



SAMET TURGUT

35909310684

Veri Madenciliği Final Projesi

PROJENİN ÖZETİ

15 Aralık 2021 ve 17 Ocak 2022 tarihleri arasında Türkiye bölgesinde “Steam” hakkında atılan 7500 adet tweet toplanmış olup, bu tweet’ler veri temizleme aşamalarından geçirildikten sonra metin madenciliği ve duygu analizi işlemlerine tabi tutulmuştur. Metin madenciliği sayesinde belirtilen tarihler aralığında bulunan tweet’ler içerisinde “Steam” kelimesi ile birlikte en fazla kullanılan diğer kelimeler ortaya çıkmıştır. Duygu analizi aşamasında ise belirtilen bu tarihler aralığında atılan tweet’lerde bulunan pozitif ve negatif kelimeler tespit edilmiştir. Cümlelerde bulunan pozitif ve negatif kelimelerden hangisi fazlalıktaysa tweet’ler bu işleme göre pozitif veya negatif cümle olarak ayrılmıştır. Cümlelerde bulunan negatif ve pozitif kelimeler eşitse veya bunlardan her ikisi de cümlelerde bulunmuyorsa ise bu cümleler de nötr olarak gruplandırılmıştır. Bu sayede bu tarihler aralığında atılan tweet’lerden yola çıkarak 7500 adet tweet içerinden kaç tanesinin “Steam” hakkında pozitif, negatif veya nötr fikre sahip olduğuna ulaşılmıştır. Bu işlemlerin tamamı R programlama diliyle yapılmıştır.

Metin Madenciliği ve Duygu Analizi Aşamasında Kullanılan Paketler

Aşağıda belirtilen paketlerin tamamı bu projede aktif durumda değildir. Bazı paketler atıl durumda kalmıştır ancak diğer popüler opsiyonları göstermek amacıyla projeye eklenmiştir.

library(ROAuth): Twitter Developer hesabının yetkilerine erişebilmek için kullanılır.

library(twitteR): Twitter üzerinden veri çekmek için gerekli olan pakettir.

library(tm): Metin madenciliği sürecinde kullanılan pakettir.

library(plyr): Veri setlerini parçalayıp üzerinde fonksiyon kullanıp tekrar birleştirmeye yarayan pakettir.

library(dplyr): Veri manipülasyonu için gerekli olan pakettir.

library(purrr): R’in fonksiyonel programlama alet takımını geliştiren pakettir.

library(RCurl): Fonksiyonlara HTTP istekleri sağlayan pakettir.

library(SnowballC): Porter’ın farklı dillerde kelime sıkma algoritmasını kullanan pakettir.

library(openssl): OpenSSL'e bağlanma paketidir.

library(wordcloud): Veri görselleştirme aşamasında kelime bulutu oluşturmak için kullanılan pakettir.

library(ggplot2): Veri görselleştirme aşamasında grafikler oluşturmak için kullanılan pakettir.

library(RColorBrewer): Veri görselleştirme aşamasında renkli görseller oluşturmak için kullanılan pakettir.

library(stringr): Hızlı ve doğru bir şekilde string uygulamaları ve string manipülasyonları sağlayan pakettir.

library(stringi): stringr'nin yaptığı işlemleri yapan, içerisinde bazı farklı fonksiyonlar bulunduran ve stringr'nin eski versiyonu olan bir pakettir.

library(tidytext): Metin madenciliği sürecinde verileri düzenlemeyi sağlayan pakettir.

library(readr): "csv", "tsv", "fwf" uzantılarına sahip verileri okumayı sağlayan pakettir.

library(hms): Zaman değişkenleri barındıran verileri şekillendirmek için kullanılan pakettir.

library(lubridate): Saat dilimi, artık gün, yaz saati uygulamaları gibi zaman değişkenleri barındıran verileri şekillendirmek için kullanılan pakettir.

library(igraph): Grafikler yaratmak, şekillendirmek ve analiz etmek için kullanılan pakettir.

library(glue): Yorumlanmış string kalıpları için kullanılan pakettir.

library(networkD3): Ağ grafikleri oluşturmak için kullanılan HTML aygıtı paketidir.

library(rtweet): Twitter'ın REST ve yayın API'na erişmek için kullanılan pakettir.

library(ggeasy): ggplot2 paketinde yapması zor olan bazı grafik işlemlerini gerçekleştirmek için kullanılan pakettir.

library(plotly): Javascript kullanarak interaktif web grafikleri oluşturmak için kullanılan pakettir.

library(magrittr): Bazı fonksiyonların daha rahat okunabilmesi için operatörler sağlayan bir pakettir.

library(tidyverse): Verileri düzenlemek için kullanılan paketleri kapsayan ana pakettir.

library(janeaustenr): Jane Austen'in yayımlanmış 6 kitabının tüm metinlerine UTF-8 formatında erişmek için kullanılan pakettir.

library(widyr): Verileri matrislere dönüştürüp tekrar düzenli forma getiren pakettir.

library(xlsx): Excel dosyalarını okumak için kullanılan pakettir.

library(readxl): Verileri Excel'den R'a aktarmak için kullanılan ve xlsx paketinden daha kolay çalışan bir pakettir.

METİN MADENCİLİĞİ

R kodlarının çalışması için öncelikle gerekli olan paketlerin bilgisayara yüklenmesi gerekiyor. Bunun için `install.packages()` komutunu kullanıyoruz.

```
1 install.packages("ROAuth")
2 install.packages("twitter")
3 install.packages("tm")
4 install.packages("plyr")
5 install.packages("dplyr")
6 install.packages("purrr")
7 install.packages("RCurl")
8 install.packages("SnowballC")
9 install.packages("openssl")
10 install.packages("wordcloud")
11 install.packages("ggplot2")
12 install.packages("RColorBrewer")
13 install.packages("stringr")
14 install.packages("textclean")
15 install.packages("tidytext")
16 install.packages("readr")
17 install.packages("xlsx")
18
19 install.packages("hms")
20 install.packages("lubridate")
21 install.packages("igraph")
22 install.packages("glue")
23 install.packages("networkD3")
24 install.packages("rtweet")
25 install.packages("ggeasy")
26 install.packages("plotly")
27 install.packages("hms")
28 install.packages("lubridate")
29 install.packages("magrittr")
30 install.packages("tidyverse")
31 install.packages("janeaustenr")
32 install.packages("widyr")
```

Gerekli olan paketler bilgisayara yüklendikten sonra ise bu paketlerin kütüphaneden çağırılması gerekiyor. Bunun için ise library() komutunu kullanıyoruz.

```
38 library(ROAuth)
39 library(twitter)
40 library(tm)
41 library(plyr)
42 library(dplyr)
43 library(purrr)
44 library(RCurl)
45 library(snowballc)
46 library(openssl)
47 library(wordcloud)
48 library(ggplot2)
49 library(RColorBrewer)
50 library(stringr)
51 library(tidytext)
52 library(readr)
53 library(stringi)
54
55 library(hms)
56 library(lubridate)
57 library(igraph)
58 library(glue)
59 library(networkD3)
60 library(rtweet)
61 library(ggeasy)
62 library(plotly)
63 library(hms)
64 library(lubridate)
65 library(magrittr)
66 library(tidyverse)
67 library(janeaustenr)
68 library(widyr)
69
70 library(xlsx)
71 library(readxl)
```

Paketleri kütüphaneden çağırdıktan sonra ise bu projenin çalışması için gerekli olan ve Twitter tarafından verilen Twitter Developer hesabına ait olan bilgileri çağırıyoruz.

```
73 options(httr_oauth_cache=T)
74
75 APIkey <- " "
76 APIsecret <- " "
77 access_token <- "{ "
78 access_secret <- " "
79
80 setup_twitter_oauth(APIkey, APIsecret, access_token, access_secret)
```

Tweet'leri Çekme

Bu işlemlerin ardından Twitter'dan tweet çekme işlemine başlayabiliriz. Türkçe dilinde içinde “steam” kelimesinin geçtiği 5000 adet tweet için arama yapıyoruz. Bu işlem sırasında otomatik olarak en fazla 1 hafta geriye gidecek şekilde veri çekiliyor. Anlık olarak çekilebilecek 5000 adet tweet bulunmuyorsa bulunan miktar kadar tweet gelecektir. Çekilen tweetleri data frame'e dönüştürüyoruz ve ardından bu data frame'i ise tweet'leri temizleme işleminde kullanacağımız yeni bir data frame'e dönüştürüyoruz. Gelen tweetleri ise UTF-8 formatına dönüştürüyoruz.

```
82
83 steam <- searchTwitter("steam", n=5000, lang = "tr")
84
85
86 steam.df <- twListToDF(steam)
87
88 length(steam)
89
90 tweet_clean <- steam.df
91
92 tweet_clean$text <- stri_enc_toutf8(tweet_clean$text)
```

O anlığına çekilen tweet miktarı için;

```
> length(steam)
[1] 1248
```

Çekilen tweet'lerin içeriğini görmek için summary(steam.df) komutu ile inceliyoruz. 16 adet sütun ve 1248 adet satırdan oluştuğunu görebiliyoruz. Ayrıca her sütunun başlığı, veri tipleri, minimum, maksimum, ortalama, 1. Çeyrek, 3. Çeyrek, medyan gibi değerleri de görebiliyoruz.

```
> summary(steam.df)
      text      favorited      favoriteCount      replyToSN      created      truncated
Length:1248   Mode :logical   Min. : 0.00   Length:1248   Min. :2022-01-09 11:31:15   Mode :logical
Class :character   FALSE:1248   1st Qu.: 0.00   Class :character   1st Qu.:2022-01-11 22:08:47   FALSE:943
Mode :character      Median : 0.00   Mode :character   Median :2022-01-13 18:23:01   TRUE :305
      Mean : 10.85
      3rd Qu.: 2.00
      Max. :3874.00
      replyToSID      id      replyToUID      statusSource      screenName      retweetCount
Length:1248      Length:1248      Length:1248      Length:1248      Length:1248      Min. : 0.00
Class :character   Class :character   Class :character   Class :character   Class :character   1st Qu.: 0.00
Mode :character     Mode :character     Mode :character     Mode :character     Mode :character     Median : 0.00
      Mean : 17.18
      3rd Qu.: 4.00
      Max. :122.00
      isRetweet      retweeted      longitude      latitude
Mode :logical      Mode :logical      Length:1248      Length:1248
FALSE:832      FALSE:1248      Class :character   Class :character
TRUE :416      Mode :character     Mode :character
```

Bu adımın ardından veri temizleme adımına geçilebilir ancak 1 ay süren bir veri toplama işleminin ardından çekilen tüm tweet'lerin tek bir dosya haline getirilip birleştirilmesinden sonra göstermek tekrarı önlemek açısından daha mantıklı olacaktır. Şimdilik tweet'leri çeker çekmez kaydettiğimizi varsayalım. Tweet'leri nereye kaydetmek istiyorsak dosya konumunu yazdıktan sonra tweet'leri hangi tarihte çektiğimizi belirten bir isimle kaydetmek oldukça faydalı olacaktır. Burada "tweet_clean" yani temizleme işlemi yapılan tweet'leri kaydettim ancak bu adımı bir sonraki aşamaya bıraktığımız için steam data frame'ini kullanarak kaydetmek daha tutarlı olacaktır.

```
154 write.xlsx(tweet_clean, "C:/Users/samet/Desktop/twitter_steam_text_mining/17.01.2022tweets.xlsx")
155
156
```

Tweet'leri Birleştirme

Artık tweet'leri Twitter'dan çekmek yerine bilgisayarımızda depoladığımız, butun-tweetler adını verdiğimiz excel dosyasından açıyoruz. Twitter Developer hesabından yaptığımız bağlantıya da ihtiyaç duymuyoruz. Bu sebeple twListToDF komutu artık burada çalışmayacak. Onun yerine as.data.frame komutunu kullanarak tweet'lerimizi tekrar data frame'e çeviriyoruz. Tweet'ler temizlendikten sonra ayrı bir data frame olacağı için tweet_clean adında bir data frame oluşturuyoruz ve UTF-8 formatına dönüştürüyoruz.

```
78 steam <- read_excel("C:/Users/samet/Desktop/twitter_steam_text_mining/butun-tweetler.xlsx")
79 steam.df <- as.data.frame(steam)
80 tweet_clean <- steam.df
81
82
83
84 tweet_clean$text <- stri_enc_toutf8(tweet_clean$text)
```

Tweet'leri temizleme işlemine retweet eden kullanıcıların başında bulunan “RT” ifadelerini temizlemekle başlıyoruz. Ardından “http” ile başlayan url linklerini, “#” ve “@” ifadelerini temizliyoruz. Daha sonra ise noktalama işaretlerini ve rakamları da tweetlerden arındırıyoruz ve tüm harfleri küçük harfe çeviriyoruz. Devamında ise ASCII formatına uygun olmayan ve alfabetik olmayan karakterleri siliyoruz.

Daha sonra bir corpus oluşturarak Turkish Stopwords ve Durak Kelimeler adındaki dosyalarımızı da projeye ekleyerek veri temizleme işleminin son adımını gerçekleştiriyoruz.

```
114 steamCorpus <- Corpus(VectorSource(tweet_clean$text))
115 inspect(steamCorpus[1:10])
```



```

123 turkish_stopwords <- read_excel("C:/Users/samet/Desktop/twitter_steam_text_mining/Turkish-Stopwords.xlsx");
124 turkish_stopwords
125 steamCorpus <- tm_map(steamCorpus, removewords, turkish_stopwords)
126
127 durakkelimeler <- read_excel("C:/Users/samet/Desktop/twitter_steam_text_mining/durakkelimeler.xlsx")
128 durakkelimeler
129 steamCorpus <- tm_map(steamCorpus, removewords, durakkelimeler)

```

```

> turkish_stopwords
# A tibble: 10,349 x 1
  STOPWORD
  <chr>
1 bir
2 ve
3 bu
4 da
5 de
6 için
7 daha
8 çok
9 gibi
10 o
# ... with 10,339 more rows

```

```

> durakkelimeler
# A tibble: 620 x 1
  a
  <chr>
1 aa
2 aaa
3 acaba
4 ad
5 ah
6 aha
7 ahada
8 aldi
9 alem
10 âlem
# ... with 610 more rows

```

Bu temizleme işlemlerinin ardından temiz verileri artık butun-tweetler-temiz olarak ayrı bir excel dosyasında depolayabiliriz.

```

169 write.xlsx(tweet_clean, "C:/Users/samet/Desktop/twitter_steam_text_mining/butun-tweetler-temiz.xlsx")|
170
171

```

Veri Görselleştirme

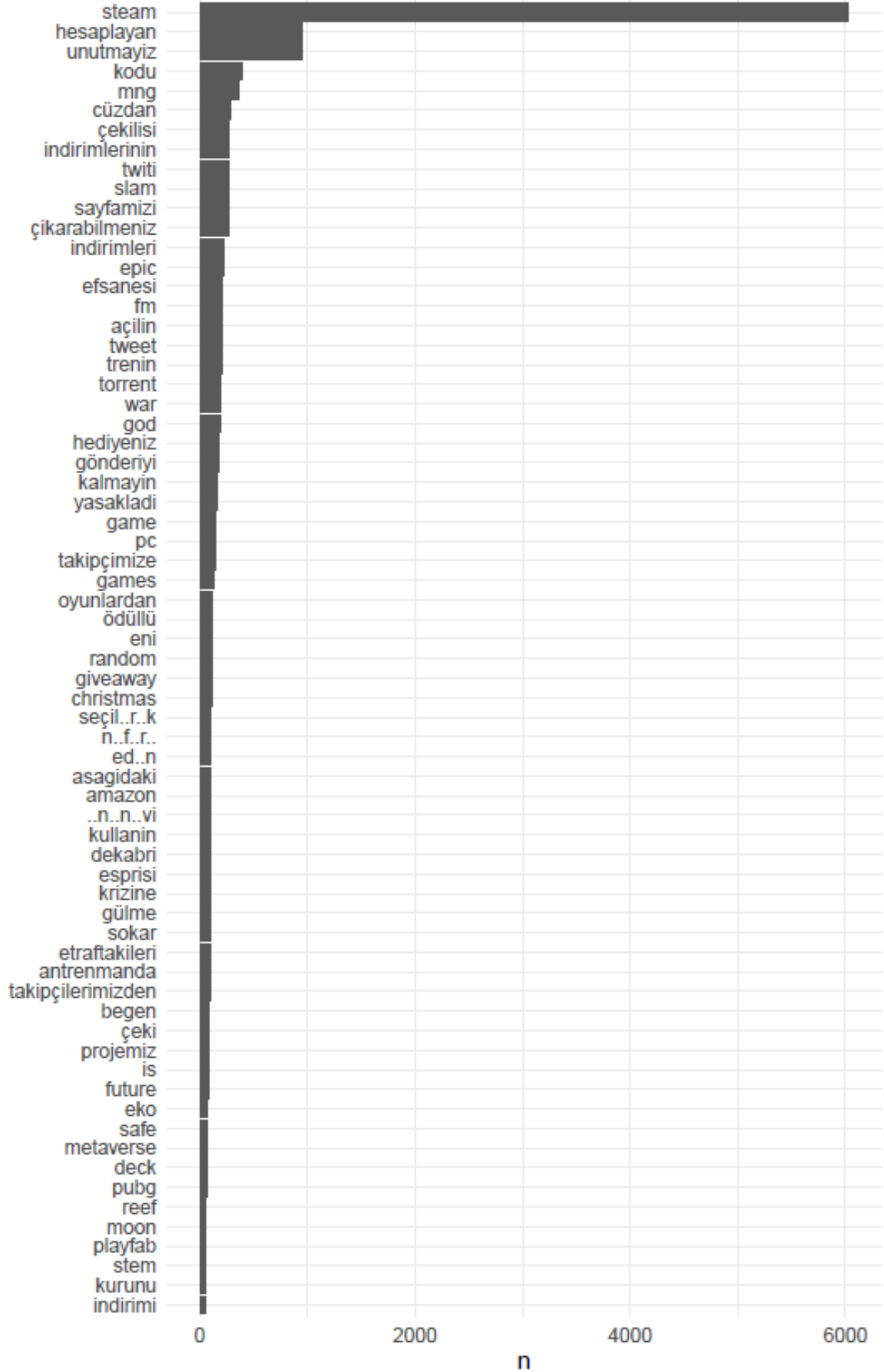
Öncelikle görselleştirilecek olan verinin nereden çekileceği işlemini yapıyoruz. Temizlenmiş olan `tweet_clean`'dan sütun başlığı `text` olan, yani tweetlerin olduğu sütunu `Turkish Stopwords`'e uygun bir şekilde oluşturuyoruz.

```
147  
148 tidy_tweets <- tweet_clean %>% select(text) %>%  
149   mutate(linenumber = row_number()) %>% unnest_tokens(word, text)  
150 tidy_tweets <- tidy_tweets %>% anti_join(turkish_stopwords, by=c("word"="STOPWORD"))  
151
```

Hemen ardından görselleştirme işlemine başlayabiliriz. Tweet'lerde 50 defadan fazla tekrar eden kelimeleri sıralı bir şekilde bar grafiğinde gösteriyoruz.

```
153 tidy_tweets %>%  
154   count(word, sort = TRUE) %>%  
155   filter(n > 50) %>%  
156   mutate(word = reorder(word, n)) %>%  
157   ggplot(aes(word, n)) +  
158   geom_col() +  
159   xlab(NULL) +  
160   coord_flip() + theme_minimal() +  
161   ggtitle("tweetlerde en çok kullanılan kelimeler")
```

tweetlerde en cok kullanılan kelimeler



[illegible]

Duygu Analizi

Duygu analizi yapmak için öncelikle farklı duygu grupları olması gerekir. Bu projede sadece pozitif, negatif ve nötr duyguları bulunmaktadır. Pozitif ve negatif kelimelerin bulunduğu txt dosyalarını projeye ekliyoruz. Özet bilgilerine summary() komutu ile erişebiliyoruz. İçeriklerinden birkaçını görmek için ise dput() komutunu kullanabiliriz.

```
178 pos.words <- scan('C:/Users/samet/Desktop/twitter_steam_text_mining/pozitifkelimeler.txt', what = 'character', comment.char = ';', skipnul = TRUE)
179 dput(pos.words[1289:1293])
180 neg.words <- scan('C:/Users/samet/Desktop/twitter_steam_text_mining/negatifkelimeler.txt', what = 'character', comment.char = ';', skipnul = TRUE)
181 dput(neg.words[1967:1971])
```

```
> summary(pos.words)
  Length      Class      Mode 
   1293  character character
```

```
> dput(pos.words[1289:1293])
c("zenginlik", "zerafet", "zevk", "zevkle", "zevкли")
```

```
> summary(neg.words)
  Length      Class      Mode 
   1981  character character
```

```
> dput(neg.words[1967:1971])
c("zindan", "ziyan", "zor", "zoraki", "zorba")
```

Ardından duygu skoru oluşturmak için bir kod bloğu yazıyoruz. Kendi içinde temizleme işlemleri barındırıyor. Duygu analizi için en önemli kısım olan `match()` komutu ise tweetlerde bulunan kelimeleri pozitif ve negatif kelimeler sözlüğümüzle kıyaslamasıdır. Ancak bu komut bize eşleşen kelimelerin pozitif ve negatif kelimeler sözlüğünde kaçınıcı kelimelere denk geldiğinin sonucunu veriyor ve cümlelere duygu analizi yapabilmek için sözlükten kaçınıcı kelimeye denk geldiğinden ziyade sözlükten bir kelimeye eşleşip eşleşmediğini bulmamız gerekir. Bu sebeple `true/false` içeren bir kod eklemesi yapıyoruz. Ardından `sum()` komutunu kullanarak bir skor formülü oluşturuyoruz. Burada `sum()` komutu `true/false` değerlerini 1 ve 0 olarak değerlendirecek. Yani pozitifle eşleşen ve negatifle eşleşmeyen bir cümle skor = 1-0 sonucundan pozitif, pozitifle eşleşmeyen ve negatifle eşleşen bir cümle skor = 0-1 sonucundan negatif, pozitifle eşleşen ve negatifle eşleşen bir cümle veya her ikisiyle de eşleşmeyen bir cümle de skor = 1-1, skor = 0-0 sonuçlarından nötr olacaktır.

```
183 score.sentiment <- function(sentences, pos.words, neg.words, .progress='none')
184 {
185   require(plyr)
186   require(stringr)
187
188   scores <- laply(sentences, function(sentence, pos.words, neg.words)
189   {
190
191     # clean up sentences with R's regex-driven global substitute, gsub() function:
192     sentence <- gsub('https://', '', sentence)
193     sentence <- gsub('http://', '', sentence)
194     sentence <- gsub('[[:graph:]]', ' ', sentence)
195     sentence <- gsub('[[:punct:]]', '', sentence)
196     sentence <- gsub('[[:cntrl:]]', '', sentence)
197     sentence <- gsub('\\d+', '', sentence)
198     sentence <- str_replace_all(sentence, "[[:graph:]]", " ")
199     # and convert to lower case:
200     sentence <- tolower(sentence)
201
202     # split into words. str_split is in the stringr package
203     word.list <- str_split(sentence, '\\s+')
204     # sometimes a list() is one level of hierarchy too much
205     words <- unlist(word.list)
206
207     # compare our words to the dictionaries of positive & negative terms
208     pos.matches <- match(words, pos.words)
209     neg.matches <- match(words, neg.words)
210
211     # match() returns the position of the matched term or NA
212     # we just want a TRUE/FALSE:
213     pos.matches <- !is.na(pos.matches)
214     neg.matches <- !is.na(neg.matches)
215
216     # TRUE/FALSE will be treated as 1/0 by sum():
217     score <- sum(pos.matches) - sum(neg.matches)
218
219     return(score)
220   }, pos.words, neg.words, .progress=.progress )
221
222   scores.df <- data.frame(score=scores, text=sentences)
223   return(scores.df)
224 }
```

Bu kod bloğunun ardından analizi hazırlamak ve duygu skor frekans tablosunu oluşturmak geliyor.

```
228  
229 analysis <- score.sentiment(tweet_clean$text, pos.words, neg.words)  
230 # sentiment score frequency table  
231 table(analysis$score)|  
232
```

Yukarıdaki skor formülüne göre bir tablo oluşturursak en düşük skor, yani en negatif yorumların skorları -4 çıkmış ve bu skora sahip 2 adet yorum bulunuyor. En yüksek, yani en pozitif yorumların skorları ise 4 çıkmış ve bu skora sahip 6 adet yorum bulunuyor. 4145 adet ise nötr yorum bulunmaktadır.

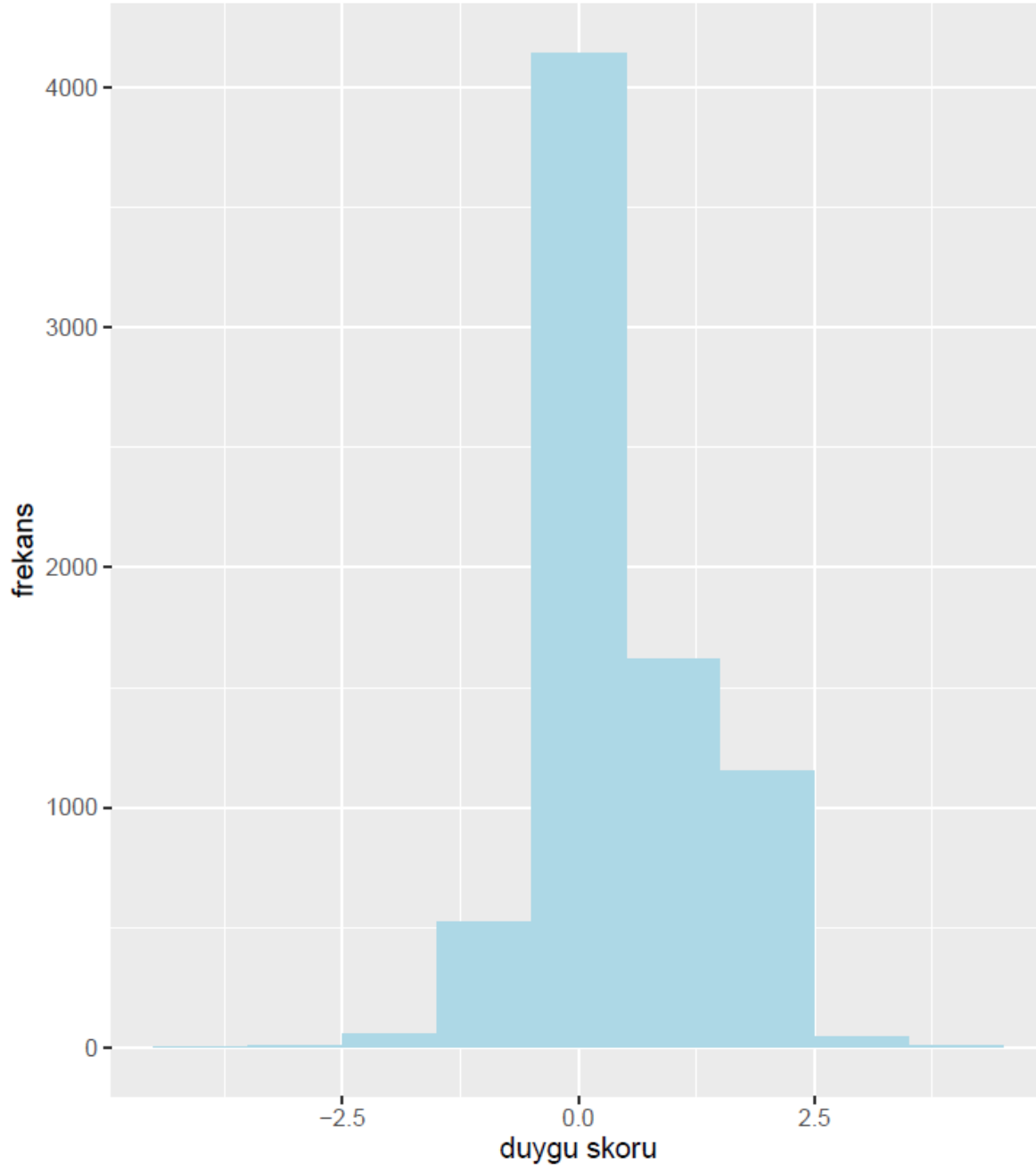
```
> table(analysis$score)
```

-4	-3	-2	-1	0	1	2	3	4
2	7	57	522	4145	1617	1152	44	6

Bu skor dağılımlarını daha iyi görmek için bir histogram grafiğinde gösterelim. Histogram grafiğini oluştururken ggplot paketini kullanıyoruz.

```
233  
234 analysis %>%  
235   ggplot(aes(x=score)) +  
236   geom_histogram(binwidth = 1, fill = "lightblue")+  
237   ylab("frekans") +  
238   xlab("duygu skoru") +  
239   ggtitle("Tweetlerin duygu skorlari dagilimi histogrami") +  
240   ggeasy::easy_center_title()  
241
```

Tweetlerin duygu skorlari dagilimi histogrami



Son olarak bu skorların 0'ın üzeri pozitif, 0'ın altı negatif ve 0'a eşit olanlarını nötr olarak gruplandırmak kalıyor. Bu aşamada artık pozitif ve negatif değerlerde 0'ın ne kadar üzerinde veya ne kadar altında olduğunun bir önemi kalmıyor. Son grafiği de ggplot paketiyle oluşturuyoruz. Artık 7500 adet tweet'in duygu tipi bar grafiği de oluşmuş oluyor.

```
242
243 neutral <- length(which(analysis$score == 0))
244 positive <- length(which(analysis$score > 0))
245 negative <- length(which(analysis$score < 0))
246 sentiment <- c("pozitif","notr","negatif")
247 Count <- c(positive,neutral,negative)
248 output <- data.frame(Sentiment,Count)
249 output$Sentiment<-factor(output$Sentiment,levels=Sentiment)
250 ggplot(output, aes(x=Sentiment,y=Count))+
251   geom_bar(stat = "identity", aes(fill = Sentiment))+
252   ggtitle("7500 tweetin duygu tipi grafiği")|
253
```

Duygu analizinin sonucunda 7500 adet tweet içerisinde büyük çoğunluğu nötr tweetlerin oluşturduğunu görmekteyiz. Yaklaşık 4250 adet nötr tweet bulunuyor. Nötr tweet'lerin ardından ise pozitif tweet'ler geliyor. Yaklaşık 2750 adet pozitif tweet bulunuyor. Son sırada ise negatif tweet'ler bulunuyor. Yaklaşık 600 adet negatif tweet bulunmaktadır. Nötr tweet'ler bütün tweet'lerin yaklaşık %56'sını, pozitif tweet'ler bütün tweet'lerin yaklaşık %36'sını ve negatif tweet'ler ise bütün tweet'lerin yaklaşık %8'ini oluşturmaktadır.

