

TEXT MINING sobre datos recolectados de Twitter.

ZEVALLOS, JOSE^{1,*}, FAZIO, BORIS^{1,**}, AND HUERTAS, ANTHONY^{1,***}

¹Escuela de Posgrado, Pontificia Universidad Católica del Perú, Lima, Perú

* Cod: 20173861

** Cod 20173896

** Cod: 20173728

Compiled May 17, 2018

Minería de datos

1. INTRODUCCIÓN

Text Mining es un mecanismo por el cual se transforman datos cualitativos e inestructurados en datos de buen uso por una computadora. Estos datos inestructurados no pueden ser medidos numéricamente; sin embargo, presentan características de color, textura y descripción textual, esta última será en la cual se verá enfocado el presente proyecto. En conclusión, text mining nos permite generar nueva información, provenientes de una amplia colección de datos, con el objetivo puedan ser examinados. Existen diversos análisis que se pueden realizar en Text mining; sin embargo, en este proyecto se harán uso del **análisis de sentimientos** y nube de palabras, también conocido como **Workcloud**.

El análisis Workcloud es una representación gráfica de palabras frecuentemente usadas en una colección de datos. El tamaño de cada palabra en el gráfico ejecutado por Workcloud es un indicador de la frecuencia de dicha palabra en el total de nuestros datos. Esto quiere decir, que mientras mayor sea la frecuencia en la que se presenta una palabra, se visualizará de mayor tamaño.

El análisis de sentimientos es una medida rescatada de los términos positivos y negativos encontrados en una descripción textual, y mediante ella permitir determinar si en general las expresiones son negativas, positivas o neutrales respecto a un tema en particular. Este análisis toma una relevancia con el criterio de comparación pues permite contrastar el entorno de sentimientos alrededor de una situación presente.

En el presente trabajo se recolectó información de la app **Twitter** utilizando el interfaz de programación provisto por la empresa, diseñando un bot mediante el lenguaje de programación Python con inclusión de la librería **tweepy** necesaria para trabajar con este tipo de datos. Esta librería de Python provee de las funcionalidades RESTfull de Twitter. Los paquetes adicionales de uso necesario en el diseño del bot, y análisis de ser conveniente, serán las siguientes:

Listing 1. Librerías de Python.

```
1 import sys,tweepy,re
2 import matplotlib.pyplot as plt
3 import dataset
4 import json
5 from textblob import TextBlob
6 from sqlalchemy.exc import ProgrammingError
7 from tweepy.streaming import StreamListener
```

El diseño de programación del bot posee dos condiciones: tener una cuenta de Twitter y una aplicación registrada en <https://apps.twitter.com>. Al hacerlo se te provee de 4 claves de acceso (*Consumer Key*, *Consumer Secret*, *Access Token* y *Access Level*), las cuales tendrán que ser almacenadas en un script alterno con el objetivo se mantenga la privacidad del investigador.

Listing 2. Función API con autentificación.

```
1 # Codigos de acceso de la App de Twitter del Usuario
2
3 # Consume:
4 CONSUMER_KEY = 'XXXXXXXXX'
5 CONSUMER_SECRET = 'XXXXXXXXX'
6
7 # Access:
8 ACCESS_TOKEN = 'XXXXXXXXX'
9 ACCESS_SECRET = 'XXXXXXXXX'
10
11 # ...
12 CONNECTION_STRING = "sqlite:///tweets.db"
13 TABLE_NAME = "tabla_nombre"
```

2. ANÁLISIS DE DATOS CON FRASES REFERENTES AL PARTIDO “BAYERN MUNICH VS REAL MADRID” DEL 25/04/18.

Se recolectaron tweets, los cuales fueron almacenados en una base de datos SQLite. Adicionalmente, por cada tweet se guardó su idioma, el id del usuario que lo realizó, y cuantos retweets se dieron. El bote ha sido diseñado mediante el lenguaje de programación Python, y el procesamiento mediante el lenguaje de programación R.

A. Diseño del bot de recolección de datos

Se programa el bot con un mecanismo de recolección de tweets, mediante cierta especificación de palabras clave (keywords), en tiempo real 5 minutos previos, durante y 5 minutos después del partido por semifinales de la Champions League 2018 entre los equipos Bayern Munich y Real Madrid.

Para ello se diseña una primera clase de clasificación y procesamiento de tweets en tiempo real

Listing 3. Bot: Clase TwitterStreamer().

```

1 from credentials import *
2
3 db = dataset.connect(CONNECTION_STRING)
4
5 class TwitterStreamer():
6     def __init__(self):
7         pass
8
9     def stream_tweets(self, fetched_tweets_filename, hash_tag_list):
10
11         # API con autentificacion y codigos de acceso
12         listener = StdOutListener(fetched_tweets_filename)
13         auth = OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
14         auth.set_access_token(ACCESS_TOKEN, ACCESS_TOKEN_SECRET)
15         stream = Stream(auth, listener)
16
17         # Filtro de captura de datos por palabras clave
18         stream.filter(track=hash_tag_list)
```

Se diseña una segunda clase, la cual permite ejecutar los tweets que se desean evaluar y ciertas especificaciones como las coordenadas, el texto, el ID, el numero de retweets y el lenguaje.

Listing 4. Bot: Clase StdOutListener().

```

1 class StdOutListener(StreamListener):
2     def __init__(self, fetched_tweets_filename):
3         self.fetched_tweets_filename = fetched_tweets_filename
4
5     def on_data(self, data):
6
7         table = db[twitter_credentials.TABLE_NAME]
8         try:
9             #print data
10            all_data = json.loads(data)
11            text = all_data['text']
12            ID = all_data["id_str"]
13            coords=all_data['coordinates']
14            retweets=all_data["retweet_count"]
15            lenguaje=all_data['lang']
16
17            table.insert(dict(
18                coordinates=coords,
19                text=text,
20                id_str=ID,
21                retweet_count=retweets,
22                lengua=lenguaje,
23            ))
24
25            with open(self.fetched_tweets_filename, 'a') as tf:
26                tf.write(data)
27            return True
28        except BaseException as e:
29            print("Error on_data %s" % str(e))
30            return True
31
32    def on_error(self, status_code):
33        if status_code == 420:
34            #returning False in on_data disconnects the stream
35            return False
36
37    def on_status(self, status):
38        all_data = json.loads(status)
39        tweet = all_data['text']
```

En esta última parte del diseño se introducen los keywords que deberán incluir los tweets, los cuales son : “Bayern Munich”, “Real Madrid”, “Cristiano Ronaldo” y “Robben”. Posteriormente el nombre del archivo de almacenamiento *tweets5.txt* que por la clase `StdOutListener` será obtenido en formato *.db*.

Listing 5. Bot: Keywords y almacenamiento.

```

1 if __name__ == '__main__':
2
3     # Authenticate using config.py and connect to Twitter Streaming API.
4     hash_tag_list = ['Bayern Munich', 'Real Madrid', 'Cristiano Ronaldo', 'Robben']
5     fetched_tweets_filename = "tweets5.txt"
6
7     twitter_streamer = TwitterStreamer()
8     twitter_streamer.stream_tweets(fetched_tweets_filename, hash_tag_list)

```

B. Convertir base de datos en formato .db (SQL) a .csv

En este proceso se desarrolla un código en R que permite convertir la base de datos en formato SQL obtenida por el bot, a un formato .csv con el objetivo puede facilitarse el análisis mediante funciones más generales. Además se adiciona un proceso de limpieza de datos en la cual se retiran emojis y tildes de las letras. Los datos son almacenados en un archivo llamado *tweets.csv*

Listing 6. Conversión de data SQL a .csv.

```

1 library(RSQLite)
2
3 # Cargar BD SQL
4 con <- dbConnect(SQLite(), dbname="tweets5.db")
5 # Ver tablas en la BD
6 sq <- dbSendQuery(con, "SELECT name FROM sqlite_master WHERE type='table'")
7 tn <- dbFetch(sq, n = -1)
8 # Ver filas en la tabla
9 sq <- dbSendQuery(con, paste0("SELECT * FROM ", tn[1,]))
10 d <- dbFetch(sq, n = -1)
11 dbDisconnect(con)
12
13 # Limpieza + Exportacion
14 d$text %<>%
15     iconv("UTF-8", "ASCII//TRANSLIT") %>%
16     ifelse(is.na(.), iconv(d$text, "latin1", "UTF-8"), .) %>%
17     stri_replace_all(regex="^[[:graph:]]", replacement=" ") %>%
18     stri_replace_all(regex="(f|ht)tp\\S+\\s*", replacement=" ") %>%
19     stri_replace_all(regex="[.,?:|]", replacement=" ") %>%
20     textclean::replace_emoji(emoji_dt = data.table(lexicon::hash_emojis[,1], y=" "))
21 write.csv(d, "tweets.csv", row.names = FALSE)

```

C. Preprocesamiento de la data

Se diseña un mecanismo de preprocesamiento de los datos previo al análisis de conglomerados y reglas de asociación que serán de interés posterior. Los resultados de este proceso desarrollado en R, los cuales serán almacenados en *02_output.Rdata*, son los siguientes:

- **resumen:** Número de tweets en total (tt), tweets que son un retweet (rt), tweets con una llamada a otro @usuario (at), tweets con al menos un hashtag (ht), y tweets con llamada @ y hashtag (ht.at), para cada idioma con >2k tweets totales.
- **rt.df:** Tabla con el número de veces que se repite cada retweet para un idioma dado, servirá para mostrar los tweets mas populares.
- **at.df, ht.df:** Tabla que muestra @usuarios y #hashtags mas populares, sin contar los retweets.
- **d.stop:** Cinco corpus (uno por idioma) sin stopwords, con la presencia de cada palabra en un tweet representada en formato matriz. Servirá para el análisis de conglomerados (tweets similares) o reglas de asociación (frases comunes).
- **d.stem:** Mecanismo similar a **d.stop** pero con reducción de las palabras a su raíz.

Listing 7. Preprocesamiento.

```

1 library(magrittr)
2 library(stringi)
3 library(textclean)
4 library(data.table)
5 library(SnowballC)
6 library(tm)
7 options(stringsAsFactors = FALSE)
8
9 # Cargar datos + Preprocesamiento general (Solo lenguas > 2k tweets)

```

```

10 d <- read.csv("tweets.csv", stringsAsFactors = FALSE)[,c(1:3,5)]
11 d <- d[d$lengua %>% is.element(c("en","es","pt","fr","de")),]
12
13 # Retirar blanco en los extremos de cada tweet + regex ( @ + # + RT .*)
14 d$text %<>% trimws
15 re.at <- "@[[:alnum:]]_"
16 re.ht <- "#[[:alpha:]]_[[:alnum:]]_"
17 re.rt <- "RT .*"
18
19 # Resúmenes
20 # Taggeo de tweets [RT|@|#]
21 # RT
22 d$text %>% stri_match_first(regex=re.rt) %>% is.na %>% not -> d$rt
23 # @
24 d$text[!d$rt] %>% stri_match_first(regex=re.at) %>% is.na %>% not -> d$at[!d$rt]
25 # #
26 d$text[!d$rt] %>% stri_match_first(regex=re.ht) %>% is.na %>% not -> d$ht[!d$rt]
27
28 # Conteo de RTs
29 d$text[d$rt] %>% table -> rt.tab
30 rt.df <- data.frame(text=names(rt.tab), n.rt=as.vector(rt.tab))
31 d.rt <- merge(d, rt.df, by = "text")[,c("text", "lengua", "n.rt")]
32 rt.df <- d.rt[!duplicated(d.rt),]
33 rm(d.rt)
34
35 # Hashtags por idioma
36 d$lengua %>% unique %>% lapply(
37   function(x){
38     d[d$lengua==x,"text"] %>% stri_extract_all(regex=re.ht) %>% unlist %>% table %>%
39       sort(decreasing = TRUE) -> ht.tab
40     data.frame(ht=names(ht.tab), n=as.vector(ht.tab), lengua=x)
41   }) %>% do.call(rbind,.) -> ht.df
42
43 # Usuarios por idioma
44 d$lengua %>% unique %>% lapply(
45   function(x){
46     d[d$lengua==x,"text"] %>% stri_extract_all(regex=re.at) %>% unlist %>% table %>%
47       sort(decreasing = TRUE) -> at.tab
48     data.frame(at=names(at.tab), n=as.vector(at.tab), lengua=x)
49   }) %>% do.call(rbind,.) -> at.df
50
51 # Limpiar palabras
52 d.tm <- d[!d$rt,c("lengua","text")]
53 d.tm$text %<>%
54   stri_replace_all(regex=re.at, replacement = " ") %>%
55   stri_replace_all(regex=re.ht, replacement = " ") %>%
56   stri_replace_all(regex="<[:alnum:][:alnum:]>", replacement = " ") %>%
57   stri_replace_all(regex="[:punct:]", replacement = " ") %>%
58   trimws %>% tolower
59 # Corpus por lenguaje
60 unique(d.tm$lengua) %>%
61   lapply(function(x) Corpus(VectorSource(d.tm[d.tm$lengua==x,"text"]))) -> d.corp
62 # Retirar stopwords
63 d.corp %>%
64   lapply(function(x) tm_map(x, removeWords, unlist(lapply(unique(d.tm$lengua),
65     stopwords)))) -> d.stop
66 # Stems
67 d.stop %>%
68   lapply(function(x) tm_map(x, stemDocument, language = "en")) %>%
69   lapply(function(x) tm_map(x, stemDocument, language = "es")) %>%
70   lapply(function(x) tm_map(x, stemDocument, language = "de")) %>%
71   lapply(function(x) tm_map(x, stemDocument, language = "fr")) %>%
72   lapply(function(x) tm_map(x, stemDocument, language = "pt")) -> d.stem
73
74 # Formar matrices
75 d.stop %<>% lapply(DocumentTermMatrix)
76 d.stem %<>% lapply(DocumentTermMatrix)
77
78 # Exportar
79 data.table(d)[, .(tt = .N, rt = sum(rt, na.rm = TRUE), at = sum(at, na.rm = TRUE),
80   ht = sum(ht, na.rm = TRUE),
81   ht.at = sum(ht|at, na.rm = TRUE))
82   , .(lengua)] -> resumen
83 resumen[, ht.at := ht + at - ht.at]
84 save(resumen, rt.df, at.df, ht.df, d.stem, d.stop, file = "02_output.Rdata")

```

| text | lengua | n.rt |
|---|--------|-------|
| <chr> | <chr> | <int> |
| 1 RT @betsattt1 Real Madrid - Bayern Muni | de | 23 |
| 2 RT @ChampionsLeague GOAL! Bayern 1-0 Real Madrid (Joshua Kimmich 28) #U~ | de | 93 |
| 3 RT @nomeserrados bayern demoniaco | de | 134 |
| 4 RT @ChampionsLeague GOAL! Bayern 1-2 Real Madrid (Marco Asensio 57) #UCL | en | 108 |
| 5 RT @ChampionsLeague Cristiano Ronaldo <c3><b0> All eyes on this season'~ | en | 130 |
| 6 RT @ChampionsLeague Marco Asensio in the knockouts for Real Madrid 2016~ | en | 172 |
| 7 RT @ZB_Media02 Tercer penalti no senalado en contra del Real Madrid min~ | es | 111 |
| 8 RT @2010MisterChip BAY 1-2 RMA (FT) - El Bayern ha perdido SEIS PARTIDO~ | es | 135 |
| 9 RT @2010MisterChip El PSG no habia perdido en toda la temporada en casa~ | es | 452 |
| 10 RT @ActuFoot_ L'Allianz Arena plus comble que jamais ! Quel stade ! | fr | 58 |
| 11 "RT @YassEncore Real vs Bayern Ribery vs Benzema les deux \" racailles ~ | fr | 64 |
| 12 RT @AmineMaTue Le Real Madrid c'est incroyable tout les matchs c'est pa~ | fr | 110 |
| 13 RT @B24pt Enquanto a maioria v<c3><83><c2><aa> o Bayern - Real Madrid D~ | pt | 70 |
| 14 RT @InfosFuteboI Marco Asensio a j<c3><83><c2><b3>ia do Real Madrid! <c~ | pt | 123 |
| 15 RT @InfosFuteboI Quando toca o hino da Champions League | pt | 157 |

(a) Retweets.

| at | n | lengua |
|---------------------|-------|--------|
| <chr> | <int> | <chr> |
| 1 @FCBayern | 59 | de |
| 2 @nomeserrados | 134 | de |
| 3 @ChampionsLeague | 141 | de |
| 4 @brfootball | 187 | en |
| 5 @FOXsoccer | 221 | en |
| 6 @ChampionsLeague | 1215 | en |
| 7 @realmadrid | 241 | es |
| 8 @Invictosomos | 365 | es |
| 9 @2010MisterChip | 1300 | es |
| 10 @YassEncore | 74 | fr |
| 11 @AmineMaTue | 110 | fr |
| 12 @ActuFoot_ | 242 | fr |
| 13 @FoxSportsBrasil | 113 | pt |
| 14 @RMadridDepre | 158 | pt |
| 15 @InfosFuteboI | 504 | pt |

(b) Usuarios.

| ht | n | lengua |
|-------------------------|-------|--------|
| <chr> | <int> | <chr> |
| 1 #Bayern | 76 | de |
| 2 #UCL | 236 | de |
| 3 #FCBRMA | 288 | de |
| 4 #BAYRMA | 412 | en |
| 5 #FCBRMA | 484 | en |
| 6 #UCL | 1505 | en |
| 7 #HalaMadrid | 229 | es |
| 8 #ChampionsLeague | 302 | es |
| 9 #UCL | 305 | es |
| 10 #UEFACHampionsLeague | 55 | fr |
| 11 #FCBRMA | 141 | fr |
| 12 #BAYRMA | 179 | fr |
| 13 #ChampionsLeague | 45 | pt |
| 14 #UCL | 78 | pt |
| 15 #HalaMadrid | 104 | pt |

(c) Hashtags.

| | | | | |
|------------|---------------|-------------|---------------|-----------------|
| \$Portuges | bayern | real | madrid | nao cristiano |
| | 2645 | 2066 | 1661 | 735 504 |
| \$Frances | bayern | real | madrid | match champions |
| | 812 | 270 | 108 | 84 67 |
| \$Ingles | bayern madrid | real munich | league | |
| | 4089 3308 | 3047 1221 | 1036 | |
| \$Espanol | madrid bayern | real munich | gol | |
| | 3175 3019 | 2740 654 | 402 | |
| \$Aleman | bayern | real madrid | munich munden | |
| | 1452 448 | 408 210 | 127 | |

(d)

| | | | | |
|------------|-------------------------|-------------------|--------|--------|
| \$Portuges | nao cristiano | ronaldo | jogo | gol |
| | 735 504 | 495 | 485 | 470 |
| \$Frances | match champions | plus | munich | league |
| | 84 67 | 60 | 56 | 44 |
| \$Ingles | munich league champions | live | watch | |
| | 1221 1036 1023 | 1018 | 489 | |
| \$Espanol | munich gol | partido champions | vamos | |
| | 654 402 | 377 362 | 251 | |
| \$Aleman | munich munden | live fur kimmich | | |
| | 210 127 | 126 76 66 | | |

(e)

Fig. S1. Resumen de datos

D. Análisis descriptivo y Workcloud

En esta subsección, se desarrollará un análisis descriptivo y un análisis de nube de palabras (Workcloud) de las palabras con mayor relevancia o popularidad respecto a los 5 idiomas referentes.

Listing 8. Inicialización.

```

1 library(ggplot2)
2 library(reshape2)
3 library(magrittr)
4 library(dplyr)
5 library(tm)
6 library(wordcloud)
7
8 # Cargar datos
9 load("@2_output.Rdata")
10 # Prepro
11 resumen$tt <- with(resumen, tt - rt - at - ht - ht.at)

```

Los análisis descriptivos serán evaluados sobre la frecuencia y proporciones de los números de tweets respecto a cada idioma. El código y resultados son los siguientes.

Listing 9. Análisis Descriptivo.

```
1 # Tipos de tweet por idioma
2 data.frame(lengua = resumen[, 1], t(apply(resumen[, -1], 1, function(x) x/sum(x)))) %>%
3   melt(id.vars = "lengua") %>% cbind(estadistico = "Proporcion") -> d1a
4 d1a$value <- d1a$value*(resumen[, -1] %>% apply(1, sum) %>% max)
5 resumen %>% melt(id.vars = "lengua") %>% cbind(estadistico = "Frecuencia") -> d1b
6
7 d1 <- rbind(d1a, d1b)
8
9 d1$lengua %<>% factor(levels = c("en", "es", "pt", "fr", "de"),
10   labels = c("Ingles", "Espanol", "Portugues", "Frances", "Aleman"))
11 d1$variable %<>% factor(levels = c("ht", "ht.at", "at", "tt", "rt"),
12   labels = c("Hashtag", "Hashtag + Usuario", "Usuario", "Tweet solo",
13     "Retweet"))
14 d1$estadistico %<>% factor(levels = c("Frecuencia", "Proporcion"))
15
16 # Plot
17 ggplot(d1, aes(x = lengua, y = value, fill = variable)) +
18   geom_bar(stat = "identity") + xlab("Idioma") + ylab("Numero de tweets") +
19   guides(fill = guide_legend(title = "Contenido de tweet")) +
20   facet_grid(. ~ estadistico, scales = "free")
```

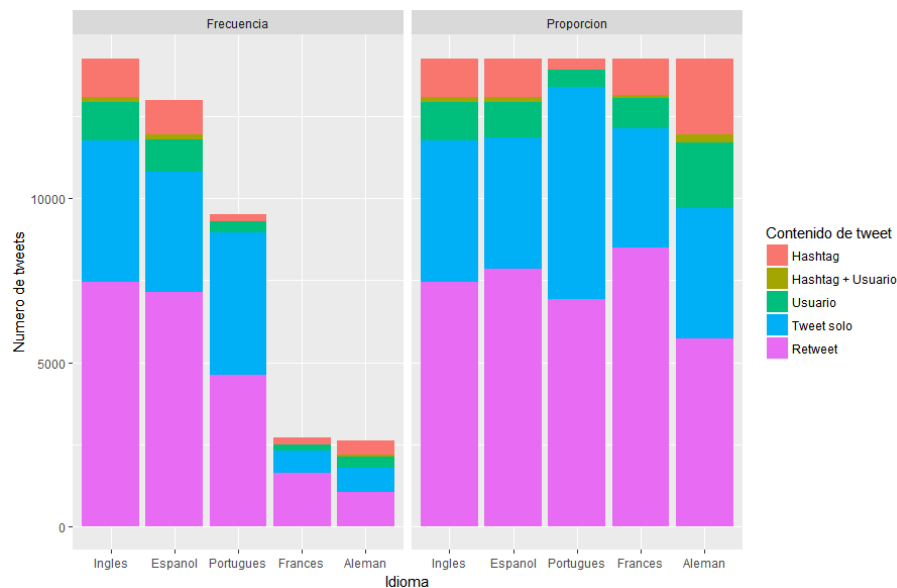


Fig. S2. Distribución de frecuencia y proporciones de los tweets por idioma.

Se diseña el código para la representación gráfica mediante WorkCloud, como se detalla a continuación

Listing 10. Análisis Workcloud.

```
1 # Palabras mas populares – sin stopwords
2 # Convertir a data.frame
3 d.stop %>% lapply(removeSparseTerms, 0.995) %>%
4   lapply(function(x) as.data.frame(as.matrix(x))) -> d.stop.s
5 d.stem %>% lapply(removeSparseTerms, 0.995) %>%
6   lapply(function(x) as.data.frame(as.matrix(x))) -> d.stem.s
7
8 # Palabras mas populares, Tabla simple
9 d.stop.s %>%
10   lapply(function(x) apply(x, 2, sum) %>%
11     as.table %>%
12     sort(decreasing = TRUE) %>%
13     head(5)) -> d.stop.top5
14 names(d.stop.top5) <- c("Portugues", "Frances", "Ingles", "Espanol", "Aleman")
15 d.stop.top5
16
17 d.stop.s %>%
18   lapply(function(x) apply(x, 2, sum) %>%
```

```

19     as.table %>%
20     sort %>%
21     head(., length(.) - 3) %>%
22     sort(decreasing = TRUE) %>%
23     head(5)) -> d.stop.top5b
24 names(d.stop.top5b) <- c("Portuges", "Frances", "Ingles", "Espanol", "Aleman")
25
26 # Wordcloud
27 d.stop.s %>%
28     lapply(function(x) apply(x, 2, sum) %>%
29         as.table %>%
30         sort %>%
31         head(., length(.) - 3)) -> d.stop.ret3
32
33 for(i in 1:length(d.stop.ret3)){
34     png(filename = paste0("descriptivo_04_wordcloud",i,".png"))
35     wordcloud(names(d.stop.ret3[[i]]),as.vector(d.stop.ret3[[i]]), scale = c(4,.5))
36     dev.off()
37 }

```



Fig. S3. Palabras más populares por cada idioma

3. ANÁLISIS DE DATOS CON TERMINO “PAOLO GUERRERO”

Los datos han sido recolectados de tweets con la frase **Paolo Guerrero** posterior al fallo del TAS del 15/05/18, organismo, que sancionó al futbolista peruano con 14 meses de suspensión. Debido a ser una información de alta relevancia en el mundo del fútbol de nuestro país y en todo el mundo se decidió a realizar un análisis de sentimientos de 1000 tweets tomados en tiempo real.

Tanto el diseño del bot que permitirá la recolección como el análisis posterior, han sido diseñados en el lenguaje de programación Python

A. Diseño del bot de recolección de datos

El mecanismo que opta el bot es el siguiente

Listing 11. Bot: Función API con autenticación.

```
1 from credentials import * Nos permite usar los codigos de acceso como variables
2
3 # Autenticacion y codigos de acceso:
4 auth = tweepy.OAuthHandler(CONSUMER_KEY, CONSUMER_SECRET)
5 auth.set_access_token(ACCESS_TOKEN, ACCESS_SECRET)
6
7 # Retorna API con autenticacion:
8 api = tweepy.API(auth)
```

Posteriormente, se recolectará los tweets que contiene el término “Paolo Guerrero”, sea hashtag o simplemente una frase de tal tipo. Para ello se diseña un mecanismo de entrada, para ingresar datos respectivos, tales como el término como la cantidad de tweets a recolectar. En particular se recolectará 1000 tweets.

Listing 12. Bot: Ingreso del término en común en los tweets y cantidad a recolección.

```
1 Termino_buscado = input("Ingrese la palabra o hashtag a buscar:")
2 numero_de_Tweets=int(input("Ingrese la cantidad de tweets que desea analizar:"))
3
4 tweets = []
5 tweets = tweepy.Cursor(api.search,q=Termino_buscado).items(numero_de_Tweets)
```

```
Ingrese la palabra o hashtag a buscar:Paolo Guerrero
Ingrese la cantidad de tweets que desea analizar:1000
```

Fig. S4. Termino “Paolo Guerrero” en 1000 tweets recolectados.

```
RT @ExpositoresTop: He añadido un vídeo a una lista de reproducción de @YouTube (https://t.co/ZPjBJFVIVB - PAOLO GUERRERO SI PUEDE JUGAR, L...
RT @americatv_peru: Primero habló Chilavert y ahora este crack https://t.co/dW82X4jR43
RT @ExpositoresTop: PAOLO GUERRERO SI PUEDE JUGAR, LA FIFA TIENE LA FACULTAD DE OTORGARLE AMNISTIA, #FifaAmnistiaAPaolo

¿Hay... https://t...
RT @o_somocurcio: Aferrarnos al FIFPRO o a la campaña FIFA AMNISTIA a Paolo Guerrero no es despertar una ilusión en el hinch a??? La FIFA no...
RT @SC_ESPN: La Federación Internacional de Futbolistas Profesionales (FIFPro) solicitó una reunión urgente con la FIFA por la sanción de 1...
RT @bieloficial: “Às pessoas que contribuíram para esta vergonhosa injustiça digo que estão me roubando o Mundial e, talvez, minha carreira...”
```

Fig. S5. Tweets Recolectados. Visualización: 6 tweets.

Es claro que debido al diseño del **Listing 4**, el estudio puede hacerse de forma general sobre cualquier término y cualquier cantidad a optar de tweets.

B. Diseño de funciones

A continuación se presentan funciones que permiten una limpieza de nuestros datos, es decir, la eliminación de caracteres innecesarios que obstaculizan el estudio; y además de una función porcentaje que permite proporcionar nuestros datos respecto al tipo de sentimientos (positivos, neutros, negativos) generados por los tweets.

Listing 13. Funciones.

```
1 def Eliminacion_caracteres(tweet):
2     return ' '.join(re.sub("([A-Za-z0-9+])|([^\0-9A-Za-z \t]) | (\w +: \ / \ / \S +)", " ",
3         tweet).split())
4
5 def porcentaje(parte, todo):
6     return 100*float(parte)/float(todo)
```


C. Análisis de sentimientos (Procesamiento)

Se lleva a cabo el análisis de sentimientos generando la polaridad respectiva de cada tweet mediante la función TextBlob

Listing 14. Proceso del Análisis de Sentimientos.

```
1 tweetText = []
2
3 positivo = negativo = neutral = polaridad = 0
4
5 for tweet in tweets:
6     print(tweet.text)
7     tweetText.append(Eliminacion_caracteres(tweet.text).encode('utf-8'))
8     analysis = TextBlob(tweet.text)
9     polaridad += analysis.sentiment.polarity
10
11     if (analysis.sentiment.polarity == 0.00):
12         neutral +=float(1)/float(3)
13         negativo +=float(1)/float(3)
14         positivo +=float(1)/float(3)
15     elif (analysis.sentiment.polarity < 0.00):
16         negativo +=1
17     elif (analysis.sentiment.polarity > 0.00):
18         positivo +=1
```

D. Almacenamiento de datos recolectados

Los datos obtenidos por el bot serán almacenados en un archivo .csv con el objetivo de posteriores estudios.

Listing 15. Almacenamiento de datos recolectados en .csv.

```
1 import csv
2 csvFile = open('datos.csv', 'a')
3 csvWriter = csv.writer(csvFile)
4 csvWriter.writerow(tweetText)
5 csvFile.close()
```

E. Análisis de sentimientos (Resultados)

Finalmente se visualizan las proporciones respectivas por cada sentimiento tipificado.

Listing 16. Resultados.

```
1 positivo = porcentaje(positivo,numero_de_Tweets)
2 negativo = porcentaje(negativo,numero_de_Tweets)
3 neutral = porcentaje(neutral,numero_de_Tweets)
4
5 positivo = format(positivo, '.2f')
6 negativo = format(negativo, '.2f')
7 neutral = format(neutral, '.2f')
8
9 print("Como el publico en Twitter reacciona a " + Termino_buscado + " mediante un analisis de " +
10      str(numero_de_Tweets)+ " Tweets.")
11
12 if (positivo > max(negativo,neutral)):
13     print("Positivo")
14 elif (negativo > max(positivo,neutral)):
15     print("Negativo")
16 elif (neutral > max(positivo,negativo)):
17     print("Neutral")
18
19 labels = ['Positivo [' + str(positivo) + '%]', 'Neutral [' + str(neutral) + '%]', 'Negativo [' +
20          str(negativo) + '%]']
21 porcentajes = [positivo,neutral,negativo]
22 colores = ['yellow','red','green']
23
24 patches, texts = plt.pie(porcentajes, colors = colores, startangle = 90)
25 plt.legend(patches, labels, loc="best")
26 plt.title("Como el publico en Twitter reacciona a " + Termino_buscado + " mediante un analisis de
27          " + str(numero_de_Tweets)+ " Tweets.")
28 plt.axis('equal')
29 plt.tight_layout()
30 plt.show()
```

Como el publico en Twitter reacciona a Paolo Guerrero mediante un analisis de 1000 Tweets.
Negativo

Fig. S6. Resultado descriptivo.

Como observamos hay una tendencia negativa obtenida de la información recolectada. Debido al contexto en evaluación esto puede deberse a que la mayoría de tweets vienen relacionados con el aspecto negativo producto de la tristeza, decepción por Paolo Guerrero, además de la critica negativa hacia los agentes que han producido este clima desfavorable en el futbol peruano.

Cómo el publico en Twitter reacciona a Paolo Guerrero mediante un analisis de 1000 Tweets.

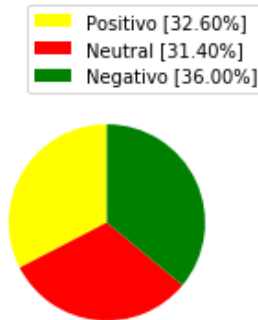


Fig. S7. Gráfico circular representativo.

Como se observa gráficamente, además de los porcentajes establecidos, que el mayor porcentaje en efecto son de sentimientos negativos; sin embargo no existe una amplia diferencia con respecto a los sentimientos positivos, estos pueden ser debido a los amplios mensajes de apoyo dedicados al futbolista. Existe además cierto porcentaje de sentimientos neutrales, estos pueden ser los relacionados con informaciones puntuales de usuarios correspondientes a noticias o los que deseen constatar informaciones en torno al tema pero sin mayor especificación de sus sentimientos.