

## Exercício 1

Considere um modelo balanceado com intercepto aleatório obtenha a estimativa de máxima verossimilhança da componente d.

**Resposta:**

O modelo balanceado com intercepto aleatório em notação matricial pode ser escrito como:

$$\mathbf{y}_i = \mathbf{X}\boldsymbol{\beta} + \mathbf{1}_n b_i + \boldsymbol{\varepsilon}_i, \quad i = 1, 2, \dots, N,$$

onde

- $\mathbf{y}_i$  é um vetor  $n \times 1$  de variáveis respostas para o i-ésimo sujeito.
- $\mathbf{X}$  é uma matriz  $n \times m$  de variáveis explicativas.
- $\boldsymbol{\beta}$  é um vetor  $m \times 1$  de efeitos fixos.
- $\boldsymbol{\varepsilon}_i$  é um vetor  $n \times 1$  de erros independentes, e assume-se que  $\boldsymbol{\varepsilon}_i \sim N(0, \sigma^2 \mathbf{I}_n)$ .
- $b_i$  é o efeito aleatório, e assume-se que  $b_i \sim N(0, \sigma^2 d)$ .
- $\mathbf{1}_n$  é um vetor unitário de dimensão  $n \times 1$ .
- $\sigma^2$  é a variância dentro do sujeito e  $d$  é a variância escalada do efeito aleatório.

O modelo também pode ser escrito como

$$\mathbf{y}_i \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma^2(\mathbf{I}_n + d\mathbf{1}_n\mathbf{1}_n')), \quad i = 1, 2, \dots, N.$$

A log-verossimilhança do modelo, ignorando a constante é

$$\ell(\boldsymbol{\theta}) = -\frac{1}{2} \left\{ Nn \log(\sigma^2) + \sum_{i=1}^N \log |\mathbf{I}_n + d\mathbf{1}_n\mathbf{1}_n'| + \sigma^{-2} \sum_{i=1}^N \mathbf{e}_i' (\mathbf{I}_n + d\mathbf{1}_n\mathbf{1}_n')^{-1} \mathbf{e}_i \right\},$$

onde  $\mathbf{e}_i = \mathbf{y}_i - \mathbf{X}\boldsymbol{\beta}$ .

Utilizando as fórmulas da redução de dimensão, temos que:

$$\begin{aligned} (\mathbf{I}_n + d\mathbf{1}_n\mathbf{1}_n')^{-1} &= \mathbf{I}_n - \mathbf{1}_n \left( \frac{1}{d} + \mathbf{1}_n' \mathbf{1}_n \right)^{-1} \mathbf{1}_n' \\ &= \mathbf{I}_n - \mathbf{1}_n \left( \frac{1}{d} + n \right)^{-1} \mathbf{1}_n' \\ &= \mathbf{I}_n - \frac{d}{1 + nd} \mathbf{1}_n\mathbf{1}_n'. \end{aligned}$$

$$|\mathbf{I}_n + d\mathbf{1}_n\mathbf{1}_n'| = |1 + dn| = 1 + dn.$$

Assim, a log-verossimilhança fica

$$\begin{aligned} \ell(\boldsymbol{\theta}) &= -\frac{1}{2} \left\{ Nn \log(\sigma^2) + \sum_{i=1}^N \log(1 + dn) + \sigma^{-2} \sum_{i=1}^N \mathbf{e}_i' \left( \mathbf{I}_n - \frac{d}{1 + nd} \mathbf{1}_n\mathbf{1}_n' \right) \mathbf{e}_i \right\} \\ &= -\frac{1}{2} \left\{ Nn \log(\sigma^2) + N \log(1 + dn) + \sigma^{-2} \sum_{i=1}^N \left[ \mathbf{e}_i' \mathbf{e}_i - \frac{d}{1 + nd} \mathbf{e}_i' \mathbf{1}_n \mathbf{1}_n' \mathbf{e}_i \right] \right\} \\ &= -\frac{1}{2} \left\{ Nn \log(\sigma^2) + N \log(1 + dn) + \sigma^{-2} \sum_{i=1}^N \left[ \mathbf{e}_i' \mathbf{e}_i - \frac{d}{1 + nd} (\mathbf{e}_i' \mathbf{1}_n)^2 \right] \right\}. \end{aligned}$$

A derivada parcial da log-verossimilhança em relação a  $d$  é:

$$\frac{\partial \ell(\boldsymbol{\theta})}{\partial d} = -\frac{1}{2} \left\{ \frac{Nn(1 + dn)\sigma^2 - \sum_{i=1}^N (\mathbf{e}'_i \mathbf{1}_n)^2}{\sigma^2(1 + dn)^2} \right\}.$$

A derivada acima é igual a zero quando

$$d = \frac{\sum_{i=1}^N (\mathbf{e}'_i \mathbf{1}_n)^2 - Nn\sigma^2}{\sigma^2 Nn^2} = \frac{\sum_{i=1}^N (\mathbf{e}'_i \mathbf{1}_n)^2}{\sigma^2 Nn^2} - \frac{1}{n},$$

assim a estimativa de máxima verossimilhança da componente  $d$  é:

$$\hat{d}_{ML} = \frac{\sum_{i=1}^N (\hat{\mathbf{e}}'_i \mathbf{1}_n)^2}{\hat{\sigma}_{ML}^2 Nn^2} - \frac{1}{n},$$

onde  $\hat{\mathbf{e}}_i = \mathbf{y}_i - \mathbf{X}\hat{\boldsymbol{\beta}}_{ML}$ , e segundo Demidenko (2013), para dados balanceados  $\hat{\boldsymbol{\beta}}_{ML} = \hat{\boldsymbol{\beta}}_{OLS} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\bar{\mathbf{y}}$ , assim, o termo  $\hat{\mathbf{e}}'_i \mathbf{1}_n$  pode ser escrito como

$$\begin{aligned} \hat{\mathbf{e}}'_i \mathbf{1}_n &= (\mathbf{y}_i - \mathbf{X}\hat{\boldsymbol{\beta}}_{OLS})' \mathbf{1}_n \\ &= (\mathbf{y}_i - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\bar{\mathbf{y}})' \mathbf{1}_n \\ &= \mathbf{y}'_i \mathbf{1}_n - \bar{\mathbf{y}}' \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbf{1}_n \\ &= \mathbf{y}'_i \mathbf{1}_n - \bar{\bar{\mathbf{y}}} \mathbf{1}_n \\ &= n(\bar{y}_i - \bar{\bar{y}}), \end{aligned}$$

onde  $\bar{\bar{\mathbf{y}}} = \bar{\mathbf{y}}' \mathbf{1}_n / n$ ,  $\bar{y}_i = \mathbf{y}'_i \mathbf{1}_n / n$ , e a prova de que  $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \mathbf{1}_n = \mathbf{1}_n$  se encontra em (Demidenko 2013, p.68).

Portanto, a estimativa de máxima verossimilhança da componente  $d$  em um modelo balanceado com intercepto aleatório é

$$\hat{d}_{ML} = \frac{\sum_{i=1}^N (\bar{y}_i - \bar{\bar{y}})^2}{\hat{\sigma}_{ML}^2 N} - \frac{1}{n},$$

## Exercício 2

Mostre que o termo do log-verossimilhança restrita

$$f(N) = -\frac{1}{2}(-m \ln \sigma^2 + \ln |\sum_{i=1}^N X_i' V_i X_i|)$$

é de ordem de  $\ln N$ . Considere:  $n_i = n$ ,  $X_i$  e  $Z_i = Z$ .

**Resposta:**

Queremos mostrar que  $\frac{f(N)}{\ln N} < M$ , com  $M \in \mathbb{R}$  e  $\forall N > N_0$ .

Temos que

$$\ln |\sum_{i=1}^N X_i' V_i X_i| = \ln |\frac{N}{N} \sum_{i=1}^N X_i' V_i X_i| = \ln(N^m |\frac{1}{N} \sum_{i=1}^N X_i' V_i X_i|) = \ln |\frac{N}{N} X' V X| + m \ln N$$

Logo,

$$f(N) = -\frac{1}{2}(-m \ln \sigma^2 + \ln |X' V X| + m \ln N)$$

só depende de  $\ln N$  quando  $\lim_{N \rightarrow \infty}$ . Logo a seguinte inequação é válida pois existe um  $M$  real tal que

$$\frac{-\frac{1}{2}(-m \ln \sigma^2 + \ln |X' V X| + m \ln N)}{\ln N} < M$$

para todo  $N > N_0$ .

### Exercício 3

Um estudo com filhotes de coruja foi conduzido com câmeras e microfones para analisar a negociação entre os irmãos definida como segue. Usando a filmagem gravada foram registrados durante intervalos de 30 segundos a cada 15 minutos o número de chamadas feitas por todos os descendentes na ausência dos pais. Para cada visita de um dos pais foi registrado o número de chamadas dos 15 minutos anteriores dividido pelo número de filhotes em cada ninho. Os dados estão no arquivo Owls. As variáveis explicativas são o sexo dos pais, tratamento de alimentos, e o tempo de chegada do pai. O tratamento de alimentos foram dois, a metade dos ninhos foram dadas presas extras (“food-satiated”), e na outra metade as presas (remanescentes) foram removidas. (“fooddeprived”). As medições ocorreram em duas noites em cada caso, e o tratamento de alimentos foi trocado na segunda noite. (Fonte Roulin and Bersier, 2007) A negociação entre irmãos (NegPerChick) pode ser transformada por  $\log_{10}(Y + 1) = \text{LogNeg}$  para ser modelada.

- Selecione um modelo misto que considere adequado para determinar os fatores que podem influenciar a variável de interesse.
- Descreva o modelo formalmente especificando as matrizes associadas aos efeitos fixos e aleatórios.
- Realize as interpretações que pertinentes.

#### item a)

**Análise descritiva** Ao todo temos 599 observações de 27 ninhos. As características observadas ao longo do experimento em todos os ninhos foram tratamento alimentar, sexo do pai, tempo de chegada do pai e tamanho da ninhada.

A característica de interesse do estudo é o número de negociação entre irmãos.

Como o número de filhotes de cada ninhada é diferente, para conseguirmos comparar as negociações entre ninhadas, dividimos o número de negociações observadas no intervalo de tempo pelo número total de filhotes do ninho. Além disso, conforme instrução do enunciado, transformamos a variável resposta NegPerChick (negociações por filhote) adicionando 1 e aplicando logaritmo de base 10 com objetivo de normalizar os dados e usar um modelo linear padrão.

Vemos que cada ninho tem pelo menos 4 observações, desta forma concluímos que as observações não são independentes e portanto, é interessante mensurar a correlação das observações de cada sujeito.

Nest	contagem
Forel	4
Sevaz	4
Chevroux	10
GDLV	10
CorcellesFavres	12
Henniez	13
Gletterens	15
LesPlanches	17
Lully	17
Rueyes	17
Jeuss	19
Trey	19
ChEsard	20
Bochet	23
StAubin	23
Murist	24

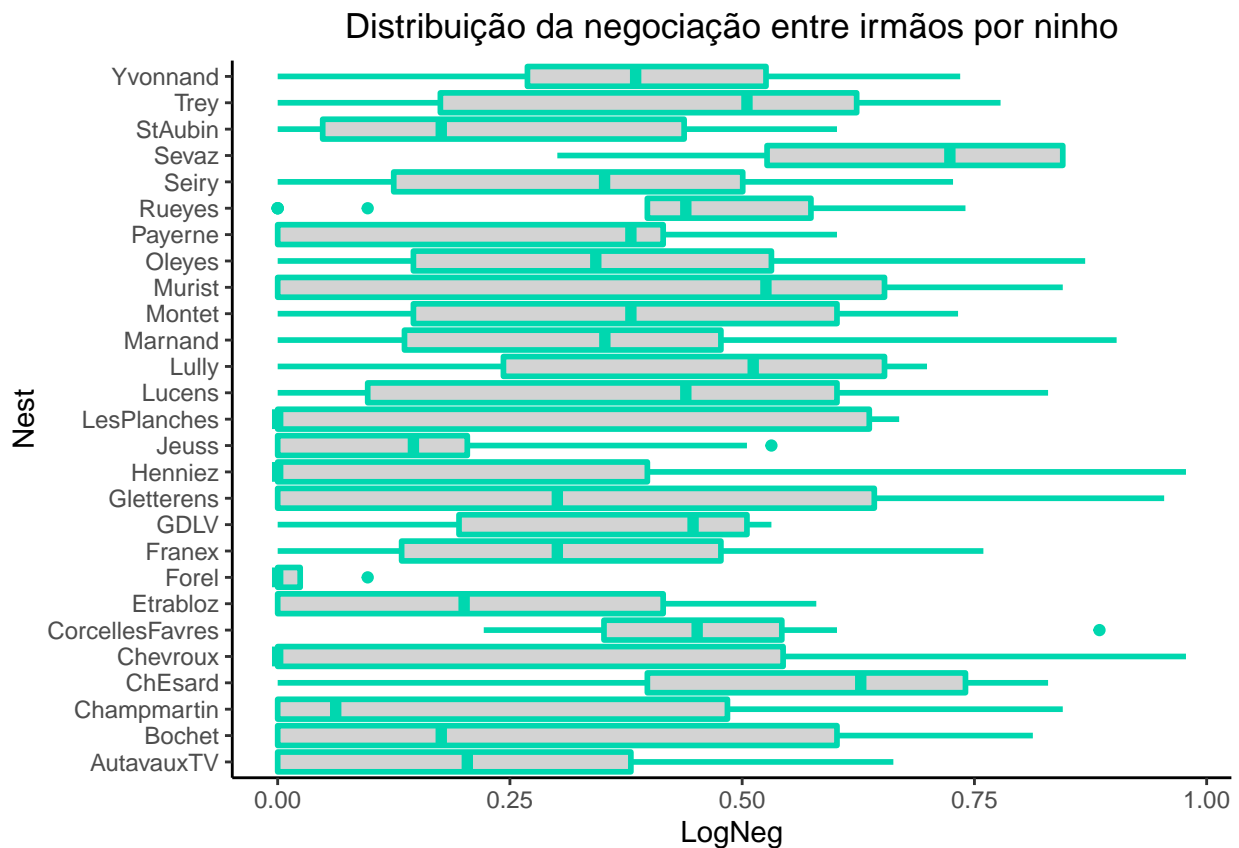
Nest	contagem
Payerne	25
Franex	26
Seiry	26
Marnand	27
AutavauxTV	28
Lucens	29
Champmartin	30
Etrabloz	34
Yvonnand	34
Montet	41
Oleyes	52

Além disso, como os ninhos estudados são uma amostra de todos os ninhos existentes, podemos estar interessados em generalizar o modelo para que ele seja capaz de descrever outros ninhos fora da amostra.

Também é importante destacar que ao considerar os ninhos como efeitos fixos, gastaríamos muitos graus de liberdade.

### Variável resposta

Abaixo fica claro que a distribuição do log das negociações entre irmãos varia entre os ninhos.

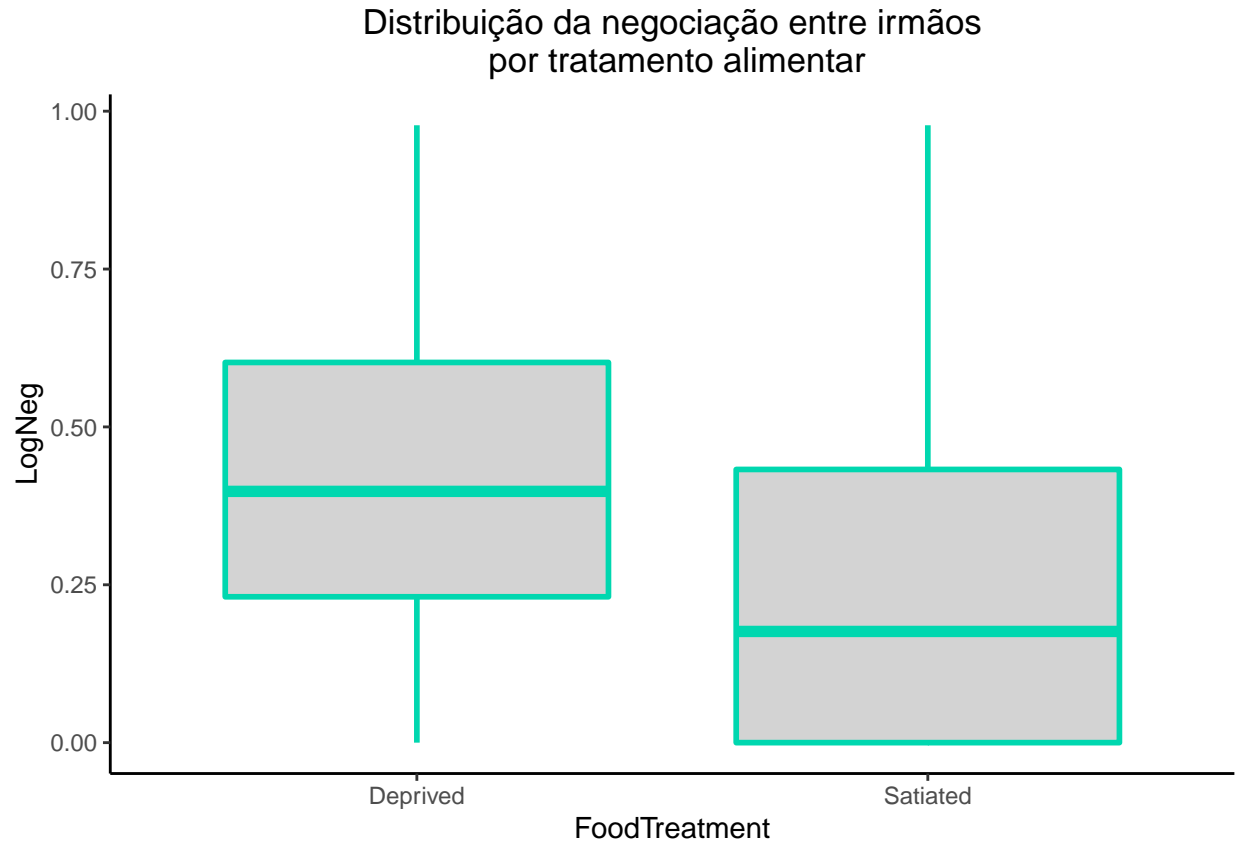


Por todo o exposto vamos construir um modelo de efeitos mistos.

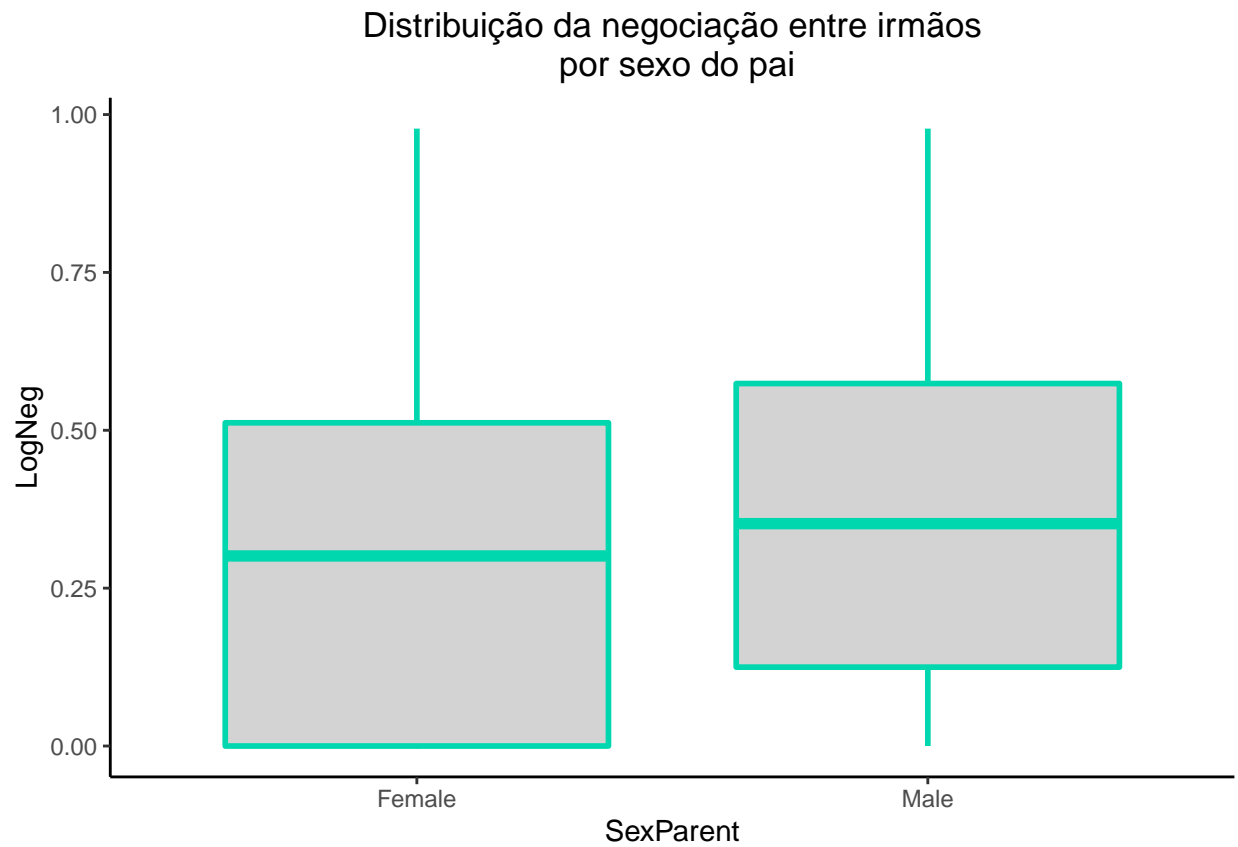
### Variáveis explicativas

Quanto às variáveis explicativas do modelo, temos disponíveis: tratamento alimentar, sexo do pai, e tempo de chegada do pai. Os demais dados disponíveis são referentes à resposta e/ou estão sendo usados para calculá-la.

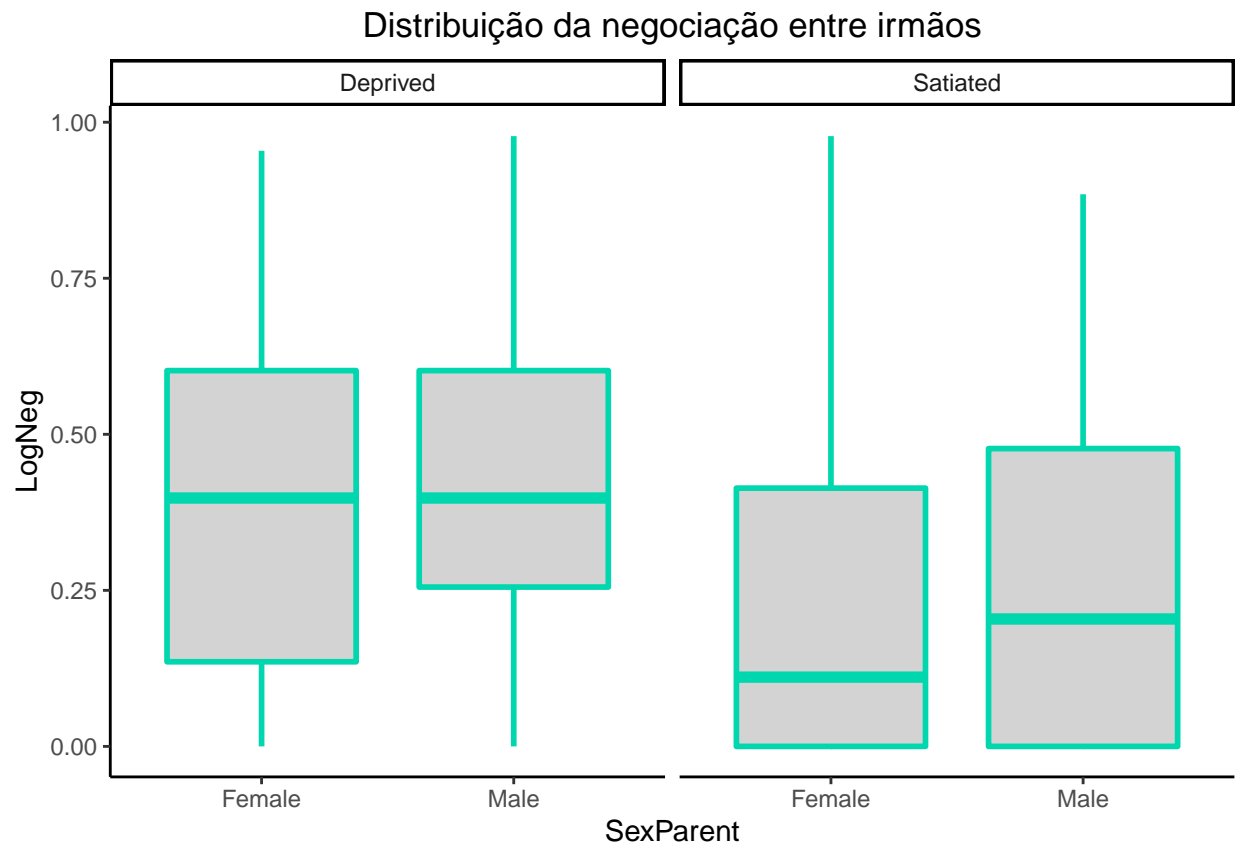
No gráfico de abaixo vemos que a distribuição do log da taxa de negociação por filhote é maior para o tratamento alimentar *Satiated*.



Avaliando agora o sexo dos pais, pela análise gráfica não identificamos diferença na negociação entre irmãos.

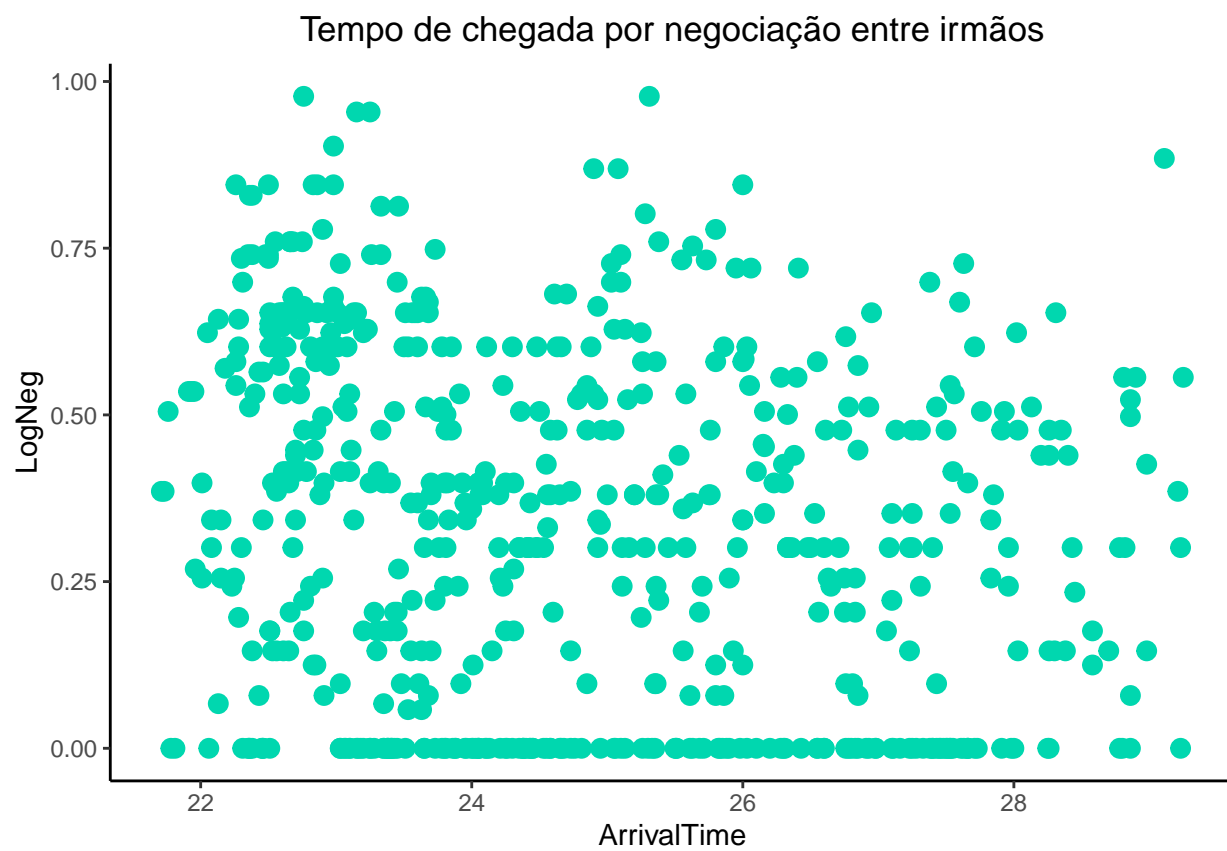


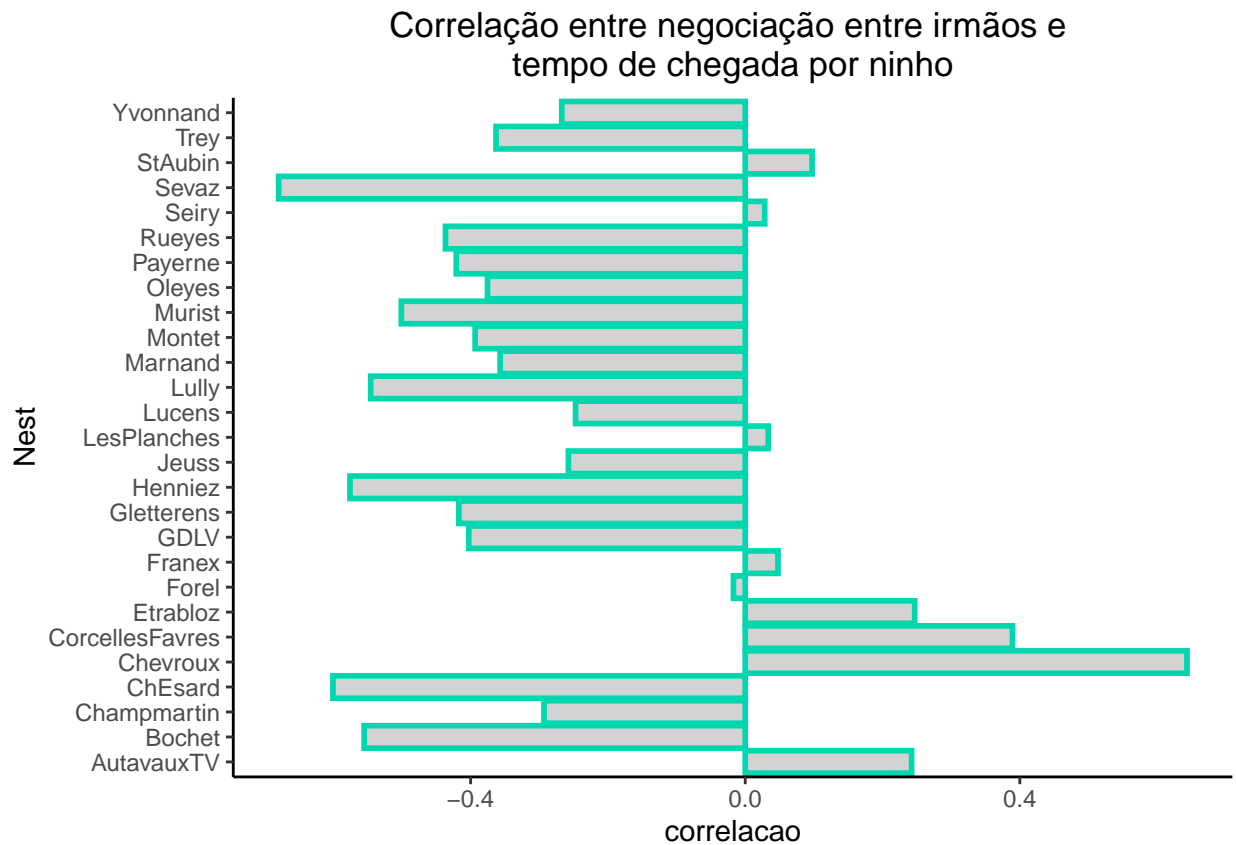
Visualizando as duas variáveis simultaneamente não encontramos evidência de interação entre elas, pois o tratamento alimentar *Deprived* mostra maior negociação entre irmãos para ambos os sexos, assim como o tratamento *Satiated* que mantém próximas as distribuições de negociação para ambos os sexos.



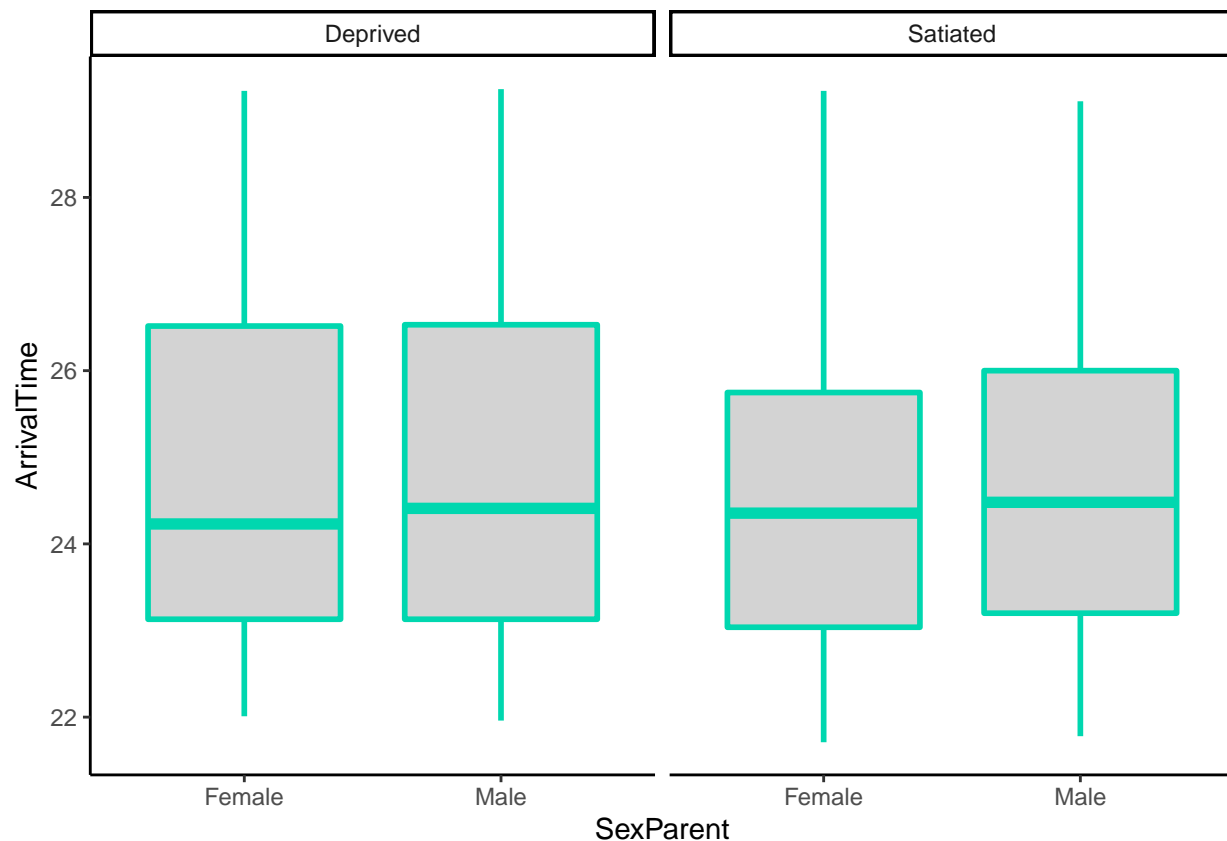
Outra variável explicativa disponível é o tempo de chegada que aparentemente não tem correlação com a variável resposta. Porém, se avaliarmos a correlação entre essas variáveis para cada um dos ninhos, vemos que o cenário muda.







Procurando interação entre as variáveis, avaliamos o tempo de chegada por tratamento alimentar e sexo dos pais e não encontramos evidências gráficas de interação. O tempo de chegada mostra mediana e dispersão parecida para ambos os sexos assim como para os dois tratamentos alimentares.



**Ajuste do modelo** Conforme evidência da análise descritiva, após ajustar o modelo, vemos que a variável *SexParent* não é significativa.

```
negociacao <- lme(LogNeg~SexParent+FoodTreatment+ArrivalTime, random= ~1| Nest, data = Owls)
summary(negociacao)
```

```
## Linear mixed-effects model fit by REML
##   Data: Owls
##       AIC      BIC    logLik
##  22.03994 48.37131 -5.019969
##
## Random effects:
## Formula: ~1 | Nest
##      (Intercept)  Residual
## StdDev:  0.09271906 0.2318242
##
## Fixed effects:  LogNeg ~ SexParent + FoodTreatment + ArrivalTime
##               Value Std.Error DF   t-value p-value
## (Intercept)    1.1748045 0.12917288 569   9.094823  0.0000
## SexParentMale    0.0202573 0.02130399 569   0.950870  0.3421
## FoodTreatmentSatiated -0.1735870 0.02001784 569  -8.671614  0.0000
## ArrivalTime     -0.0312510 0.00512116 569  -6.102322  0.0000
## Correlation:
##               (Intr) SxPrnM FdTrtS
## SexParentMale    -0.055
```

```
## FoodTreatmentSatiated -0.115 0.068
## ArrivalTime          -0.979 -0.051 0.035
##
## Standardized Within-Group Residuals:
##      Min      Q1      Med      Q3      Max
## -2.23169999 -0.77923134 -0.08570653 0.71095687 3.29621690
##
## Number of Observations: 599
## Number of Groups: 27
```

Outro teste que pode ser feito para verificar a necessidade de inserir a variável *SexParent* no modelo é o teste de verossimilhança, porém, para usá-lo precisamos que o método de estimativa dos parâmetros seja máxima verossimilhança.

Então ajustamos dois modelos semelhantes (pelo método da verossimilhança), um deles com as variáveis *SexParent*, *FoodTreatment* e *ArrivalTime* e outro apenas com *FoodTreatment* e *ArrivalTime* e aplicamos o teste da razão de verossimilhança por meio do comando *anova*. Conforme o esperado, o resultado obtido foi que a variável *SexParent* não melhora o ajuste e pode ser desconsiderada do modelo.

```
com_SexParent <- lme(LogNeg~SexParent+FoodTreatment+ArrivalTime,
                     random= ~1| Nest,
                     method = "ML",
                     data = Owls)
sem_SexParent <- lme(LogNeg~FoodTreatment+ArrivalTime,
                     random= ~1| Nest,
                     method = "ML",
                     data = Owls)
anova(sem_SexParent, com_SexParent)
```

```
##           Model df      AIC      BIC  logLik  Test  L.Ratio p-value
## sem_SexParent   1  5 -5.545145 16.43116 7.772572
## com_SexParent   2  6 -4.476920 21.89465 8.238460 1 vs 2 0.9317755 0.3344
```

O **modelo final** escolhido para explicar o log da taxa de negociações por filhote é um modelo misto, cujos parâmetros são estimados pelo método da máxima verossimilhança restrita porque ele retorna estimativas não viciadas para os componentes da variância.

As variáveis escolhidas para este modelo são tratamento alimentar e tempo de chegada.

```
sem_SexParent <- lme(LogNeg~FoodTreatment+ArrivalTime, random= ~1| Nest, data = Owls)
sem_SexParent
```

```
## Linear mixed-effects model fit by REML
##   Data: Owls
##   Log-restricted-likelihood: -2.536915
##   Fixed: LogNeg ~ FoodTreatment + ArrivalTime
##           (Intercept) FoodTreatmentSatiated      ArrivalTime
##           1.18213859      -0.17507539      -0.03102135
##
## Random effects:
##   Formula: ~1 | Nest
##           (Intercept) Residual
## StdDev: 0.09468769 0.2316398
##
## Number of Observations: 599
## Number of Groups: 27
```

item b)

O modelo escolhido no item anterior pode ser formalmente escrito conforme segue:

$$y = X\beta + Zb + \epsilon$$

onde

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Vetor coluna de variável resposta de dimensão  $N \times 1$ .  $N$  é o número total de observações.

$$N = \sum_{i=1}^{27} n_i = 599$$

.

A Matrix design dos efeitos fixos do modelo com as covariáveis *FoodTreatment* e *ArrivalTime* tem dimensão  $N \times 3$  e pode ser escrita como

$$X = \begin{bmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{1N} & x_{2N} \end{bmatrix}$$

em que  $x_{11}, x_{12}, \dots, x_{1N}$  são os valores observados da covariável *FoodTreatment* para cada um dos  $N$  ninhos, sendo 0 se Deprived e 1 se Satiated. As observações  $x_{21}, x_{22}, \dots, x_{2N}$  são os valores da covariável *ArrivalTime* para cada um dos  $N$  ninhos. O Vetor de parâmetros  $\beta$  tem dimensão  $3 \times 1$  e é dado por

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

com  $\beta_0$  sendo o intercepto,  $\beta_1$  o efeito do tratamento *Satiated* e  $\beta_2$  o impacto da covariável *ArrivalTime*.

Matriz dos efeitos aleatórios tem dimensão  $N \times 27$ , e é denotada por

$$Z = \begin{bmatrix} Z_1 & 0 & \dots & 0 \\ 0 & Z_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & Z_{27} \end{bmatrix}$$

em que cada  $Z_i$ , com  $i = 1, \dots, 27$ , é um vetor de uns, com dimensão  $n_i \times 1$ . Como estamos assumindo apenas o intercepto aleatório temos que para cada ninho  $j$   $b_j \sim N(0, \sigma_b)$ . O vetor de efeitos aleatórios de dimensão  $27 \times 1$  pode ser escrito como

$$b = \begin{bmatrix} b_1 \\ \vdots \\ b_{27} \end{bmatrix}$$

Por fim, o vetor de erros de dimensão  $N \times 1$

$$\epsilon = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_N \end{bmatrix}$$

além disso,

$$b \sim N(0, \sigma^2 D)$$

e

$$\epsilon \sim N(0, \sigma^2 I)$$

em que  $D$  é uma matriz diagonal de dimensão  $27 \times 27$ . Como estamos trabalhando com modelos mistos, podemos escrever um modelo para cada sujeito, no caso do exercício, um modelo para cada ninho. Para escrever o modelo com os dados do exercício, um dos menores ninhos, o **Forel** foi escolhido.

$$y_i = X_i \beta + Z_i b_i + \epsilon_i$$

Para  $i = 8$ , temos duas variáveis explicativas e 4 observações para este sujeito e mantivemos o intercepto no modelo, a matriz  $X$  de design tem dimensão  $4 \times 3$ .

$$X_1 = \begin{bmatrix} 1 & 1 & 23.25 \\ 1 & 0 & 23.92 \\ 1 & 1 & 24.26 \\ 1 & 1 & 24.30 \end{bmatrix}$$

Nest	FoodTreatment	SexParent	ArrivalTime	LogNeg
Forel	Satiated	Male	23.25	0.00000
Forel	Deprived	Male	23.92	0.09691
Forel	Satiated	Male	24.26	0.00000
Forel	Satiated	Male	24.30	0.00000

O vetor de respostas  $y_8^T = (0, 0.096, 0, 0)$

---



---

0.00000  
0.09691  
0.00000  
0.00000

---

Como temos duas variáveis explicativas e 4 observações para este sujeito e mantivemos o intercepto no modelo, a matriz  $X$  de design tem dimensão  $4 \times 3$ .

intercepto	FoodTreatmentSatiated	ArrivalTime
1	1	23.25
1	0	23.92
1	1	24.26
1	1	24.30

O vetor de parâmetros dos efeitos fixos tem dimensão 3x1. As estimativas dos parâmetros são apresentadas na tabela abaixo.

(Intercept)	1.1821386
FoodTreatmentSatiated	-0.1750754
ArrivalTime	-0.0310214

No modelo escolhido, definimos o intercepto aleatório, então a matriz  $Z$  de efeitos aleatórios é uma matriz diagonal em blocos, onde temos 1 nas linhas e as colunas preenchidas são as colunas referentes ao ninho  $i$ , no caso *Forel*  $i = 8$ .

$$Z_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

E, por fim, temos a estimativa do parâmetro de variância do efeito aleatório *Forel* ( $\sigma_b^2$ ).

$$-0.1138491$$

#### item c)

Estamos usando o modelo para entender o comportamento da taxa de negociação entre irmãos dadas algumas características observadas de cada ninho. São elas, o próprio ninho, o tratamento alimentar oferecido no momento da observação e o tempo de chegada dos pais.

```
## Linear mixed-effects model fit by REML
##   Data: Owls
##   Log-restricted-likelihood: -2.536915
##   Fixed: LogNeg ~ FoodTreatment + ArrivalTime
##           (Intercept) FoodTreatmentSatiated      ArrivalTime
##           1.18213859      -0.17507539      -0.03102135
##
## Random effects:
##   Formula: ~1 | Nest
##           (Intercept) Residual
##   StdDev:  0.09468769 0.2316398
##
## Number of Observations: 599
## Number of Groups: 27
```

O intercepto do modelo nos diz que a taxa de negociação entre irmãos média, sabendo que o tratamento alimentar é *Deprived*, é 1.1821386.

O tratamento alimentar *Satiated* reduz a taxa de negociação em 0.1750754.

E a cada unidade de tempo acrescida no valor médio de tempo de chegada, temos um decréscimo da taxa de negociação de 0.0310214