

Análise de dados (EDA)

Cristian Villegas

2023-05-18

Contents

1	Leitura de dados	1
2	Alguns resumos dos dados	2
3	Gráficos de interesse	3
4	Transformando os dados	4
4.1	Gráficos antes de transformar dados	4
4.2	Transformando dados (seguindo o feito pelo Jalmar)	8
5	Gráficos de perfis	9
5.1	Sexo Female == 0	9
5.2	Sexo Male == 1	9
6	Ajuste de modelos	10
6.1	cloglog	10
6.2	logit	12
6.3	probit	14
6.4	cauchit	16
7	GEE	18
7.1	independence	18
7.2	exchangeable	19
7.3	unstructured	20

1 Leitura de dados

```
library(tidyverse)
library(hnp)

dados <- read.csv("Arthritis.txt", sep="",
                  stringsAsFactors=TRUE)

dados$id <- 1:nrow(dados)

dados <- tibble(dados)
dados

## # A tibble: 51 x 9
##   Sex      Age Group Week0 Week1 Week5 Week9 Week13   id
```

```
##      <fct> <int> <fct> <int> <int> <int> <int> <int> <int>
## 1 M      48 A      1      1      1      1      1      1
## 2 M      29 A      1      1      1      1      1      2
## 3 M      59 P      1      1      1      1      1      3
## 4 F      56 P      1      1      1      1      1      4
## 5 M      33 P      1      1      1      1      1      5
## 6 M      61 P      1      1      0      1      1      6
## 7 M      63 A      0      0      1     NA     NA      7
## 8 M      57 P      1      0      1      1      1      8
## 9 M      47 P      1      1      1      0      1      9
## 10 F     42 A      0      0      1     NA      0     10
## # ... with 41 more rows
```

2 Alguns resumos dos dados

```
dados %>%
  group_by(Sex) %>%
  summarise( n = n())
```

```
## # A tibble: 2 x 2
##   Sex      n
##   <fct> <int>
## 1 F      13
## 2 M      38
```

```
dados %>%
  group_by(Sex) %>%
  summarise(media_Sex = mean(Age))
```

```
## # A tibble: 2 x 2
##   Sex  media_Sex
##   <fct>    <dbl>
## 1 F      51.8
## 2 M      50.2
```

```
dados %>%
  group_by(Group) %>%
  summarise( n = n())
```

```
## # A tibble: 2 x 2
##   Group      n
##   <fct> <int>
## 1 A      27
## 2 P      24
```

```
dados %>%
  group_by(Group) %>%
  summarise( media_Age = mean(Age))
```

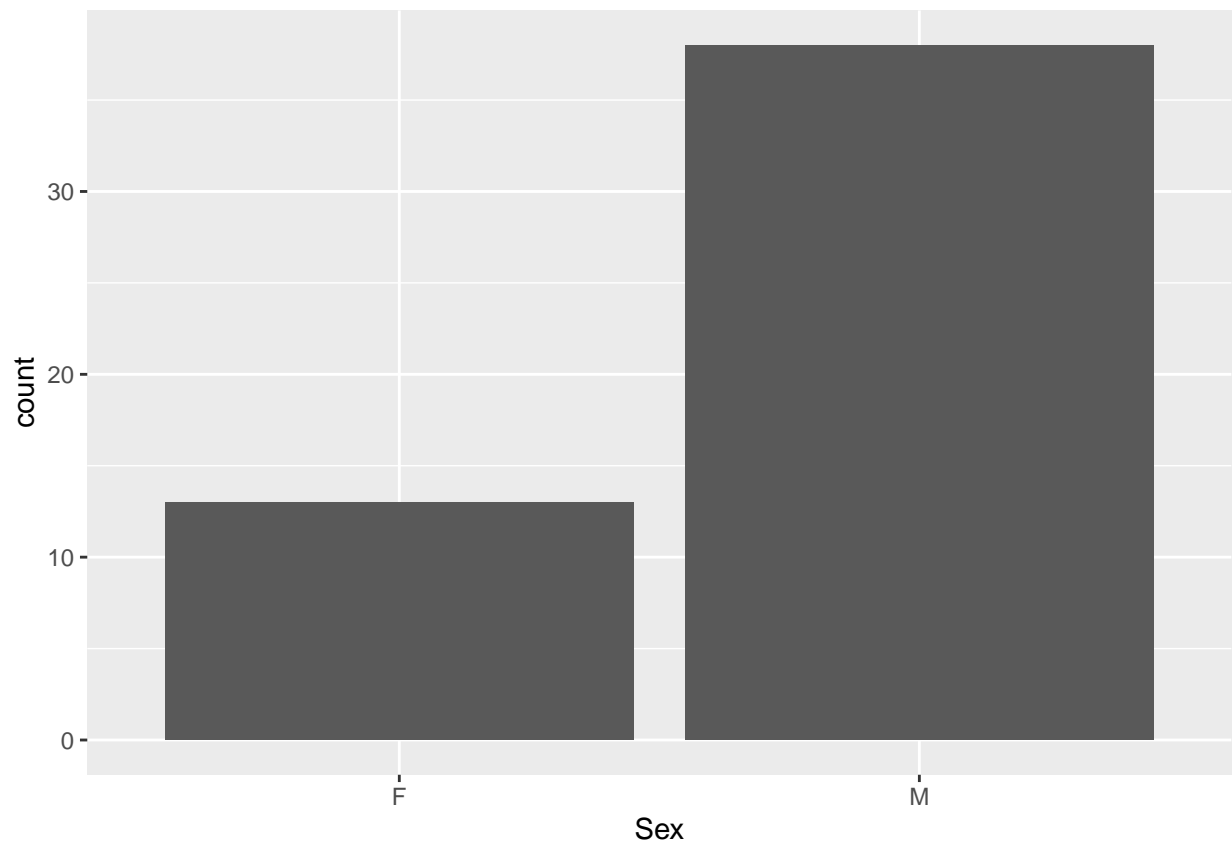
```
## # A tibble: 2 x 2
##   Group media_Age
##   <fct>    <dbl>
## 1 A      51.0
## 2 P      50.2
```

```
dados %>%
  group_by(Sex, Group) %>%
  summarise( n = n())
```

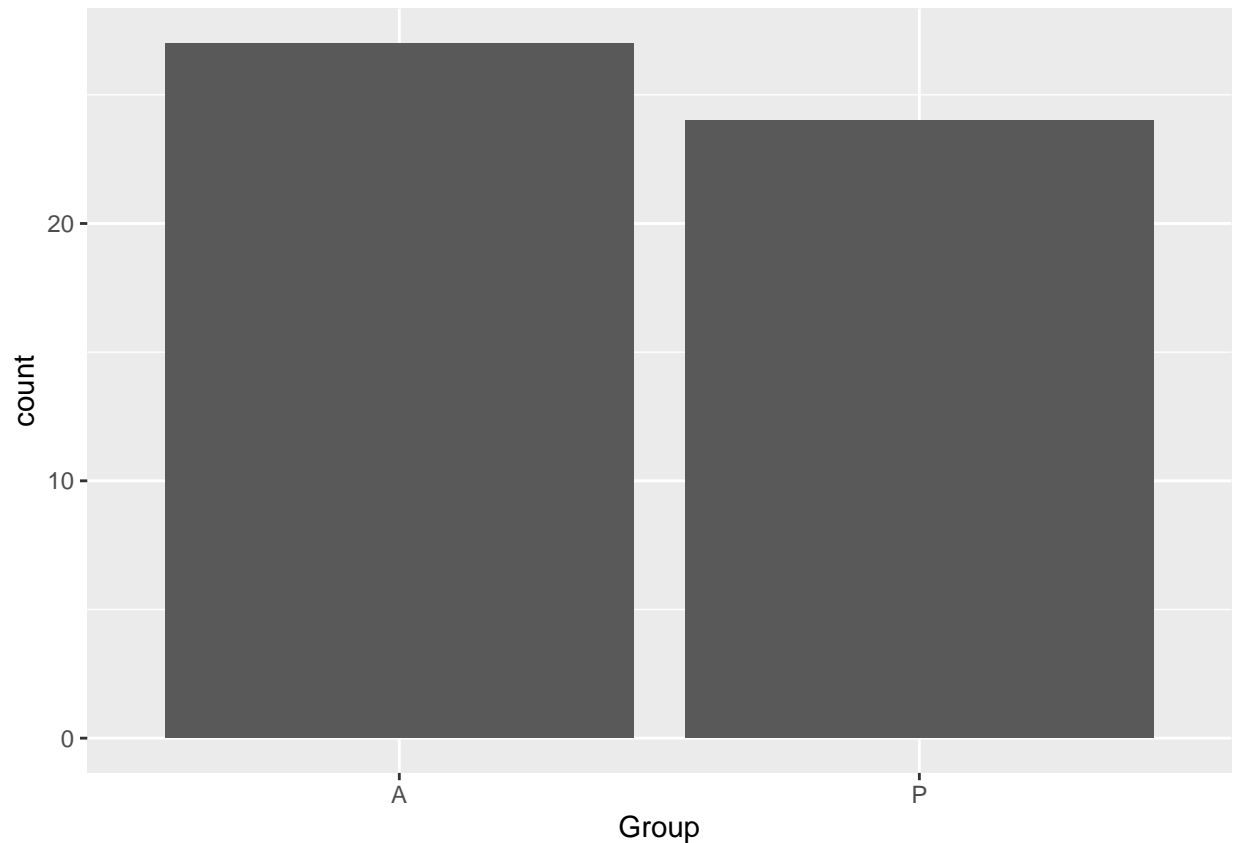
```
## # A tibble: 4 x 3
## # Groups:   Sex [2]
##   Sex   Group     n
##   <fct> <fct> <int>
## 1 F     A       7
## 2 F     P       6
## 3 M     A      20
## 4 M     P      18
```

3 Gráficos de interesse

```
ggplot(dados, aes(x = Sex)) +
  geom_bar()
```



```
ggplot(dados, aes(x = Group)) +
  geom_bar()
```

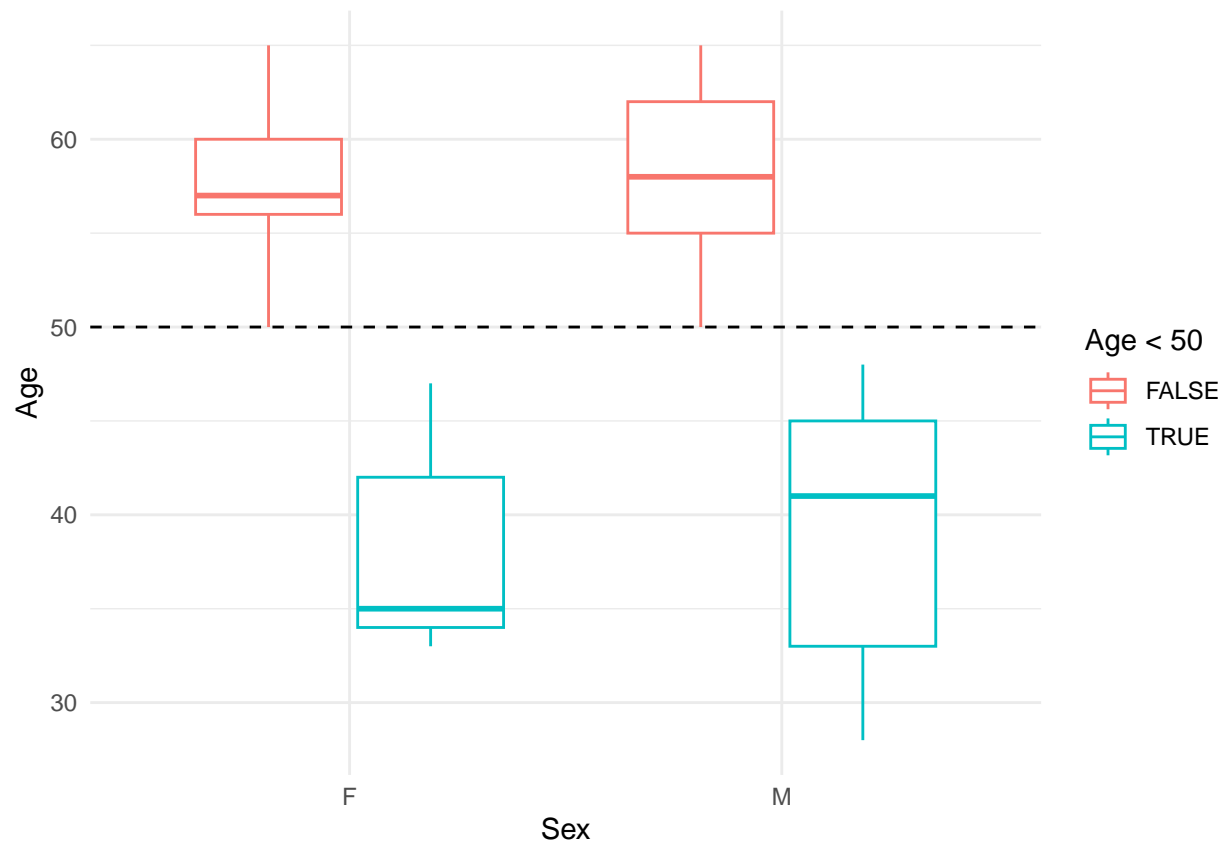


4 Transformando os dados

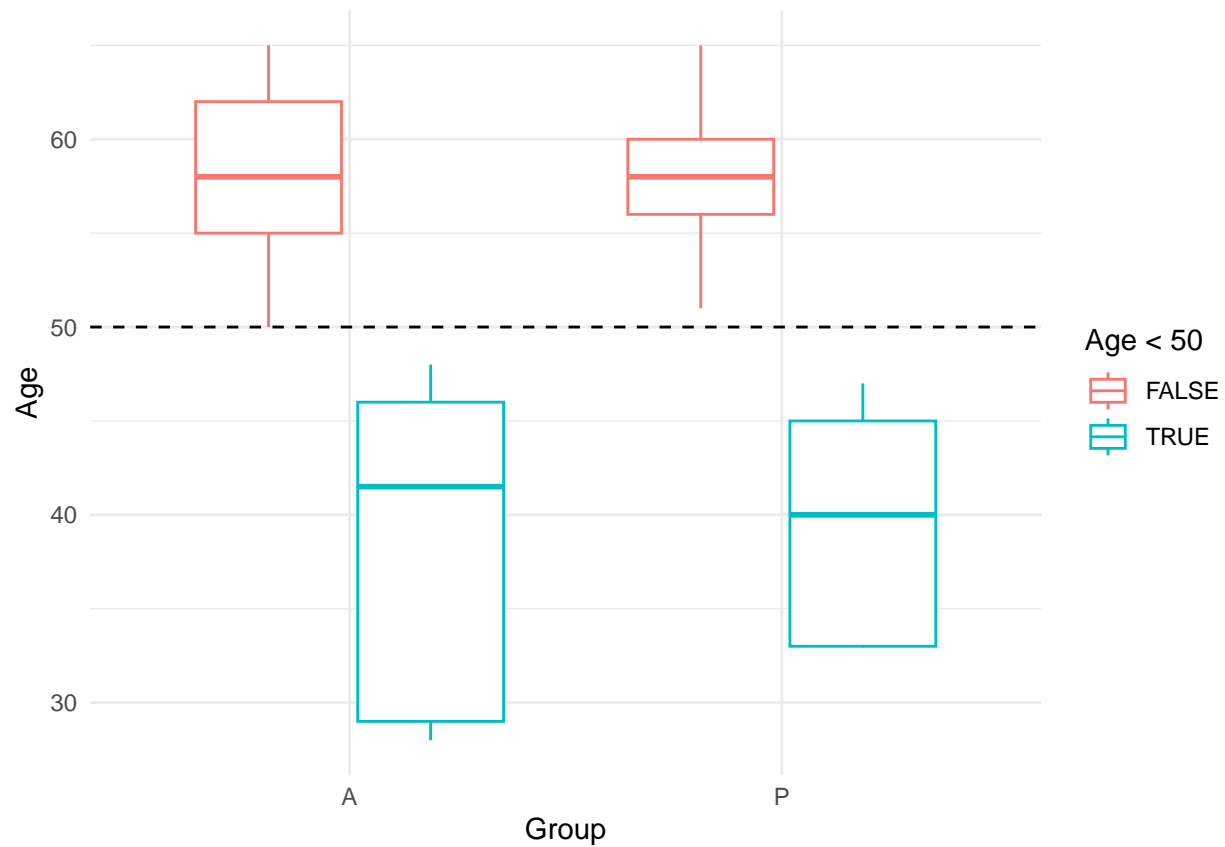
```
dados_longos<- dados %>%
  pivot_longer(
    cols = starts_with("Week"),
    names_to = "week",
    names_prefix = "Week",
    values_to = "Y",
    values_drop_na = TRUE
  )
```

4.1 Gráficos antes de transformar dados

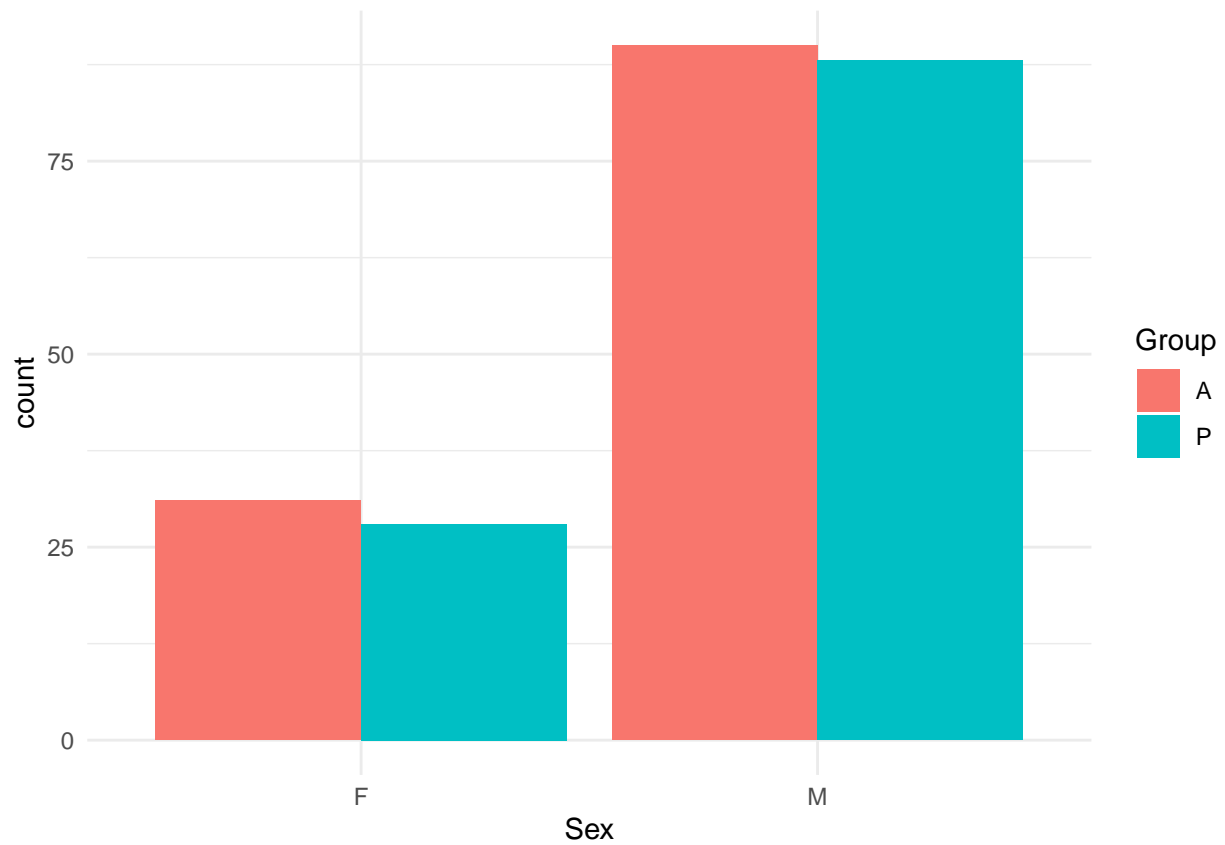
```
ggplot(dados_longos, aes(Sex, Age, col = Age < 50)) +
  geom_boxplot()+
  geom_hline(yintercept = 50, col = "black", linetype = 2)+
  theme_minimal()
```



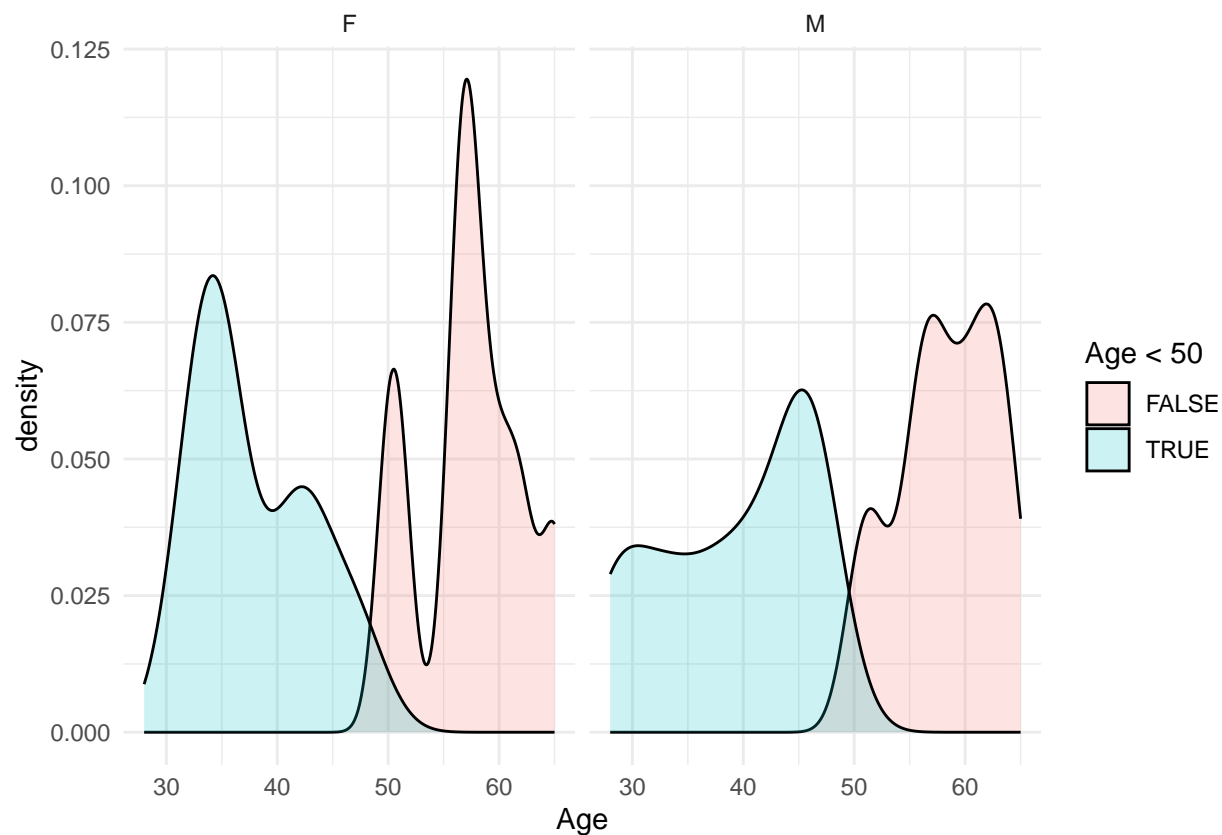
```
ggplot(dados_longos, aes(Group, Age, col = Age < 50)) +  
  geom_boxplot()+  
  geom_hline(yintercept = 50, col = "black", linetype = 2)+  
  theme_minimal()
```



```
ggplot(dados_longos, aes(x = Sex, fill = Group)) +  
  geom_bar(position=position_dodge()) +  
  theme_minimal()
```



```
ggplot(dados_longos, aes(Age, fill = Age < 50)) +  
  #geom_histogram(fill = "yellow",  
  #               aes(y = after_stat(density)), bins=6)+  
  geom_density(alpha=0.2)+  
  facet_wrap(~Sex)+  
  theme_minimal()
```



4.2 Transformando dados (segundo o feito pelo Jalmar)

```
dados_longos$Sex<- recode_factor(dados_longos$Sex, `F` = "0", `M` = "1")

dados_longos$Age<- factor(case_when(dados_longos$Age <50 ~ 1,
  dados_longos$Age >=50 ~ 0, .default = dados_longos$Age),
  levels = c(0, 1))

dados_longos$Group<- recode_factor(dados_longos$Group, `P` = "0", `A` = "1")

dados_longos$week<- factor(dados_longos$week,
  levels = c(0, 1, 5, 9, 13) )

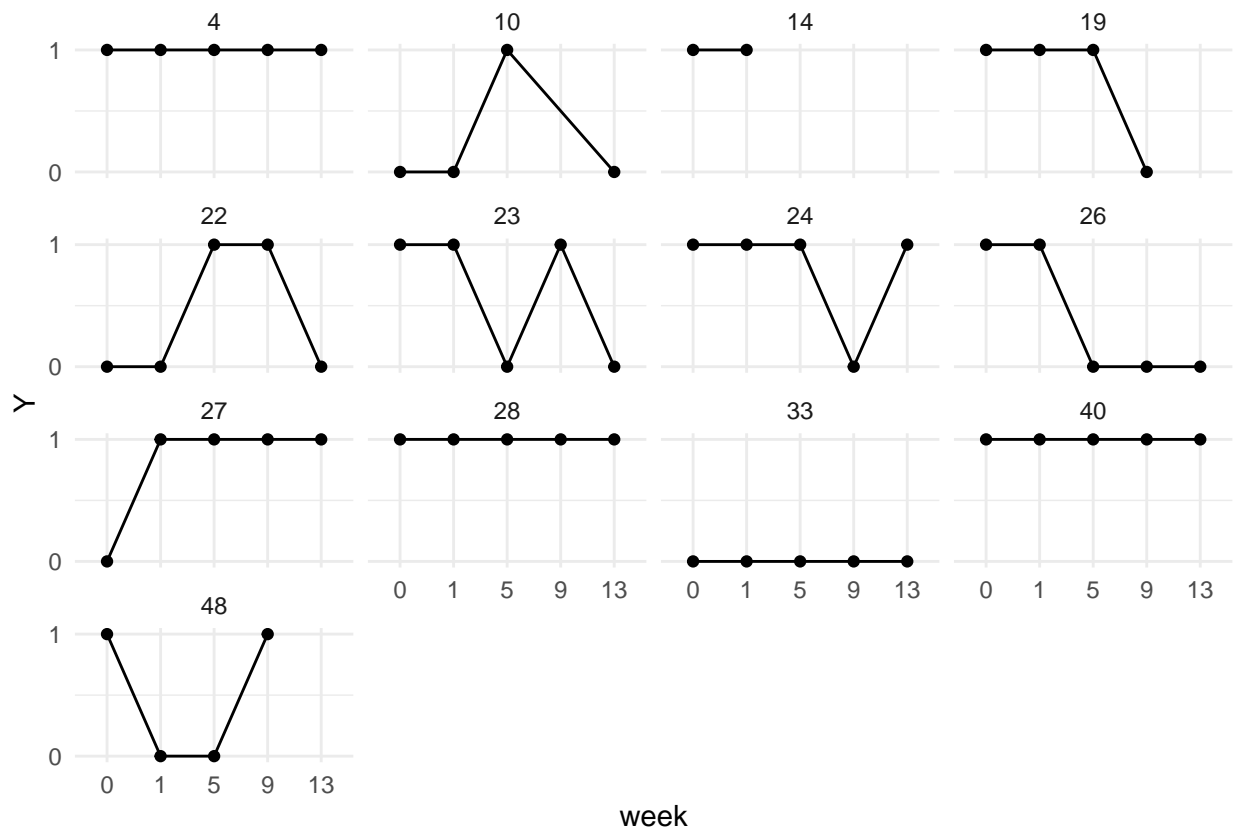
dados_longos %>%
  group_by(week) %>%
  summarise( n = n())
```

```
## # A tibble: 5 x 2
##   week      n
##   <fct> <int>
## 1 0       51
## 2 1       51
## 3 5       48
## 4 9       45
## 5 13      42
```


5 Gráficos de perfis

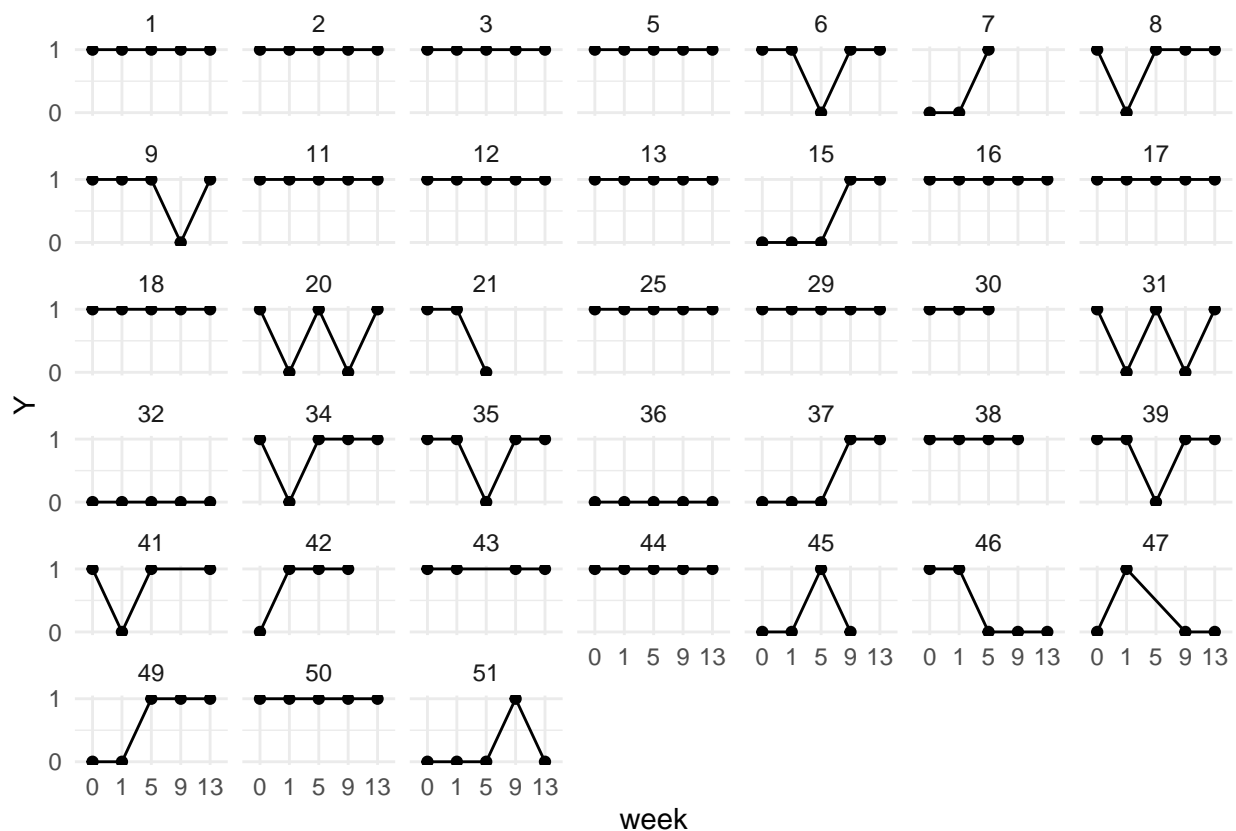
5.1 Sexo Female == 0

```
dados_longos %>% filter(Sex == "0") %>%  
  ggplot(aes(week, Y, group = id)) +  
  geom_point()+  
  geom_line()+  
  theme_minimal()+  
  scale_y_continuous(breaks = c(0,1))+  
  facet_wrap(~id)
```



5.2 Sexo Male == 1

```
dados_longos %>% filter(Sex == "1") %>%  
  ggplot(aes(week, Y, group = id)) +  
  geom_point()+  
  geom_line()+  
  theme_minimal()+  
  scale_y_continuous(breaks = c(0,1))+  
  facet_wrap(~id)
```



6 Ajuste de modelos

6.1 cloglog

```
modelo_cloglog<- glm(Y ~ Sex +
  Age +
  Group +
  as.numeric(week),
  family = binomial(link = "cloglog"),
  data= dados_longos)
modelo_cloglog$family
```

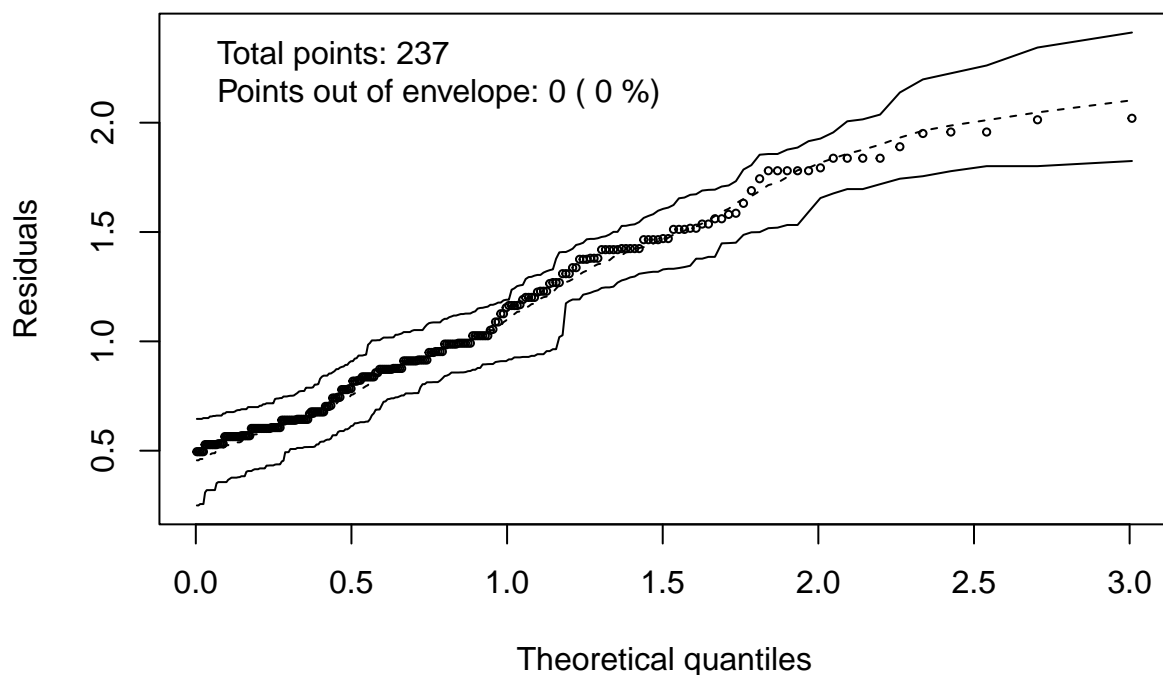
```
##
## Family: binomial
## Link function: cloglog
```

```
summary(modelo_cloglog)
```

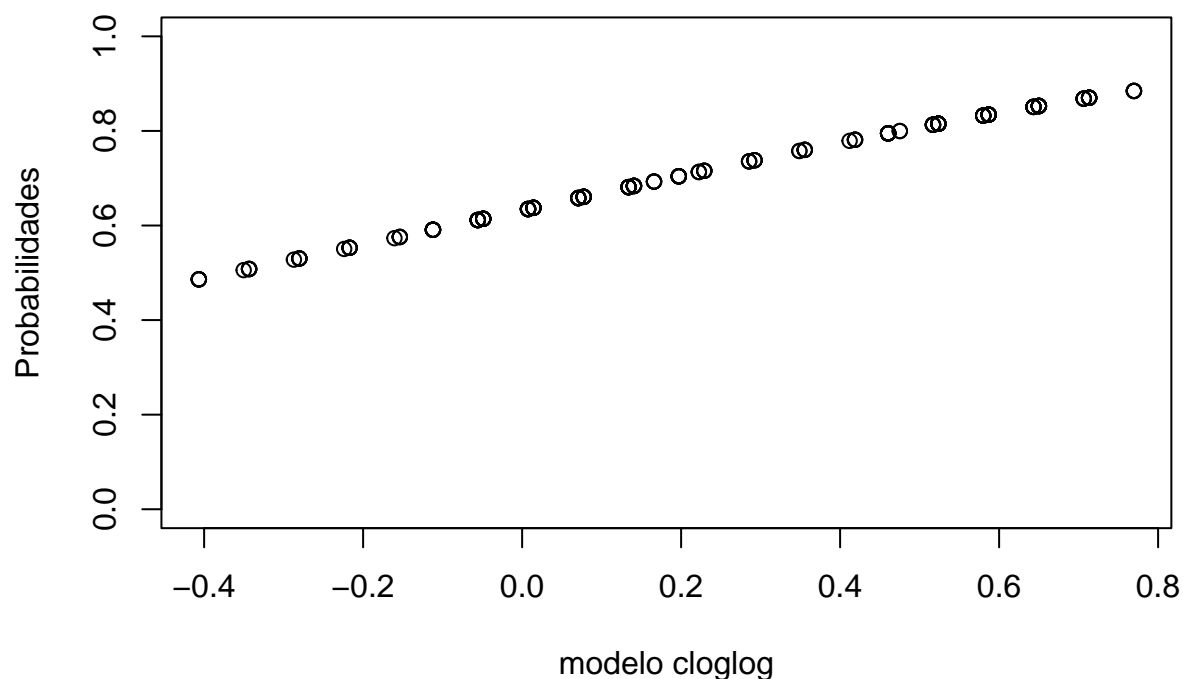
```
##
## Call:
## glm(formula = Y ~ Sex + Age + Group + as.numeric(week), family = binomial(link = "cloglog"),
##      data = dados_longos)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.0202  -1.2688   0.6394   0.8721   1.2010
```

```
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.46985    0.27391  -1.715 0.086281 .
## Sex1         0.29448    0.19859   1.483 0.138115
## Age1         0.05639    0.17402   0.324 0.745898
## Group1       0.57257    0.16874   3.393 0.000691 ***
## as.numeric(week) 0.06321    0.05979   1.057 0.290470
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 282.26  on 236  degrees of freedom
## Residual deviance: 268.18  on 232  degrees of freedom
## AIC: 278.18
##
## Number of Fisher Scoring iterations: 5
hnp(modelo_cloglog, print.on = TRUE)

## Binomial model
```



```
plot(predict.glm(modelo_cloglog, type="response")~predict.glm(modelo_cloglog, type="link"),
      ylab = "Probabilidades",
      xlab = "modelo cloglog",
      ylim=c(0,1))
```



6.2 logit

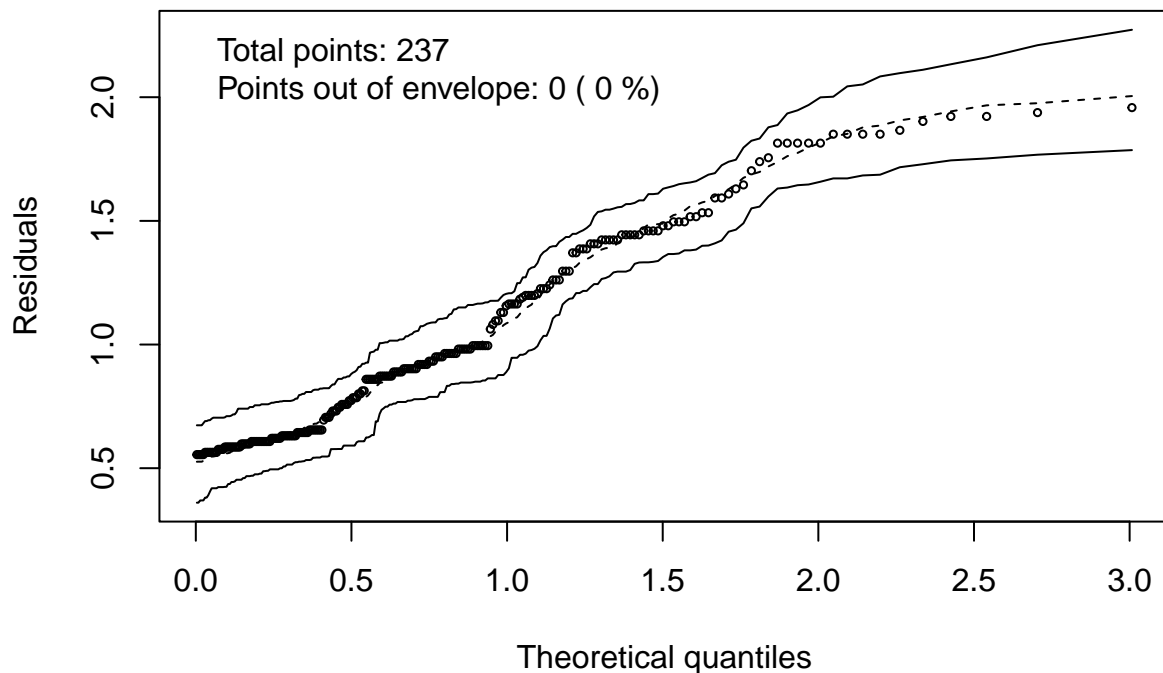
```
modelo_logit<- glm(Y ~ Sex +
  Age +
  Group +
  as.numeric(week),
  family = binomial(link = "logit"),
  data= dados_longos)
modelo_logit$family
```

```
##
## Family: binomial
## Link function: logit
summary(modelo_logit)
```

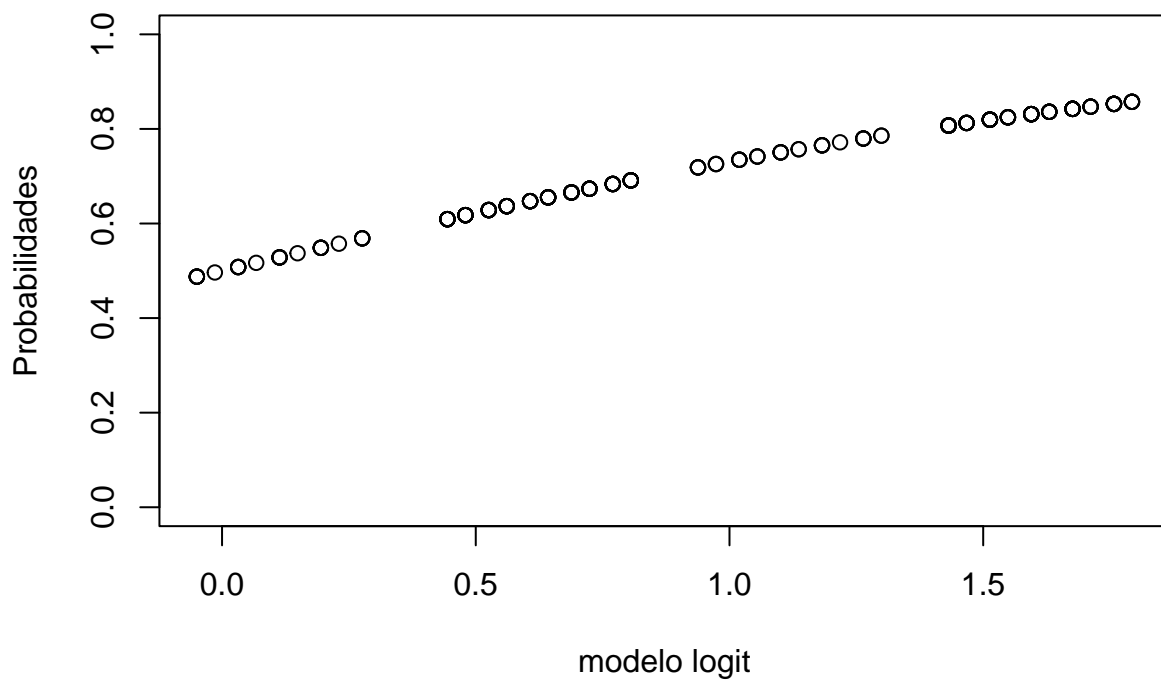
```
##
## Call:
## glm(formula = Y ~ Sex + Age + Group + as.numeric(week), family = binomial(link = "logit"),
##      data = dados_longos)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9579  -1.2612   0.6309   0.8723   1.1986
##
## Coefficients:
```

```
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -0.13114    0.44776  -0.293  0.76962
## Sex1           0.49379    0.33640   1.468  0.14214
## Age1           0.03550    0.31166   0.114  0.90930
## Group1         0.98760    0.30382   3.251  0.00115 **
## as.numeric(week) 0.08148    0.10635   0.766  0.44360
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 282.26  on 236  degrees of freedom
## Residual deviance: 268.79  on 232  degrees of freedom
## AIC: 278.79
##
## Number of Fisher Scoring iterations: 4
hnp(modelo_logit, print.on = TRUE)

## Binomial model
```



```
plot(predict.glm(modelo_logit, type="response")~predict.glm(modelo_logit, type="link"),
      ylab = "Probabilidades",
      xlab = "modelo logit",
      ylim=c(0,1))
```



6.3 probit

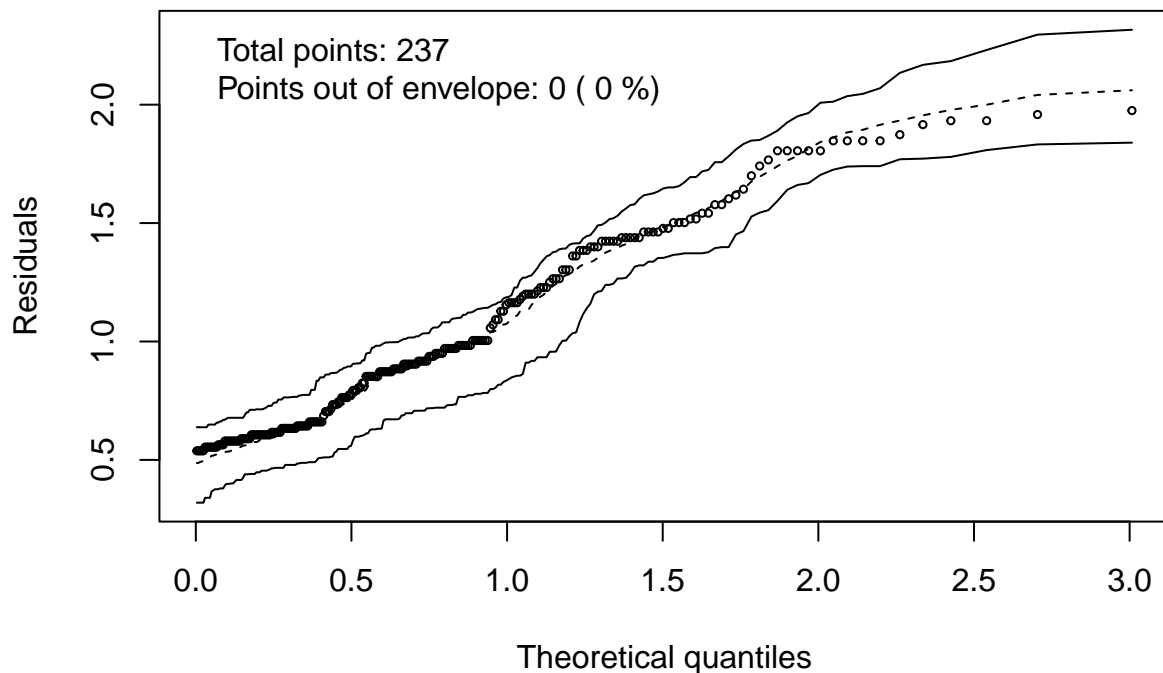
```
modelo_probit<- glm(Y ~ Sex +
  Age +
  Group +
  as.numeric(week),
  family = binomial(link = "probit"),
  data= dados_longos)
modelo_probit$family
```

```
##
## Family: binomial
## Link function: probit
summary(modelo_probit)
```

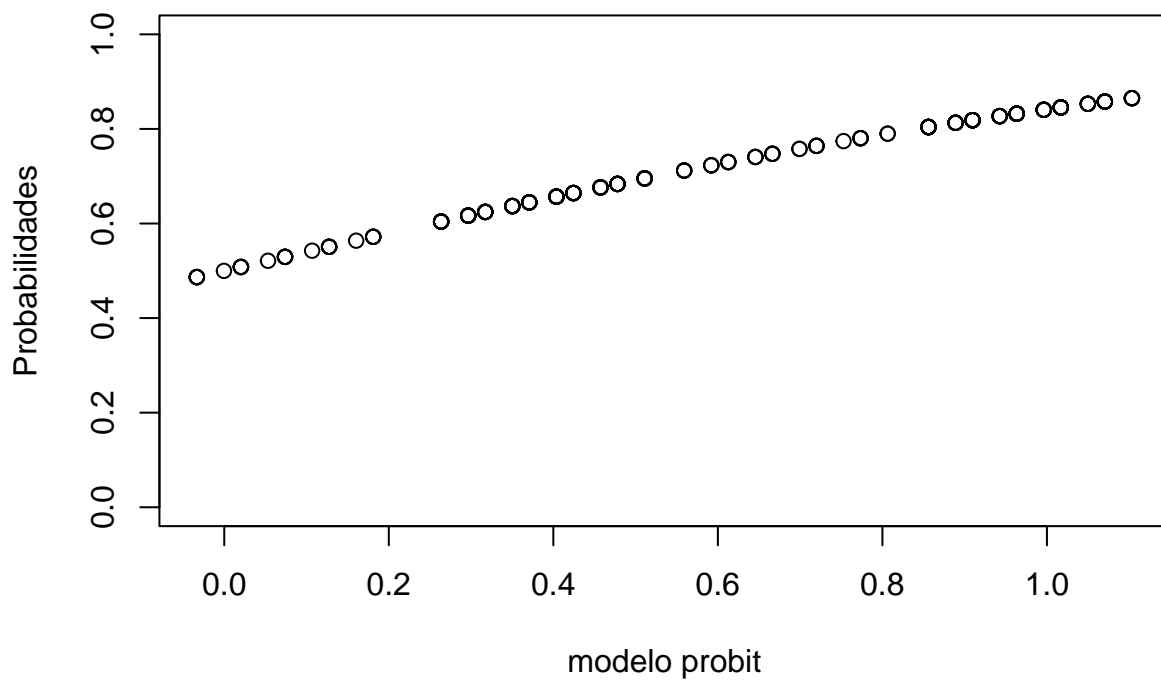
```
##
## Call:
## glm(formula = Y ~ Sex + Age + Group + as.numeric(week), family = binomial(link = "probit"),
##      data = dados_longos)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.9750  -1.2650   0.6329   0.8721   1.2001
##
## Coefficients:
```

```
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -0.08693    0.27115  -0.321 0.748510
## Sex1          0.29695    0.20235   1.468 0.142233
## Age1          0.03294    0.18530   0.178 0.858905
## Group1        0.59240    0.17866   3.316 0.000914 ***
## as.numeric(week) 0.05359    0.06329   0.847 0.397121
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 282.26  on 236  degrees of freedom
## Residual deviance: 268.63  on 232  degrees of freedom
## AIC: 278.63
##
## Number of Fisher Scoring iterations: 4
hnp(modelo_probit, print.on = TRUE)

## Binomial model
```



```
plot(predict.glm(modelo_probit, type="response")~predict.glm(modelo_probit, type="link"),
      ylab = "Probabilidades",
      xlab = "modelo probit",
      ylim=c(0,1))
```



6.4 cauchit

```

modelo_cauchit<- glm(Y ~ Sex +
  Age +
  Group +
  as.numeric(week),
  family = binomial(link = "cauchit"),
  data= dados_longos)
modelo_cauchit$family

##
## Family: binomial
## Link function: cauchit
summary(modelo_cauchit)

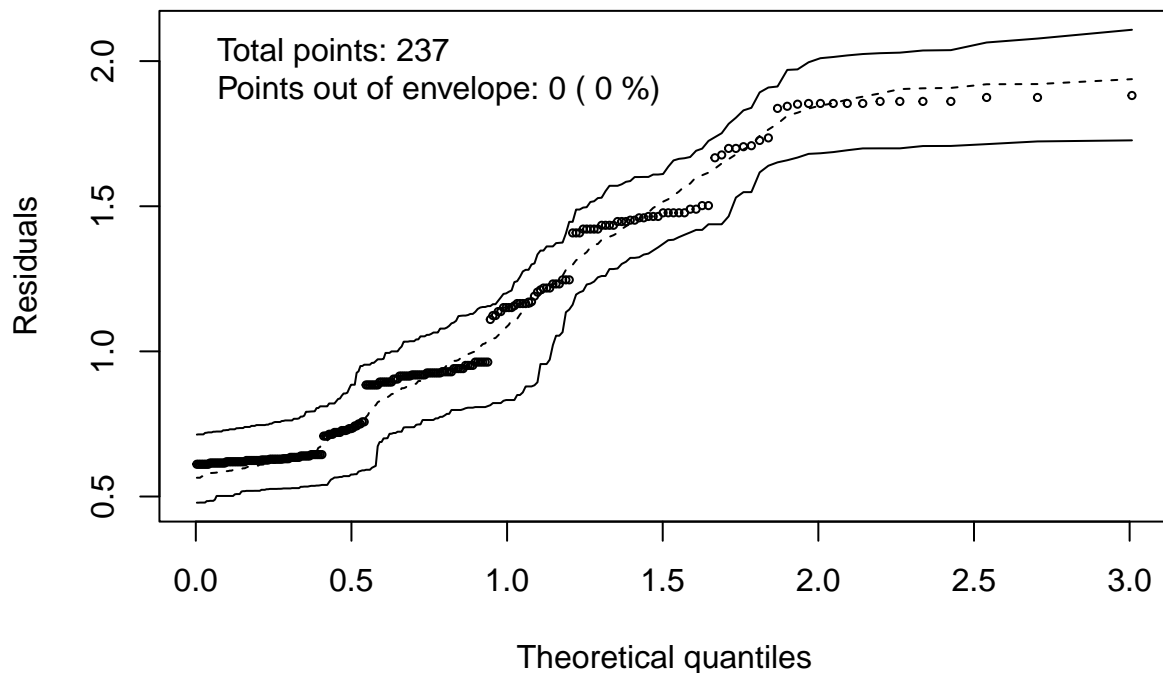
##
## Call:
## glm(formula = Y ~ Sex + Age + Group + as.numeric(week), family = binomial(link = "cauchit"),
##      data = dados_longos)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8815  -1.2323   0.6286   0.8943   1.2118
##
## Coefficients:

```

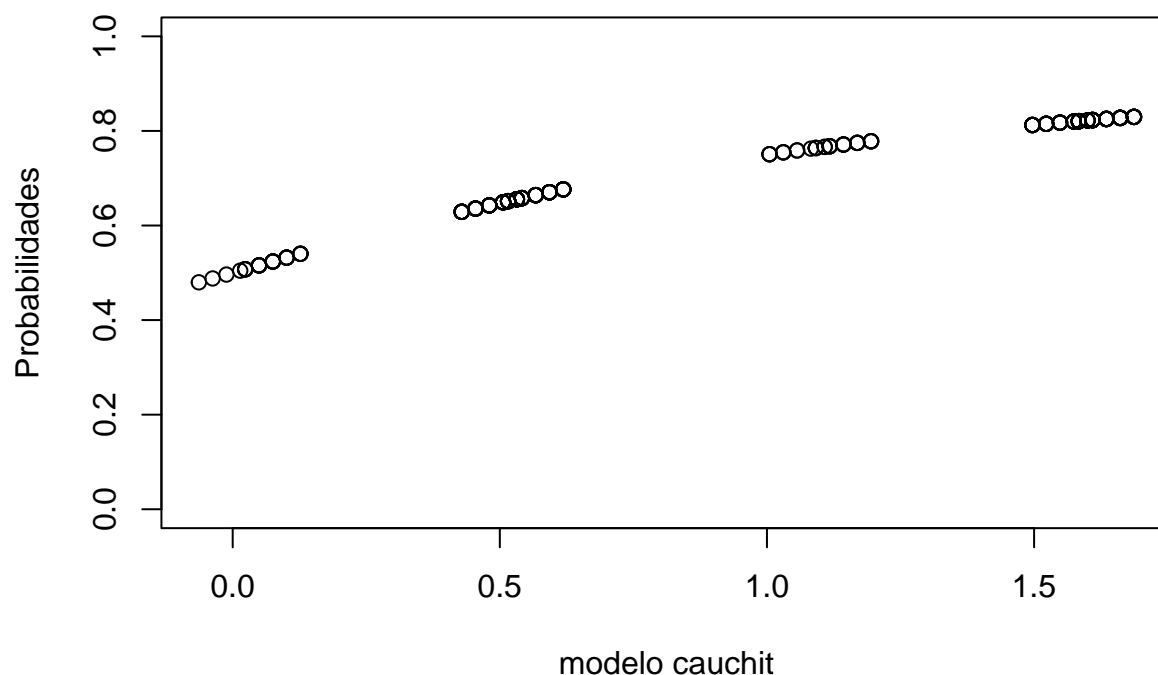


```
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -0.002309   0.420176  -0.005   0.9956
## Sex1          0.492029   0.334590   1.471   0.1414
## Age1         -0.086834   0.321406  -0.270   0.7870
## Group1        1.068118   0.383931   2.782   0.0054 **
## as.numeric(week) 0.025811   0.107747   0.240   0.8107
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 282.26  on 236  degrees of freedom
## Residual deviance: 269.49  on 232  degrees of freedom
## AIC: 279.49
##
## Number of Fisher Scoring iterations: 6
hnp(modelo_cauchit, print.on = TRUE)

## Binomial model
```



```
plot(predict.glm(modelo_cauchit, type="response")~predict.glm(modelo_cauchit, type="link"),
     ylab = "Probabilidades",
     xlab = "modelo cauchit",
     ylim=c(0,1))
```



7 GEE

7.1 independence

```
library(gee)
modelo_gee_1 <- gee(Y ~ Sex + Age + Group + as.numeric(week),
  data = dados_longos,
  id = id,
  family = binomial(link = "cloglog"),
  corstr = "independence")

##      (Intercept)           Sex1           Age1           Group1
##      -0.46985374      0.29447851      0.05639087      0.57257244
## as.numeric(week)
##      0.06320692

summary(modelo_gee_1)

##
## GEE:  GENERALIZED LINEAR MODELS FOR DEPENDENT DATA
## gee S-function, version 4.13 modified 98/01/27 (1998)
##
## Model:
## Link:                      Cloglog
## Variance to Mean Relation: Binomial
## Correlation Structure:      Independent
```

```
##
## Call:
## gee(formula = Y ~ Sex + Age + Group + as.numeric(week), id = id,
##      data = dados_longos, family = binomial(link = "cloglog"),
##      corstr = "independence")
##
## Summary of Residuals:
##      Min      1Q   Median      3Q      Max
## -0.8700482 -0.5528727  0.1848709  0.3163075  0.5138260
##
##
## Coefficients:
##              Estimate Naive S.E.   Naive z Robust S.E.   Robust z
## (Intercept)   -0.46987242  0.2760654 -1.702033  0.32613496 -1.4407300
## Sex1           0.29447626  0.2001501  1.471277  0.27987872  1.0521567
## Age1           0.05640169  0.1753833  0.321591  0.24953232  0.2260296
## Group1         0.57257352  0.1700686  3.366721  0.24311673  2.3551383
## as.numeric(week) 0.06321248  0.0602626  1.048950  0.05275103  1.1983173
##
## Estimated Scale Parameter:  1.015735
## Number of Iterations:  1
##
## Working Correlation
##      [,1] [,2] [,3] [,4] [,5]
## [1,]    1    0    0    0    0
## [2,]    0    1    0    0    0
## [3,]    0    0    1    0    0
## [4,]    0    0    0    1    0
## [5,]    0    0    0    0    1
```

7.2 exchangeable

```
modelo_gee_2 <- gee(Y ~ Sex + Age + Group + as.numeric(week),
  data = dados_longos,
  id = id,
  family = binomial(link = "cloglog"),
  corstr = "exchangeable")
```

```
##      (Intercept)      Sex1      Age1      Group1
##      -0.46985374    0.29447851    0.05639087    0.57257244
## as.numeric(week)
##      0.06320692
```

```
summary(modelo_gee_2)
```

```
##
## GEE:  GENERALIZED LINEAR MODELS FOR DEPENDENT DATA
## gee S-function, version 4.13 modified 98/01/27 (1998)
##
## Model:
## Link:                      Cloglog
## Variance to Mean Relation: Binomial
## Correlation Structure:     Exchangeable
##
## Call:
```

```
## gee(formula = Y ~ Sex + Age + Group + as.numeric(week), id = id,
##     data = dados_longos, family = binomial(link = "cloglog"),
##     corstr = "exchangeable")
##
## Summary of Residuals:
##      Min      1Q   Median      3Q      Max
## -0.8599897 -0.5530691  0.1870639  0.3276707  0.5032505
##
##
## Coefficients:
##              Estimate Naive S.E.   Naive z Robust S.E.   Robust z
## (Intercept)    -0.42904864 0.33365054 -1.2859222  0.32034671 -1.3393259
## Sex1           0.27285705 0.28277327  0.9649323  0.27849472  0.9797567
## Age1           0.08134966 0.24774357  0.3283623  0.24724658  0.3290224
## Group1         0.53829103 0.24201844  2.2241736  0.24310512  2.2142316
## as.numeric(week) 0.05314300 0.05066053  1.0490021  0.05083011  1.0455023
##
## Estimated Scale Parameter: 0.9965314
## Number of Iterations: 4
##
## Working Correlation
##      [,1]      [,2]      [,3]      [,4]      [,5]
## [1,] 1.0000000 0.2907233 0.2907233 0.2907233 0.2907233
## [2,] 0.2907233 1.0000000 0.2907233 0.2907233 0.2907233
## [3,] 0.2907233 0.2907233 1.0000000 0.2907233 0.2907233
## [4,] 0.2907233 0.2907233 0.2907233 1.0000000 0.2907233
## [5,] 0.2907233 0.2907233 0.2907233 0.2907233 1.0000000
```

7.3 unstructured

```
modelo_gee_3 <- gee(Y ~ Sex + Age + Group + as.numeric(week),
  data = dados_longos,
  id = id,
  family = binomial(link = "cloglog"),
  corstr = "unstructured")
```

```
##      (Intercept)      Sex1      Age1      Group1
##      -0.46985374      0.29447851      0.05639087      0.57257244
## as.numeric(week)
##      0.06320692
```

```
summary(modelo_gee_3)
```

```
##
## GEE: GENERALIZED LINEAR MODELS FOR DEPENDENT DATA
## gee S-function, version 4.13 modified 98/01/27 (1998)
##
## Model:
## Link:                      Cloglog
## Variance to Mean Relation: Binomial
## Correlation Structure:     Unstructured
##
## Call:
## gee(formula = Y ~ Sex + Age + Group + as.numeric(week), id = id,
##     data = dados_longos, family = binomial(link = "cloglog"),
```

```

##      corstr = "unstructured")
##
## Summary of Residuals:
##      Min      1Q      Median      3Q      Max
## -0.8703654 -0.5230487  0.1650142  0.3603214  0.5205362
##
##
## Coefficients:
##              Estimate Naive S.E.      Naive z Robust S.E.      Robust z
## (Intercept)   -0.468235192  0.3394253 -1.37949401  0.32501737 -1.44064666
## Sex1           0.279293127  0.2731156  1.02261888  0.27175387  1.02774297
## Age1           0.007621797  0.2397362  0.03179243  0.24396340  0.03124156
## Group1         0.693890430  0.2333050  2.97417687  0.24223920  2.86448444
## as.numeric(week) 0.041897668  0.0552891  0.75779257  0.04990972  0.83946913
##
## Estimated Scale Parameter:  1.028027
## Number of Iterations:  7
##
## Working Correlation
##      [,1]      [,2]      [,3]      [,4]      [,5]
## [1,] 1.0000000  0.64699683  0.19879899  0.3105900  0.2345909
## [2,] 0.6469968  1.00000000  0.02545539  0.2919917  0.1593287
## [3,] 0.1987990  0.02545539  1.00000000  0.1337962  0.2364074
## [4,] 0.3105900  0.29199173  0.13379622  1.0000000  0.1851979
## [5,] 0.2345909  0.15932868  0.23640737  0.1851979  1.0000000

```