

# Proyecto Final: Clasificación de tweets basado en el sentimiento.

Antonio Godoy, Diego Mendoza, Jesús Campos

**Resumen**—En los últimos años, el análisis de sentimiento ha sido una de las áreas que más fuerza ha tomado, debido a la producción masiva de comentarios o publicaciones por parte de usuarios en redes sociales, blogs, y otros sitios web. Asimismo, ha representado una fuente importante de información para las grandes empresas y organizaciones, puesto que representa el punto de partida para conocer lo que los consumidores están buscando, así como para conocer su reputación frente a ellos. Este documento busca abordar una comparación entre el Galaxy Note 8, el iPhone X, y el Pixel 2, a través de un proyecto de clasificación de publicaciones en Twitter (a partir de ahora tweets), basándose en un criterio de sentimiento y emociones.

**Index Terms**—R, RStudio, Aprendizaje máquina, Clasificador Bayesiano, Análisis de Sentimientos, Minería de Opinión, Minería de Datos, Twitter.

## 1. INTRODUCCIÓN

EN este documento se muestran los resultados obtenidos en el proyecto final elaborado para el curso de "Aprendizaje de Máquina", impartida por el Dr. Ulises Orozco Rosas en CETYS Universidad Campus Tijuana. El objetivo principal de éste reporte es el de dar a conocer el desenlace de la aplicación del proyecto final, que consiste en la clasificación de tweets sobre el iPhone X, el Galaxy Note 8, y el Pixel 2, de acuerdo con la aplicación de un análisis de sentimiento. Asimismo, se muestran los requerimientos y su implementación en el lenguaje de programación R.

## 2. DESCRIPCIÓN DEL PROBLEMA

El problema a tratar específicamente es la clasificación de tweets en la red social Twitter, de acuerdo a un criterio basado en el sentimiento. El proyecto se enfoca en el "data mining" dentro de Twitter. La idea principal es la de buscar tweets que contengan información sobre el iPhone X, el Galaxy Note 8, y el Pixel 2, y clasificarlos de acuerdo a tres categorías de sentimiento, positivo, neutral, y negativo. El modelo de clasificación de tweets se haría de acuerdo a las palabras contenidas en el mismo. Cada palabra tendrá una categoría de subjetividad asignada (positivo, negativo o neutral), dándole como resultado un valor numérico representando su grado de positividad o negatividad. Asimismo, se tendrá la posibilidad de agrupar tweets de acuerdo a la emoción que los caracteriza. La máquina, dado un conjunto de tweets y librerías de apoyo para el análisis de los mismos, agrupará dichos tweets en las categorías ya mencionadas.

## 3. OBJETIVOS

Los objetivos planteados para este proyecto se muestran a continuación:

- Crear una aplicación de Twitter y autenticarlo mediante código fuente en el lenguaje de programación R.

*Antonio Godoy, Diego Mendoza, Jesús Campos, son estudiantes de Licenciatura en Ingeniería en Ciencias Computacionales en CETYS Universidad Campus Tijuana.*

- Extraer tweets desde la red social de Twitter.
- Depurar el contenido de los tweets, eliminando caracteres extraños y palabras irrelevantes en el análisis.
- Obtener las palabras con mayor frecuencia de aparición en los tweets.
- Representar gráficamente las palabras con mayor frecuencia (nube de palabras, gráficas de barras, etc.).
- Ligar un conjunto de datos que contengan las palabras positivas y negativas, con las cuales se le dará valor a los tweets.
- Crear un método/función que consiga asociar el contenido de los tweets con las palabras del conjunto de datos correspondiente.
- Normalizar los valores o resultados obtenidos de los tweets a una "calificación" estandarizada.
- Ligar la "calificación" con la clasificación de sentimientos (Positivo, neutro, negativo).
- Clasificar los tweets de acuerdo al sentimiento predominante, utilizando una función apoyada en el clasificador Bayesiano.
- Clasificar los tweets de acuerdo a la emoción predominante (miedo, enojo, alegría, etc.), utilizando el clasificador mencionado.
- De acuerdo a una cierta palabra, o a un grupo de palabras, y a las categorías mencionadas anteriormente, proponer soluciones a problemas encontrados por las palabras negativas. Por ejemplo, sugerir una solución a la opinión pública que cree que el iPhone X es muy caro, o que ventajas/desventajas presenta sobre otros dispositivos.

## 4. MARCO TEÓRICO

- **API (Application Programming Interface):** Una API es un marco que facilita la creación de servicios HTTP disponibles para una amplia variedad de clientes, entre los que se incluyen exploradores y dispositivos móviles.
- **Análisis de Sentimiento:** Es el proceso donde la máquina por medio de algoritmos de clasificación y regresión

(para este caso), determina si una frase contiene alguna opinión, positiva, negativa o neutral.

- **Minería de Opinión:** Serie de aplicaciones de técnicas del procesamiento del lenguaje natural, lingüística computacional y minería de textos, que tienen como objetivo la extracción de información subjetiva a partir de contenidos generados por los usuarios, ejemplo: Facebook, Twitter, etc.
- **Minería de Datos:** Es el conjunto de técnicas y tecnologías que permiten explorar grandes bases de datos, de manera automática o semiautomática, con el objetivo de encontrar patrones repetitivos, tendencias o reglas que expliquen el comportamiento de los datos en un determinado contexto.
- **Nube de Palabras:** Es la representación visual de las palabras e oraciones que conforman un texto, donde el tamaño es mayor para las palabras buscadas o aparecidas con mas frecuencia.
- **Procesamiento de Lenguaje Natural (NLP):** es el campo que combina las tecnologías de la ciencia computacional (como la inteligencia artificial, el aprendizaje automático o la inferencia estadística) con la lingüística aplicada, con el objetivo de hacer posible la comprensión y el procesamiento asistidos por ordenador de información expresada en lenguaje humano para determinadas tareas, como la traducción automática, los sistemas de diálogo interactivos, el análisis de opiniones, etc.
- **Clasificador Bayesiano:** El algoritmo de clasificación Naive Bayes es un clasificador probabilístico. Se basa en modelos de probabilidades que incorporan fuertes suposiciones de independencia. Las suposiciones de independencia a menos no tienen ningún efecto sobre la realidad. Por lo tanto, se consideran ingenuas (naive en inglés).

## 5. METODOLOGÍA DE LA INVESTIGACIÓN

Como se mencionó con anterioridad, el problema radica en la clasificación de tweets encontrados de acuerdo a un análisis de sentimiento, obteniendo como resultado tres agrupaciones de tweets, positivos, neutrales, y negativos, asimismo, se realizará una clasificación de acuerdo a la emoción que predomine en dicho tweet. Los pasos a seguir en nuestra investigación y a través de los cuales se obtendrá este resultado, se muestran a continuación:

- **Conexión con Twitter:** Lo primero que se busca es lograr establecer una conexión entre la aplicación de Twitter y el ambiente de RStudio a través de librerías, realizando el proceso de autenticación solicitado por la red social.
- **Búsqueda paramétrica:** Una vez lograda la conexión entre R y Twitter, será necesario realizar una búsqueda paramétrica, en donde los parámetros a recibir serán las palabras o temas de interés, así como una cantidad determinada de tweets. Sobre éstas palabras se buscarán tweets en la red social.
- **Extracción de tweets:** De acuerdo a los parámetros especificados en el paso anterior, mediante las librerías pertinentes se hará uso de una función para extraer los

tweets solicitados desde la red social y se almacenarán en una variable de especificada en R.

- **Depuración de tweets:** Una vez extraídos los tweets, será necesario limpiarlos. En este proceso se busca eliminar todo tipo de palabras que no tenga importancia, conocidas como "stop words". Asimismo, se eliminarán palabras que no sean de nuestro interés en la búsqueda, o bien, que sean redundantes, tales como la misma palabra buscada, o palabras relacionadas de manera obvia a ella.
- **Diccionario de palabras:** Previo al análisis y valoración de los tweets extraídos, será necesario conocer una lista o diccionario de palabras sobre las cuales será evaluado. Estas palabras podrán ser extraídas a través de alguna librería para R, de algún repositorio de internet, o bien, formulada por parte del usuario.
- **Valoración de tweets:** A través de una lectura del tweet, se hará un mapeado entre las palabras leídas en la publicación y aquellas que coincidan de los diccionarios. Una vez terminada la lectura de la publicación, se obtendrá una valoración general y estandarizada del mismo.
- **Determinar subjetividad:** Al obtener su valoración, se buscará conocer el grado de subjetividad de tweet, este grado estará dado por una categoría, a la cual pertenecerá el tweet. Dicha categoría tiene como opciones los tweets muy positivos, positivos, neutrales, negativos, y muy negativos.
- **Clasificación de tweets:** Una vez realizado lo anterior, los grupos de tweets se verán reducidos a tres categorías, en las cuales el tweet podrá ser ubicado. Estas categorías serán las siguientes: "Tweets Positivos", "Tweets Neutros", "Tweets Negativos". Cada categoría contendrá un rango de valores estimado a través de los cuales la máquina podrá predecir a cuál pertenecerá el tweet. De la misma manera que en la valoración de tweets, los intervalos de las categorías estarán en función de un valor determinado por la puntuación del tweet.
- **Graficación de resultados:** Al haber conseguido clasificar los tweets, se buscará representar los resultados de la clasificación a través de gráficos que permitan analizar la relación de los dispositivos con las salidas.
- **Presentación de palabras representativas de cada clasificación:** Al haber clasificado correctamente los tweets, se buscará que la máquina extraiga las palabras más frecuentes en la búsqueda, así como las más significativas dentro de cada categoría de emociones, con la finalidad de que el usuario conozca la relación entre dichas palabras y el sentimiento u opiniones de los usuarios de la red social.

## 6. DESCRIPCIÓN DE LA IMPLEMENTACIÓN

### 6.1. Creación de Aplicación y Autenticación

El primer paso antes de la codificación, es la creación de una aplicación en Twitter. Se crea una cuenta de usuario de Twitter (en caso de no tener una) y se accede a la página web [www.app.twitter.com](http://www.app.twitter.com), se iniciará sesión con la cuenta creada, y se creará una app. Una vez creada la app de Twitter, se generarán las siguientes claves, *api key*, *api secret*, *access*

*token*, y *access token secret*, las cuales serán utilizadas para la autorización. La autenticación requerirá que el usuario ingrese en consola una clave dinámica generada por twitter.

## 6.2. Extracción de Tweets

Una vez autenticada la aplicación, el siguiente paso es extraer los tweets utilizando la función *searchTwitter()*. Esta función lleva como argumentos la palabra a buscar, el número de tweets a recabar, el lenguaje, el tipo de tweet a recabar (popular, reciente, etc.), un rango de fechas, entre otros. No obstante, puesto que nuestro objetivo es el de recabar información sobre el Iphone X y el Galaxy Note 8 para compararlos, bastará con hacer una búsqueda para cada uno de esos productos, desde mediados de 2017.

## 6.3. Limpieza de Tweets

Al extraer los tweets, éstos traen consigo demasiados caracteres (como los emoticones, los *hashtags*, enlaces, etc), así como palabras irrelevantes a nuestra búsqueda. Es por eso que mediante una función creada llamada *Clean()*, se depurará el texto de cada uno de los tweets sobre ambos dispositivos, eliminando cualquier palabra y caracter innecesario.

## 6.4. Adjuntar los Conjuntos de Datos de Palabras

Al mismo tiempo, teniendo listos los tweets (únicamente con palabras significativas), es posible comenzar con el proceso de puntuación. El primer paso para lograrlo es introduciendo un listado de palabras buenas y malas, sobre las cuales será puntuado el tweet. Las listas utilizadas en este proyecto se encuentran referenciadas en la sección de bibliografías.

## 6.5. Creación de las Funciones de Polaridad y Emociones

Al lograr vincular los listados de palabras con el programa, es posible comenzar con la clasificación de tweets. Sin embargo, para llegar a este paso es necesario crear funciones que se encarguen de asignar una puntuación a cada uno de los tweets encontrados.

## 6.6. Clasificación de Tweets

Una vez creadas las funciones de polaridad y de emociones, el siguiente paso es clasificar los tweets. Esta clasificación, a su vez, se logrará mediante otra función que logre vincular las puntuaciones de cada tweet con la polaridad o emoción uw mejor se ajuste.

## 6.7. Graficación de Resultados

Al haber clasificado los tweets, el último paso consta de mostrar gráficamente los resultados, de manera que éstos sean fáciles de entender e interpretar por las personas. Por lo tanto, para la elaboración de esta tarea, se utilizarán gráficos simples, tales como las gráficas de barras, las gráficas de pastel, entre otros.

## 6.8. Creación de Nube de Palabras

Una vez que los tweets han sido depurados y que solo contienen palabras significativas, es posible utilizarlos para realizar una nube de palabras con la ayuda de la librería *wordcloud*. Esta librería permite crear una nube de palabras que contiene las palabras con mayor frecuencia de aparición en los tweets ya clasificados, sobre cada dispositivo, dando una pauta para conocer qué es lo que opinan los usuarios de Twitter acerca de ellos.

## 7. CONCLUSIÓN

De acuerdo a lo mencionado a lo largo de este documento, se puede inferir que el Análisis de Sentimiento y la Minería de Datos son herramientas muy poderosas actualmente. A través de la implementación de ellas en el lenguaje de programación R, el equipo ha sido capaz de analizar las opiniones de los usuarios de twitter sobre los distintos dispositivos propuestos (Apple iPhone X, Samsung Galaxy Note 8 y Google Pixel 2). Esto ha servido para darnos una idea de la perspectiva de las personas sobre cada uno de estos terminales, así como para proponer soluciones que consigan erradicar cualquier sentimiento negativo o de rechazo hacia un dispositivo en específico. Cabe destacar la importancia y la extensa área de aplicación que tienen estas herramientas, pues pueden ser utilizadas en casi cualquier rubro, siempre y cuando exista un tráfico de datos considerable para realizar el análisis. Estas herramientas, además, permiten llegar a soluciones de una manera efectiva y rápida.

Es importante considerar que el curso de aprendizaje de máquina fue totalmente beneficioso para el desarrollo de este proyecto, pues gracias al contenido de la materia le fue posible al equipo comprender la estructura del problema, así como la metodología a seguir para lograr el objetivo, que era el análisis de sentimiento en los tweets relacionados al iPhone X, Note 8 y Pixel 2. Se espera que se le de un seguimiento a estas herramientas y lograr, a través de ellas, formular soluciones a problemas de gran importancia para nuestra generación, así como aprovechar al máximo la gran cantidad de datos que tenemos actualmente. De esta manera, no solo se aplicarán los conocimientos aprendidos en el curso de aprendizaje de máquina, sino que también se llevarán a cabo en proyectos aplicados en situaciones reales y con datos reales.

## 8. RESULTADOS DE LA IMPLEMENTACIÓN

A continuación, se muestran los resultados obtenidos tras la implementación del proyecto.

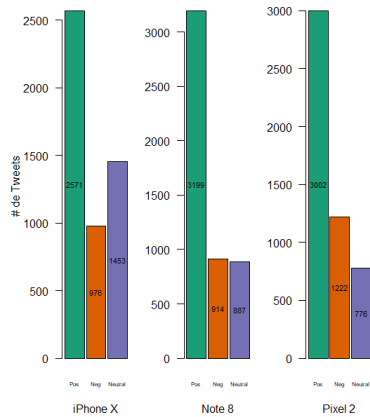


Figura 1. Tweets categorizados por subjetividad en cada dispositivo.

De acuerdo a la clasificación de tweets por subjetividad en cada dispositivo, se obtuvo la Figura 1. Por los números indicados en las gráficas, se puede observar cómo es que de un total de 5000 tweets recabados por cada dispositivo, el iPhone X resultó ser el que más tweets neutros obtuvo, sin embargo, fue el que consiguió la menor cantidad de tweets positivos. Por otro lado, el Note 8 obtuvo la menor cantidad de tweets negativos, y fue el dispositivo con mayor número de tweets positivos. Lo que realmente sorprende es que el Pixel 2 haya conseguido superar al iPhone X en cuanto a tweets positivos, pues fue el que mayor consiguió, aunque también consiguió la mayor cantidad de tweets negativos y la menor cantidad de tweets neutros.

Comentado lo anterior, se puede inferir que, gracias a que las marcas de Apple y Samsung son muy reconocidas por sus smartphones, las críticas hacia ellos han sido un poco más regulares, inclinándose hacia el lado positivo, caso contrario con el Pixel 2, que acaba de introducirse a esta rama y cuyos tweets negativos representan más de la tercera parte de los positivos, siendo muy irregular.

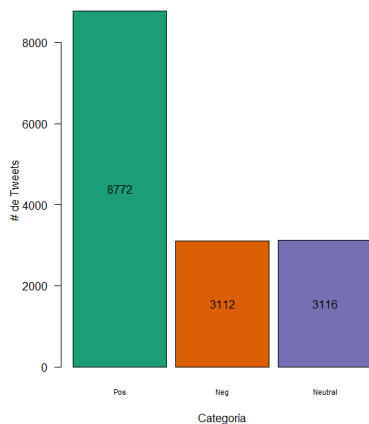


Figura 2. Tweets categorizados por subjetividad de todos los dispositivos.

Con respecto a los datos recabados por la Figura 1,

es posible construir un gráfico que represente de manera general la totalidad de tweets recibidos por categoría de subjetividad. Como era de esperarse, los tweets positivos resultaron ser la mayor parte de los totales recabados, sin embargo, los tweets negativos se encuentran a la par con los tweets neutrales. Esta elevada cantidad de tweets negativos puede deberse a diferentes factores, tales como el costo del iPhone X, Note 8, y Pixel 2, las explosiones de las baterías en smartphones de Samsung, o la inexperiencia de Google en la industria de la telefonía con el Pixel 2.

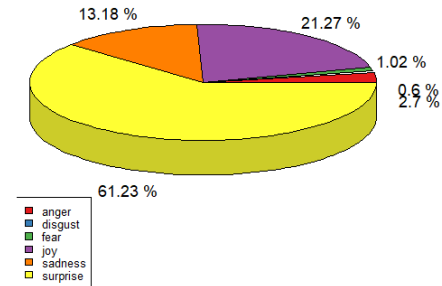


Figura 3. Distribución de las emociones en el iPhone X.

En la Figura 3, se busca analizar la distribución de los tweets en las diferentes categorías de emociones. Se puede inferir, por obviedad, que la gran mayoría de los tweets recaen en la emoción de *sorpres*a, pues el iPhone X llegó para revolucionar en Apple con todas sus novedades, además, se trataba de nada más y nada menos que del celular con el que la compañía festejaría el décimo aniversario del iPhone. Cabe destacar que existe una parte sumamente considerable que se categoriza triste con el iPhone X. Esto puede deberse a su elevado costo de lanzamiento, a su poca disponibilidad, o bien, a su diseño.

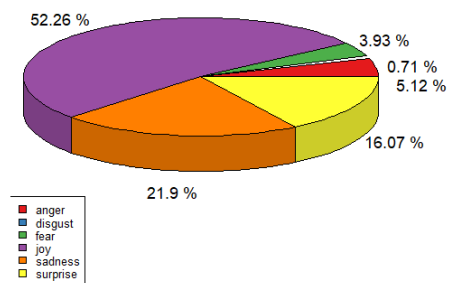


Figura 4. Distribución de las emociones en el Note 8.

En cuanto al Galaxy Note 8, la gran mayoría de los tweets se encuentran concentrados en la emoción de *alegría*, mientras que casi la cuarta parte del total recabado se categorizó como *triste*. Los principales puntos a considerar para la tristeza de los usuarios puede ser el alto costo impuesto por Samsung al dispositivo, la fragilidad de la pantalla, o bien, por el diseño presentado. Por otro lado, la alegría de los usuarios podría radicar en la gran potencia y funcionalidad que permite el dispositivo, así como en el gran tamaño de la pantalla presentado en este modelo.







respecto a los diseños presentados, así como a la autonomía de la batería y al tamaño de la pantalla, sin embargo, puede entenderse por *buzzing* a la gran espera que han tenido por los dispositivos, sintiendo ese miedo de que no sea de su agrado, o bien, que se tenga el temor de fallos de audio relacionados con zumbidos.

## 9. REQUERIMIENTOS

Para la correcta implementación del proyecto en el lenguaje de programación R, se necesitarán las siguientes librerías y archivos listados a continuación:

- **ROAuth:** Librería utilizada para realizar una autenticación con la aplicación de Twitter.
- **twitterR:** Librería requerida para extraer los tweets.
- **base64enc:** Librería para la codificación en base64.
- **httr:** Librería necesaria para realizar una conexión vía URL y HTTPS.
- **devtools:** Librería que permite instalar librerías de R desde un repositorio de GitHub, entre otras funciones.
- **tm:** Librería enfocada en la implementación de *text mining*, requerida para realizar la limpieza de los tweets.
- **wordcloud:** Librería que permite crear una nube de palabras con mayor frecuencia de aparición en un conjunto de datos.
- **RColorBrewer:** Librería que permite utilizar paletas de colores en la graficación de nubes de palabras.
- **RTextTools:** Librería utilizada como base para la función creada *createMatrix()* (*create\_matrix* de *RTextTools*).
- **Opinion/sentiment Lexicon:** Obtenido de <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html#lexicon>. (Link de descarga)
- **Funciones bayesianas:** Funciones implementadas en el lenguaje de programación R, necesarias para realizar la clasificación de tweets de acuerdo a la emoción correspondiente, así como su clasificación de acuerdo a la polaridad.
- **Relación de palabras con emociones y con subjetividad:** Base de datos que contiene palabras relacionadas con una emoción y subjetividad específica, requerida por las funciones bayesianas descritas anteriormente. Últimos dos requerimientos obtenidos de: <https://cran.r-project.org/src/contrib/Archive/sentiment/>. (Link de descarga)

## 10. CÓDIGO DE IMPLEMENTACIÓN EN R

Para fines prácticos y de estética del documento, se han omitido los comentarios del código. De cualquier manera, en el repositorio del proyecto (listado en la sección de Enlace del Video de Resultados) se encuentra el código de manera íntegra.

```
1 #-----INSTALAR LIBRERIAS-----#
2 # install.packages("ROAuth")
3 # install.packages("twitterR")
4 # install.packages("base64enc")
5 # install.packages("httr")
6 # install.packages("devtools")
7 # install.packages("tm")
8 # install.packages("wordcloud")
9 # install.packages("RColorBrewer")
10 # install.packages("RTextTools")

11 #-----CARGAR LIBRERIAS-----#
12 library(ROAuth)
13 library(twitterR)
14 library(base64enc)
15 library(httr)
```

```
16 library(devtools)
17 library(tm)
18 library(wordcloud)
19 library(RColorBrewer)
20 library(RTextTools)

21 #-----AUTENTICACION DE TWITTER-----#
22 api_key = "Y5fbQA5lJodmk0E4q4c1DYbXD"
23 api_secret = "f4OQTcDmCg56EiH5aSZZ3P0HrEnZT0QHj0KbqBmZyRWmNkVL"
24 access_token = "919998032938012672-6Gf05cGcS1SZacw19QwXw8VWuIKDXDe"
25 access_token_secret = "YW0mWl4B3UTxOeYmVypwVhQyzlqygfO2cN2ySd1RZk4"
26 request_url = "https://api.twitter.com/oauth/request_token"
27 access_url = "https://api.twitter.com/oauth/access_token"
28 auth_url = "https://api.twitter.com/oauth/authorize"
29
30 setup_twitter_oauth(api_key, api_secret, access_token, access_token_secret)
31
32 credential = OAuthFactory$new(consumerKey = api_key, consumerSecret =
33   api_secret, requestURL = request_url, accessURL = access_url, authURL =
34   auth_url)
35
36 credential$handshake(cainfo = system.file("CurlSSL", "cacert.pem", package =
37   "RCurl"))
38
39 #-----EXTRACCION DE TWEETS-----#
40 iphoneTweets = searchTwitter("iPhone X", n=5000, lang="en", since = "2017-06-01"
41 )
42 noteTweets = searchTwitter("Galaxy Note 8", n=5000, lang="en", since =
43   "2017-06-01")
44 pixelTweets = searchTwitter("Google+Pixel 2", n=5000, lang="en", since =
45   "2017-06-01")
46
47 iphoneTweets.df = twListToDF(iphoneTweets)
48 noteTweets.df = twListToDF(noteTweets)
49
50 #-----LIMPIEZA DE TWEETS-----#
51 tweets.clean = function(tweetsList){
52   cleanedTweets = apply(tweetsList, function(row) iconv(row, "latin1", "ASCII",
53     sub=""))
54   cleanedTweets = gsub("(f|ht)tp(s)?/(.*)" [a-z]+", "", cleanedTweets)
55   tweetsCorpus = Corpus(VectorSource(cleanedTweets[names(cleanedTweets)]))
56   tweetsCorpus = tm_map(tweetsCorpus, function(tweet) gsub("http[^:space:]]*",
57     "", tweet))
58   tweetsCorpus = tm_map(tweetsCorpus, function(tweet) gsub("\\xed[\\[:space:]]*",
59     "", tweet))
60   tweetsCorpus = tm_map(tweetsCorpus, function(tweet) gsub("/[\\[:space:]]*",
61     "", tweet))
62   tweetsCorpus = tm_map(tweetsCorpus, function(tweet) gsub("@[\\[:space:]]*",
63     "", tweet))
64   tweetsCorpus = tm_map(tweetsCorpus, removePunctuation)
65   tweetsCorpus = tm_map(tweetsCorpus, content_transformer(tolower))
66   tweetsCorpus = tm_map(tweetsCorpus, removeWords, c(stopwords("english"), "\n",
67     "rt"))
68   tweetsCorpus = tm_map(tweetsCorpus, removeNumbers)
69   searchedWords = c("iphone x", "iphonex", "iphone", "apple", "ios", "galaxy",
70     "galaxynote", "note", "note 8", "samsung", "android", "google", "pixel",
71     "pixel 2", "pixel2", "phone", "smartphone", "cellphone", "giveaway",
72     "international", "tweet")
73
74   tweetsCorpus = tm_map(tweetsCorpus, removeWords, searchedWords)
75   tweetsCorpus = tm_map(tweetsCorpus, stripWhitespace)
76   cleanedTweets = data.frame(text = apply(tweetsCorpus, as.character),
77     stringsAsFactors = FALSE)
78   return(cleanedTweets)
79 }
80 iphoneCleanedTweets = tweets.clean(iphoneTweets.text)
81 noteCleanedTweets = tweets.clean(noteTweets.text)
82 pixelCleanedTweets = tweets.clean(pixelTweets.text)
83
84 #-----ADJUNTAR CONJUNTO DE DATOS DE PALABRAS-----#
85 posWords = scan("./posWords.txt", what='character', comment.char = ';')
86 negWords = scan("./negWords.txt", what='character', comment.char = ';')
87
88 #-----FUNCIONES DE PUNTUACION POR PALABRAS POSITIVAS Y NEGATIVAS-----#
89 getScores = function(tweets, pos.words, neg.words){
90   results.df = data.frame(matrix(nrow=0, ncol=5))
91
92   for(i in 1:length(tweets)){
93     tweet = tweets[i]
94     words = strsplit(tweet, ' ')
95     words = unlist(words)
96     words = words[!words %in% c("", " ", "NA")]
97     pos.match = match(words, posWords)
98     pos.match = !is.na(pos.match)
99     neg.match = match(words, negWords)
100    neg.match = !is.na(neg.match)
101    totalPos = sum(pos.match)
102    totalNeg = sum(neg.match)
103    score = totalPos - totalNeg
104
105    if(score <= -3){
106      category = "very negative"
107    }
108    else if(score > -3 && score < 0){
109      category = "negative"
110    }
111    else if(score == 0){
112      category = "neutral"
113    }
114    else if(score > 0 && score <= 2){
115      category = "positive"
116    }
117    else{
118      category = "very positive"
119    }
120    if(length(words) != 0L){
121      results.df[i,] = c(tweet, totalPos, totalNeg, score, category)
122    }
123  }
124
125  colnames(results.df) = c("tweet", "pos", "neg", "score", "category")
126  return(results.df)
127 }
```

```

110 iphone.scores = getScores(iphoneCleanedTweets$text, posWords, negWords)
111 note.scores = getScores(noteCleanedTweets$text, posWords, negWords)
112 pixel.scores = getScores(pixelCleanedTweets$text, posWords, negWords)

113 #-----FUNCIONES DE POLARIDAD Y EMOCIONES-----#
114 createMatrix = function(textColumns, language="english", minDocFreq=1,
115   minWordLength=3, weighting=weightTf){
116   control = list(language=language, minDocFreq=minDocFreq, minWordLength=
117     minWordLength, weighting=weighting)
118   trainingColumn = apply(as.matrix(textColumns), 1, paste, collapse=" ")
119   trainingColumn = supply(as.vector(trainingColumn, mode="character"), iconv, to
120     ="UTF8", sub="byte")
121   corpus = Corpus(VectorSource(trainingColumn), readerControl=list(language=
122     language))
123   matrix = DocumentTermMatrix(corpus, control=control)
124   return(matrix)
125 }
126
127 bayesianEmotions = function(tweets, prior=1.0){
128   matrix = createMatrix(tweets)
129   lexicon = read.csv(file = "emotions.csv", header=FALSE, sep=",")
130   counts = list(anger=length(which(lexicon[,2]=="anger")), disgust=length(
131     which(lexicon[,2]=="disgust")), fear=length(which(lexicon[,2]=="fear")),
132     joy=length(which(lexicon[,2]=="joy")), sadness=length(which(lexicon
133       [,2]=="sadness")), surprise=length(which(lexicon[,2]=="surprise")),
134     total=nrow(lexicon))
135   documents = c()
136   for(i in 1:nrow(matrix)){
137     scores = list(anger=0, disgust=0, fear=0, joy=0, sadness=0, surprise=0)
138     doc = matrix[i,]
139     words = findFreqTerms(doc, lowfreq=1)
140     for(word in words){
141       for(key in names(scores)){
142         emotions = lexicon[which(lexicon[,2]==key),]
143         index = pmatch(word, emotions[,1], nomatch=0)
144         if(index>0){
145           entry = emotions[index,]
146           category = as.character(entry[[2]])
147           count = counts[[category]]
148           score = 1.0
149           score = abs(log(score*prior/count))
150           scores[[category]] = scores[[category]]+score
151         }
152       }
153     }
154     best_fit = names(scores)[which.max(unlist(scores))]
155     if(best_fit == "disgust" && as.numeric(unlist(scores[2]))-3.09234 < .01){
156       best_fit = NA
157     }
158     documents = rbind(documents, c(scores$anger, scores$disgust, scores$fear,
159       scores$joy, scores$sadness, scores$surprise, best_fit))
160   }
161   colnames(documents) = c("Anger", "Disgust", "Fear", "Joy", "Sadness", "Surprise",
162     "Emotion")
163   return(documents)
164 }
165
166 bayesianPolarity = function(tweets, pStrong=0.5, pWeak=1.0, prior=1.0){
167   matrix = createMatrix(tweets)
168   lexicon = read.csv(file = "subjectivity.csv", header=FALSE, sep=",")
169   counts = list(positive=length(which(lexicon[,3]=="positive")), negative=
170     length(which(lexicon[,3]=="negative")), total=nrow(lexicon))
171   documents = c()
172   for(i in 1:nrow(matrix)){
173     scores = list(positive=0, negative=0)
174     doc = matrix[i,]
175     words = findFreqTerms(doc, lowfreq=1)
176     for(word in words){
177       index = pmatch(word, lexicon[,1], nomatch=0)
178       if(index>0){
179         entry = lexicon[index,]
180         polarity = as.character(entry[[2]])
181         category = as.character(entry[[3]])
182         count = counts[[category]]
183         score = pWeak
184         score = abs(log(score*prior/count))
185         scores[[category]] = scores[[category]]+score
186       }
187     }
188     best_fit = names(scores)[which.max(unlist(scores))]
189     ratio = as.integer(abs(scores$positive/scores$negative))
190     if(!length(ratio)==0){
191       if(ratio==1){
192         best_fit="neutral"
193       }
194     }
195     documents = rbind(documents, c(scores$positive, scores$negative, abs(
196       scores$positive/scores$negative), best_fit))
197   }
198   colnames(documents) = c("Pos", "Neg", "Pos/Neg", "Polarity")
199   return(documents)
200 }
201
202 emotions = c("anger", "disgust", "fear", "joy", "sadness", "surprise")
203 iphoneEmotions.class = bayesianEmotions(iphone.scores$tweet)
204
205 noteEmotions.class = bayesianEmotions(note.scores$tweet)
206 pixelEmotions.class = bayesianEmotions(pixel.scores$tweet)
207
208 iphoneEmotions = iphoneEmotions.class[,7]
209 noteEmotions = noteEmotions.class[,7]
210 pixelEmotions = pixelEmotions.class[,7]
211
212 iphoneEmotions[is.na(iphoneEmotions)] = "unknown"
213 noteEmotions[is.na(noteEmotions)] = "unknown"
214 pixelEmotions[is.na(pixelEmotions)] = "unknown"
215
216 iphonePolarity.class = bayesianPolarity(iphone.scores$tweet)
217 notePolarity.class = bayesianPolarity(note.scores$tweet)
218 pixelPolarity.class = bayesianPolarity(pixel.scores$tweet)
219
220 iphonePolarity = iphonePolarity.class[,4]
221 notePolarity = notePolarity.class[,4]
222 pixelPolarity = pixelPolarity.class[,4]
223
224 iphonePolarity = data.frame(text=iphone.scores$tweet, emotion=iphoneEmotions,
225   polarity=iphonePolarity, stringAsFactors=FALSE)
226 notePolarity = data.frame(text=note.scores$tweet, emotion=noteEmotions,
227   polarity=notePolarity, stringAsFactors=FALSE)
228 pixelPolarity = data.frame(text=pixel.scores$tweet, emotion=pixelEmotions,
229   polarity=pixelPolarity, stringAsFactors=FALSE)
230
231 #-----GRAFICACION DE RESULTADOS-----#
232
233 #-----Subjetividad por dispositivo-----#
234
235 iphone.sub = data.frame(Pos = length(which(iphonePolarity.class=="positive")),
236   Neg = length(which(iphonePolarity.class=="negative")), Neutral =
237   length(which(iphonePolarity.class=="neutral")))
238 note.sub = data.frame(Pos = length(which(notePolarity.class=="positive")),
239   Neg = length(which(notePolarity.class=="negative")), Neutral =
240   length(which(notePolarity.class=="neutral")))
241 pixel.sub = data.frame(Pos = length(which(pixelPolarity.class=="positive")),
242   Neg = length(which(pixelPolarity.class=="negative")), Neutral =
243   length(which(pixelPolarity.class=="neutral")))
244
245 par(fig = c(0.05, 0.35, 0, 0.95), mar = c(5,4,3,0))
246 barplot(as.numeric(c(iphone.sub$Pos, iphone.sub$Neg, iphone.sub$Neutral)),
247   names.arg = names(iphone.sub), space=c(0.1,1), cex.names = 0.5, xlab="
248   iPhone X", ylab="# de Tweets", las = 1, col = brewer.pal(3, "Dark2"))
249 text(c(1.5, 2.6, 3.7), c(iphone.sub$Pos/2, iphone.sub$Neg/2, iphone.sub$Neutral/2),
250   c(iphone.sub$Pos, iphone.sub$Neg, iphone.sub$Neutral), cex = 0.7)
251
252 par(fig = c(0.35, 0.65, 0, 0.95), mar = c(5,4,3,0), new = TRUE)
253 barplot(as.numeric(c(note.sub$Pos, note.sub$Neg, note.sub$Neutral)), names.arg
254   = names(note.sub), space=c(0.1,0), cex.names = 0.5, xlab="Note 8", las
255   = 1, col = brewer.pal(3, "Dark2"))
256 text(c(0.5, 1.6, 2.7), c(note.sub$Pos/2, note.sub$Neg/2, note.sub$Neutral/2),
257   c(note.sub$Pos, note.sub$Neg, note.sub$Neutral), cex = 0.7)
258
259 par(fig = c(0.65, 0.95, 0, 0.95), mar = c(5,4,3,0), new = TRUE)
260 barplot(as.numeric(c(pixel.sub$Pos, pixel.sub$Neg, pixel.sub$Neutral)), names.
261   arg = names(pixel.sub), space=c(0.1,0), cex.names = 0.5, xlab="Pixel 2",
262   las = 1, col = brewer.pal(3, "Dark2"))
263 text(c(0.5, 1.6, 2.7), c(pixel.sub$Pos/2, pixel.sub$Neg/2, pixel.sub$Neutral/2),
264   c(pixel.sub$Pos, pixel.sub$Neg, pixel.sub$Neutral), cex = 0.7)
265
266 title(main = "Tweets Categorizados por Subjetividad por Dispositivo", outer =
267   TRUE, line = -2)
268
269 #-----Emociones por dispositivo-----#
270
271 mixed.sub = data.frame(Pos = length(which(iphonePolarity.class=="positive"))
272   + length(which(notePolarity.class=="positive")) + length(which(
273     pixelPolarity.class=="positive"))),
274   Neg = length(which(iphonePolarity.class=="negative"))
275   + length(which(notePolarity.class=="negative")) + length(which(
276     pixelPolarity.class=="negative"))),
277   Neutral = length(which(iphonePolarity.class=="neutral")
278   + length(which(notePolarity.class=="neutral")) + length(which(
279     pixelPolarity.class=="neutral"))))
280
281 barplot(as.numeric(c(mixed.sub$Pos, mixed.sub$Neg, mixed.sub$Neutral)), names.
282   arg = names(mixed.sub), space=c(0.1,1), cex.names = 0.7, xlab="
283   Categoria", ylab="# de Tweets", las = 1, col = brewer.pal(3, "Dark2"))
284 text(c(1.5, 2.6, 3.7), c(mixed.sub$Pos/2, mixed.sub$Neg/2, mixed.sub$Neutral/2),
285   c(mixed.sub$Pos, mixed.sub$Neg, mixed.sub$Neutral), cex = 1)
286
287 title(main = "Tweets Categorizados por Subjetividad de Todos los Dispositivos",
288   outer = TRUE, line = -2)
289
290 #-----Emociones por dispositivo-----#
291
292 iphoneAnger = sum(length(which(iphoneEmotions == "anger")))
293 iphoneDisgust = sum(length(which(iphoneEmotions == "disgust")))
294 iphoneFear = sum(length(which(iphoneEmotions == "fear")))
295 iphoneJoy = sum(length(which(iphoneEmotions == "joy")))
296 iphoneSadness = sum(length(which(iphoneEmotions == "sadness")))
297 iphoneSurprise = sum(length(which(iphoneEmotions == "surprise")))
298 emotionsIphone = c(iphoneAnger, iphoneDisgust, iphoneFear, iphoneJoy,
299   iphoneSadness, iphoneSurprise)
300 total = sum(iphoneAnger, iphoneDisgust, iphoneFear, iphoneJoy, iphoneSadness,
301   iphoneSurprise)
302 percentages = c(paste(round(iphoneAnger/total*100,2), "%"), paste(round(
303   iphoneDisgust/total*100,2), "%"), paste(round(iphoneFear/total*100,2), "%"),
304   paste(round(iphoneJoy/total*100,2), "%"), paste(round(iphoneSadness/
305     total*100,2), "%"), paste(round(iphoneSurprise/total*100,2), "%"))
306
307 pie3D(x=emotionsIphone, labels = percentages, col = brewer.pal(6, "Set1"), main=
308   "Distribucion de Emociones en el iPhone X (%)", labelcex = 1.1)
309 legend("bottomleft", legend = emotions, fill = brewer.pal(6, "Set1"), cex =
310   0.55)
311
312 noteAnger = sum(length(which(noteEmotions == "anger")))
313 noteDisgust = sum(length(which(noteEmotions == "disgust")))
314 noteFear = sum(length(which(noteEmotions == "fear")))
315 noteJoy = sum(length(which(noteEmotions == "joy")))
316 noteSadness = sum(length(which(noteEmotions == "sadness")))
317 noteSurprise = sum(length(which(noteEmotions == "surprise")))
318 emotionsNote = c(noteAnger, noteDisgust, noteFear, noteJoy, noteSadness,
319   noteSurprise)
320 total = sum(noteAnger, noteDisgust, noteFear, noteJoy, noteSadness, noteSurprise)
321 percentages = c(paste(round(noteAnger/total*100,2), "%"), paste(round(
322   noteDisgust/total*100,2), "%"), paste(round(noteFear/total*100,2), "%"),
323   paste(round(noteJoy/total*100,2), "%"), paste(round(noteSadness/total

```



```

285     *100,2),"%"),paste(round(noteSurprise/total*100,2),"%"))
286 pie3D(emotionsNote, labels = percentages, col = brewer.pal(6,'Set1'),main =
287   "Distribucion de Emociones en el Note 8 (%)", labelcex = 1.1)
288 legend("bottomleft", legend = emotions, fill = brewer.pal(6,'Set1'), cex =
289   .5)
290 pixelAnger = sum(length(which(pixelEmotions == "anger")))
291 pixelDisgust = sum(length(which(pixelEmotions == "disgust")))
292 pixelFear = sum(length(which(pixelEmotions == "fear")))
293 pixelJoy = sum(length(which(pixelEmotions == "joy")))
294 pixelSadness = sum(length(which(pixelEmotions == "sadness")))
295 pixelSurprise = sum(length(which(pixelEmotions == "surprise")))
296 emotionsPixel = c(pixelAnger, pixelDisgust, pixelFear, pixelJoy, pixelSadness,
297   pixelSurprise)
298 total = sum(pixelAnger, pixelDisgust, pixelFear, pixelJoy, pixelSadness,
299   pixelSurprise)
300 percentages = c(paste(round(pixelAnger/total*100,2),"%"),paste(round(
301   pixelDisgust/total*100,2),"%"),paste(round(pixelFear/total*100,2),"%"),
302   paste(round(pixelJoy/total*100,2),"%"),paste(round(pixelSadness/total
303   *100,2),"%"),paste(round(pixelSurprise/total*100,2),"%"))
304
305 pie3D(emotionsPixel, labels = percentages, col = brewer.pal(6,'Set1'), main =
306   "Distribucion de Emociones en el Pixel 2 (%)", labelcex = 1.1)
307 legend("bottomleft", legend = emotions, fill = brewer.pal(6,'Set1'), cex =
308   .5)
309
310 #-----Emociones por todos los dispositivos-----#
311 mixPos = sum(length(which(iphonePolarity.class == "positive")),length(which(
312   notePolarity.class == "positive")))
313 mixNeg = sum(length(which(iphonePolarity.class == "negative")),length(which(
314   notePolarity.class == "negative")))
315 mixNeutral = sum(length(which(iphonePolarity.class == "neutral")),length(which(
316   notePolarity.class == "neutral")))
317 devices = c(mixPos,mixNeg,mixNeutral)
318 total = sum(devices)
319 percentages = c(paste(round(mixPos*100/total,2),"%"), paste(round(mixNeg*100/
320   total,2),"%"),paste(round(mixNeutral*100/total,2),"%"))
321
322 pie3D(devices, labels = percentages, col = brewer.pal(3,'Dark2'), main = "
323   Distribucion de Emociones en Todos los Dispositivos (%)", labelcex =
324   1.1)
325 legend("bottomleft",legend = c("Positive","Negative","Neutral"), fill =
326   brewer.pal(3,'Dark2'), cex = .8)
327
328 #-----Puntuaciones de tweets por todos los dispositivos-----#
329 getAccumulated = function(list){
330   accum = c()
331   sum = 0
332   for(i in 1:length(list)){
333     if(!is.na(list[i])){
334       sum = sum+as.numeric(list[i])
335       accum = c(accum,sum)
336     }
337   }
338   return(accum)
339 }
340
341 iphoneAccum = getAccumulated(iphone.scores$score)
342 noteAccum = getAccumulated(note.scores$score)
343 pixelAccum = getAccumulated(pixel.scores$score)
344
345 plot(iphoneAccum, type = "l", ylim = c(min(iphoneAccum,noteAccum,pixelAccum),
346   max(iphoneAccum,noteAccum,pixelAccum)), col = brewer.pal(8,'Set1')[3],
347   xlab = "# de Tweets", ylab = "Puntuacion Acumulada", main = "
348   Puntuacion Historica Acumulada por Dispositivo")
349 lines(noteAccum, type = "l", col = brewer.pal(8,'Set1')[4])
350 lines(pixelAccum, type = "l", col = brewer.pal(8,'Set1')[1])
351 legend(0,1850,c("iPhone X","Note 8","Pixel 2"),lty=c(1),lwd=c(2.5),col = c(
352   brewer.pal(8,'Set1')[3],brewer.pal(8,'Set1')[4],brewer.pal(8,'Set1')
353   [1]))
354
355 #-----Nube de palabras por emociones-----#
356 getTdm = function(tweets,tweetsEmotions){
357   wcEmotions = levels(factor(emotions))
358   nEmotions = length(wcEmotions)
359   wcEmotions.docs = rep("", nEmotions)
360
361   for(i in 1:nEmotions)
362   {
363     tmp = tweets$text[tweetsEmotions == wcEmotions[i]]
364     wcEmotions.docs[i] = paste(tmp, collapse=' ')
365   }
366
367   corpus = Corpus(VectorSource(wcEmotions.docs))
368   tdm = TermDocumentMatrix(corpus)
369   tdm = as.matrix(tdm)
370   colnames(tdm) = wcEmotions
371   return(tdm)
372 }
373
374 iphoneCorpus = getTdm(iphoneCleanedTweets,iphoneEmotions)
375 noteCorpus = getTdm(noteCleanedTweets,noteEmotions)
376 pixelCorpus = getTdm(pixelCleanedTweets,pixelEmotions)
377
378 mixedTweets = rbind(iphoneCleanedTweets,noteCleanedTweets,pixelCleanedTweets)
379 mixedEmotions = rbind(iphoneEmotions,noteEmotions,pixelEmotions)
380 mixedCorpus = getTdm(mixedTweets,mixedEmotions)
381
382 comparison.cloud(iphoneCorpus, random.order = FALSE, max.words = 1000, rot.
383   per=.15, colors = brewer.pal(emotions, 'Dark2'), scale = c(4,0.5), title.
384   size = 1)
385 title(main = "Nube de Palabras de Emociones del iPhone X", outer = TRUE, line
386   = -0.7)
387
388 comparison.cloud(noteCorpus, random.order = FALSE, max.words = 1000, rot.per
389   =.15, colors = brewer.pal(emotions, 'Dark2'), scale = c(4,0.5), title.
390   size = 1)
391 title(main = "Nube de Palabras de Emociones del Note 8", outer = TRUE, line =
392   -0.7)

```

```

368 comparison.cloud(pixelCorpus, random.order = FALSE, max.words = 1000, rot.per
369   =.15, colors = brewer.pal(emotions, 'Dark2'), scale = c(4,0.5), title.
370   size = 1)
371 title(main = "Nube de Palabras de Emociones del Pixel 2", outer = TRUE, line
372   = -0.7)
373
374 comparison.cloud(mixedCorpus, random.order = FALSE, max.words = 1000, rot.per
375   =.15, colors = brewer.pal(emotions, 'Dark2'), scale = c(4,0.5), title.
376   size = 1)
377 title(main = "Nube de Palabras de Emociones de Todos los Dispositivos", outer
378   = TRUE, line = -0.7)
379
380 #-----Frecuencia de palabras-----#
381 termsFreq = rowSums(as.matrix(mixedCorpus))
382 mostFreq = subset(termsFreq,termsFreq>=50)
383 df = data.frame(term = names(mostFreq), freq=mostFreq)
384
385 par(mai=c(1,1.5,1,1))
386 barplot(mostFreq, xlim=c(0,round(max(mostFreq)/100,0)*100), col = c(brewer.
387   pal(name = 'Dark2',n = 8),brewer.pal(name = 'Paired',n = 12)), horiz =
388   TRUE, las=1, cex.names = 0.5, space = c(0.5,0))
389
390 par(mar = c(5,7,4,2) + 0.1)
391 title(main = "Palabras Mas Frecuentes por Todos los Dispositivos", ylab = "
392   Palabras", xlab = "Frecuencia")

```

## 11. ENLACE DEL VIDEO DE RESULTADOS

En el siguiente enlace se encuentra el video de demostración del código de implementación y resultados parciales del proyecto (video). Además, se incluye el enlace al repositorio del proyecto para mayor información. (repositorio).

## REFERENCIAS

- [1] ASP.NET Web API, Msdn.microsoft.com, 2017. [En línea]. Disponible en: [https://msdn.microsoft.com/es-es/library/hh833994\(v=vs.108\).aspx](https://msdn.microsoft.com/es-es/library/hh833994(v=vs.108).aspx).
- [2] Business Intelligence, Sinnexus, 2016. [En línea]. Disponible en: [http://www.sinnexus.com/business\\_intelligence/datamining.aspx](http://www.sinnexus.com/business_intelligence/datamining.aspx)
- [3] Entendiendo el análisis de sentimiento: qué es y para qué se usa, brandwatch, 2015. [En línea]. Disponible en: <https://www.brandwatch.com/es/2015/02/analisis-de-sentimiento/>
- [4] H. Fujita, M. Ali, A. Selamat, J. Sasaki and M. Kurematsu, Trends in applied knowledge-based systems and data science. Springer, 2016, pp. 269-279.
- [5] IBM Knowledge Center, Ibm.com, 2017. [En línea]. Disponible en : [https://www.ibm.com/support/knowledgecenter/es/SSEPGG\\_9.7.0/com.ibm.im.overview.doc/c\\_naive\\_bayes\\_classification.html](https://www.ibm.com/support/knowledgecenter/es/SSEPGG_9.7.0/com.ibm.im.overview.doc/c_naive_bayes_classification.html).
- [6] Liu, B., Hu, M. and Cheng, J. Opinion Observer: Analyzing and Comparing Opinions on the Web. Chiba, 2005 [En línea]. Disponible en: <https://www.cs.uic.edu/liub/publications/www05-p536.pdf>
- [7] Minería de Opiniones, BrainSINS, 2016.[En línea]. Disponible en: [https://www.brainsins.com/es/blog/mineria-opiniones/3555/modeler\\_mainhelp\\_client\\_ddita/clementine/svm\\_howwork.html](https://www.brainsins.com/es/blog/mineria-opiniones/3555/modeler_mainhelp_client_ddita/clementine/svm_howwork.html)
- [8] Procesamiento del lenguaje natural, Vicomtech.org, 2017. [En línea]. Disponible: <http://www.vicomtech.org/t4/e11/procesamiento-del-lenguaje-natural>.
- [9] Y. Zhao, R. and Data Mining: Examples and Case Studies, 1st ed. Elsevier, 2015.