

Text/Opinion Mining no Twitter: Prisao do Presidente da Dolly

Rafael Bicudo

10 de maio de 2018

Projeto 1 - Analise de Sentimentos em Redes Sociais - Dolly

O objetivo deste trabalho e capturar dados da rede social Twitter e realizar analises de texto e sentimento. Escolheu-se 10/05/2018 por ser a data da prisao preventiva do dono da empresa de refrigerantes 'Dolly' (Laerte Codonho), e em que a palavra "dolly" manteve-se como um dos assuntos mais citados na plataforma (fonte: <https://trends24.in/brazil/>), portanto sendo um bom referencial para a proposta do projeto. Para sua execucao, diversos pacotes devem ser instalados e carregados.

Todo o projeto sera descrito de acordo com suas etapas. Primeiro, apos a coleta dos tweets e limpeza, far-se-a uma analise preliminar dos termos e suas relacoes estatisticas para, em seguida, prosseguir com a analise de sentimento.

```
library(rtweet)
library(dplyr)
library(knitr)
library(rmarkdown)
```

Etapa 1 - Autenticacao

Abaixo, montam-se as variaveis contendo as chaves necessarias a autenticacao e conexao. Lembre-se que precisa ter uma conta criada no Twitter e se criar uma aplicacao. Para fins de privacidade, os dados referentes a essas chaves foram omitidos.

```
# Criando autenticacao no Twitter
app_rbr <- "xxxxx"
consumer_key_rbr <- "xxxxx"
consumer_secret_rbr <- "xxxxx"
```

Etapa 2 - Conexao

Aqui, testa-se a conexão e se capturam os tweets. Quanto maior sua amostra, mais precisa sera a analise, assim a quantidade escolhida foi de 1000 tweets e, da mesma forma como descrito no inicio do trabalho, o objeto de procura 'dolly'. (Obs.: Como a Rest API do Twitter possui algumas limitacoes, das quais nao reaver dados com mais de 2 semanas, após se fazer a pesquisa, os dados dos tweets brutos foram salvos no arquivo 'tweetsSearch_dolly.csv' para se manter a reproducibilidade)

```
# Realizando a coleta dos tweets (10/05/18)
tema <- "Dolly"
qtd_tweets <- 1000
lingua <- "pt"
#tweetsSearch <- search_tweets(tema, n = qtd_tweets, lang = lingua, include_rts = FALSE)
tweetsSearch <- as.tbl(read.csv('tweetsSearch_dolly.csv', stringsAsFactors = F))
# Selecionando o vetor de textos para analise e criando um indicador para cada tweet
tweets_text <- tweetsSearch %>%
```

```
select(text) %>%
mutate(tweet_id = row_number()) %>%
select(tweet_id, text)
```

Etapa 3 - Tratamento dos dados coletados

Aqui, instalam-se os pacotes tidytext, stringr e stopwords para text mining. Começando pela decodificação de alguns encodings comuns de arquivos retirados diretamente da web, parte-se para a transformação dos tweets coletados em tokens, processo que transforma nosso dataset de tweets em um com 1 termo por observação/linha, e realiza algumas etapas de limpeza; para, na sequência, fazerem-se as transformações de limpeza remanescentes específicas (nomes de usuários, alguns resíduos, dentre outros) e, através de análises suscetivas dos termos mais frequentes, a criação das stopwords. Por fim acontecem duas visualizações para ilustrar os termos mais usados: um gráfico de barras para contagem, e uma nuvem de palavras (wordcloud), para ilustrar os termos com maior frequência de acordo com seu tamanho e cor.

Etapa 2: Tratamento dos dados

```
# Ativação dos pacotes necessários ao Text Mining e manipulação dos dados
library(tidytext)
library(stringr)
library(stopwords)
library(ggplot2)
library(wordcloud)
library(RColorBrewer)

# Remoção de problemas de encoding comuns ao fazer text mining pela internet

tweets_text$text <- tweets_text$text %>%
  iconv(to = "ASCII//TRANSLIT") %>% # Corrige encoding do texto do post (ao menos uma parte deles)
  iconv(sub="", 'UTF-8', 'ASCII') # Remove emojis

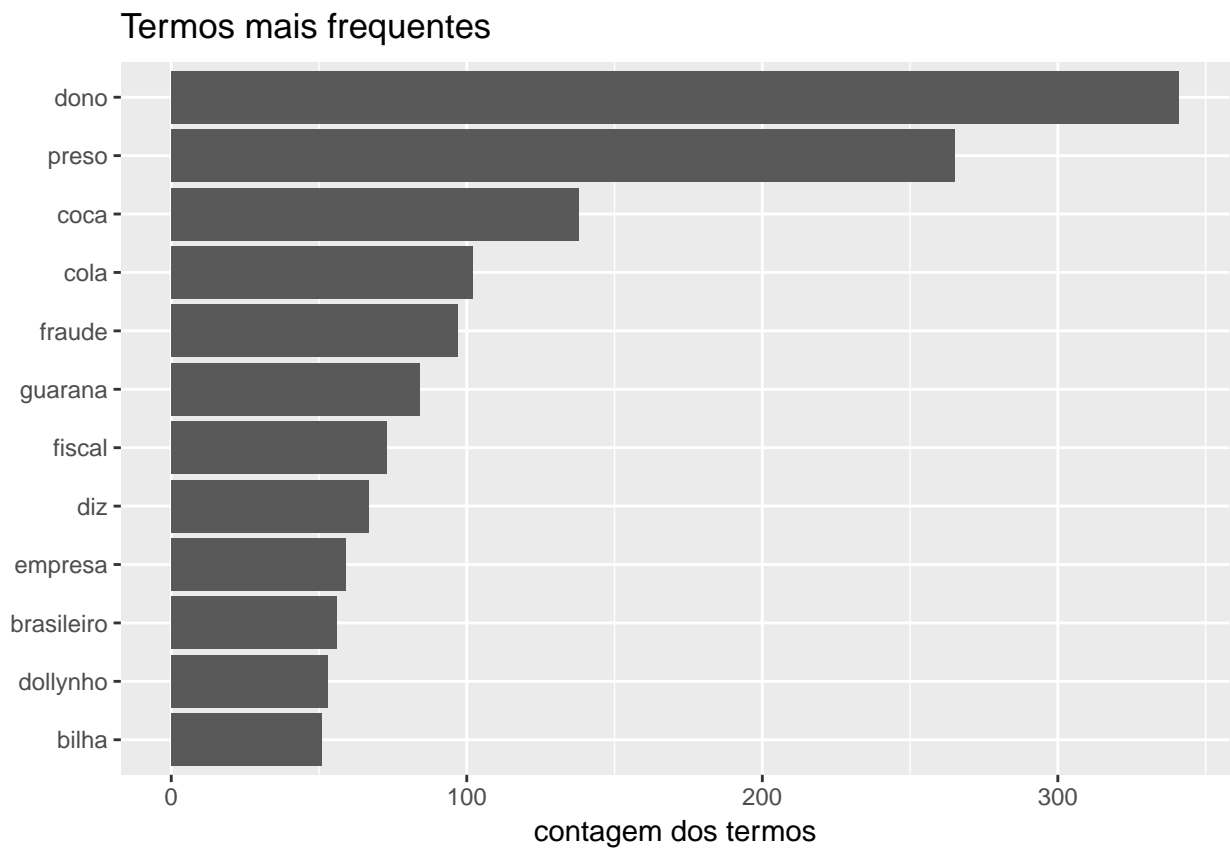
# Criação dos tokens e limpeza final dos dados

tweets_tidy <- tweets_text %>%
  unnest_tokens(word, text) %>% # Transformando texto em tokens
  filter(!word %in% c("https", "t.co", "amp", 'http', 'htt', 'bmr9ju8t8y', 'fwplgwXu4o', 'f0', 'es', '009',
    !word %in% tolower(tweetsSearch$screen_name), # Removendo nome de usuários
    !grepl("^\\d+$", word)) # Removendo números

# Definição das stopwords
stopwords_assunto <- c('nao', 'sao', 'pra', 'q', 'dolly', 'r', 'u', 'refrigerante', 'refrigerantes')
stopwords_tweets <- data_frame(word = c(stopwords('portuguese'), stopwords_assunto))

# Termos mais frequentes
tweets_tidy %>%
  anti_join(stopwords_tweets) %>%
  count(word, sort = TRUE) %>%
  filter(n > 50) %>%
  mutate(word = reorder(word, n)) %>%
  ggplot(aes(word, n)) +
  geom_col() +
  xlab(NULL) +
  coord_flip() +
```

```
labs(title = 'Termos mais frequentes', y = 'contagem dos termos')
```



```
# Gerando uma nuvem palavras
tweets_tidy %>%
  anti_join(stopwords_tweets) %>%
  count(word) %>%
  with(wordcloud(word, n, min.freq = 5,
                 max.words = 200, random.order = FALSE, rot.per = 0.35,
                 colors = brewer.pal(8, "Dark2")))
```



Atraves da observacao dos grafico, percebe-se a quantidade alta de termos relacionados as noticias em pauta do dia (dono, preso, fraude, fiscal, coca, cola), confirmando a suposicao inicial de que alavancagem do assunto (dolly) em relacao aos termos mais usados deveu-se a prisao do dono da marca e a grande repercussao publica do evento.

Etapa 4 - Associações estatísticas entre a frequência das palavras por termo e dendograma

Cria-se, agora, com o auxílio do pacote `tm`, gera-se uma TDM, ou Matriz de Termos por Documento, com o intuito de se verificar relações estatísticas em relação a quantidade de termos em cada tweet, terminando com a visualização de um processo de clusterização hierárquica (dendograma).

```
# Carregado pacotes necessario
library(tm)

# Criacao do Corpus e TDM para analise estatistica
tweets_tdm <- tweets_tidy %>%
  anti_join(stopwords_tweets) %>%
  group_by(tweet_id, word) %>%
  summarise(count = n()) %>%
  cast_tdm(document = tweet_id, term = word, value = count)

# Encontrando as palavras que aparecem com maior frequencia por termo
findFreqTerms(tweets_tdm, lowfreq = 90)
```

```
## [1] "coca"    "cola"    "fraude"  "preso"   "dono"
```

```

# Buscando associacoes de palavras intra-termos
findAssocs(tweets_tdm, c('dono'), c(0.30))

## $dono
##  preso fiscal fraude
##  0.51  0.35  0.31

# Removendo termos esparsos (nao utilizados frequentemente)
tweets_maisUsados <- removeSparseTerms(tweets_tdm, sparse = 0.95)

# Matriz de distancias
tweets_df = as.data.frame(as.matrix(tweets_maisUsados))
tweets_distancias <- dist(tweets_df)

# Realizando a clusterizacao hierarquica dos dados
tweets_hclust <- hclust(d = tweets_distancias)

# Criando o dendograma (verificando como as palavras se agrupam)
plot(tweets_hclust)

# Verificando os grupos
cutree(tweets_hclust, k = 5)

##      coca      cola      bilha      fraude      preso      dono
##      1         1         2         2         3         4
## empresa  fiscal brasileiro  guarana      diz
##      2         2         5         5         2

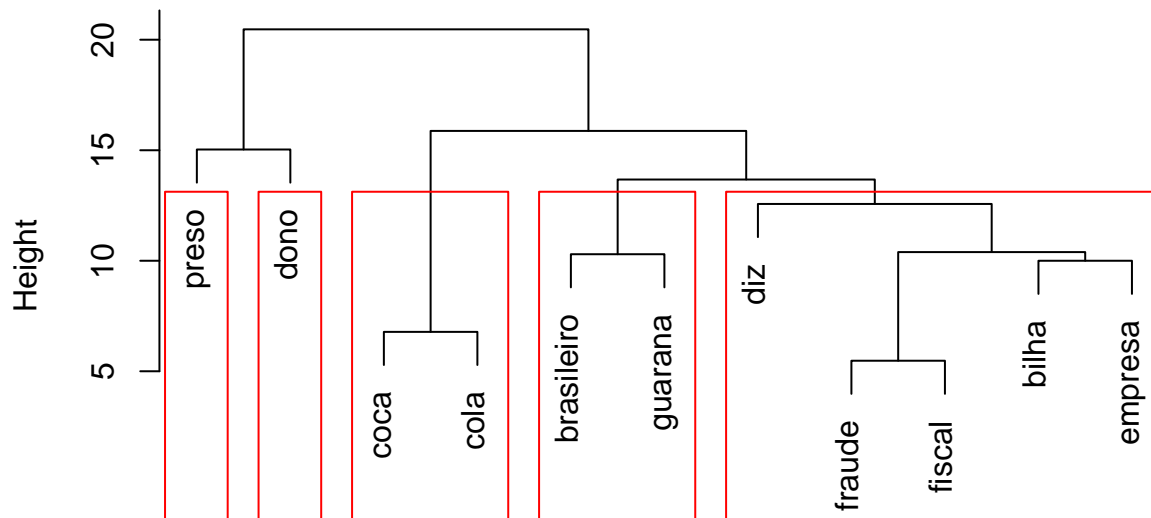
```

```

# Visualizando os grupos de palavras no dendograma
rect.hclust(tweets_hclust, k = 5, border = "red")

```

Cluster Dendrogram



```
tweets_distancias
hclust (*, "complete")
```

Um fato interessante é a altíssima frequência dos termos “coca” e “cola”, um forte indicador da alta absorção e aceitação popular da foto do momento da prisão (em que o acusado segura um cartaz com os dizeres “Preso pela Coca Cola”). Outro fato de importância é a visualização da clusterização hierárquica das frequências: tirando os tweets com maior enfoque no produto e na foto supracitada, os 2 maiores supostos grupos focam-se no fato da prisão em si, e no motivo da acusação (fraude fiscal).

Etapa 5 - Análise de Sentimento

Por último, pode-se proceder à análise de sentimento. A partir do pacote `lexiconPT`, desenvolvido por um brasileiro e reunindo os principais léxicos nacionais, gera-se uma lista de palavras positivas e negativas a partir da escolha de um dos integrantes do pacote. Após se atribuir a cada termo uma pontuação positiva, negativa ou neutra (1, -1, 0, respectivamente), realiza-se uma comparação da distribuição do sentimento predominante por tweet, e uma comparação da contagem dos termos mais usados por sentimento.

```
## Análise de Sentimento (Opinion Mining)

# Carregando pacotes necessários
# devtools::install_github("sillasgonzaga/lexiconPT")
library(lexiconPT)
library(tidyr)

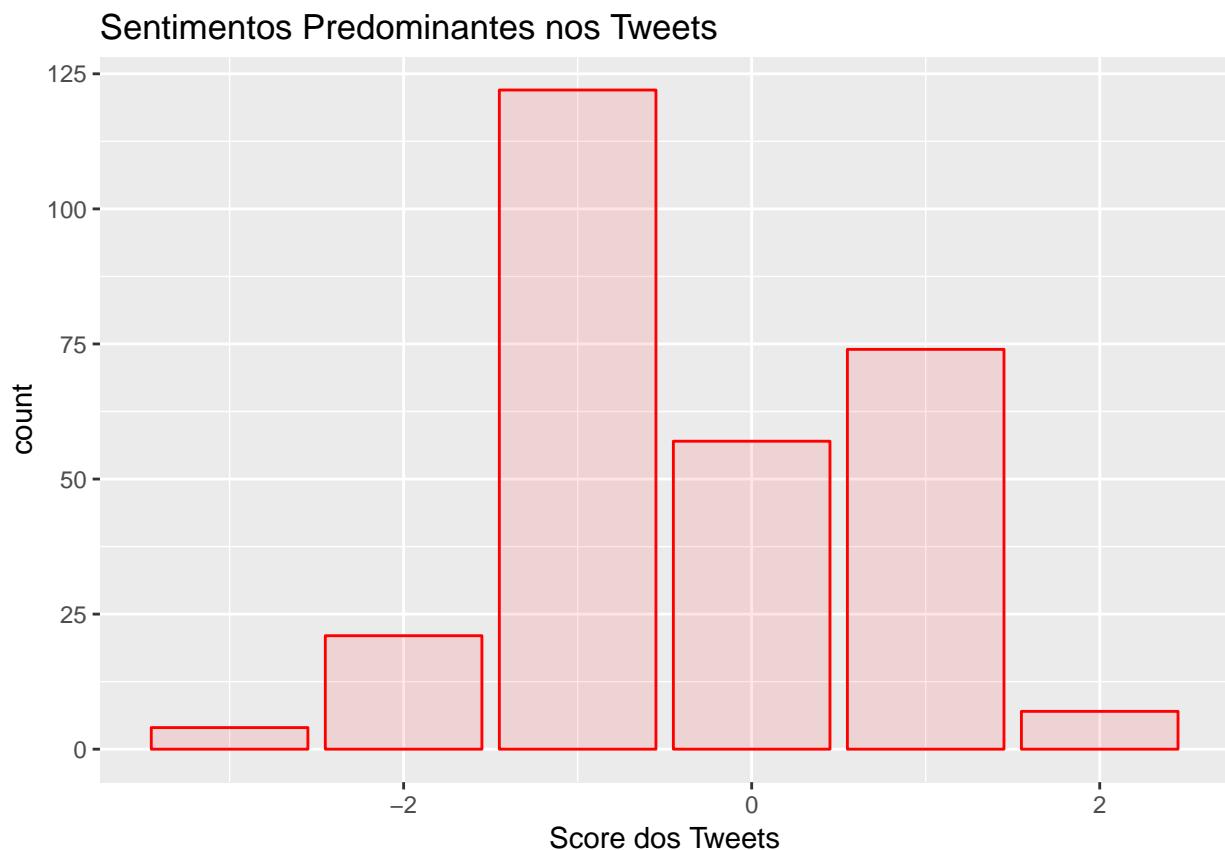
# Carregando léxico escolhido
data("sentiLex_lem_PT02")
sentiLex <- as.tbl(sentiLex_lem_PT02)
sentimentos <- sentiLex %>%
  select(word = term, polarity)
```

```

# Calculando o sentimento de cada palavra
tweets_sentimentos <- tweets_tidy %>%
  filter(!word %in% c('preso', 'fraude', 'acusado', 'imposto', 'acusado')) %>% # excluem-se essas palavras
  inner_join(sentimentos)

# Plotando a distribuicao do sentimento agregado dos tweets
tweets_sentimentos %>%
  group_by(tweet_id) %>%
  summarise(sentimento_tweet = sum(polarity)) %>%
  ggplot(aes(sentimento_tweet)) +
  geom_bar(fill = 'red', colour = 'red', alpha = .1) +
  labs(title = 'Sentimentos Predominantes nos Tweets', x = 'Score dos Tweets')

```

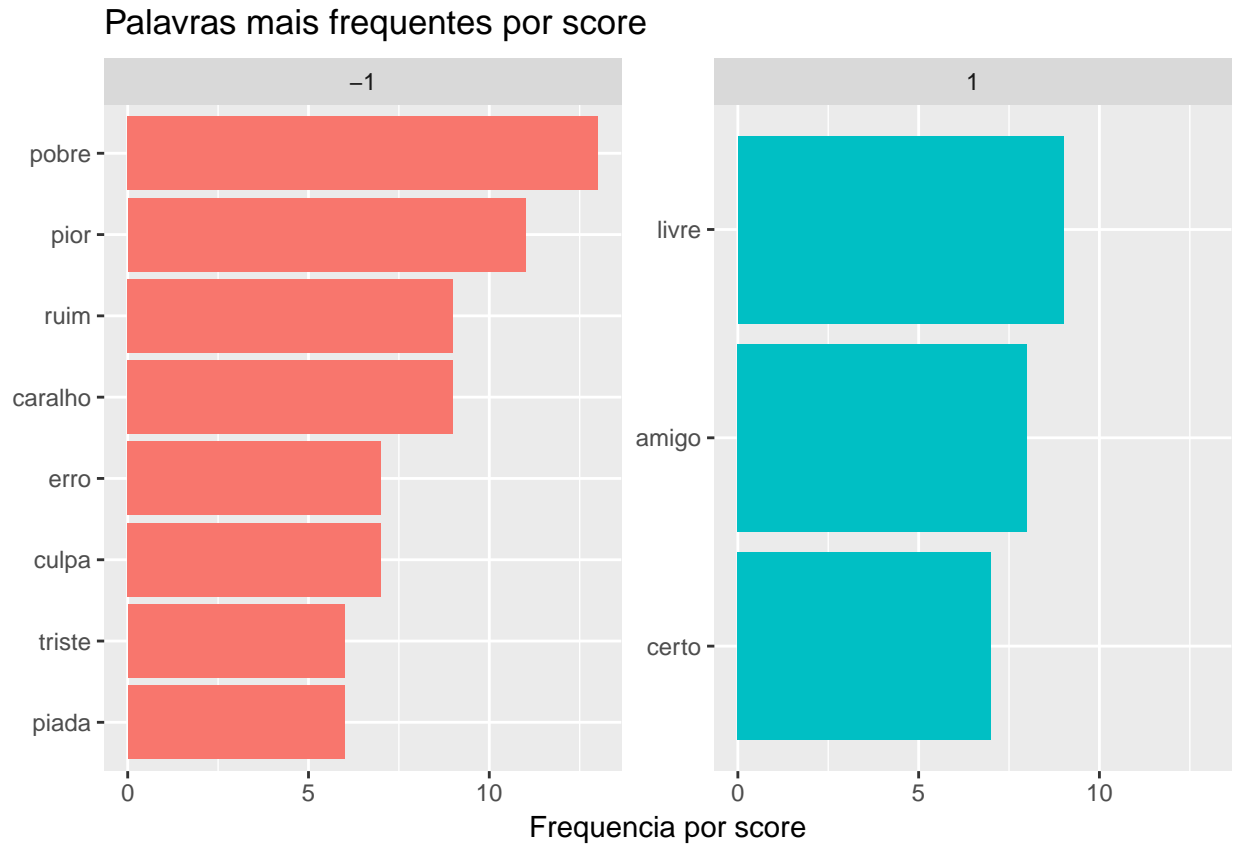


```

# Analisando as palavras mais frequentes por sentimento
tweets_sentimentos %>%
  count(word, polarity, sort = TRUE) %>%
  filter(n > 5, !polarity == 0) %>%
  mutate(word = reorder(word, n)) %>%
  ggplot(aes(word, n, fill = as.factor(polarity))) +
  geom_col() +
  xlab(NULL) +
  coord_flip() +
  facet_wrap(~ polarity, scales = "free_y") +
  labs(title = 'Palavras mais frequentes por score',
       y = "Frequencia por score",
       x = NULL) +

```

```
coord_flip() +  
theme(legend.position = "none")
```



Por fim, estas ultimas analises graficas permitem concluir por uma tendencia maior a aceitacao negativa do publico com a prisao, provavelmente por possiveis temores em relacao ao futuro da marca assim como se pode ver nas palavras mais frequentes por score: um uso muito alto de termos como “pobre”, “triste”, “ruim”, dentre outros, indicando uma predominancia do sentimento de infelicidade e descontentamento com a acao.

Fim

Obrigado!