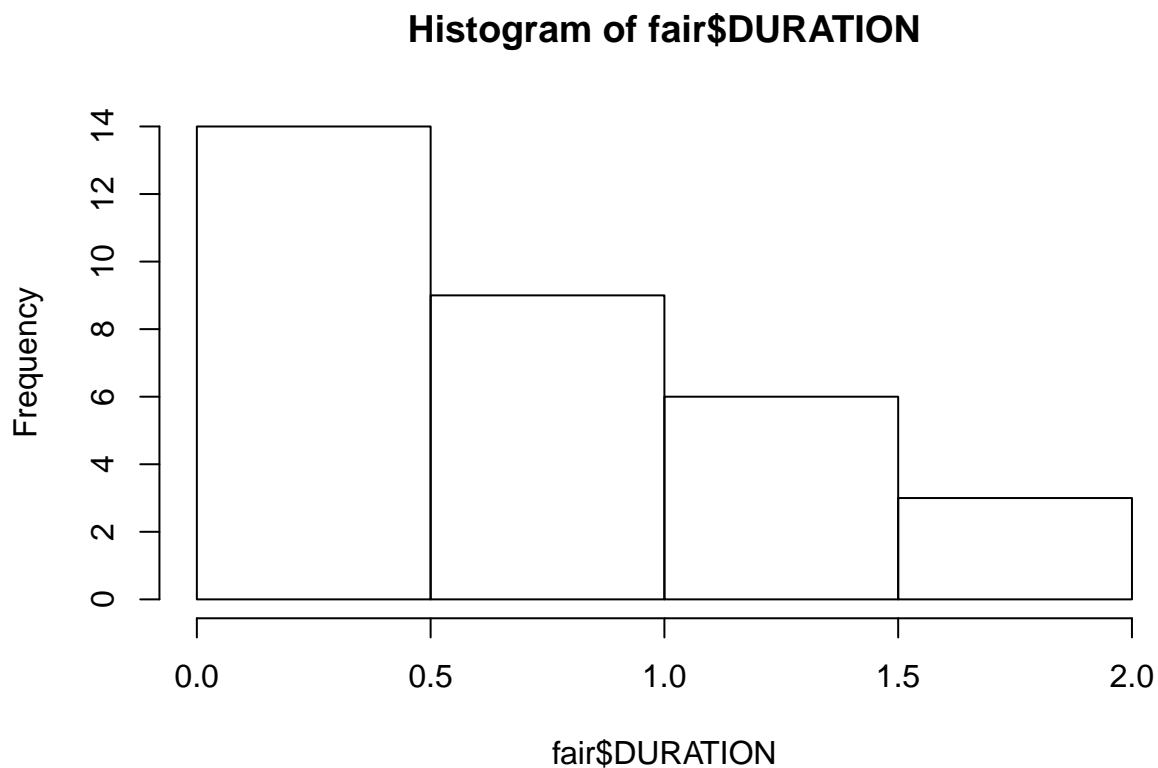
A variável duration apresenta os valores de 2 a 0. Com o comando `fivenum()` possível obter uma análise descritiva das variáveis:

```
fivenum(fair$DURATION)
```

```
## [1] 0.00 0.00 1.00 1.25 2.00
```

O valor mínimo e do primeiro quartil é 0, a mediana é 1, o 3º quartil é 1.25 e o valor máximo da distribuição é 2.

```
hist(fair$DURATION)
```



Pelo histograma, vemos que a maior parte das observações são 0 e 1 (25)

War

A variável war é uma variável binária que é melhor observada por meio do comando `fable`:

```
fable(fair$WAR)
```

```
##    0    1  
##  
##  29    3
```

é possível verificar que 29 casos possuem valor 0 e 3 casos apresentam valor 1

Growth

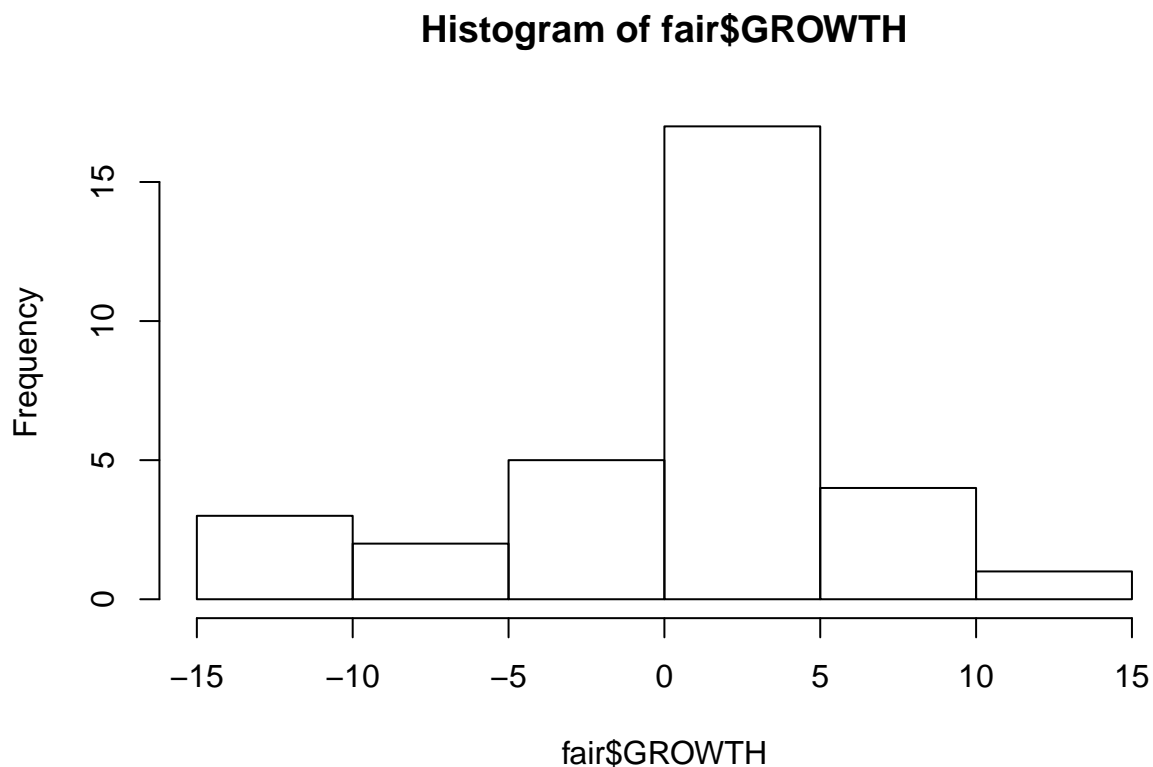
A variável growth contém informações acerca da variação de crescimento do PIB por ano:

```
fivenum(fair$GROWTH)
```

```
## [1] -14.557 -1.923  2.245  4.095 11.677
```

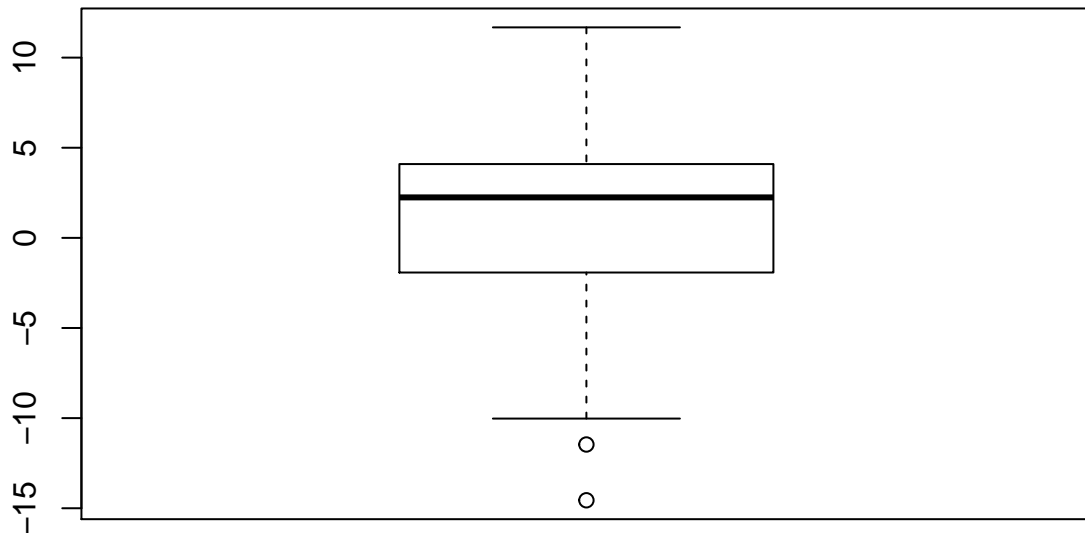
A descrição por meio do comando `fivenum` mostra que o valor mínimo da distribuição foi uma contração de 14.56 no PIB, enquanto que o 1º quartil apresenta um valor de -1.92, a mediana é de 2.24, o 3º quartil é 4.1 e o valor máximo da distribuição é de 11.68. Por meio de um histograma é possível uma melhor visualização da distribuição.

```
hist(fair$GROWTH)
```



Com o histograma, percebemos que boa parte da distribui??o est? localizada no centro, com casos entre -5 e 5. Com um box plot, podemos identificar a presen?a ou n?o de outliers na distribui??o.

```
boxplot(fair$GROWTH)
```



Por meio do Boxplot, podemos verificar que existem dois valores que s?o outliers: -14.55 e -11.46.

Inflation

Vamos verificar as estat?sticas principais da vari?vel inflation:

```
fivenum(fair$INFLATION)
```

```
## [1] 0.0000 1.3545 2.1590 3.3715 7.9260
```

O valor m?nimo da distribui??o ? 0, o 1? quartil ? 1.35, a mediana ? de 2.16, o 3? quartil ? de 3.38 e o valor m?ximo ? de 7.93.

Goodnews

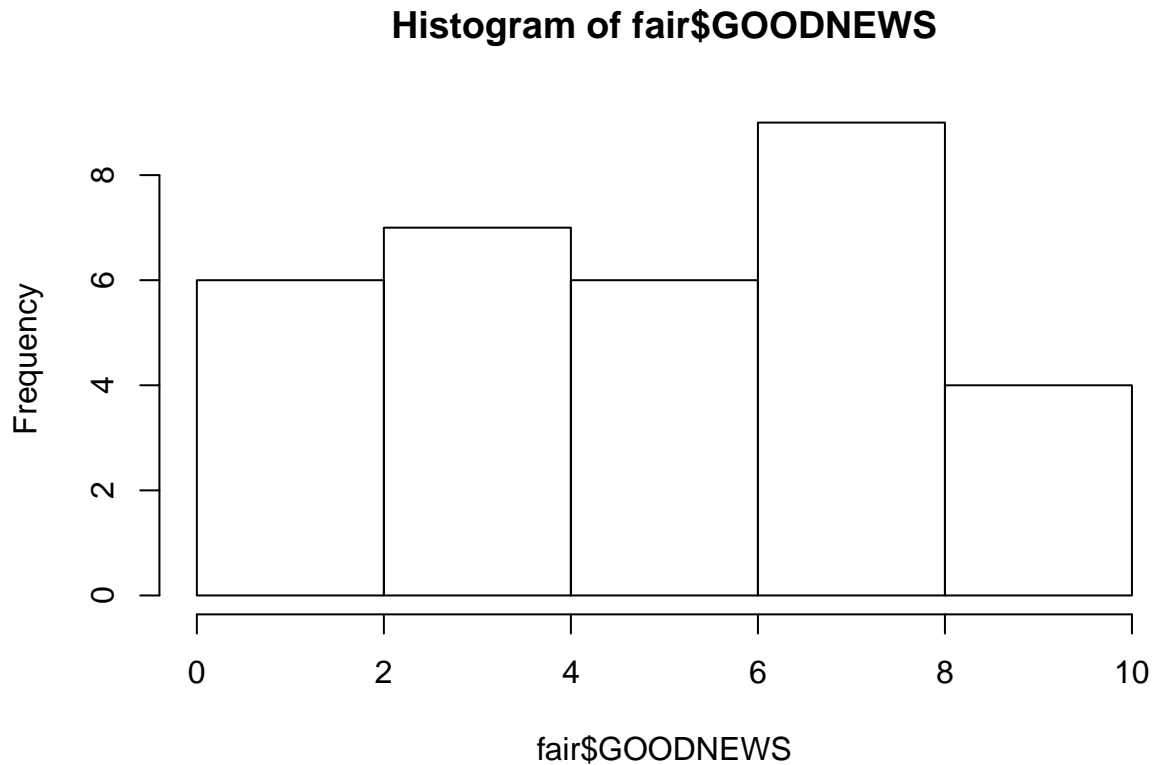
Abaixo vemos os cinco valores da vari?vel goodnews

```
fivenum(fair$GOODNEWS)
```

```
## [1] 0.0 3.5 5.0 7.5 10.0
```

O valor m nimo da vari vel   0, o 1  quartil   3.5, a mediana   5, o 3  quartil   7.5 e o valor m ximo   10.

```
hist(fair$GOODNEWS)
```



O histograma mostra que a distribui  o dos valores   bem balanceada.

b) Modelo de regress o

Para executar o modelo, com a vari vel Vote como VD e Growth como VI, vamos utilizar o seguinte comando:

```
Linear <- lm(VOTE ~ GROWTH, data = fair)
summary(Linear)
```

```
##
## Call:
## lm(formula = VOTE ~ GROWTH, data = fair)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.2487 -3.3330 -0.4282  3.1425  9.7286
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   51.8598     0.8817  58.821  < 2e-16 ***
## GROWTH         0.6536     0.1607   4.068 0.000316 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.955 on 30 degrees of freedom
## Multiple R-squared:  0.3555, Adjusted R-squared:  0.3341
## F-statistic: 16.55 on 1 and 30 DF,  p-value: 0.0003165
```

Resultado da regressão

Como é possível observar no resultado, a variação de uma unidade na VI (Growth) leva a um aumento de 0.65 na VD (Votes). Nesse caso, o modelo de regressão utilizado é bivariado, dado a presença de apenas duas variáveis no modelo.

Em relação aos resíduos do modelo, o valor mínimo é de -8.25, o 1º quartil é de -3.33, a mediana é de -0.423, o 3º quartil é de 3.14 e o valor máximo é de 9.72.

Quando a VI assume o valor de 0, espera-se que o valor da VD seja de 51.86 (Valor do intercepto). Já em relação a VI, além do valor do coeficiente, o erro padrão foi de 0.16, o teste-f que testa a hipótese nula de que não há relação entre as variáveis apresentou um p-valor menor que 0.001. Ou seja, rejeita-se a hipótese nula de que não há relação entre as variáveis.

Capacidade explicativa do modelo

No que se refere a capacidade explicativa do modelo, o R² ajustado foi de 0.33, significando que o modelo consegue explicar 33% da variação da variável dependente. O erro padrão do resíduo foi de 4.955.

Ajuste do modelo

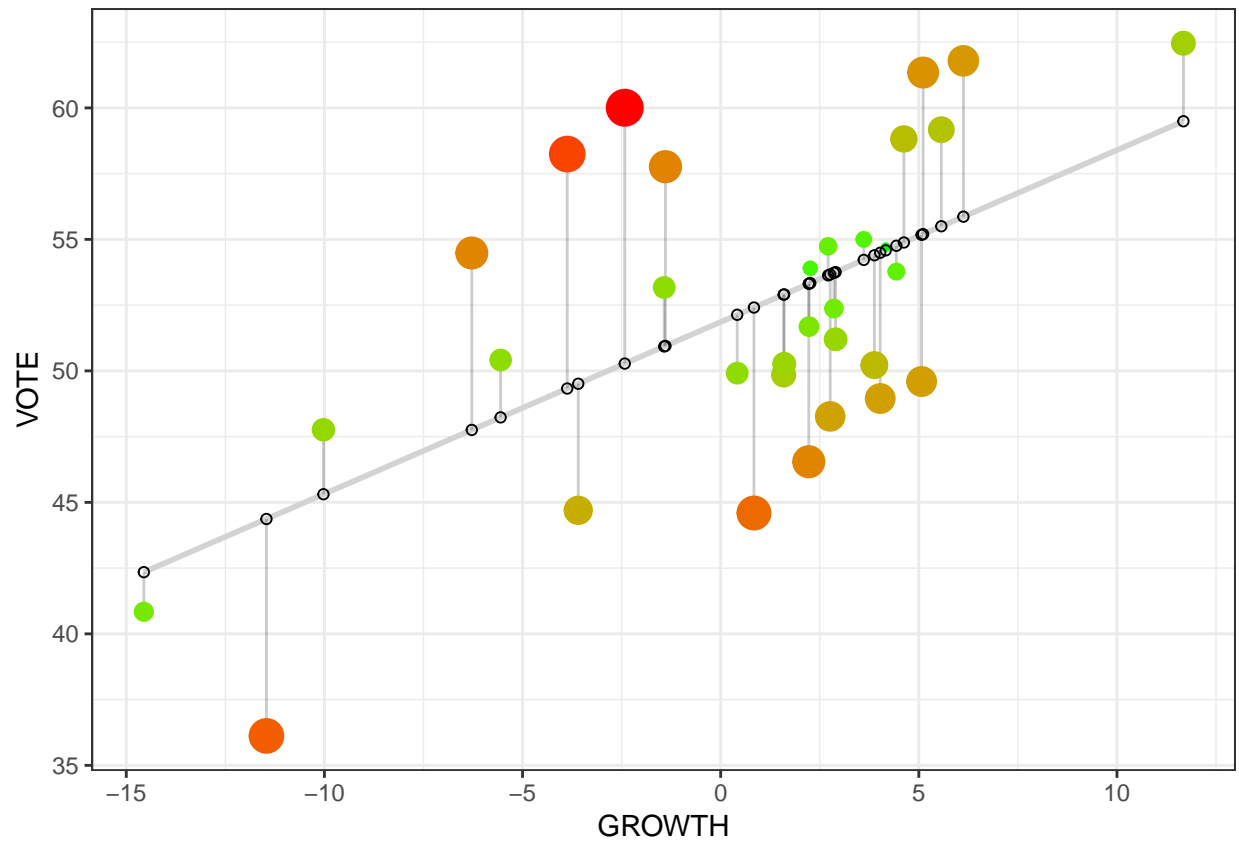
Já no tocante ao ajuste, podemos analisar alguns gráficos e verificar a adequabilidade do modelo:

```
library(ggplot2)
```

```
## Warning: package 'ggplot2' was built under R version 3.5.3
```

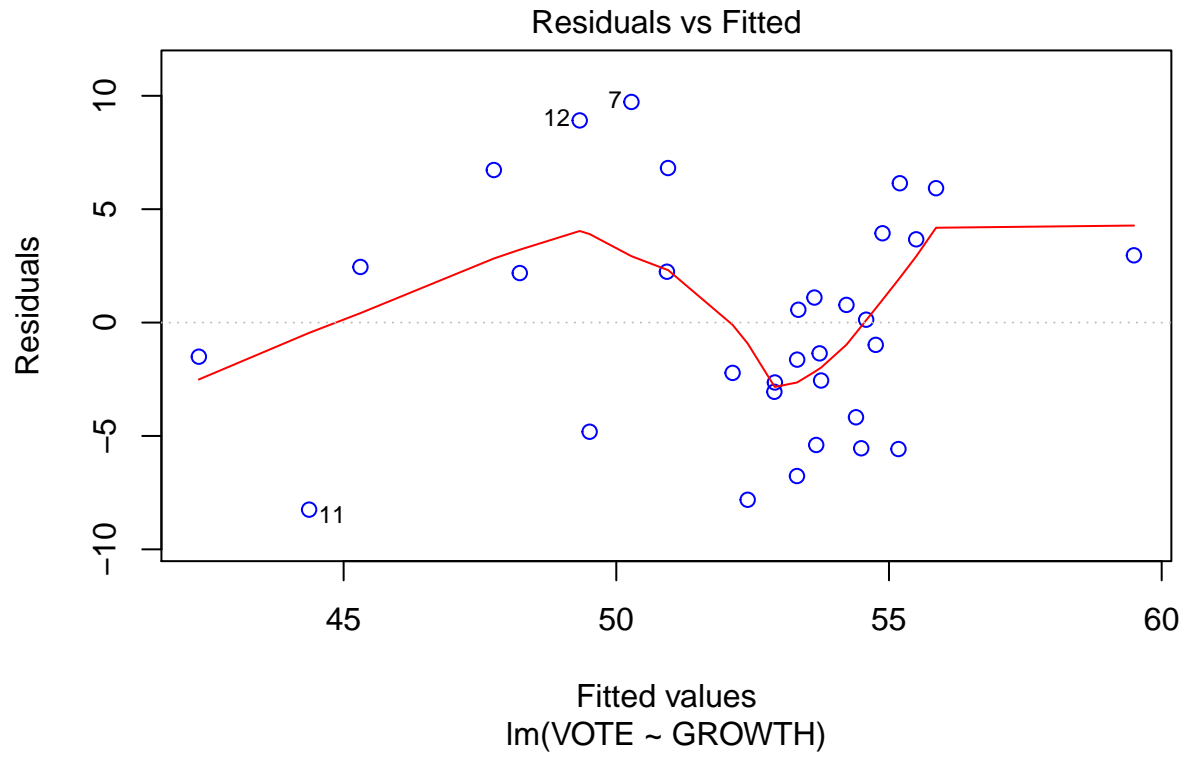
O gráfico abaixo apresenta o tamanho do resíduo por meio do tamanho do ponto e cor (quanto mais vermelho mais longe e quanto mais verde menor o resíduo). O tamanho do resíduo é a distância entre o ponto e a linha de regressão.

```
fair$predicted <- predict(Linear)
fair$residuals <- residuals(Linear)
ggplot(fair, aes(x = GROWTH, y = VOTE)) +
  geom_smooth(method = "lm", se = FALSE, color = "lightgrey") +
  geom_segment(aes(xend = GROWTH, yend = predicted), alpha = .2) +
  geom_point(aes(color = abs(residuals), size = abs(residuals))) +
  scale_color_continuous(low = "green", high = "red") +
  guides(color = FALSE, size = FALSE) +
  geom_point(aes(y = predicted), shape = 1) +
  theme_bw()
```

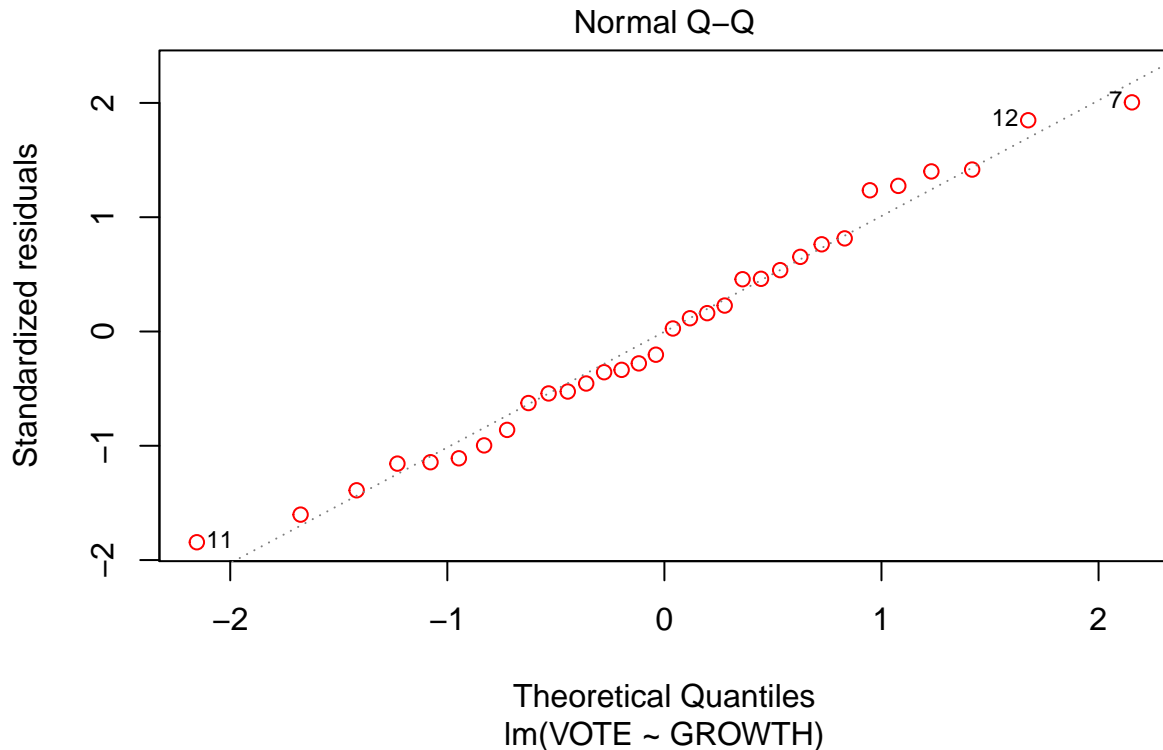
Podemos por meio do gráfico observar a distribuição dos resíduos do modelo. Por meio do gráfico de resíduos, é possível verificar a adequabilidade dos dados. Dado que o resíduo é aquela parte do modelo que não é explicada, espera-se que este não possua nem um padrão. Em suma, a distribuição do resíduo deve ser aleatória.

```
plot(Linear, which=1, col=c("blue"))
```



Quando observamos o gráfico, é possível verificar uma certa concentração de observações na parte um pouco a direita do centro do gráfico, o que pode apresentar uma não normalidade dos dados. Por isso será executado um gráfico de quantis (quantile-quantile plot). Caso os resíduos sigam uma linha no gráfico, é uma boa indicação de uma distribuição normal.

```
plot(Linear, which=2, col=c("red"))
```



Observa-se que os res?duos seguem a linha da regress?o de maneira bastante alinhada.

Com base nessas an?lises, ? poss?vel concluir que o modelo ? adequado.

c) Modelo de regress?o com mais de uma VI

Para executar o modelo de regress?o com VOTES como VD com Growth e Goodnews como VI, executamos o seguinte comando:

```
Linear_2 <- lm(VOTE ~ GROWTH + GOODNEWS, data = fair)
summary(Linear)
```

```
##
## Call:
## lm(formula = VOTE ~ GROWTH, data = fair)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.2487 -3.3330 -0.4282  3.1425  9.7286
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   51.8598     0.8817  58.821  < 2e-16 ***
## GROWTH         0.6536     0.1607   4.068 0.000316 ***
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.955 on 30 degrees of freedom
## Multiple R-squared:  0.3555, Adjusted R-squared:  0.3341
## F-statistic: 16.55 on 1 and 30 DF,  p-value: 0.0003165
```

Resultados do modelo

A variável adicionada ao modelo é a de GoodNews, relacionada a quantidade de notícias boas na área econômica.

Ao analisar o modelo, percebemos que ambas as variáveis independentes apresentam um efeito sobre a variável dependente. Entretanto, a variável Growth apresenta um P-valor menor que 0.001 enquanto que a variável Goodnews é apresenta um efeito significativo sobre a VD quando consideramos um P-valor menor que 0.05. Mesmo assim, ambas as variáveis estão abaixo do patamar de 0.05, permitindo rejeitar a hipótese nula de que não há relação entre as variáveis. Observando o coeficiente da variável Growth, podemos interpretar que, o aumento de 1% na variável do PIB leva a um aumento de 0.57 na VD (mantendo todo o resto constante). Já em relação a GOODNEWS, o aumento em 25% de notícias econômicas positivas, leva ao aumento de 0.71 na VD (mantendo todo o resto constante).

Os valores do resíduo mostram um resíduo mínimo de -8.3125, com o 1º quartil igual a -3.9191, Mediana de 0.4876, 3º quartil de 3.05 e valor máximo de 9.68.

Capacidade explicativa do modelo

No que se refere a capacidade explicativa, o modelo apresenta um R² ajustado de 0.43, significando uma explicação de 43% de variação na VD. O erro padrão do resíduo foi de 4.596.

RMSE

Para analisar o RMSE (Raiz quadrada do erro médio) será necessário primeiro abrir o pacote `sjstats`

```
library("sjstats")
```

```
## Warning: package 'sjstats' was built under R version 3.5.3
```

Assim é possível verificar o valor da RMSE:

```
rmse(Linear)
```

```
## [1] 4.797286
```

O RMSE é calculado para mensurar a diferença entre os valores preditos pelo modelo e os valores observados. Também é conhecido como erro padrão do modelo e é calculado elevando ao quadrado cada erro do modelo, somando-os, dividindo pelo número de casos e obtendo a raiz quadrada. Nesse caso, o valor obtido foi de 4.37. O RMSE tem a mesma unidade que a VD e quanto menor, melhor.

Comparando resultados dos modelos

Para poder comparar os resultados desse modelo com o anterior, é necessário padronizar os coeficientes de ambos os modelos. Para isso será utilizado o pacote `lm.beta`:

```
library('lm.beta')
```

```
## Warning: package 'lm.beta' was built under R version 3.5.2
```

```
lm.beta(Linear)
```

```
##  
## Call:  
## lm(formula = VOTE ~ GROWTH, data = fair)  
##  
## Standardized Coefficients:  
## (Intercept)      GROWTH  
##  0.0000000    0.5962706
```

```
lm.beta(Linear_2)
```

```
##  
## Call:  
## lm(formula = VOTE ~ GROWTH + GOODNEWS, data = fair)  
##  
## Standardized Coefficients:  
## (Intercept)      GROWTH      GOODNEWS  
##  0.0000000    0.5227285    0.3373268
```

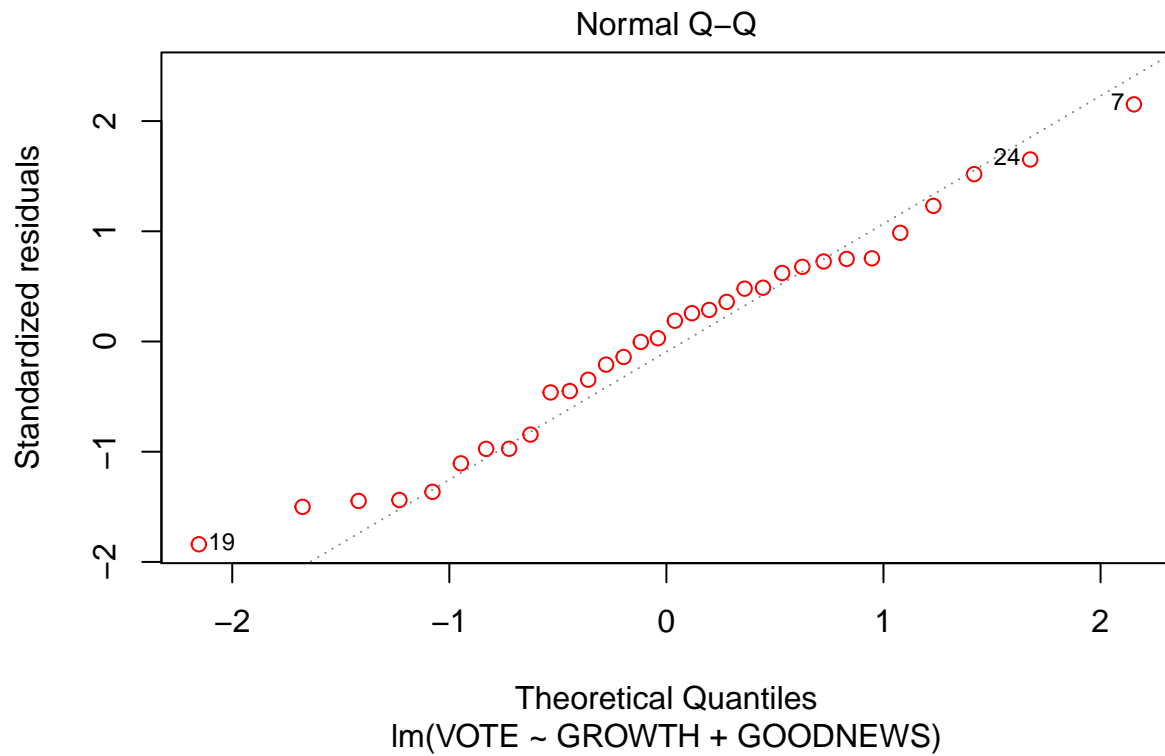
Por meio dos coeficientes padronizados, é possível comparar os modelos. No primeiro modelo executado, verifica-se que para o aumento de 1 desvio padrão da VI Growth, espera-se o aumento de .60 desvio padrão na VD.

Já no segundo modelo, o aumento de 1 desvio padrão da VI Growth leva a um aumento de 0.52 desvio padrão na VD. Ou seja, houve uma redução no efeito da VI Growth do modelo 1 para o modelo 2.

Ajuste do modelo

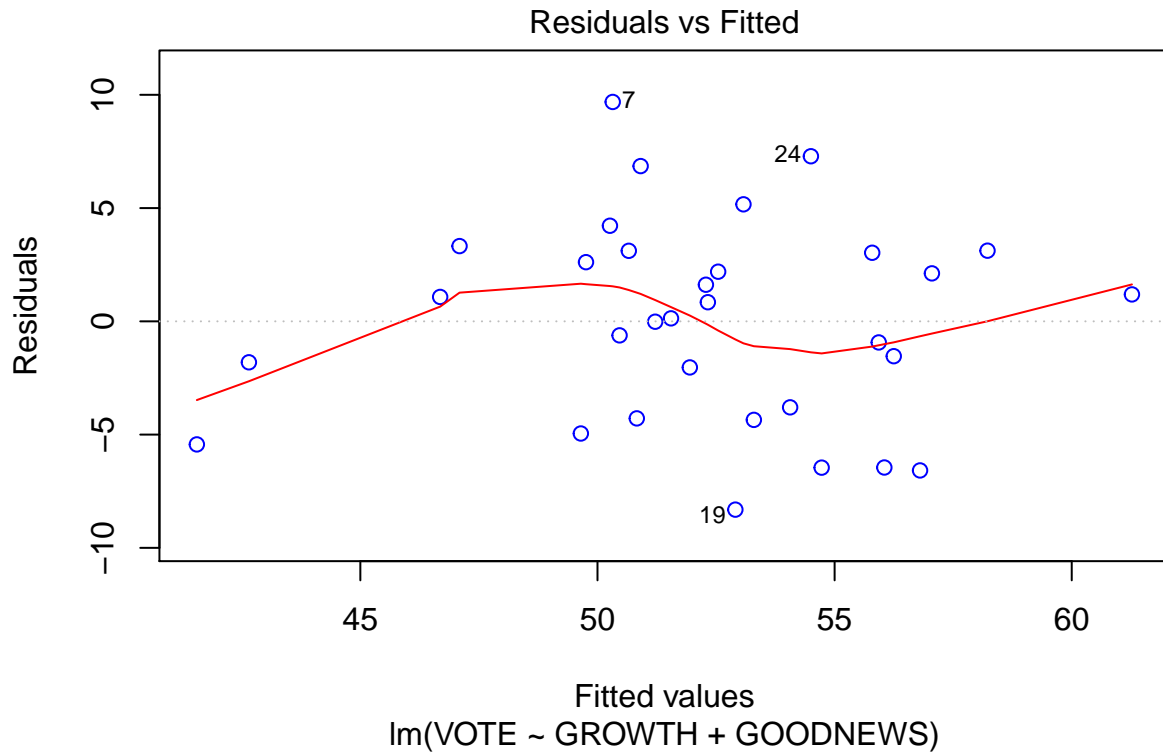
Para verificar o ajuste do modelo, será analisado os resíduos graficamente.

```
plot(Linear_2, which=2, col=c("red"))
```



O gráfico Q-Q mostra um ajuste adequado do modelo dado a tendência dos pontos em relação a linha de regressão.

```
plot(Linear_2, which=1, col=c("blue"))
```



J? o gr?fico dos res?duos apresenta valores tanto acima quanto abaixo da linha, representando tamb?m um bom ajuste do modelo.

Por fim, o teste de Shapiro permite verificar a normalidade da distribui??o:

```
res <- residuals(Linear_2)

shapiro.test(res)

##
##  Shapiro-Wilk normality test
##
## data:  res
## W = 0.97866, p-value = 0.7598
```

Os resultados do teste apontam um p-valor de 0.76, levando a n?o rejeitar a hip?tese nula de que os dados testados n?o est?o normalmente distribu?dos.