

Modelos para dados positivos



Universidade de Brasília
Instituto de Ciências Exatas - Departamento de Estatística
Modelos Lineares Generalizados

Alunos responsáveis:

Álvaro Jeronimo da Silva Kothe - 17/0004694
Marcos Augusto Daza Barbosa - 17/0017834
Mathews de Noronha Silveira Lisboa - 17/0042324
William Edward Rappel de Amorim - 17/0047385

Professor:
Guilherme Souza Rodrigues

29 de Outubro de 2019

Construção do modelo

Para a construção do modelo, primeiramente foram modificados os formatos de algumas variáveis, como tirar a ',' de milhar, o '\$', separou-se as datas recebidas em ano, mês e dia da semana. Nas colunas em que havia '%', considerou-se os valores como numéricos, compreendidos na escala de zero a um. Valores faltantes em variáveis categóricas foram substituídos por uma nova categoria chamada de '*Unknown*'. A coluna *Amenities* foi separada em várias variáveis *dummys*, identificando se cada item identificado em *amenities* estava presente na moradia ou não.

Em seguida, foi realizada uma seleção de variáveis, com o objetivo de obter um modelo prévio, utilizando o método *forward*.

Após esta etapa, preencheu-se valores faltantes utilizando a informação do *host*. Caso não desse para ser estimado, foi imputada a média. Na maior parte, foi imputado considerando o seguinte critério:

- Se a variável for limitada inferiormente e superiormente, foi utilizada a regressão *beta*.
- Se a variável for de contagem, foi utilizada a regressão de Poisson.

As colunas *host_response_rate* e as que começavam com *reviews_scores* foram estimadas com regressão *beta* e transformadas para a escala original. As colunas *bathrooms*, *bedrooms*, *beds*, *reviews_per_month*, *days_first_review*, *days_last_review* foram estimadas utilizando regressão Poisson.

Foi criada uma *dummy* para cada coluna estimada identificando se teve que imputar naquela coluna ou não, e, na parte de seleção de variáveis, será verificado se ela é necessária.

Após isso, foi realizada uma seleção de variáveis utilizando o método *Stepwise*, com uma pequena amostra escolhida aleatoriamente do banco de dados com valores imputados. Em seguida, foram feitas pequenas alterações para chegar no modelo abaixo, com *AIC* com o valor de 7.486.994.

$$\begin{aligned}
\ln(\mu_{price}) = & \beta_0 + \beta_{host_response_timewithinanehour}X_1 + \beta_{host_response_timewithinafewhours}X_2 \\
& + \beta_{host_response_timewithinaday}X_3 + \beta_{host_response_timeafewdaysormore}X_4 \\
& + \beta_{host_is_superhost}X_5 + \beta_{host_total_listings_count}X_6 \\
& + \beta_{host_has_profile_pic}X_7 + \beta_{neighbourhood_group_cleansedBeaconHill}X_8 \\
& + \beta_{neighbourhood_group_cleansedCapitolHill}X_9 \\
& + \beta_{neighbourhood_group_cleansedCascade}X_{10} \\
& + \beta_{neighbourhood_group_cleansedCentralArea}X_{11} \\
& + \beta_{neighbourhood_group_cleansedDelridge}X_{12} \\
& + \beta_{neighbourhood_group_cleansedDowntown}X_{13} \\
& + \beta_{neighbourhood_group_cleansedInterbay}X_{14} \\
& + \beta_{neighbourhood_group_cleansedLakeCity}X_{15} \\
& + \beta_{neighbourhood_group_cleansedMagnolia}X_{16} \\
& + \beta_{neighbourhood_group_cleansedNorthgate}X_{17} \\
& + \beta_{neighbourhood_group_cleansedOtherneighborhoods}X_{18} \\
& + \beta_{neighbourhood_group_cleansedQueenAnne}X_{19} \\
& + \beta_{neighbourhood_group_cleansedRainierValley}X_{20} \\
& + \beta_{neighbourhood_group_cleansedSewardPark}X_{21} \\
& + \beta_{neighbourhood_group_cleansedUniversityDistrict}X_{22} \\
& + \beta_{neighbourhood_group_cleansedWestSeattle}X_{23} + \beta_{latitude}X_{24} \\
& + \beta_{is_location_exact}X_{25} + \beta_{property_typeBed\&Breakfast}X_{26} \\
& + \beta_{property_typeBoat}X_{27} + \beta_{property_typeBungalow}X_{28} + \beta_{property_typeCabin}X_{29} \\
& + \beta_{property_typeCamper/RV}X_{30} + \beta_{property_typeChalet}X_{31} \\
& + \beta_{property_typeCondominium}X_{32} + \beta_{property_typeDorm}X_{33} \\
& + \beta_{property_typeHouse}X_{34} + \beta_{property_typeLoft}X_{35} + \beta_{property_typeOther}X_{36} \\
& + \beta_{property_typeTent}X_{37} + \beta_{property_typeTownhouse}X_{38} \\
& + \beta_{property_typeTreehouse}X_{39} + \beta_{property_typeUnknown}X_{40} + \beta_{property_typeYurt}X_{41} \\
& + \beta_{room_typePrivateroom}X_{42} + \beta_{room_typeSharedroom}X_{43} + \beta_{accommodates}X_{44} \\
& + \beta_{bathrooms}X_{45} + \beta_{bedrooms}X_{46} + \beta_{security_deposit}X_{47} + \beta_{cleaning_fee}X_{48} \\
& + \beta_{guests_included}X_{49} + \beta_{number_of_reviews}X_{50} + \beta_{review_scores_rating}X_{51} \\
& + \beta_{review_scores_accuracy}X_{52} + \beta_{review_scores_location}X_{53} + \beta_{review_scores_value}X_{54} \\
& + \beta_{cancellation_policymoderate}X_{55} + \beta_{cancellation_policystRICT}X_{56} \\
& + \beta_{calculated_host_listings_count}X_{57} + \beta_{reviews_per_month}X_{58} + \beta_{monthfev}X_{59} \\
& + \beta_{monthmar}X_{60} + \beta_{monthabr}X_{61} + \beta_{monthmai}X_{62} + \beta_{monthjun}X_{63} \\
& + \beta_{monthjul}X_{64} + \beta_{monthago}X_{65} + \beta_{monthset}X_{66} + \beta_{monthout}X_{67} \\
& + \beta_{monthnov}X_{68} + \beta_{monthdez}X_{69} + \beta_{week_dayseg}X_{70} + \beta_{week_dayter}X_{71} \\
& + \beta_{week_dayqua}X_{72} + \beta_{week_dayqui}X_{73} + \beta_{week_daysex}X_{74} + \beta_{week_daysáb}X_{75} \\
& + \beta_{days_first_review}X_{76} + \beta_{days_last_review}X_{77} + \beta_{Internet}X_{78} + \beta_{TV}X_{79} \\
& + \beta_{WirelessInternet}X_{80} + \beta_{CableTV}X_{81} + \beta_{ElevatorinBuilding}X_{82} + \beta_{longitude}X_{83}
\end{aligned}$$

a) Ajuste um MLG para estimar a probabilidade de um dado apartamento (isso é, com um nível fixo das covariáveis) ser anunciado por um valor superior a \$200. Qual é o erro quadrático médio do modelo no conjunto de teste?

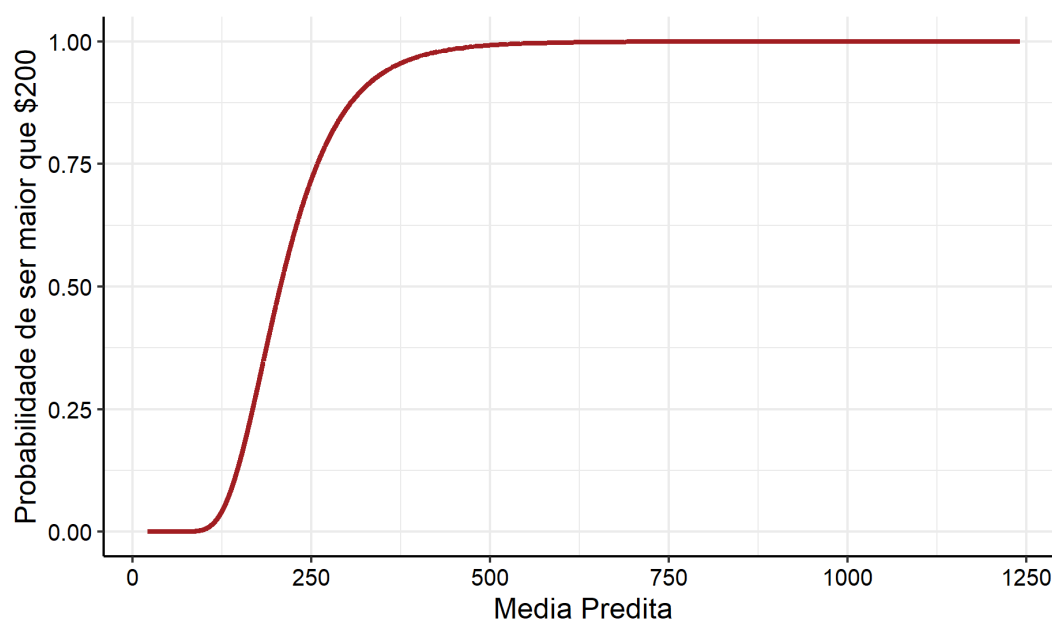
Para estimar a probabilidade de um dado apartamento ser anunciado por um valor superior a \$200, será utilizada como distribuição da variável resposta a distribuição Gama, com função de ligação igual ao logaritmo natural. Utilizando a parametrização dos slides, tem-se que:

- $\alpha = \phi$
- $\beta = \frac{\phi}{\mu}$

Dessa forma, para estimar essas probabilidades, basta utilizar o valor estimado da média e da precisão (ϕ) para obter a distribuição estimada para cada observação e, para cada uma delas, calcular a cauda à direita, a partir do valor de preço igual a \$200.

```
1 phi <- 1/MASS::gamma.dispersion(final) # estimativa da precisao
2 alpha <- phi
3 beta <- phi/predict(final, type = "response")
4 prob.200 <- 1 - pgamma(200, shape = alpha, rate = beta) # calculo
  das probabilidades
5 mu.hat <- predict(final, type = "response") # valores preditos
6 amostra <- sample(1:length(mu.hat), 10000) # amostra de pontos para
  o grafico
7 df10000 <- data.frame(mu.hat = mu.hat[amostra], prob.200 = prob
  .200[amostra])
8 ggplot(df10000, aes(x=mu.hat, y=prob.200)) +
9   geom_line(colour="#A11D21", size=1) +
10   labs(x="Media Predita", y="Probabilidade de ser maior que $200")
  +
11   theme_bw() +
12   theme(axis.title.y=element_text(colour="black", size=12),
13         axis.title.x = element_text(colour="black", size=12),
14         axis.text = element_text(colour = "black", size=9.5),
15         panel.border = element_blank(),
16         axis.line = element_line(colour = "black"))
17 ggsave("prob_200.png", width = 158, height = 93, units = "mm")
```

Figura 1: Média Predita x Probabilidade Estimada do Preço Superar \$200



Abaixo, segue um quadro contendo os erros quadráticos médios do modelo final avaliado no banco de dados de treinamento e no de teste.

Banco de Dados	Erro Quadrático Médio
Treinamento	62,2379
Teste	62,4896

Por meio do quadro acima, percebe-se que os erros quadráticos médios de ambos os bancos são muito próximos, o que indica uma boa capacidade de previsão do preço, a partir das covariáveis mantidas no modelo.

b) Considerando fixas as demais covariáveis, qual é o ponto ou região da cidade com maiores valores anunciados?

Considerando fixas as demais covariáveis e levando em consideração o modelo ajustado, percebe-se que a latitude e a longitude possuem efeito β negativo; portanto, o incremento multiplicativo é menor que 1. Logo, para maximizar o preço, deve-se minimizar a longitude e a latitude. Observando o banco de dados, encontra-se que os menores valores para latitude e longitude são, respectivamente, 47,50728 e -122,41722. Sendo assim, a região que possui os maiores valores anunciados é Cascade.

Também é possível verificar isso observando o incremento multiplicativo β da variável "neighbourhood_group_cleansed" para o fator Cascade, pois possui o maior valor de β . Sendo assim, pode-se dizer que Cascade é a região com maior preço médio anunciado.

c) Descreva a influência da data no preço dos alugueis

É esperado que o mercado de alugueis de imóveis, tal qual é o *Airbnb*, seja fortemente influenciado pela data em que se planeja alugar o imóvel, visto que é uma procura sazonal.

Variável	Coefficiente
fev	0,00989
mar	0,03563
abr	0,05623
mai	0,08402
jun	0,12360
jul	0,13380
ago	0,12560
set	0,08305
out	0,04348
nov	0,03161
dez	0,03378

Ao observar os meses, pode-se notar que, em relação a janeiro, os meses com maiores valores de preços anunciados são Junho, Julho e Agosto, o que é coerente com o período de férias escolares de Verão nos EUA. Dessa forma, espera-se que nesses meses o preço seja maior.

Variável	Coefficiente
Segunda	-0,003849
Terça	-0,006828
Quarta	-0,007522
Quinta	-0,000819
Sexta	0,045930
Sábado	0,048040

Ao observar os dias da semana, pode-se notar que, comparado a domingo, sexta e sábado são os dias em que os preços são maiores, enquanto que, de segunda a quinta, os preços são menores. Este fato condiz com a ideia de que imóveis são mais procurados em finais de semana.

d) Quais são, em ordem decrescente, os cinco principais fatores na determinação do preço das unidades ofertadas?

Os cinco principais fatores, em ordem decrescente, na determinação do preço das unidades são longitude, latitude, "room_typeShared room", "property_typeTent" e "property_typeBoat". Essa constatação foi baseada nos maiores valores absolutos dos coeficientes que estão no quadro abaixo.

Variável	Coeficiente
latitude	-1,2722705
longitude	-0,9326195
room_typeShared room	-0,8230671
property_typeDorm	-0,7131058
property_typeTent	-0,6134797

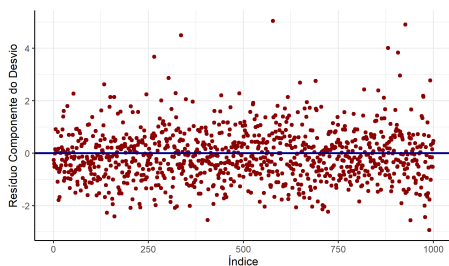
e) De que modo o preço é impactado pelo aumento das notas dos reviews?

Considerando o modelo ajustado, pode-se observar que o valor β para a variável "reviews_score_rating" indica um efeito positivo de aproximadamente $6,7 \times 10^{-3}$. Sendo assim, espera-se que tenha um incremento multiplicativo de $e^{6,7 \times 10^{-3}}$ no valor do preço final. Ou seja, o efeito das reviews não é tão impactante no preço anunciado do imóvel quanto se espera de um mecanismo de avaliação do imóvel.

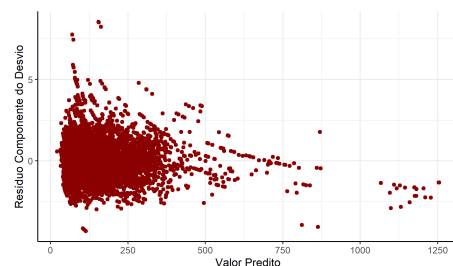
Talvez, a falta de um efeito mais significativo seja resultante do número de valores não observados que as variáveis de "reviews" têm. Foram mais de 140 mil valores faltantes.

f) Enumere as suposições feitas pelo modelo, comente a razoabilidade de cada uma delas e indique formas de melhorar o modelo (não é preciso implementar as estratégias elencadas).

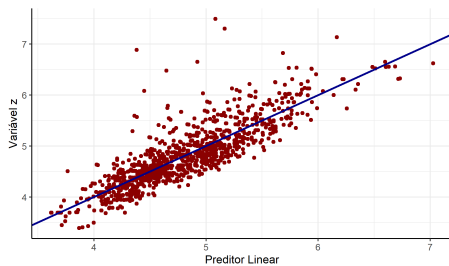
- Independência das observações
- A distribuição da variável "price" é uma Gama
- A forma do preditor linear é adequada
- A função de ligação é adequada



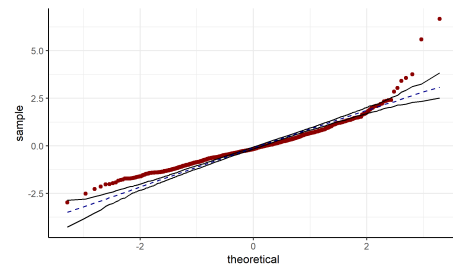
(a) RCD vs. Índice



(b) RCD vs. Valores Preditos



(c) Variável Z vs. Preditor Linear



(d) Gráfico de Envelope

Observando a figura (a), pode-se notar que a suposição de independência é razoável.

Voltando a atenção para a figura (b), pode-se notar que não há muitos valores discrepantes, e que o maior ponto de alavanca encontrado foi o de valor 0,00697 aproximadamente.

Ao observar o gráfico da figura (c), pode-se notar que a forma do preditor linear é razoavelmente adequada, visto que os pontos estão seguindo a reta. Porém, pode-se notar um certo afastamento dessa forma linear para alguns pontos da variável Z.

Em relação ao gráfico da figura (d), observa-se o gráfico de envelope, que é crucial para identificar a adequabilidade do modelo. Observando o gráfico, pode-se notar que o modelo não parece estar totalmente adequado, levando em consideração que muitos pontos observados não estão dentro do intervalo criado por simulações do modelo.

Ou seja, não consegue-se simular, a partir do modelo, um comportamento similar ao comportamento dos dados observados, principalmente para os primeiros quantis e para valores mais à direita, ou seja, para os últimos quantis.

Para melhorar o modelo, poderia ser considerada uma outra distribuição para o preço. Além disso, seria possível utilizar modelos mistos, considerando o efeito do bairro, ou outros estratos. Como haviam muitas variáveis de data no banco, também poderia ser considerado um modelo de séries temporais.

A Código R

```
1 # Arrumando o banco de dados -----
2
3 ### Carregando os pacotes exigidos
4 library(readr)
5 require(tidyverse)
6 set.seed(22019)
7 setwd("C:\\Users\\alvar\\Dropbox\\6Semestre\\Modelos_Lineares_
   Generalizados")
8
9 ### Carregando os dados
10 listings <- read_csv("listings.csv")
11 calendar <- read_csv("calendar.csv")
12 ### Limpando a base
13 calendar <- calendar %>%
14   filter(!is.na(price)) %>% # filtrando os imóveis disponíveis
15   mutate(price=str_replace(.$price, "\\$", "") %>%
16     str_replace("\\,", ",") %>%
17     as.numeric()) # transformando os dados de preços em
   numérico
18 ### Dividindo o banco de dados em treino e teste
19 n.obs <- nrow(calendar)
20 treino <- sample(n.obs, n.obs*.8)
21 teste <- setdiff(1:n.obs, treino)
22 calendar.treino <- calendar[treino, ]
23 calendar.teste <- calendar[teste, ]
24
25 ### Definindo as informações a serem consideradas no modelo
26 covariaveis <- c("date", "host_since", "host_location", "host_
   response_time",
27   "host_response_rate", "host_acceptance_rate", "
   host_is_superhost",
28   "host_total_listings_count", "host_has_profile_pic",
29   "host_identity_verified",
30   "neighbourhood_group_cleansed", "zipcode", "
   latitude", "longitude", "is_location_exact",
31   "property_type", "room_type", "accommodates", "
   bathrooms", "bedrooms", "beds",
32   "bed_type", "amenities", "square_feet", "security_
   deposit", "cleaning_fee",
33   "guests_included", "extra_people", "minimum_nights",
34   "maximum_nights", "number_of_reviews",
35   "first_review", "last_review", "review_scores_
   rating", "review_scores_accuracy",
36   "review_scores_cleanliness", "review_scores_
   checkin", "review_scores_communication",
37   "review_scores_location", "review_scores_value", "
   requires_license", "instant_bookable",
38   "cancellation_policy", "require_guest_profile_
   picture", "require_guest_phone_verification",
39   "calculated_host_listings_count", "reviews_per_
   month")
40
41 # Início do nosso trabalho
42 listings2 = listings %>% dplyr::select(id, covariaveis[-1])
43 calendar2 = calendar %>% dplyr::select(-available)
44 df = left_join(calendar2, listings2, by = c("listing_id" = "id"))
```

```

43 df = df[,-1]
44 df.treino = df[treino, ]
45 df.teste = df[teste, ]
46 dir.create('D:/AirBnB/cache')
47 lapply(ls(), function(i) saveRDS(get(i), paste0('D:/AirBnB/cache/',
48 i, '.rds'))))
49 rm(list = ls()[!ls() %in% c("df.treino", "df.teste")])
50 gc()
51 # Descritiva Preço
52 summary(df.treino$price)
53 hist(df.treino$price) # Escolhemos a Gama
54
55 # Arrumando as variáveis
56 require(lubridate)
57 df.treino2 <- df.treino %>%
58   mutate(year = factor(year(date)),
59           month = month(date, label = TRUE),
60           week_day = wday(date, label = TRUE),
61           days_as_host = as.numeric(date - host_since),
62           host_response_time = factor(host_response_time,
63                                       levels = c("N/A",
64                                                  "within an hour",
65                                                  "within a few hours
66",
67                                                  "within a day",
68                                                  "a few days or more
69",
70                                                  "Unknown",
71                                                  "within an hour",
72                                                  "within a few hours
73",
74                                                  "within a day",
75                                                  "a few days or more
76",
77                                                  ordered = TRUE),
78           host_response_rate = as.numeric(sub("%","",host_response_
79 rate))/100,
80           host_acceptance_rate = as.numeric(sub("%","",host_
81 acceptance_rate))/100,
82           property_type = ifelse(is.na(property_type), "Unknown",
83 property_type),
84           room_type = factor(room_type),
85           bed_type = factor(bed_type),
86           security_deposit = as.numeric(gsub(",|[$]","",security_
87 deposit)),
88           cleaning_fee = as.numeric(gsub(",|[$]","",cleaning_fee)),
89           extra_people = as.numeric(gsub(",|[$]","",extra_people)),
90           days_first_review = as.numeric(date - first_review),
91           days_last_review = as.numeric(date - last_review),
92           cancellation_policy = factor(cancellation_policy)) %>%
93   dplyr::select(-date,-host_since,-host_location,-zipcode,-first_
94 review,-last_review,-requires_license)
95
96 amen = str_split(str_remove_all(df.treino2$amenities, '[{}\\"]'),
97                  ', ', simplify = T) %>% as.data.frame() %>%
98   gather(V, word) %>% filter(str_length(word) > 0) %>% distinct(
99   word) %>% na.omit() %>% as_vector()

```

```

90 names(amen) <- NULL
91
92 for(i in amen){
93   df.treino2 = df.treino2 %>%
94     mutate(!i := str_detect(amenities, i))
95 }
96
97 df.treino2 <- dplyr::select(df.treino2, -amenities)
98
99 df.treino3 = df.treino2 %>%
100   dplyr::select(-square_feet, -host_acceptance_rate) %>%
101   replace_na(list(cleaning_fee = 0,
102                  security_deposit = 0))
103 covars <- names(sort(colMeans(is.na(df.treino3))[colMeans(is.na(df.
104   treino3)) > 0],
105                 decreasing = T))
106 df.treino3 %>% filter_at(vars(covars), all_vars(!is.na(.))) %>%
107   nrow()
108 df.treino4 = df.treino3 %>% filter_at(vars(covars), all_vars(!is.na
109   (.))) %>%
110   mutate_if(is.logical, list(~as.numeric(.)))
111 rm(list = ls()[!ls()%in%c("df.treino4")])
112 gc()
113
114 mod3 <- glm(price ~ . - 'Dog(s)' - 'Cat(s)' - 'Other pet(s)', family
115   = Gamma(link="log"), data = df.treino4)
116 summary(mod3)
117
118 require(MASS)
119 mod0 = glm(price ~ 1, family = Gamma(link="log"), data = df.treino4
120   )
121
122 ultimo_modelo <- glm(price ~ accommodates + room_type +
123   neighbourhood_group_cleansed +
124   bedrooms + cleaning_fee + month + bathrooms
125   + days_last_review +
126   host_is_superhost + review_scores_location +
127   reviews_per_month +
128   property_type + 'Elevator in Building' +
129   week_day + calculated_host_listings_count +
130   latitude + bed_type + guests_included +
131   longitude + 'Cable TV' +
132   review_scores_value + review_scores_rating +
133   host_has_profile_pic +
134   Doorman + 'Laptop Friendly Workspace' +
135   Shampoo + 'Suitable for Events' +
136   'Smoking Allowed' + review_scores_accuracy +
137   host_total_listings_count +
138   '24-Hour Check-in' + 'Pets live on this
139   property' + 'Indoor Fireplace' +
140   'Air Conditioning' + host_identity_verified
141   + 'Fire Extinguisher' +
142   security_deposit + 'Wireless Internet' +
143   Internet + maximum_nights +
144   Iron + Hangers + is_location_exact + 'Hot

```

```

132   Tub' + host_response_time +
133       'Wheelchair Accessible' + Kitchen + require_
134   guest_profile_picture +
135       'Safety Card' + beds + 'Lock on Bedroom Door
136   ' + 'Hair Dryer' +
137       'Washer / Dryer' + Heating + cancellation_
138   policy + 'Smoke Detector' +
139       Breakfast + days_first_review + number_of_
140   reviews + minimum_nights +
141       days_as_host + instant_bookable + TV + host_
142   response_rate +
143       'Family/Kid Friendly' + 'Carbon Monoxide
144   Detector' + review_scores_checkin +
145       'Buzzer/Wireless Intercom' + extra_people +
146   Essentials +
147       review_scores_cleanliness + Pool + year +
148   review_scores_communication +
149       'First Aid Kit' + 'Free Parking on Premises'
150   + require_guest_phone_verification +
151       'Pets Allowed', family = Gamma(link="log"),
152   data = df.treino4)
153
154 mod_step = stepAIC(ultimo_modelo, direction = "forward", scope =
155   list(upper = mod3, lower = mod0))
156
157 saveRDS(mod_step, 'D:/AirBnB/cache/modelo_stepAIC.rds')
158 saveRDS(mod3, 'D:/AirBnB/cache/modelo_3.rds')
159
160 rm(mod_step, mod3, df.treino4)
161 gc()
162 # Imputação -----
163
164 require(tidyverse)
165 require(lubridate)
166 df.teste <- readRDS('D:/AirBnB/cache/df.teste.rds')
167 df.teste2 <- df.teste %>%
168   mutate(year = factor(year(date)),
169     month = month(date, label = TRUE),
170     week_day = wday(date, label = TRUE),
171     days_as_host = as.numeric(date - host_since),
172     host_response_time = factor(host_response_time,
173       levels = c("N/A",
174         "within an hour",
175         "within a few hours
176
177     ",
178         "within a day",
179         "a few days or more
180
181     "),
182     labels = c("Unknown",
183       "within an hour",
184       "within a few hours
185
186     ",
187       "within a day",
188       "a few days or more
189
190     "),
191     ordered = TRUE),

```

```

174     host_response_rate = as.numeric(sub("%","",host_response_
175     rate))/100,
176     host_acceptance_rate = as.numeric(sub("%","",host_
177     acceptance_rate))/100,
178     property_type = ifelse(is.na(property_type), "Unknown",
179     property_type),
180     room_type = factor(room_type),
181     bed_type = factor(bed_type),
182     security_deposit = as.numeric(gsub(",|[$]","",security_
183     deposit)),
184     cleaning_fee = as.numeric(gsub(",|[$]","",cleaning_fee)),
185     extra_people = as.numeric(gsub(",|[$]","",extra_people)),
186     days_first_review = as.numeric(date - first_review),
187     days_last_review = as.numeric(date - last_review),
188     cancellation_policy = factor(cancellation_policy)) %>%
189     dplyr::select(-date,-host_since,-host_location,-zipcode,-first_
190     review,-last_review,-requires_license)
191
192 amen <- str_split(str_remove_all(df.teste2$amenities, '[{}\\"]'),
193     ', ', simplify = T) %>% as.data.frame() %>%
194     gather(V, word) %>% filter(str_length(word) > 0) %>% distinct(
195     word) %>% na.omit() %>% as_vector()
196
197 names(amen) <- NULL
198
199 for(i in amen){
200     df.teste2 <- df.teste2 %>%
201     mutate(!i := str_detect(amenities, i))
202 }
203
204 df.teste2 <- dplyr::select(df.teste2, -amenities)
205
206 saveRDS(df.teste2, 'D:/AirBnB/cache/df.teste2.rds')
207
208 df.treino3 <- readRDS('D:/AirBnB/cache/df.treino3.rds')
209 df.teste2 <- readRDS('D:/AirBnB/cache/df.teste2.rds')
210 teste <- df.teste2
211
212 # Colunas numericas faltando
213
214 colSums(is.na(df.treino3))[colMeans(is.na(df.treino3))>0] # Remove
215 os 67 errados
216
217 df_correto <- df.treino3 %>% filter(!is.na(days_as_host))
218
219 colSums(is.na(df_correto))[colMeans(is.na(df_correto))>0] # Todos
220 que tinham 67 desapareceram
221
222 col_incomp <- names(colMeans(is.na(df_correto))>0)
223
224 col_comp <- names(df_correto)[colMeans(is.na(df_correto))==0][-1] #
225 usar as colunas completas para prever as incompletas
226 retirando price
227
228 fo_host <- paste0(' ', str_subset(col_comp, 'host')[-1], ' ',
229     collapse = '+') # Prever apenas para as colunas de host, não
230 converge se usar todas
231 # as colunas completas

```



```

220
221
222 # Métodos escolhidos -----
223 # host_response_rate, reviews_scores: Regressão beta
224 # bathrooms, bedrooms, beds, reviews_per_month, days_first_review,
    days_last_review: regressão poisson
225
226 # regressão beta -----
227
228 require(betareg)
229 # host_response_rate
230 temp <- df_correto
231 temp$host_response_rate[temp$host_response_rate==1] <- .99
232 temp$host_response_rate[temp$host_response_rate==0] <- .01
233
234
235 hrr_betareg_mod <- betareg(as.formula(paste('host_response_rate ~',
    fo_host))), data = temp)
236 hrr_predicted <- predict(hrr_betareg_mod, newdata = df_correto)
237 df_correto <- df_correto %>%
238   mutate(host_response_rate_informed = as.numeric(!is.na(host_
    response_rate)),
239     host_response_rate = coalesce(host_response_rate, hrr_
    predicted))
240
241 test_pred <- predict(hrr_betareg_mod, newdata = teste)
242 teste <- teste %>%
243   mutate(host_response_rate_informed = as.numeric(!is.na(host_
    response_rate)),
244     host_response_rate = coalesce(host_response_rate, test_
    pred, mean(host_response_rate, na.rm = T)))
245
246 saveRDS(hrr_betareg_mod, 'D:/AirBnB/cache/host_response_rate_
    betareg.rds')
247
248 # review_scores_rating escala de 0 a 100
249
250 temp <- df_correto
251 temp <- temp %>%
252   mutate_at(vars(starts_with('review_scores_')), list(~./max(., na.
    rm = T))) %>%
253   mutate_at(vars(starts_with('review_scores_')), list(~case_when(
254     .==1~.99,
255     .==0~.01,
256     TRUE~.
257   )))
258
259 rev_sc_rating <- betareg(as.formula(paste('review_scores_rating ~',
    fo_host))), data = temp)
260 rev_sc_predicted <- predict(rev_sc_rating, newdata = df_correto)
261
262 df_correto <- df_correto %>%
263   mutate(review_scores_rating_informed = as.numeric(!is.na(review_
    scores_rating)),
264     review_scores_rating = coalesce(review_scores_rating, rev_
    sc_predicted*100))
265
266 test_pred <- predict(rev_sc_rating, newdata = teste)

```

```

267 teste <- teste %>%
268   mutate(review_scores_rating_informed = as.numeric(!is.na(review_
269     scores_rating)),
270     review_scores_rating = coalesce(review_scores_rating, test
271       _pred*100, mean(review_scores_rating, na.rm = T)))
272 rm(rev_sc_rating, rev_sc_predicted, hrr_betareg_mod, hrr_predicted)
273 gc()
274 # Os outros reviews_scores possui escala de 0 a 10
275 col_reviews <- str_subset(col_incomp, 'review_scores_(?!rating)')
276 for(j in col_reviews){
277   # j= col_reviews[1]
278   j <- as.name(j)
279
280   rev_sc_reg <- betareg(as.formula(paste(j, '~', fo_host)), data =
281     temp)
282   rev_sc_predicted <- predict(rev_sc_reg, newdata = df_correto)
283
284   df_correto <- df_correto %>%
285     mutate(!paste0(j, '_informed') := as.numeric(!is.na(!j)),
286       !!j := coalesce(!j, rev_sc_predicted*10))
287
288   test_pred <- predict(rev_sc_reg, newdata = teste)
289   teste <- teste %>%
290     mutate(!paste0(j, '_informed') := as.numeric(!is.na(!j)),
291       !!j := coalesce(!j, test_pred*10, mean(!j, na.rm = T))
292     )
293
294   saveRDS(rev_sc_reg, paste0('D:/AirBnB/cache/', j, '_betareg.rds')
295     )
296   cat('writted ', paste0('D:/AirBnB/cache/', j, '_betareg.rds'), '\
297     n')
298 }
299 rm(temp, col_reviews); gc()
300 # Regressão Poisson -----
301
302 poi_reg_var <- unlist(str_split('bathrooms,bedrooms,beds,reviews_
303   per_month,days_first_review,days_last_review', ','))
304
305 for(j in poi_reg_var){
306   # j= poi_reg_var[2]
307   j <- as.name(j)
308
309   poi_reg <- glm(as.formula(paste(j, '~', fo_host)), data = df_
310     correto, family = poisson('log'))
311   poi_predicted <- predict(rev_sc_reg, newdata = df_correto, type=
312     'response')
313
314   df_correto <- df_correto %>%
315     mutate(!paste0(j, '_informed') := as.numeric(!is.na(!j)),
316       !!j := coalesce(!j, rev_sc_predicted))
317
318   test_pred <- predict(poi_reg, newdata = teste)
319   teste <- teste %>%
320     mutate(!paste0(j, '_informed') := as.numeric(!is.na(!j)),
321       !!j := coalesce(!j, test_pred, mean(!j, na.rm = T)))
322 }

```

```

316 saveRDS(rev_sc_reg, paste0('D:/AirBnB/cache/', j, '_poisson_reg.
317 rds'))
318 cat('writted ', paste0('D:/AirBnB/cache/', j, '_poisson_reg.rds'),
    '\n')
319 }
320 teste <- teste %>% replace_na(list(cleaning_fee = 0,
321                                   security_deposit = 0))
322 saveRDS(teste, 'D:/AirBnB/cache/teste_imputed.rds')
323 saveRDS(df_correto, 'D:/AirBnB/cache/df_imputed.rds')
324
325 rm(list = ls()); gc()
326
327 # Construção de Modelo -----
328
329 require(tidyverse)
330 df_imputed <- readRDS('D:/AirBnB/cache/df_imputed.rds')
331 df_imputed <- df_imputed %>% mutate_if(is.factor, list(~factor(.,
    ordered = F))) %>%
332   mutate_if(is.logical, list(~as.numeric(.)))
333 # Modelo escolhido pelo método automático, com algumas alterações
334 fit <- glm(price ~ host_response_time + host_is_superhost + host_
    total_listings_count +
335             host_has_profile_pic + neighbourhood_group_cleansed +
    latitude +
336             is_location_exact + property_type + room_type +
    accommodates +
337             bathrooms + bedrooms + security_deposit + cleaning_fee
    +
338             guests_included + number_of_reviews + review_scores_
    rating +
339             review_scores_accuracy + review_scores_location +
    review_scores_value +
340             cancellation_policy + calculated_host_listings_count +
    reviews_per_month +
341             month + week_day + days_first_review + days_last_
    review +
342             Internet + TV + 'Wireless Internet' + 'Cable TV' +
    'Elevator in Building' + longitude, family = Gamma('
343 log'), data = df_imputed)
344
345 summary(fit)
346
347 teste_imputed <- readRDS('D:/AirBnB/cache/teste_imputed.rds')
348 teste_imputed <- teste_imputed %>%
349   replace_na(list(host_is_superhost = T, # imputado pela Moda
350                   host_response_time = 'Unknown', # Nova categoria
351                   host_has_profile_pic = T)) %>% # Moda
352   mutate(host_total_listings_count = ifelse(is.na(host_total_
    listings_count), mean(host_total_listings_count, na.rm = T),
353                                             host_total_listings_
    count)) %>% # host_total_listings_count imputado pela média
354   mutate_if(is.factor, list(~factor(., ordered = F))) %>%
355   mutate_if(is.logical, list(~as.numeric(.)))
356
357
358
359

```

```

360 predito_teste <- predict(fit, newdata = teste_imputed, type = '
    response') # Predizendo
361
362 # Calculando o RMSE
363 (RMSE_treino <- sqrt(mean((df_imputed$price - fit$fitted.values)^2)
    )) # EQM treino
364 (RMSE_teste <- sqrt(mean((teste_imputed$price - predito_teste)^2)))
    # EQM teste
365 cbind(RMSE_treino, RMSE_teste)
366
367
368 # Gráficos de Diagnóstico
369
370 h <- hatvalues(fit)
371
372
373
374 ro <- residuals.glm(fit, type = 'deviance')
375
376 fi <- MASS::gamma.shape(fit)$alpha
377
378
379 RcD <- ro*sqrt(fi/(1 - h))
380
381
382
383 require(ggplot2)
384
385 ggplot(data.frame(x = 1:1000, y = RcD[1:1000]), aes(x, y)) + geom_
    point(color = 'darkred') +
386   geom_hline(yintercept = 0, color = 'darkblue', size = 1.0) +
387   labs(x = 'Índice', y = 'Resíduo Componente do Desvio') +
388   theme_bw() +
389   theme(panel.border = element_blank(),
390         axis.line = element_line('black'))
391
392
393 ggplot(data.frame(x = fit$fitted.values[1:10000], y = RcD[1:10000])
    , aes(x, y)) + geom_point(color = 'darkred') +
394   labs(x = 'Valor Predito', y = 'Resíduo Componente do Desvio') +
395   theme_bw() +
396   theme(panel.border = element_blank(),
397         axis.line = element_line('black'))
398
399 w <- fit$weights
400 eta <- predict(fit)
401 z <- eta + resid(fit, type="pearson")/sqrt(w)
402 ggplot(data.frame(x = eta[1:1000], y = z[1:1000]), aes(x, y)) +
    geom_point(color = 'darkred') +
403   labs(x = 'Preditor Linear', y = 'Variável z') +
404   theme_bw() +
405   theme(panel.border = element_blank(),
406         axis.line = element_line('black')) +
407   geom_abline(intercept = 0, slope = 1, color = 'darkblue', size =
    1.0)
408
409
410 # Gráfico de Envelope -----

```

```

411 set.seed(22019)
412 ind = sample(1:747566 , 1000)
413 df_graf = df_imputed[ind, ]
414 fit_graf = glm(price ~ host_response_time + host_is_superhost +
415               host_total_listings_count +
416               host_has_profile_pic + neighbourhood_group_
417               cleansed + latitude +
418               is_location_exact + property_type + room_type +
419               accommodates +
420               bathrooms + bedrooms + security_deposit + cleaning
421               _fee +
422               guests_included + number_of_reviews + review_
423               scores_rating +
424               review_scores_accuracy + review_scores_location +
425               review_scores_value +
426               cancellation_policy + calculated_host_listings_
427               count + reviews_per_month +
428               month + week_day + days_first_review + days_last_
429               review +
430               Internet + TV + 'Wireless Internet' + 'Cable TV' +
431               'Elevator in Building' + longitude, family = Gamma
432               ('log'), data = df_graf)
433 envelgama <- function(fit.model,ligacao,estphi){
434   # fit.model: objeto com o resultado do ajuste do MLG obtido
435   #          atravs da funcao glm
436   # ligacao: funcao de ligacao (mesmo nome usado pela funcao glm (colocar
437   #          entre aspas))
438   # estphi: mtodo de estimacao para o parametro phi
439   #          1 - Mxima verossimilhana
440   #          2 - Mtodo dos Momentos
441   #          3 - Default do R
442
443   #par(mfrow=c(1,1))
444   X <- model.matrix(fit.model)
445   n <- nrow(X)
446   p <- ncol(X)
447   h <- hatvalues(fit.model)
448   ro <- resid(fit.model,type="response")
449   #fi <- (n-p)/sum((ro/(fitted(fit.model)))^ 2)
450   #library(MASS)
451   #fi <- gamma.shape(fit.model)$alpha
452   if (estphi == 1)
453   {
454     library(MASS)
455     fi <- gamma.shape(fit.model)$alpha
456   }
457   else if (estphi == 2)
458   {
459     #ro <- resid(fit.model,type="response")
460     fi <- (n-p)/sum((ro/(fitted(fit.model)))^ 2)}
461   else if (estphi == 3)
462   {
463     fi <- 1/summary(fit.model)$dispersion
464   }
465   td <- resid(fit.model,type="deviance")*sqrt(fi/(1-h))
466   #
467   e <- matrix(0,n,100)
468   #

```

```

458 for(i in 1:100){
459   resp <- rgamma(n,fi)
460   resp <- (fitted(fit.model)/fi)*resp
461   fit <- glm(resp ~ X, family=Gamma(link=ligacao))
462   h <- hatvalues(fit)
463   if (h==1){h = 0.99}
464   ro <- resid(fit,type="response")
465   #phi <- (n-p)/sum((ro/(fitted(fit)))^ 2)
466   #library(MASS)
467   #fi <- gamma.shape(fit.model)$alpha
468   if (estphi == 1)
469   {
470     library(MASS)
471     phi <- gamma.shape(fit.model)$alpha
472   }
473   else if (estphi == 2)
474   {
475     #ro <- resid(fit.model,type="response")
476     phi <-(n-p)/sum((ro/(fitted(fit.model)))^ 2)}
477   else if (estphi == 3)
478   {
479     phi <- 1/summary(fit.model)$dispersion
480   }
481
482   e[,i] <- sort(resid(fit,type="deviance")*sqrt(phi/(1-h)))}
483 #
484 e1 <- numeric(n)
485 e2 <- numeric(n)
486 #
487 for(i in 1:n){
488   eo <- sort(e[i,])
489   e1[i] <- (eo[2]+eo[3])/2
490   e2[i] <- (eo[97]+eo[98])/2}
491 #
492 med <- apply(e,1,mean)
493 faixa <- range(td,e1,e2)
494 #par(pty="s")
495 qqnorm(td, xlab="Percentil da N(0,1)",
496         ylab="Resduo Componente do Desvio", ylim=faixa, pch=16,
497         main="",cex=.5,cex.axis=1.1,cex.lab=1.1)
498 par(new=T)
499 qqnorm(e1,axes=F,xlab="",ylab="",type="l",ylim=faixa,lty=1, main=
500        "")
501 par(new=T)
502 qqnorm(e2,axes=F,xlab="",ylab="", type="l",ylim=faixa,lty=1, main=
503        "")
504 par(new=T)
505 qqnorm(med,axes=F,xlab="", ylab="", type="l",ylim=faixa,lty=2,
506        main="")
507 #-----#
508 td <-<- td
509 med <-<- med
510 e1 <-<- e1
511 e2 <-<- e2
512 }
513
514 envelgama(fit_graf, 'log', 3)

```

```

512 ggplot(data.frame(td, med, e1, e2)) +
513   stat_qq(aes(sample = td), col = 'darkred') +
514   stat_qq(aes(sample = e1), geom = 'line') +
515   stat_qq(aes(sample = e2), geom = 'line') +
516   stat_qq(aes(sample = med), geom = 'line', linetype = 2, color = '
    darkblue') +
517   theme_bw() +
518   theme(panel.border = element_blank(),
519         axis.line = element_line('black'))

```