





INFOTEC CENTRO DE INVESTIGACIÓN E  
INNOVACIÓN EN TECNOLOGÍAS DE LA  
INFORMACIÓN Y COMUNICACIÓN

DIRECCIÓN ADJUNTA DE INNOVACIÓN Y  
CONOCIMIENTO  
GERENCIA DE CAPITAL HUMANO  
POSGRADOS

# “CUANTIFICACIÓN Y PREDICCIÓN DE POSIBLES CASOS DE DEPRESIÓN EN USUARIOS DE TWITTER UTILIZANDO MACHINE LEARNING”

REPORTE ANALÍTICO DE EXPERIENCIA LABORAL  
Que para obtener el grado de MAESTRO EN  
CIENCIA DE DATOS E INFORMACIÓN

Presenta:

**Enrique Isidoro Vázquez Ramos**

Asesores:

**Dr. Sabino Miranda Jiménez**  
**Dr. Eric Sadit Téllez Avila**

Aguascalientes, noviembre, 2021.

# Autorización de Impresión



## AUTORIZACIÓN DE IMPRESIÓN Y NO ADEUDO EN BIBLIOTECA

### Maestría en Ciencia de Datos e Información, MCDI.

Ciudad de México, 1 de junio de 2022.

La Gerencia de Capital Humano / Gerencia de Investigación hacen constar que el trabajo de titulación intitulado:

**"Cuantificación y predicción de posibles casos de depresión en usuarios de Twitter utilizando machine learning"**,

Desarrollado por el alumno: **Enrique Isidoro Vázquez Ramos** y bajo la asesoría del **Dr. Sabino Miranda Jiménez** y el **Dr. Eric Sadit Téllez Avila**; cumple con el formato de Biblioteca. Por lo cual, se expide la presente autorización para impresión del proyecto terminal al que se ha hecho mención.

Asimismo, se hace constar que no debe material de la biblioteca de INFOTEC.

Vo. Bo.

**Lic. Juan Ramón Abarca Damián**  
Coordinador de Biblioteca

**Anexar a la presente autorización al inicio de la versión impresa del trabajo referido que ampara la misma.**

## Agradecimientos

Agradezco a Dios,  
sin él nada sería posible.

Al amor de mi vida y esposa Izadora,  
que siempre me apoya en los momentos más difíciles.

A mi papá, mamá y hermanas y familia en general,  
espero que estén orgullosos.

A la sociedad en general,  
esperando que en algún momento esta tesis les sea de ayuda.

Agradezco a mi amigo Víctor Silva Cuevas,  
compañero de la maestría por su paciencia y ayuda.

También a todos mis profesores  
por su gran trabajo y por su tiempo.

## Tabla de contenido

Introducción .....	1
Capítulo 1 - Planteamiento.....	1
1.1 Motivación .....	1
1.2 Problemática .....	2
1.3 Objetivos .....	3
1.3.1 Objetivo General.....	3
1.3.2  Objetivos Específicos.....	3
1.4 Contribución .....	3
Capítulo 2 - Marco Teórico.....	5
2.1 Aprendizaje Supervisado.....	5
2.2 Partición de los Datos y Métricas .....	6
2.2.1 Entrenamiento - Prueba .....	6
2.2.2 Métricas de desempeño .....	7
2.3 Algoritmos de Clasificación.....	8
2.3.1 Soporte Lineal para la Clasificación de Vectores .....	9
2.3.2 Clasificador de Vecinos más Cercanos .....	9
2.3.3 Tf-Idf.....	10
2.3.4 Clasificador Naive Bayes para modelos multinomiales .....	11
2.3.5 MicroTC - TextModel.....	12
2.4 Estado del Arte.....	12
2.4.1 Aprendizaje Automático .....	12
2.4.2 Minería de Texto.....	14
2.4.3 Clasificación de pacientes con depresión usando machine learning .....	15
2.4.4 Análisis de Texto .....	16

2.5 Representaciones semánticas / Vectores de palabras (word embeddings) .	17
2.5.1 FastText .....	19
2.5.2 InferSent .....	20
Capítulo 3 – Metodología .....	22
3.1 Descripción General .....	22
3.2 Conjunto de Datos .....	22
3.3 Preprocesamiento .....	23
3.4 Procedimiento .....	23
3.5 Ambiente de Desarrollo .....	24
3.5.1 Servidor utilizado para correr el código .....	24
3.5.2 Tecnologías y librerías utilizadas .....	25
Capítulo 4 - Análisis y Evaluación.....	27
4.1 Estructura de los Datos por Usuario .....	27
4.2 Estructura de los Datos por Tweet con y sin preprocesamiento .....	28
4.3 Experimentos y Resultados .....	29
4.3.1 Métricas por usuarios con preprocesamiento.....	30
4.3.2 Métricas por Usuarios sin preprocesamiento .....	31
4.3.3 Métricas por Tweet con preprocesamiento .....	32
4.3.4 Métricas por Tweet sin preprocesamiento.....	33
4.4 Discusión de los resultados obtenidos.....	33
4.5 MicroTC Parámetros:.....	34
Capítulo 5 – Conclusiones y Trabajo Futuro .....	38
5.1 Conclusiones.....	38
5.2 Trabajo Futuro.....	40

Bibliografía .....	41
ANEXO I: Repositorio del código de Python utilizado en este proyecto .....	46
ANEXO II: Códigos individuales de Python utilizados en este proyecto .....	47

## Índice de figuras

Figura 1. Ejemplo bidimensional de un problema de clasificación. ....	5
Figura 2. Ejemplo bidimensional de una representación vectorial .....	18
Figura 3. Arquitectura general del proyecto .....	25



## Índice de cuadros

Cuadro 1. Distribución de las clases. ....	27
Cuadro 2. Distribución de las clases con preprocesamiento .....	28
Cuadro 3. Distribución de las clases sin preprocesamiento .....	28
Cuadro 4. Desempeños por usuario con preprocesamiento .....	30
Cuadro 5. Desempeños por usuario sin preprocesamiento.....	31
Cuadro 6. Desempeños por tweet con preprocesamiento.....	32
Cuadro 7. Desempeños por tweet sin preprocesamiento.....	33
Cuadro 8. Enfoques y Modelos.....	34
Cuadro 9. Mejores Resultados.....	34
Cuadro 10. Parámetros MicroTc.....	35

## Glosario

### “A”

**Análisis sintáctico:** Es el análisis de las funciones sintácticas o relaciones de concordancia y jerarquía que guardan las palabras cuando se agrupan entre sí en forma de sintagmas, oraciones simples y oraciones compuestas de proposiciones.

### “L”

**Lematización:** Es un proceso lingüístico que consiste en, dada una forma flexionada, hallar el lema correspondiente. El lema es la forma que por convenio se acepta como representante de todas las formas flexionadas de una misma palabra.

### “M”

**Machine learning:** Es el subcampo de las ciencias de la computación y una rama de la inteligencia artificial, cuyo objetivo es desarrollar técnicas que permitan que las computadoras aprendan.

### “N”

**Normalización:** Segmentación de las palabras, formato de las palabras y segmentación de las oraciones en el texto.

### “T”

**Tokenización:** Una instancia de un tipo en un texto dado.

# Capítulo 1 - Planteamiento

The background of the page is a light gray technical drawing. It features several large gears on the left side, with various lines, dashed lines, and arrows indicating mechanical components and movement. The drawing is partially obscured by the chapter title.

## Introducción

En el presente trabajo de investigación se abordarán conceptos enfocados a la creación y aplicación de un modelo de machine learning<sup>1</sup> alimentado por un corpus de tweets que dará como resultado una herramienta para poder predecir posibles casos de depresión en personas. Durante el transcurso del proyecto se estarán detallando las herramientas, teorías, análisis, procesos y resultados llevados a cabo para la realización de dicho modelo. En resumen, el trabajo se divide en 5 capítulos, en el primero se explica el planteamiento en general, la motivación y problemática, los objetivos y la contribución que se aportará al lector.

El segundo capítulo muestra los conceptos y definiciones básicas que se deben tener para la construcción del modelo de machine learning, aquí mismo se muestran algunos estudios previos que diversos investigadores han realizado similares al que se aplica en esta investigación. En el tercer capítulo se habla acerca de la metodología que se aplicará, haciendo énfasis en los datos, herramientas y procedimiento.

El cuarto capítulo define la estructura de los datos y da a conocer los resultados obtenidos de los experimentos realizados. Finalmente, en el quinto y último capítulo se podrán leer las conclusiones del proyecto, así como las áreas de oportunidad que surgieron al completarlo.

## Capítulo 1 - Planteamiento

### 1.1 Motivación

El problema social referente a que las personas sufran depresión es un tema preocupante en nuestra sociedad actual, la cifra de personas que la padecen crece cada vez más. La vida acelerada, competitiva y estresante (entre más factores) que se vive hoy en día ha dado paso a que muchas personas caigan en depresión. La gente

---

<sup>1</sup> Un modelo de machine learning es la salida de información que se genera cuando entrena su algoritmo de machine learning con datos. Después del entrenamiento, al proporcionar un modelo con una entrada, se le dará una salida (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009).

encuentra diversas válvulas de escape para desahogarse, entre ellas están las redes sociales, y dentro de las redes sociales más populares se encuentra Twitter. Es aquí donde el presente proyecto entra, haciendo uso de herramientas de inteligencia artificial y de un corpus (textos de Twitter) clasificado por expertos donde se identifican personas con depresión y sin depresión, se desarrollará una propuesta para predecir posibles casos de depresión en usuarios de Twitter localizados en México.

## **1.2 Problemática**

Los gobiernos actuales, así como diferentes organizaciones alrededor del mundo están desarrollando propuestas. Un ejemplo de lo anterior se puede consultar en (Hosseini, Hassan y Rostami, 2013) que habla sobre salud mental, para poder prevenir y/o detectar casos de depresión en las personas.

El tema de la depresión está muy ligado al del suicidio (Kulasinghe, Jayasinghe, Rathnayaka, Karunarathne, Suranjini Silva y Anuradha Jayakodi, 2019) las personas que atentan contra su vida generalmente sufren algún tipo de depresión y/o problemas que creen no tienen solución. Es posible que no cuenten con personas cercanas a ellas y se sientan solos, tratando de desahogarse haciéndose notar en sus perfiles de redes sociales.

Dado el enorme abanico de herramientas y tecnologías que se tiene actualmente, como son las redes sociales (Twitter, Facebook, etc.) y la inteligencia artificial (machine learning), muchas propuestas para apoyar en el estudio de este problema se han desarrollado. Los resultados han sido variados (Joulin, Grave, Bojanowski y Mikolov, 2017). (Conneau, Kiela, Schwenk, Barrault y Bordes, 2017) como desarrollo de léxicos y predicción de desempeño. Se utilizarán técnicas de machine learning y Procesamiento de Lenguaje Natural (PLN) sobre datos de redes sociales para detectar a usuarios con posibles casos de depresión.

## **1.3 Objetivos**

### **1.3.1 Objetivo General**

Proponer un modelo de clasificación automática<sup>2</sup> basado en aprendizaje computacional que sea capaz de detectar posibles casos de depresión entre los usuarios de Twitter. En el caso particular de este proyecto, el enfoque será trabajar sólo en lo que se refiere a México.

### **1.3.2 Objetivos Específicos**

- Obtención y caracterización de la base de conocimiento.
- Creación y prueba de modelos.
- Selección del modelo que se desempeñe mejor para su aplicación en un conjunto de datos grande de tweets geolocalizados en México.
- Análisis de los resultados.

## **1.4 Contribución**

Este proyecto aportará al lector una herramienta para poder pronosticar posibles casos de depresión en las personas, lo anterior se logra haciendo uso de un algoritmo de machine learning. El algoritmo cuenta con varias implementaciones las cuales se explicarán más adelante en esta investigación. Los resultados, así como el algoritmo generado en este trabajo estarán disponibles a través de un repositorio de Github.

---

<sup>2</sup> Consiste en utilizar técnicas de Inteligencia Artificial sobre un conjunto de elementos para ordenarlos por clases o categorías (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009).

# Capítulo 2 - Marco Teórico



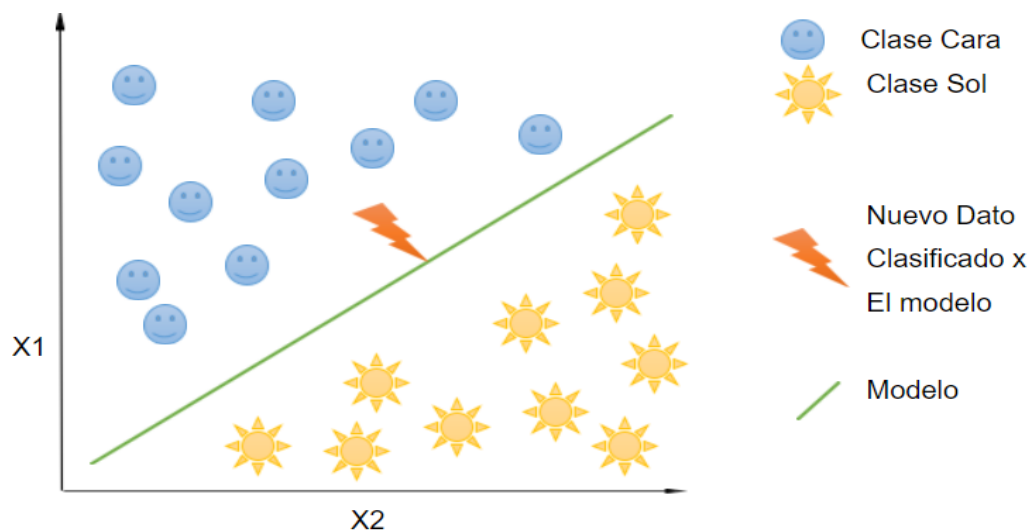
## Capítulo 2 - Marco Teórico

A continuación, se abordarán los conceptos de las herramientas, librerías, y tecnologías usadas en este proyecto de investigación.

### 2.1 Aprendizaje Supervisado

Esta técnica es una rama de la inteligencia artificial usada para entrenar o crear modelos a partir de datos numéricos o textos que han sido previamente etiquetados, estas etiquetas son realizadas manualmente y por lo general son nombradas como las clases del problema, el aprendizaje supervisado cuenta con una gran variedad de aplicaciones como el uso de las redes sociales para monitorear la salud mental (McClellan, Ali, Mutter, Kroutil y Landwehr, 2016) por mencionar alguna.

En (Román, 2019) se expone que un problema de clasificación<sup>3</sup> es una subcategoría del aprendizaje supervisado en la que el objetivo es predecir las etiquetas de clase categóricas de las nuevas instancias, basándonos en observaciones pasadas.



**Figura 1.** Ejemplo bidimensional de un problema de clasificación.

Fuente: Elaboración propia.

<sup>3</sup> Se caracterizan por tener una variable cualitativa 'Y' como respuesta.



Las clases identificarán a cada individuo de los datos, por ejemplo, supongamos que los datos son extraídos de un corpus de tweets que cuenta con textos escritos por usuarios seleccionados al azar, como es de esperarse habrá gran variedad de combinaciones de unidades léxicas con diversos significados, una de estas combinaciones podría ser el texto “¿jestoy muy feliiiiiz!? #2021” que, siguiendo la definición de aprendizaje supervisado se clasificaría manualmente con la clase “Positivo”. Entre más grande sea el corpus mejor serán los resultados.

La finalidad de esta técnica es, que una vez entrenado los datos con algoritmos de machine learning, se podrán clasificar nuevos datos sin la necesidad de hacerlo manualmente.

## **2.2 Partición de los Datos y Métricas**

Dentro del aprendizaje supervisado se requiere de la realización de pruebas a los modelos creados para de esta forma medir su eficiencia, en este proyecto se aplicará la estrategia de particionar los datos en dos<sup>4</sup>, con el primero se entrenarán los modelos de machine learning y con el segundo se realizarán pruebas para generar las métricas de desempeño (véase 2.2.2).

### **2.2.1 Entrenamiento - Prueba**

El conjunto de datos manejado en este proyecto proviene del INAOE, fue en este instituto donde se etiquetó a los usuarios de Twitter. La forma en que lo hicieron fue la siguiente: los usuarios etiquetados con depresión auto declararon esta condición y los usuarios etiquetados sin depresión son cuentas genéricas. Cabe señalar que en este proceso no intervinieron especialistas en el campo de la salud mental<sup>5</sup>. Los datos se encuentran en formato json mostrando la siguiente estructura (cada línea del archivo json contiene 3 campos: user, klass y tweets):

- User. Hay 317 usuarios (se debe tener en cuenta que las clases no están balanceadas, véase 4.1), 92 considerados con depresión y 225 sin depresión.

---

<sup>4</sup> Conjunto de entrenamiento y conjunto de prueba.

<sup>5</sup> Psiquiatras y psicólogos.

- **Klass.** El número 1 representa a los considerados con depresión y el número 0 a los considerados sin depresión.
- **Tweets.** El campo tweets es una lista por usuario de tamaño no constante donde cada elemento es un tweet que contiene todos los datos de éste: fecha, id, texto, fecha de creación, etc. En este proyecto se utilizará sólo la variable texto.

La partición de los datos tanto por usuarios como por tweets se realizará dividiendo la información (arreglos) en subconjuntos de pruebas y trenes aleatorios, quedando los porcentajes de la siguiente manera:

- **Entrenamiento.** Tendrá el 70% de los datos.
- **Prueba.** Tendrá el 30 %.

### 2.2.2 Métricas de desempeño

Las métricas que se utilizarán para calificar los resultados serán las siguientes:

- **Precision.** La precisión es la relación entre las observaciones positivas predichas correctamente<sup>6</sup> y el total de observaciones positivas predichas<sup>7</sup> (Renuka, 2016).

$$precision = \frac{TP_{(k)}}{TP_k + FP_k}$$

- **Recall.** Es la proporción de observaciones positivas predichas correctamente a todas las observaciones en la clase real (Renuka, 2016).

$$recall = \frac{TP_{(k)}}{TP_k + FN_k}$$

---

<sup>6</sup> Son los aciertos que el modelo obtuvo.

<sup>7</sup> Es el total de todas las predicciones.

- *F1 Score Measure*. Es el promedio ponderado de Precision y Recall (Renuka, 2016).

$$F1 = \frac{2 * \text{precisión}_{(k)} * \text{recall}_{(k)}}{\text{precisión}_{(k)} + \text{recall}_{(k)}}$$

$$\text{micro } F1 = \frac{2 * \text{precisión} * \text{recall}}{\text{precisión} + \text{recall}}$$

$$(\text{macro o ponderado}) F1 = \sum_{K \in \Gamma} W_k * F1_k$$

$$W_k = \left\{ \begin{array}{l} \text{macro} \rightarrow 1/\Gamma \\ \text{ponderado} \rightarrow \frac{|\{u \in D | u=k\}|}{|D|} \end{array} \right\}$$

- *Accuracy*. Es simplemente una relación entre la observación predicha correctamente y el total de ellas (Renuka, 2016).

$$\text{accuracy} = \frac{TP+TN}{\#muestras}$$

Donde:

TP = número de verdaderos positivos.

TN = verdaderos negativos.

FP = número de falsos positivos.

FN = número de falsos negativos.

k = clase.

## 2.3 Algoritmos de Clasificación

Como se mencionó anteriormente la estrategia en este proyecto será utilizar algoritmos de clasificación supervisada (véase 2.1), pues se conoce la cantidad y el tipo de las clases asociadas a los textos que se trabajarán. En los puntos siguientes se

definen los algoritmos utilizados en este proyecto que serán requeridos para la clasificación automática de los tweets.

### 2.3.1 Soporte Lineal para la Clasificación de Vectores

LinearSVC (Linear Support Vector Classification) (Pedregosa, 2011) es una implementación de clasificación de vectores de soporte para el caso de un núcleo lineal<sup>8</sup>.

Como otros clasificadores, LinearSVC toma como entrada dos matrices<sup>9</sup>: una matriz 'X' de forma (n\_samples, n\_features) que contiene las muestras de entrenamiento<sup>10</sup> y una matriz 'Y' de etiquetas de clase<sup>11</sup> (cadenas o enteros), de forma (n\_samples).

### 2.3.2 Clasificador de Vecinos más Cercanos

La clasificación basada en los vecinos más cercanos es un tipo de aprendizaje basado en instancias o aprendizaje no generalizador<sup>12</sup>. La clasificación se calcula a partir de un voto de mayoría simple de los vecinos más cercanos de cada punto: a un punto de consulta se le asigna la clase de datos que tiene más representantes dentro de los vecinos más cercanos del punto (Pedregosa, 2011).

El nombre también se refiere a la clase de algoritmos llamados "no paramétricos". Estos son algoritmos basados en instancia, que se caracterizan por memorizar el conjunto de datos de entrenamiento, y el aprendizaje perezoso es un caso particular de estos algoritmos, asociados con coste computacional cero durante el aprendizaje (Román, 2019). El proceso que sigue el algoritmo es (Román, 2019):

- Escoge el número de k y la distancia.
- Encuentra el k vecino más cercano de la muestra que se pretende clasificar.
- Asigna la etiqueta de clase por votación mayoritaria.

---

<sup>8</sup> Es el núcleo de la aplicación lineal definida por la matriz (Pedregosa, 2011).

<sup>9</sup> Tablas bidimensionales de números en cantidades abstractas que pueden sumarse y multiplicarse (Pedregosa, 2011).

<sup>10</sup> Datos que serán utilizados para el entrenamiento (Pedregosa, 2011).

<sup>11</sup> Etiqueta que se asigna a cada uno de los datos (Pedregosa, 2011).

<sup>12</sup> No intenta construir un modelo interno general, sino que simplemente almacena instancias de los datos de capacitación (Pedregosa, 2011).

El algoritmo encuentra las  $k^{13}$  muestras que son más cercanas al punto que se quiere clasificar, basando sus predicciones en la distancia métrica (Román, 2019). La principal ventaja es que se adapta a los nuevos datos de entrenamiento, al ser un algoritmo basado en la memoria. La desventaja es que el coste computacional se incrementa linealmente con el tamaño de los datos de entrenamiento (Román, 2019).

### 2.3.3 Tf-Idf

Esta es una estrategia específica de tokenización, conteo y normalización, a la que se llama la bolsa de palabras o la representación de "bolsa de n-gramas". Los documentos se describen por las ocurrencias de las palabras mientras se ignora por completo la información de posición relativa de las palabras en el documento (Pedregosa, 2011).

En un corpus<sup>14</sup> de texto grande, algunas palabras estarán muy presentes, por lo tanto, llevarán muy poca información significativa sobre el contenido real del documento. Si tuviéramos que alimentar los datos de conteo directo directamente a un clasificador, esos términos muy frecuentes sombrearían las frecuencias de términos más raros aún más interesantes. Tf significa frecuencia de término, mientras que Tf-Idf significa frecuencia de término por frecuencia inversa de documento.

Originalmente, este era un esquema de ponderación de términos desarrollado para la recuperación de información que también ha encontrado un buen uso en la clasificación y agrupación de documentos (Pedregosa, 2011).

$$tfidf(t, d, D) = tf(t, d) * idf(t, D)$$

$$tf(t, d) = \log \log (1 + freq(t, d))$$

$$idf(t, D) = \log \left( \frac{N}{|(d \in D: t \in d)|} \right)$$

---

<sup>13</sup> Variable que representa a cualquier valor posible.

<sup>14</sup> Es el total de los datos.

Donde:

t = término

d = documento

N = número total de documentos en el corpus

$d \in D$ :  $t \in d$  = número de documentos en los que aparece el término "t"

### 2.3.4 Clasificador Naive Bayes para modelos multinomiales

MultinomialNB implementa el algoritmo Bayes ingenuo para datos distribuidos multinomialmente, y es una de las dos variantes Bayes ingenuas clásicas utilizadas en la clasificación de texto. La distribución está parametrizada por vectores (Pedregosa, 2011).

En este algoritmo se asume que las variables predictoras<sup>15</sup> son independientes entre sí. En otras palabras, que la presencia de una cierta característica en un conjunto de datos no está en absoluto relacionada con la presencia de cualquier otra característica (Román, 2019).

Lo consigue proporcionando una forma de calcular la probabilidad<sup>16</sup> 'posterior' de que ocurra un cierto evento  $A$ , dadas algunas probabilidades de eventos 'anteriores' (Román, 2019). Los pasos generales para un algoritmo Naive Bayes supervisado son: (Román, 2019):

- Convertir el conjunto de datos en una tabla de frecuencias<sup>17</sup>.
- Crear una tabla de probabilidad calculando las correspondientes a que ocurran los diversos eventos.
- La ecuación Naive Bayes se usa para calcular la probabilidad posterior de cada clase.
- La clase con la probabilidad posterior más alta es el resultado de la predicción.

---

<sup>15</sup> Son las variables de interés en un experimento (Román, 2019).

<sup>16</sup> Es una medida del grado de certidumbre de que un suceso pueda ocurrir (Román, 2019).

<sup>17</sup> Agrupación de datos en categorías mutuamente excluyentes que indican el número de observaciones en cada categoría (Román, 2019).

### 2.3.5 MicroTC - TextModel

MicroTC (Tellez, Moctezuma, Miranda y Graff, 2018) es intencionalmente simple, por lo que solo se implementó una pequeña cantidad de funciones. Sin embargo, utiliza algunas herramientas complejas de numpy y scikit-learn. El número de dependencias está limitado y cumplido por casi cualquier distribución de Scientific Python, por ejemplo, Anaconda.

Está diseñado para abordar los problemas de clasificación de texto de forma agnóstica<sup>18</sup>, ya que es independiente del dominio y del lenguaje.

El punto de entrada del núcleo de  $\mu$ TC es `microtc.textmodel.TextModel`, que puede verse como una función con la forma  $\mathbf{m}(\text{texto}) \rightarrow \mathbf{R}^d$  donde ' $d$ ' es el vocabulario, es decir, la dimensión del espacio vectorial. Como puede verse, ' $m$ ' puede usarse para transformar un texto en un vector y, en consecuencia, puede usarse para transformar un conjunto de entrenamiento de pares, texto y etiqueta, en un conjunto de entrenamiento de pares, vectores y etiqueta, que puede ser utilizado directamente por cualquier algoritmo de aprendizaje supervisado para obtener un clasificador de texto (Tellez, Moctezuma, Miranda y Graff, 2018).

## 2.4 Estado del Arte

Es importante resaltar que una de las principales etapas que debe desarrollarse en un proyecto de titulación es la realización de su Estado del Arte. Ya que no sólo muestra un panorama de cómo se ha trabajado anteriormente el tema de este proyecto también ofrece información sobre los avances y las estrategias que se esperan en un futuro.

### 2.4.1 Aprendizaje Automático

Existe una gran cantidad de trabajos enfocados en el aprendizaje automático (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009) hablan sobre como aprender de los datos, de cómo en un escenario típico, se tiene una medida de resultado, generalmente cuantitativa o categórica, que se desea predecir basado en un conjunto de

---

<sup>18</sup> Considera que la veracidad de ciertas afirmaciones son desconocidas.

características. Tienen un conjunto de datos de entrenamiento, en el que observaron el resultado y las mediciones de características para un conjunto de objetos. Usando esos datos construyeron un modelo de predicción, que permite predecir el resultado para nuevos objetos.

En (Pratama y Purwarianti, 2017) trabajan con técnicas de aprendizaje automático en sus dos variantes: supervisado y no supervisado. Ven al público de Bandung (ciudad indonesia) como un usuario activo de las redes sociales, en especial de Twitter, para que los ciudadanos denuncien sus quejas. Para gestionar las quejas denunciadas hacen uso del aprendizaje automático y así detectar los temas de los tweets. El enfoque de aprendizaje supervisado lo implementaron para clasificar el tema de los tweets en función de las quejas, mientras que el enfoque de aprendizaje no supervisado se utiliza para agrupar los tweets de quejas en función de la similitud de la información detallada contenida en las quejas. Su enfoque de aprendizaje supervisado lo evaluaron utilizando accuracy, precisión, recall y medida F1 (véase 2.2.2). Con su enfoque de aprendizaje no supervisado, utilizaron el valor del índice de agrupación para evaluar los grupos de temas detectados.

En (Larsen, Boonstra, Batterham, O’Dea, Paris y Christensen, 2015) describen que la mayoría de los datos de investigación sobre problemas de salud mental se recopilan a través de encuestas, esto trae como consecuencia que la información recolectada no proporciona una visión en tiempo real<sup>19</sup> de los individuos encuestados. Mencionan que una solución para obtener información en tiempo real es extraer ésta desde las redes sociales, pues los individuos participan en ellas continuamente. Describieron un sistema que llamaron “We feel” donde analizan variaciones globales y regionales de la expresión emocional. Su fuente de información fue Twitter, obteniendo así en un periodo de 12 semanas información en forma de textos y mediante el uso de PCA reconocer patrones que clasificaron como positivos y negativos.

En (Chopade, Edwards, Khan, Andrade y Pu, 2019). presentan un marco basado en el aprendizaje automático para la identificación de evidencias en las habilidades cognitivas y de aprendizaje socioemocional en la resolución colaborativa de

---

<sup>19</sup> Acontecimientos que están sucediendo en el momento presente.



problemas de la diádica. En otras palabras, demostrar las habilidades y la dinámica del trabajo en equipo como comportamientos verbales y no verbales y cómo éstos se pueden capturar y analizar mediante la recopilación pasiva de datos. Para lo anterior, analizaron datos de juegos multimodales y desarrollaron modelos de habilidades y técnicas para su medición, utilizaron técnicas de procesamiento de lenguaje natural (PLN) como la inserción de oraciones y la bolsa de palabras obteniendo resultados que muestran el desempeño de los jugadores clasificándolos como exitosos y no exitosos.

#### **2.4.2 Minería de Texto**

Dada la enorme cantidad de información<sup>20</sup> que actualmente se genera en el mundo, es necesario contar con técnicas avanzadas de minería de texto, así lo hacen notar (Hogenboom, Frasinca, Kaymak y Jong, 2011) al mencionar que la utilización de información extraída en los procesos de toma de decisiones<sup>21</sup> se convierte en un problema cada vez más urgente y difícil. Un problema omnipresente es el hecho de que la mayoría de los datos inicialmente no están estructurados y se describen usando un lenguaje natural, comprensible para el ser humano, que hace que los datos estén limitados en el grado en que son interpretados por las máquinas. Este problema frustra la automatización de, por ejemplo, los procesos de recuperación de información vital (IR) y extracción de información (IE), utilizados para la toma de decisiones.

En (Tekiner, Tsuruoka, Tsujii y Ananiadou, 2009) mencionan que las tareas de etiquetar datos de secuencia, la fragmentación y el reconocimiento de entidades nombradas son una de las tareas más importantes dentro de la minería de textos.

Cuando se trata de grandes cantidades de datos, Text Mining se refiere al aprendizaje de información a partir de texto preprocesado.

En (Agnihotri, Verma y Tripathi, 2014) analizan que, debido al uso intensivo de dispositivos electrónicos en la actualidad, la mayor parte de la información está

---

<sup>20</sup> Grandes cantidades de textos que se pueden recolectar de todas las redes sociales públicas.

<sup>21</sup> Conjunto de estrategias, aplicaciones, datos, productos, tecnologías y arquitectura técnicas, los cuales están enfocados a la administración y creación de conocimiento sobre el medio, a través del análisis de los datos existentes (Hogenboom, Frasinca, Kaymak y Jong, 2011).

almacenada en formato digital y gran parte de ella se encuentra en forma de texto. Mencionan que extraer el conocimiento, como la búsqueda de patrones o agrupaciones de palabras similares, es uno de los temas más importantes en la actualidad<sup>22</sup>. El documento se centra en extraer la información más importante de los textos, utilizaron el lenguaje de programación R para sus experimentos, su algoritmo utiliza la similitud coseno para medir la distancia entre palabras y de esa manera agruparlas usando k-medias.

Por medio de la minería de texto, a menudo utilizando el procesamiento del lenguaje natural, la información se extrae de textos de diversas fuentes, como noticias, mensajes, tweets o blogs, y se representa y almacena de forma estructurada.

#### **2.4.3 Clasificación de pacientes con depresión usando machine learning**

Existe una amplia variedad de aplicaciones<sup>23</sup> donde se puede utilizar machine learning, como es conocido, en el presente proyecto se aplicará esta técnica para detectar a usuarios de Twitter con posible depresión. Este mismo argumento también lo plantean (Hosseinifard, Hassan y Rostami, 2013) pues estudiaron el análisis no lineal de la señal EEG (electroencefalograma) para discriminar pacientes con depresión y controles normales. Cuarenta y cinco pacientes deprimidos no medicados y cuarenta y cinco sujetos normales participaron en este estudio.

Para discriminar los dos grupos<sup>24</sup>, utilizaron el vecino K más cercano, el análisis discriminante lineal y la regresión logística<sup>25</sup> como los clasificadores. Todas las características no lineales y el clasificador LR lograron una precisión de clasificación del 90%. Su estudio muestra que el análisis no lineal de EEG puede ser un método útil para discriminar pacientes deprimidos y sujetos normales.

En (Hooda, Saxena, Madhulika y Yadav, 2017) abordan el tema de que la depresión es un problema crítico que puede afectar a las personas de diversas formas. Su objetivo es predecir a aquellas personas que ni siquiera saben que sufren del padecimiento. Para lo anterior usaron tres modelos: uso de clasificadores de

---

<sup>22</sup> Los patrones y similitudes de palabras dan un contexto natural al problema a resolver.

<sup>23</sup> Seguridad de datos, comercio financiero, cuidado de la salud, marketing personalizado, etc.

<sup>24</sup> Pacientes con depresión y controles normales (Hosseinifard, Hassan y Rostami, 2013).

<sup>25</sup> Clasificadores utilizados por los investigadores (Hosseinifard, Hassan y Rostami, 2013).

aprendizaje automático y WEKA<sup>26</sup>, uso de métodos de imágenes y aprendizaje automático y uso de factores de riesgo.

#### **2.4.4 Análisis de Texto**

En (Hakdağlı, Özcan y Oğul, 2018) mencionan que, con la tecnología actual en desarrollo, el acceso de las personas a la información<sup>27</sup> y su producción ha alcanzado un nivel muy rápido. Esta información generada y obtenida se crea instantáneamente, se ingresa en sistemas de datos y se actualiza. Describen que estas fuentes de transmisión de datos se pueden transformar en valiosos resultados de análisis cuando se manejan con métodos específicos. Determinaron un campo de datos de texto para realizar análisis de datos generados instantáneamente y utilizaron Twitter, la plataforma más rica para datos de texto instantáneos. Twitter genera instantáneamente una variedad de datos en grandes cantidades y los presenta como código abierto. Utilizaron el entorno de análisis de transmisión de Apache Spark para analizar los recursos, el análisis de la situación lo realizaron utilizando la máquina de vectores de soporte, los árboles de decisión y los algoritmos de regresión logística de Apache Spark<sup>28</sup>. De esta manera, con la información generada instantáneamente y enfocando estos recursos y estrategias al análisis de sentimientos se pueden detectar posibles casos de depresión en usuarios de Twitter en tiempo real y de esta forma reducir el tiempo de diagnóstico.

En (Xiaomeng, Chong y Zipei, 2020) analizan que, con el desarrollo de la tecnología de comunicación móvil, el aumento de las redes sociales en línea brinda oportunidades para que las personas expresen sus puntos de vista<sup>29</sup>, lo que juega un papel importante en el análisis y monitoreo de la opinión pública. El tema principal de su estudio es el análisis emocional de los usuarios de Twitter en inglés. Aplican tecnología Reptile para obtener y limpiar tweets relacionados con China, luego de

---

<sup>26</sup> Plataforma de software para el aprendizaje automático y la minería de datos escrito en Java y desarrollado en la Universidad de Waikato (Hooda, Saxena, Madhulika y Yadav, 2017).

<sup>27</sup> Internet, medios de comunicación, globalización, etc.

<sup>28</sup> Motor ultrarrápido para el almacenamiento, procesamiento y análisis de grandes volúmenes de datos.

<sup>29</sup> Cualquier persona con acceso a internet puede expresar su opinión en cualquier red social.

etiquetar datos referidos al diccionario de palabras fuertes que se compone de frecuencias estadísticas de palabras, usan GloVe<sup>30</sup> para generar vectores de palabras y usan redes neuronales convolucionales<sup>31</sup> (CNN) para clasificar tweets automáticamente.

Así establecen un modelo de CNN preciso para analizar los tweets populares relacionados con China de 2014 a 2018. Sus resultados muestran que China ocupó una posición importante en la política mundial en 2018 con grandes comentarios negativos en Twitter. El número de simpatizantes mayor se dio en 2014, pero tiene un descenso repentino en 2015 y se mantiene en torno al 45% de 2016 a 2018. Al igual que en la anterior investigación (Hakdağlı, Özcan y Oğul, 2018), si se cambia el enfoque hacia el análisis de sentimientos enfocado a la depresión, éste estudio puede arrojar nuevos resultados y nuevas métricas para la posible detección de personas con depresión.

## **2.5 Representaciones semánticas / Vectores de palabras (word embeddings)**

El lenguaje natural<sup>32</sup> que usamos los humanos no puede ser procesado literalmente por las máquinas, dado que los textos trabajados en este proyecto están realizados por personas y escritos en lenguaje natural, se deben buscar estrategias que le faciliten a la máquina la manipulación e interpretación de éstos.

Word Embedding es una técnica utilizada para modelar textos escritos en lenguaje natural, su finalidad consiste en crear vectores numéricos que representen a los textos relacionándolos semánticamente. Es muy eficiente al momento de entrenar grandes cantidades de datos por tal motivo es utilizado en este proyecto.

Cada vector captura la información semántica<sup>33</sup> de la palabra o frase asociada según convenga y es visto en un plano multidimensional, es decir, la longitud

---

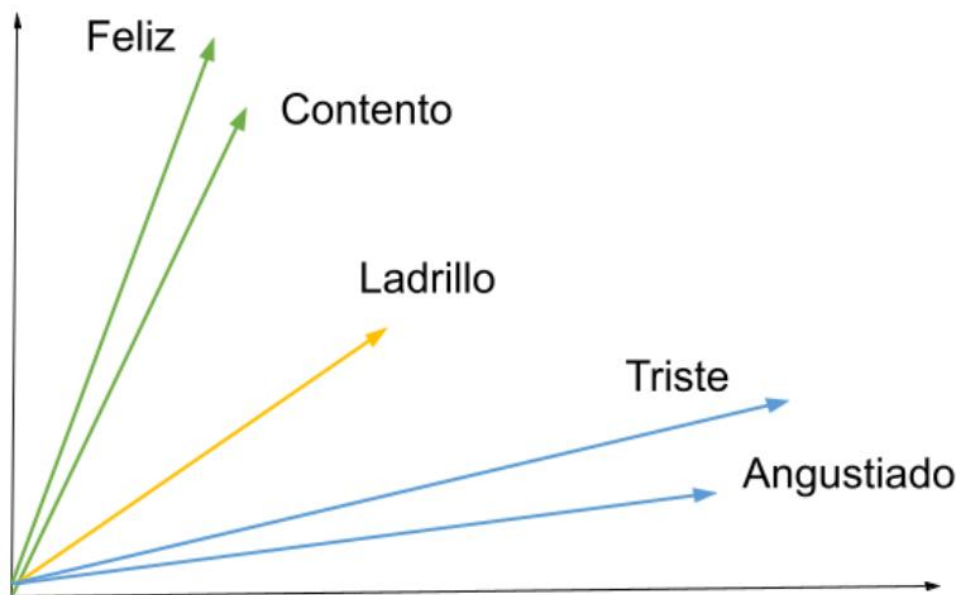
<sup>30</sup> Algoritmo de aprendizaje no supervisado para obtener representaciones vectoriales de palabras (Xiaomeng, Chong y Zipei, 2020).

<sup>31</sup> Tipo de red neuronal artificial donde las “neuronas” corresponden a campos receptivos de una manera muy similar a las neuronas en la corteza visual primaria (V1) de un cerebro biológico (Xiaomeng, Chong y Zipei, 2020).

<sup>32</sup> Lenguaje utilizado por los humanos de manera natural (español, inglés, francés, etc.).

<sup>33</sup> Significado de las palabras.

y la orientación que tiene cada vector puede determinar el contexto de su palabra o frase asociada, esto se da gracias a que, si dos o más vectores presentan comportamientos parecidos, las palabras o frases que representan cuentan con significados similares.



**Figura 2.** *Ejemplo bidimensional de una representación vectorial*

Fuente: Elaboración propia.

La figura 2 muestra un ejemplo de cómo los vectores pueden relacionarse entre sí dada su longitud y dirección, se aprecia que las palabras “feliz” y “contento” tienen similitud, pero están muy alejadas de las palabras “triste” y “angustiado” pues su contexto es diferente, en cambio el vector que representa la palabra “ladrillo” permanece neutral a los dos contextos anteriores.

### 2.5.1 FastText

FastText es una biblioteca de código abierto<sup>34</sup>, gratuita y liviana que permite a los usuarios aprender representaciones de texto y clasificadores de texto. Funciona en hardware estándar y genérico (Joulin, Grave, Bojanowski y Mikolov, 2017).

Para calcular vectores de palabras<sup>35</sup>, se necesita un gran corpus de texto. Dependiendo del corpus, los vectores de palabras capturan información diferente (Joulin, Grave, Bojanowski y Mikolov, 2017).

FastText proporciona dos modelos para calcular representaciones de palabras: skipgram y cbow ('bolsa de palabras continúa'). El modelo de skipgram aprende a predecir una palabra objetivo gracias a una palabra cercana. Por otro lado, el modelo cbow predice la palabra objetivo según su contexto. El contexto se representa como una bolsa de palabras contenidas en una ventana de tamaño fijo alrededor de la palabra de destino (Joulin, Grave, Bojanowski y Mikolov, 2017).

Los parámetros<sup>36</sup> más importantes del modelo son su dimensión y el rango de tamaño de las subpalabras. La dimensión controla el tamaño de los vectores, cuanto más grandes son, más información pueden capturar, pero requiere que se aprendan más datos. Pero, si son demasiado grandes, son más difíciles y lentos de entrenar. De forma predeterminada, se usan 100 dimensiones, pero cualquier valor en el rango de 100 a 300 es igual de popular. Las subpalabras son todas las subcadenas contenidas en una palabra entre el tamaño mínimo y el tamaño máximo. Por defecto, se toman todas las subpalabras entre 3 y 6 caracteres, pero otro rango podría ser más apropiado para diferentes idiomas (Joulin, Grave, Bojanowski y Mikolov, 2017).

Una forma sencilla de comprobar la calidad de un vector de palabra es mirar a sus vecinos más cercanos. Esto da una intuición del tipo de información semántica

---

<sup>34</sup> Software cuyo código fuente y otros derechos que normalmente son exclusivos para quienes poseen los derechos de autor, son publicados bajo una licencia de código abierto o forman parte del dominio público (Joulin, Grave, Bojanowski y Mikolov, 2017).

<sup>35</sup> Listas que contienen a un determinado número de palabras (Joulin, Grave, Bojanowski y Mikolov, 2017).

<sup>36</sup> Elemento de un sistema que es útil o crítico al identificar el sistema o al evaluar su rendimiento, estado, condición, etc. (Joulin, Grave, Bojanowski y Mikolov, 2017).

que los vectores son capaces de capturar (Joulin, Grave, Bojanowski y Mikolov, 2017). Por ejemplo:

[(0.891, 'tomate'), (0.889, 'espinacas'), (0.892, 'zanahoria')]

### **2.5.2 InferSent**

InferSent es un método de incrustación de oraciones que proporciona representaciones semánticas para oraciones. Está capacitado en datos de inferencia de lenguaje natural y se generaliza bien para muchas tareas diferentes (Conneau, Kiela, Schwenk, Barrault y Bordes, 2017).

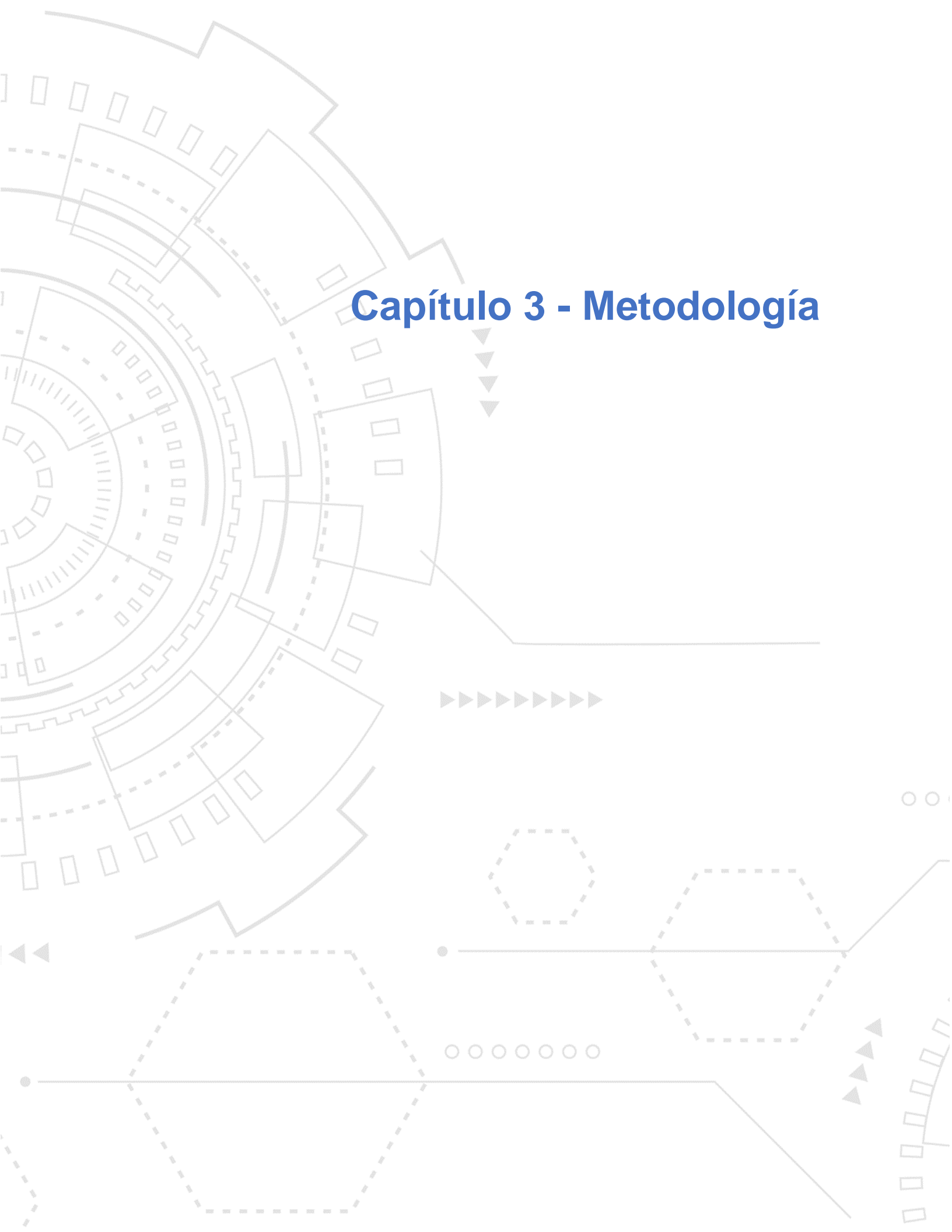
Para su funcionamiento se requieren descargar los vectores GloVe o fastText. Se debe tener en cuenta que `inferred1` se entrena con GloVe e `inferred2` se entrena con fastText. Éste último también elimina el relleno de ceros con agrupación máxima (Conneau, Kiela, Schwenk, Barrault y Bordes, 2017).

Inferred genera una matriz con  $n$  vectores de dimensión 4096. La velocidad es de alrededor de 1000 frases por segundo con un tamaño de lote de 128 en una sola GPU<sup>37</sup>.

---

<sup>37</sup> Unidad de procesamiento de gráficos.

# Capítulo 3 - Metodología





## Capítulo 3 – Metodología

### 3.1 Descripción General

A continuación, se detallará el procedimiento<sup>38</sup> llevado a cabo en la investigación, el ambiente de desarrollo, el servidor<sup>39</sup> utilizado, el ecosistema donde se trabajará y finalmente las tecnologías y librerías utilizadas.

### 3.2 Conjunto de Datos

Para la elaboración de los experimentos propuestos se hizo uso de un corpus compuesto de 1,046,582 tweets realizados por 317 usuarios, éstos últimos están clasificados con dos etiquetas: “con depresión” (1) y “sin depresión” (0) por un grupo de expertos en el tema (INAOE: “usuario\_depresion.json”). El análisis se divide en dos grandes grupos:

- **Por Usuario.** A cada uno de los 317 usuarios se le concatenan todos sus tweets creando así un único texto por usuario, el cual, contiene todos los textos que el usuario escribió, este texto en consecuencia tiene una dimensión relativamente grande (el promedio de tweets escritos por usuario es de 3,302), pero en total son pocos textos generados (317 textos). Como la clasificación de los expertos se hizo por usuario (INAOE: “usuario\_depresion.json”), entonces la clasificación de los textos es la misma.
- **Por Tweet.** En este caso a cada tweet se le considera como un texto independiente, la única relación que tiene con los usuarios es la clasificación, es decir, la clasificación de cada tweet es la misma a la que tiene el usuario que lo escribió. Cabe mencionar que para este análisis las dimensiones de los textos son pequeñas, pero el total de textos generados es muy alto, pues son 1,046,582 tweets en total.

---

<sup>38</sup> Conjunto de acciones que tienen que realizarse todas igualmente, para obtener los mismos resultados bajo las mismas circunstancias (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009).

<sup>39</sup> Sistema que proporciona recursos, datos, servicios o programas a otros ordenadores, conocidos como clientes, a través de una red (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009).

### 3.3 Preprocesamiento

El preprocesamiento que se le aplicó a los textos fue el siguiente: se eliminaron los acentos, signos, puntuaciones, tildes, dígitos, letras repetidas (oooyeeeeee => oye); se transformaron todos los caracteres a minúscula, se eliminaron las siguientes palabras puntuales: http, rt. Finalmente se aplicó una eliminación de stopwords (véase Glosario) en español. Cabe mencionar que se omitió la limpieza de emojis y la aplicación de stemming, lo anterior para no perder peso en las emociones de los usuarios.

### 3.4 Procedimiento

Con la finalidad de seguir una secuencia estructurada de los pasos a desarrollar en este proyecto de investigación y de llevar un orden en los experimentos a realizar se propone la siguiente línea de trabajo:

- Extraer y leer el corpus (INAOE: "usuario\_depresion. json") mediante el uso de la librería de python. (Tellez, Moctezuma, Miranda y Graff, 2018).
- Modelar el problema de identificación de posible depresión, como uno de aprendizaje supervisado.
- Se cuenta con un corpus anotado con posibles usuarios con depresión en investigaciones paralelas por investigadores nacionales.
- Preparación y procesamiento del corpus para su explotación; preparación y procesamiento de los datos a ser anotados de manera automática.
- Preparación y procesamiento del corpus para su explotación; preparación y procesamiento de los datos a ser anotados de manera automática.
- Creación y selección de modelos para la solución del problema; se usará validación cruzada. (Tellez, Moctezuma, Miranda y Graff, 2018), (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009), (Hogenboom, Frasinca, Kaymak y Jong, 2011).
- Análisis de los modelos (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009), los modelos utilizados en este proyecto (Soporte Lineal para la Clasificación de Vectores, Clasificador de Vecinos más Cercanos, Tf-Idf, Clasificador Naive Bayes para modelos multinomiales y microTC – TextModel) son tan

solo una pequeña fracción de los que están actualmente disponibles queda a criterio del lector utilizar el modelo que más le convenga (véase 2.3).

- Etiquetado masivo y análisis de los resultados. Esto con el propósito de crear un dataset funcional para su correcto uso en el aprendizaje automático que se le aplicará, así mismo, al tener la información correctamente etiquetada se podrán obtener resultados sin sesgos (Hastie, Trevor, Tibshirani, Robert y Friedman, 2009), (Hogenboom, Frasinca, Kaymak y Jong, 2011).

Con la anterior secuencia se espera tener un orden en el proyecto y evitar así retrabajos y confusiones, para de esta forma presentar resultados lo más claros posibles sin dificultades de apreciación.

## 3.5 Ambiente de Desarrollo

### 3.5.1 Servidor utilizado para correr el código

El servidor utilizado para la ejecución del código fue proporcionado por INFOTEC<sup>40</sup> campus Aguascalientes<sup>41</sup>. Este proyecto se llevó a cabo haciendo uso de la terminal en un sistema operativo Ubuntu, en éste se corrieron los scripts de Python correspondientes a esta investigación. Las especificaciones del servidor son las siguientes:

- Servidor Intel® Xeon®<sup>42</sup>
- CPU<sup>43</sup> E5-2680 v4, 2.40 GHz<sup>44</sup>
- Sistema Operativo Ubuntu<sup>45</sup> 16.04.2 LTS
- 256 GB de memoria Ram<sup>46</sup>.
- 24 procesadores<sup>47</sup>

---

<sup>40</sup> Centro de Investigación e Innovación en Tecnologías de la Información y Comunicación.

<sup>41</sup> Ciudad de la zona central de México.

<sup>42</sup> Xeon es una familia de microprocesadores Intel para servidores.

<sup>43</sup> Unidad Central de Procesamiento, procesa todas las instrucciones del hardware y software.

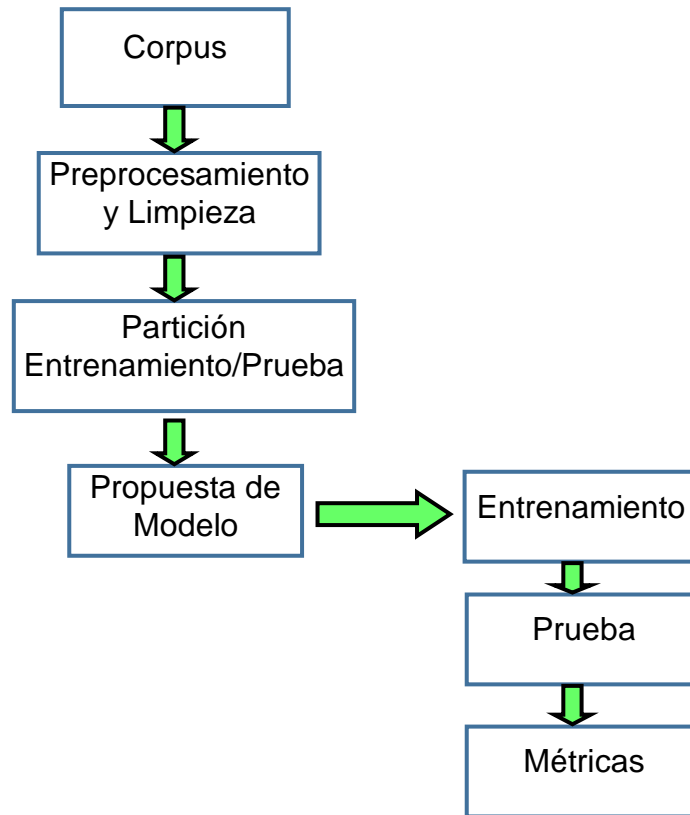
<sup>44</sup> Mide la frecuencia de un chip electrónico, el más conocido es el procesador.

<sup>45</sup> Sistema operativo de software libre y código abierto.

<sup>46</sup> Donde se almacenan de forma temporal los datos de los programas que son utilizados en tiempo real.

<sup>47</sup> CPU.

A continuación, se muestra un diagrama que resume la arquitectura del proyecto:



**Figura 3.** *Arquitectura general del proyecto*

Fuente: Elaboración propia.

### **3.5.2 Tecnologías y librerías utilizadas.**

El lenguaje de programación usado en este proyecto es Python (versión 3.7.6), el cual cuenta con librerías instaladas por default que son utilizadas en este proyecto. Sin embargo, es necesario la instalación de otras para el correcto funcionamiento del código aplicado en este trabajo, las librerías que se deben instalar son:

- conda 4.8.3
- scikit-learn 0.24.2
- microtc 2.2.5
- fasttext 0.9.2
- torch 0.4.1
- infersent

# Capítulo 4 - Análisis y Evaluación

The background features a complex technical illustration. On the left, there are several interlocking gears of different sizes, some with dashed outlines. A network of solid and dashed lines, along with various geometric shapes like rectangles and hexagons, extends across the page. Some lines have arrowheads pointing in different directions, suggesting a flow or process. The overall style is clean and technical, typical of a professional report or textbook cover.

## Capítulo 4 - Análisis y Evaluación

Este capítulo muestra los experimentos realizados, los resultados logrados y la configuración detallada de cada uno de éstos. En los siguientes puntos se detalla la estructura, partición, modelos realizados y los resultados de los experimentos.

### 4.1 Estructura de los Datos por Usuario

El siguiente cuadro detalla la estructura de los datos por Usuario, así como su partición de entrenamiento y prueba.

Estructura	Clases	Textos	Usuarios	FastText	InferSent
Completo	NEG	225	225	225	225
	POS	92	92	92	92
Entrenamiento	NEG	159	159	159	159
	POS	62	62	62	62
Prueba	NEG	66	66	66	66
	POS	30	30	30	30

**Cuadro 1.** *Distribución de las clases.*

Fuente: Elaboración propia.

Como se puede observar en el cuadro anterior no existe un balance entre las clases, con dominio claro en la clase NEG, es decir, en los usuarios sin depresión. Esto se puede explicar debido a que de manera natural la mayoría de la gente y por consecuencia la mayoría de las personas que escriben en Twitter no sufren de depresión<sup>48</sup>.

La partición de los datos fue de 70% para el conjunto de entrenamiento<sup>49</sup> y el 30% para el conjunto de prueba<sup>50</sup>. La estructura de las particiones se puede observar en el cuadro 1.

---

<sup>48</sup> Existen más personas sin depresión que con ella (Fernández 2019).

<sup>49</sup> Datos que se usarán para entrenar el modelo de machine learning.

<sup>50</sup> Datos que se usarán para probar el modelo de machine learning.

## 4.2 Estructura de los Datos por Tweet con y sin preprocesamiento

Para el análisis de los datos por tweet se encontró que debido a la alta cantidad de textos no se pudieron generar los experimentos con FastText y con InferSent (aún con el uso del Clúster), la solución fue aplicar un preprocesamiento de limpieza a los textos (véase 3.3) para disminuir sus dimensiones, de esta manera la experimentación con FastText se pudo lograr, aunque con InferSent persistió el problema de no poder generar los experimentos por la alta cantidad de textos.

Estructura	Clases	Textos	FastText
Completo	NEG	745,783	745,783
	POS	300,799	300,799
Entrenamiento	NEG	522,048	522,048
	POS	210,559	210,559
Prueba	NEG	223,735	223,735
	POS	90,240	90,240

**Cuadro 2.** *Distribución de las clases con preprocesamiento*

Fuente: Elaboración propia.

Estructura	Clases	Textos
Completo	NEG	746,271
	POS	300,923
Entrenamiento	NEG	522,421
	POS	210,614
Prueba	NEG	223,850
	POS	90,309

**Cuadro 3.** *Distribución de las clases sin preprocesamiento*

Fuente: Elaboración propia.

Así que se generaron dos corpus por tweet, uno con preprocesamiento (véase 3.3) donde se incluye la experimentación con FastText (cuadro 2) y el otro sin éste, donde no se toma en cuenta FastText ni InferSent (cuadro 3).

### 4.3 Experimentos y Resultados

En el proceso de experimentación se usaron varias combinaciones de modelos<sup>51</sup> para de esta forma tener una visión más amplia<sup>52</sup> de los resultados que se pueden presentar en problemas de clasificación enfocados a la detección de depresión. En los siguientes puntos están los cuadros con los resultados arrojados por los experimentos, las métricas mostradas son precisión, recall y F1 en sus tres formas (selección automática de parámetros, macro y micro) y accuracy. Para una mejor redacción se abreviará la selección automática de parámetros como **SAP** y la selección editada de parámetros como **SEP**. La partición seleccionada fue con una proporción de 70% para entrenamiento y 30% para prueba.

Todos los experimentos fueron realizados en el lenguaje de programación Python y ejecutados en el clúster (véase 3.5.1). Como se mencionó anteriormente, el uso de Word embeddings en este proyecto se llevó a cabo con FastText y con InferSent, los vectores resultantes fueron de dimensión 300 en el caso de FastText y de dimensión 4096 en el caso de InferSent. Los vectores fueron entrenados de la siguiente manera:

- FastText (véase 2.5.1): crawl-300d-2M-subword.bin
- InferSent1 (véase 2.5.2): GloVe.840B.300d.txt

Parámetros del modelo de InferSent:

- params\_model (véase 2.5.2) = {'bsize': 64, 'word\_emb\_dim': 300, 'enc\_lstm\_dim': 2048, 'pool\_type': 'max', 'dpout\_model': 0.0, 'version': model\_version}

---

<sup>51</sup> Véase cuadro 8.

<sup>52</sup> Al contar con varias combinaciones de modelos se tiene un mayor rango de resultados que pueden compararse.



### 4.3.1 Métricas por usuarios con preprocesamiento

Modelo	Métrica / Test	SAP	Macro	Micro
TfIdf / MultinomialNB	Precision	0.72	0.36	0.72
	Recall	1	0.50	0.72
	F1	0.84	0.42	0.72
	Accuracy	0.72	-	-
MicroTc / SAP	Precision	0.78	0.89	0.80
	Recall	1	0.65	0.80
	F1	0.88	0.67	0.80
	Accuracy	0.80	-	-
MicroTc / SEP	Precision	0.76	0.88	0.78
	Recall	1	0.66	0.78
	F1	0.86	0.67	0.78
	Accuracy	0.78	-	-
SVM / FastText	Precision	0.72	0.36	0.72
	Recall	1	0.50	0.72
	F1	0.84	0.42	0.72
	Accuracy	0.72	-	-
SVM / InferSent	Precision	0.90	0.92	0.91
	Recall	0.99	0.85	0.91
	F1	0.94	0.87	0.91
	Accuracy	0.91	-	-
KNN / FastText	Precision	0.87	0.82	0.84
	Recall	0.93	0.78	0.84
	F1	0.90	0.79	0.84
	Accuracy	0.84	-	-
KNN / InferSent	Precision	0.80	0.70	0.76
	Recall	0.90	0.65	0.76
	F1	0.84	0.67	0.76
	Accuracy	0.76	-	-

**Cuadro 4.** Desempeños por usuario con preprocesamiento

Fuente: Elaboración propia.

En el cuadro 4 se puede observar que se corrieron 7 diferentes modelos donde la combinación SVM / Infsent obtuvo los mejores resultados, las métricas más bajas son para la combinación de SVM / FastText y para TfIdf. En esta etapa, todos los experimentos se corrieron sin problema en el clúster.

#### 4.3.2 Métricas por Usuarios sin preprocesamiento

Modelo	Métrica / Test	SAP	Macro	Micro
MicroTc / SAP	Precision	0.85	0.92	0.88
	Recall	1	0.80	0.88
	F1	0.92	0.83	0.88
	Accuracy	0.88	-	-
MicroTc / SEP	Precision	0.83	0.91	0.85
	Recall	1	0.77	0.85
	F1	0.90	0.80	0.85
	Accuracy	0.85	-	-

**Cuadro 5.** *Desempeños por usuario sin preprocesamiento*

Fuente: Elaboración propia

En el cuadro 5 se aprecia que se corrió solo 1 modelo que fue microTc, esto debido a que en la prueba anterior (véase 4.3.1) se corrieron todos los modelos con un preprocesamiento (véase 3.3) pero en este caso no existió tal limpieza, así que se probó el modelo de MicroTc en una forma pura con dos configuraciones de parámetros diferentes: SAP y SEP (véase 4.5) para de esta forma poder comparar resultados. Los mejores resultados se obtuvieron con la configuración SAP. En esta etapa, los experimentos se corrieron sin problema en el clúster, dando prioridad al modelo microTc y observando cómo se comporta cuando se utilizan sus parámetros automáticos y cuando se les edita sin un preprocesamiento previo.

Lo anterior pues microTc cuenta con su propia limpieza (véase 4.5 microTc Parámetros), dado esto, en el actual proyecto se puede observar que las métricas arrojadas por microTc con una limpieza tradicional antes de aplicar la suya (cuadro 4) son menores a las arrojadas solo aplicando la limpieza de microTC.

### 4.3.3 Métricas por Tweet con preprocesamiento

Modelo	Métrica / Test	SAP	Macro	Micro
SVM / FastText	Precision	0.72	0.61	0.71
	Recall	0.97	0.52	0.71
	F1	0.83	0.48	0.71
	Accuracy	0.71	-	-
MicroTc / SAP	Precision	0.89	0.87	0.88
	Recall	0.95	0.83	0.88
	F1	0.92	0.85	0.88
	Accuracy	0.88	-	-
MicroTc / SEP	Precision	0.89	0.86	0.88
	Recall	0.94	0.83	0.88
	F1	0.91	0.84	0.88
	Accuracy	0.88	-	-

**Cuadro 6.** *Desempeños por tweet con preprocesamiento*

Fuente: Elaboración propia

El cuadro 6 muestra sólo 3 modelos, esto por el hecho de que el clúster no corrió las demás pruebas, la causa fue el gran volumen de textos que se manejaron al trabajar por tweet y no por usuario. Sólo hay resultados para SVM / FastText y microTC, siendo este último en su versión SAP el que obtuvo mejores métricas.

#### 4.3.4 Métricas por Tweet sin preprocesamiento

Modelo	Métrica / Test	SAP	Macro	Micro
MicroTc / SAP	Precision	0.84	0.77	0.81
	Recall	0.90	0.74	0.81
	F1	0.87	0.75	0.81
	Accuracy	0.81	-	-
MicroTc / SEP	Precision	0.84	0.79	0.82
	Recall	0.92	0.74	0.82
	F1	0.88	0.76	0.82
	Accuracy	0.82	-	-

**Cuadro 7.** *Desempeños por tweet sin preprocesamiento*

Fuente: Elaboración propia

En el cuadro 7 al igual que en el cuadro 5 sólo se corrió microTc, esto para poder ver cómo se comporta este modelo con sus propios parámetros sin que haya un preprocesamiento previo en los datos que trabajará, al contrario de lo que sucedió antes, en esta prueba, se lograron mejores resultados cuando se realizó un preprocesamiento previo al que microTc realiza, mejorando de esta forma sus métricas.

#### 4.4 Discusión de los resultados obtenidos

El actual proyecto aplicó varios experimentos con diferentes combinaciones de clasificadores, embeddings y modelos para la resolución de un problema de clasificación aplicando machine learning.

La investigación está basada en dos enfoques que a su vez están ramificados en múltiples experimentos todos con el mismo fin pero que dieron resultados diferentes. Los dos enfoques y sus experimentos correspondientes se muestran en los cuadros 8 y 9 respectivamente.

Por Usuario		Por Tweet	
Con Prepro-cesamiento	Sin Preprocesa- miento	Con Preprocesa- miento	Sin Preprocesa- miento
MicroTC SAP	MicroTC SAP	MicroTC SAP	MicroTC SAP
MicroTC SEP	MicroTC SEP	MicroTC SEP	MicroTC SEP
SVM/FastText	*	SVM/FastText	*
TfIdf	*	*	*
SVM/Infer- Sent	*	*	*
KNN/FastText	*	*	*
KNN/Infer- Sent	*	*	*

**Cuadro 8. Enfoques y Modelos**

Fuente: Elaboración propia

\* No se pudieron correr debido a su alto consumo de recursos computacionales.

Del análisis realizado los mejores resultados fueron:

Por Usuario				
Modelo	Métrica / Test	Resultado	Macro	Micro
SVM/Infer- Sent con pre- procesa- miento	Precision	0.90	0.92	0.91
	Recall	0.99	0.85	0.91
	F1	0.94	0.87	0.91
	Accuracy	0.91	-	-

Por Tweet				
Modelo	Métrica / Test	Resultado	Macro	Micro
MicroTC SAP con preproce- samiento	Precision	0.89	0.87	0.88
	Recall	0.95	0.83	0.88
	F1	0.92	0.85	0.88
	Accuracy	0.88	-	-

**Cuadro 9. Mejores Resultados**

Fuente: Elaboración propia

#### 4.5 MicroTC Parámetros:

Los siguientes parámetros de MicroTC fueron los que se utilizaron en los experimentos SAP y SEP correspondientemente.

MicroTc SEP		MicroTc SAP	
docs	None	docs	None
text	'text'	text	'text'
num_option	'delete'	num_option	'group'
usr_option	'delete'	usr_option	'group'
url_option	'delete'	url_option	'group'
emo_option	'group'	emo_option	'group'
hashteg_option	'delete'	hashteg_option	'none'
ent_option	'none'	ent_option	'none'
lc	True	lc	True
del_dup	True	del_dup	True
del_punc	True	del_punc	False
del_diac	True	del_diac	True
token_list	[[2,1], - 1,3,4]	token_list	[-1]
token_min_filter	0	token_min_filter	0
token_max_filter	1	token_max_filter	1
select_ent	False	select_ent	False
select_suff	False	select_suff	False
select_conn	False	select_conn	False
weighting	'tfidf'	weighting	'tfidf'

**Cuadro 10.** *Parámetros MicroTc*

Fuente: Elaboración propia

Para una mejor redacción se abreviará la selección automática de parámetros como SAP y SEP. A continuación, se muestran los parámetros que contiene el clasificador microTC y su significado:

- **docs:** Debe ser una lista que contendrá el corpus.
- **text** : En el caso de que el corpus sea un diccionario, el nombre del campo será la clave del texto .
- **num\_option:** Transformación de números, posibles valores: none, group y delete.
- **usr\_option:** Transformación de los usuarios, posibles valores: none, group y delete.
- **url\_option:** Transformación de las url, posibles valores, posibles valores: none, group y delete.

- **emo\_option**: Transformaciones de emojis y emoticones, posibles valores: none, group y delete.
- **hashtag\_option**: Transformaciones de hashtags, posibles valores: none, group y delete.
- **ent\_option**: Transformación de entidades, posibles valores: none, group y delete.
- **lc**: Minúsculas
- **del\_dup**: Remover símbolos duplicados.
- **del\_punc**: Elimina los símbolos de puntuación.
- **del\_diac**: Elimina símbolos diacríticos.
- **token\_list**: Tokens<sup>53</sup> > 0 qgrams<sup>54</sup> < 0 word-grams<sup>55</sup>.
- **token\_min\_filter**: Conserva los tokens que aparecen más veces que el parámetro (utilizado en la clase de ponderación).
- **token\_max\_filter**: Conserva los tokens que aparecen menos veces que el parámetro (utilizado en la clase de ponderación).
- **select\_ent**: Booleano.
- **select\_suff**: Booleano.
- **select\_conn**: Booleano.
- **weighting**: Esquema de ponderación, posibles valores: tfidf, tf y entropy.

---

<sup>53</sup> Conjunto de elementos seleccionados y separadas del texto original (Tellez, Moctezuma, Miranda y Graff, 2018).

<sup>54</sup> Conjunto de elementos seleccionados y separados por 'q' posiciones del texto original (Tellez, Moctezuma, Miranda y Graff, 2018).

<sup>55</sup> Conjunto de palabras seleccionadas y separadas del texto original (Tellez, Moctezuma, Miranda y Graff, 2018).

The background of the page is a light gray technical drawing. It features several large gears of different sizes and orientations, some with dashed lines indicating hidden parts. There are also various geometric shapes like rectangles and hexagons, and lines representing mechanical components and paths. The overall style is clean and technical.

## Capítulo 5 – Conclusiones y Trabajo Futuro



## Capítulo 5 – Conclusiones y Trabajo Futuro

### 5.1 Conclusiones

En la investigación realizada se cumplieron los objetivos planteados logrando crear una herramienta de machine learning capaz de detectar posibles casos de depresión entre los usuarios de Twitter.

El proyecto trabajó dos tipos de pruebas: una por usuario y otra por tweet, en cada una se aplicaron experimentos y se obtuvo igual cantidad de resultados. La cantidad de datos por usuario eran pocos, pero con textos de dimensiones enormes y la cantidad de datos por tweet eran muchos, pero con textos cortos.

Los experimentos realizados mostraron que el modelo SVM/InferSent por usuario con preprocesamiento y el modelo microTC SAP por tweet con preprocesamiento fueron los que mejores métricas obtuvieron:

Por Usuario				
Modelo	Métrica / Test	Resultado	Macro	Micro
SVM / Infersent con preprocesamiento	Precision	0.90	0.92	0.91
	Recall	0.99	0.85	0.91
	F1	0.94	0.87	0.91
	Accuracy	0.91	-	-

Por Tweet				
Modelo	Métrica / Test	Resultado	Macro	Micro
MicroTC SAP con preprocesamiento	Precision	0.89	0.87	0.88
	Recall	0.95	0.83	0.88
	F1	0.92	0.85	0.88
	Accuracy	0.88	-	-

Cabe mencionar que los resultados de los modelos restantes son cercanos a los dos antes mencionados, teniendo generalmente la métrica más alta en el recall. Es importante destacar que microTc se comportó de manera robusta en todos los experimentos, arrojando no sólo buenos resultados en las métricas, sino también fue el único modelo que se pudo correr en todas las pruebas, mostrando que es un clasificador confiable y robusto, como se puede observar en la siguiente tabla:

Por Usuario		Por Tweet	
Con Preprocesamiento	Sin Preprocesamiento	Con Preprocesamiento	Sin Preprocesamiento
MicroTC SAP	MicroTC SAP	MicroTC SAP	MicroTC SAP
MicroTC SEP	MicroTC SEP	microTC SEP	MicroTC SEP
SVM/FastText	*	SVM/FastText	*
TfIdf	*	*	*
SVM/InferSent	*	*	*
KNN/FastText	*	*	*
KNN/InferSent	*	*	*

Un punto interesante descubierto en los experimentos fue el hecho de que, por usuario, microTc se comportó de mejor manera cuando realizó su propia limpieza sin aplicarle una previa. Y por tweet, microTc se comportó de mejor manera cuando se realizó un preprocesamiento previo al de su propia limpieza. Como aporte final el corpus utilizado, así como los experimentos realizados facilitan nuevos resultados al análisis de predicciones en problemas de salud mental.

## 5.2 Trabajo Futuro

Algunas propuestas de posibles trabajos futuros para continuar la presente investigación son:

- Realizar el análisis de los datos por tweet con word embedding aplicando InferSent. Lo anterior ya que en el presente proyecto se propuso este análisis, pero la capacidad del clúster y los recursos de cómputo disponibles entre otros factores no dejaron realizar el experimento.
- Se puede mejorar el corpus del proyecto, involucrando a expertos en salud mental al momento de clasificar a los usuarios, ya que en el presente proyecto la clasificación con depresión fue auto declarada por éstos.
- En el actual proyecto se trabajó con un solo corpus el cual estaba enfocado en la clasificación de personas con depresión. Pero existen más desórdenes mentales que aquejan a las personas y que igualmente se pueden analizar con los métodos propuestos en esta investigación.
- Visualizar de manera georreferenciada a los usuarios clasificados con y sin depresión en un mapa, para de esta manera poder identificar las zonas con más densidad de población que cuenten con población afectada por este u otro problema mental, llámense, por ejemplo, zonas rurales, urbanas, nuevos asentamientos, etcétera. El código y el trabajo propuesto en este proyecto será publicado y puesto al alcance del público para que de esta manera ayude a trabajos futuros.

## Bibliografía

Agnihotri, D., Verma, K., & Tripathi, P. (2014). Pattern and Cluster Mining on Text Data. *IEEE Xplore*. <https://ieeexplore.ieee.org/document/6821432>

Chopade, P., Edwards, D., Khan, S., Andrade, A., & Pu, S. (2019). CPSX: Using AI-Machine Learning for Mapping Human-Human Interaction and Measurement of CPS Teamwork Skills. *IEEE Xplore*. <https://ieeexplore.ieee.org/abstract/document/9032906>

Conneau, A., Kiela, D., Schwenk, H., Barrault, L., & Bordes, A. (2017). Supervised Learning of Universal Sentence Representations from Natural Language Inference Data. *ACL Anthology, D17 (1070)*, 670-680. <https://aclanthology.org/D17-1070/>

Fernández, A. (2019). DE CADA 100 MEXICANOS, 15 PADECEN DEPRESIÓN. *Boletín UNAM-DGCS*, 455. [https://www.dgcs.unam.mx/boletin/bdboletin/2019\\_455.html](https://www.dgcs.unam.mx/boletin/bdboletin/2019_455.html)

Hakdağlı, Ö., Özcan, C., & Oğul, I. (2018). Stream text data analysis on twitter using apache spark streaming. *IEEE Xplore*. <https://ieeexplore.ieee.org/document/8404540>.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Second Edition*. Berlin: Springer.

Hogenboom, F., Frasinca, F., Kaymak, U., & De Jong, F. (2011). An overview of event extraction from text. *CEUR Workshop Proceedings, 779*, 48-57. <http://ceur-ws.org/Vol-779/>

Hooda, M., Saxena, A., Madhulika, D., & Yadav, B. (2017). A Study and Comparison of Prediction Algorithms for Depression Detection among Millennials: A Machine Learning Approach. *IEEE Xplore*. <https://ieeexplore.ieee.org/document/8455078>

Hosseinifard, B., Hassan, M., & Rostami, R. (2013). Classifying depression patients and normal subjects using machine learning techniques and nonlinear features from EEG signal. *Computer Methods and Programs in Biomedicine*, 109 (3), 339-345. <https://www.sciencedirect.com/science/article/abs/pii/S0169260712002507>

Joulin, A., Grave, E., Bojanowski, P., & Mikolov, T. (2017). Bag of Tricks for Efficient Text Classification. *ACL Anthology*, E17 (2068), 427-431. <https://aclanthology.org/E17-2068/>

Kulasinghe, S., Jayasinghe, A., Rathnayaka, R., Karunarathne, P., Silva, S., & Jayakodi, A. (2019). AI Based Depression and Suicide Prevention System. *IEEE Xplore*. <https://ieeexplore.ieee.org/document/9103411>

Larsen, M., Boonstra, T., Batterham, P., O'Dea, B., Paris, C., & Christensen, H. (2015). We Feel: Mapping Emotion on Twitter. *IEEE Xplore*. <https://ieeexplore.ieee.org/document/7042256>

McClellan, C., Ali, M., Mutter, R., Kroutil, L., & Landwehr, J. (2016). Using social media to monitor mental health discussions - evidence from Twitter. *Journal of the American Medical Informatics Association*, 24 (3), 496-502.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel V., Thirion, B., Grisel, O., ... Duchesnay, É. (2011). Scikit-learn: Machine Learning in Python. *JMLR*, 12 (85), 2825-2830. <https://jmlr.csail.mit.edu/papers/v12/pedregosa11a.html>

Pratama, T., & Purwarianti, A. (2017). Topic classification and clustering on Indonesian complaint tweets for bandung government using supervised and unsupervised learning. *IEEE Xplore*. <https://ieeexplore.ieee.org/document/8090981>

Renuka, J. (2016). Accuracy, Precision, Recall & F1 Score: Interpretation of Performance Measures. *Exsilio Solutions*. <https://blog.exsilio.com/all/accuracy-precision-recall-f1-score-interpretation-of-performance-measures/>

Roman, V. (2019). Algoritmos Naive Bayes: Fundamentos e Implementación. *Medium*. <https://medium.com/datos-y-ciencia/algoritmos-naive-bayes-fudamentos-e-implementaci%C3%B3n-4bcb24b307f>

Roman, V. (2019). Aprendizaje Supervizado: Introducción a la Clasificación y Principales Algoritmos. *Medium*. <https://medium.com/datos-y-ciencia/aprendizaje-supervisado-introducci%C3%B3n-a-la-clasificaci%C3%B3n-y-principales-algoritmos-dadee99c9407>

Tekiner, F., Tsuruoka, Y., Tsujii, J., & Ananiadou, S. (2009). Highly scalable Text Mining - parallel tagging application. *IEEE Xplore*. <https://ieeexplore.ieee.org/abstract/document/5379432/references#references>

Tellez, E., Moctezuma, D., Miranda, S., & Graff, M. (2018). An automated text categorization framework based on hyperparameter optimization. *Knowledge-Based Systems*, 149, 110-123. <https://www.sciencedirect.com/science/article/abs/pii/S0950705118301217>

Xiaomeng, Z., Chong, Z., & Zipei, Y. (2020). Emotional Tendency Analysis of Twitter Texts Based on CNNs Model. *IEEE Xplore*. <https://ieeexplore.ieee.org/document/9150239>

The background features a complex, light gray geometric pattern. On the left, there are several interlocking gear-like shapes with various teeth and internal structures. A network of solid and dashed lines crisscrosses the page, some ending in small circles or dots. Several sets of arrows are scattered throughout: three pointing down, a row of eight pointing right, two pointing left, and a cluster of four pointing up and right. The overall aesthetic is technical and architectural.

# ANEXOS



## ANEXO I: Repositorio del código de Python utilizado en este proyecto

### Liga:

<https://github.com/mlCorpusDepresion/mlPythonDepresionCode>

### Hipervínculo:

mlCorpusDepresion/mlPythonDepresionCode: This project shows the code used to predict possible cases of depression in Twitter users, a tagged corpus was used. The steps to follow were: preprocessing, cleaning, modeling, training, testing and results. (github.com)

## ANEXO II: Códigos individuales de Python utilizados en este proyecto

**Código de Python utilizado para el análisis por tweet utilizando microTC (SEP y SAP). Sin limpieza previa.**

```
1 from sklearn.model_selection import train_test_split
2 import numpy as np
3 import tarfile
4 from microtc.utils import tweet_iterator
5
6 tar = tarfile.open("usuarios_depresion.json.tar.gz")
7 tar.extractall()
8 tar.close()
9 data = [i for i in tweet_iterator('usuarios_depresion.json')]
10
11 ktu,klases,tweets,users = [],[],[],[]
12 for i in data:
13     for j in i['tweets']:
14         ktu.append([j['text'],int(i['klass']),i['user']])
15
16 for i in ktu:
17     tweets.append(i[0])
18     klases.append(i[1])
19     users.append(i[2])
20
21 import pandas as pd
22 dataE = pd.DataFrame({'Textos':tweets, 'Klases':klases,'Usuarios':users})
23
24 # Train Test
25 XE = dataE
26 yE = dataE['Klases']
27 XE_data_train, XE_data_test, yE_train, yE_test = train_test_split(XE, yE,
28 test_size=0.3)
29 print('se creo el train y el test')
30
31 print(XE.groupby(['Klases']).count())
32 print(XE_data_train.groupby(['Klases']).count())
33 print(XE_data_test.groupby(['Klases']).count())
34
35 yE_train = list(yE_train)
36 yE_test = list(yE_test)
37
38 XE_train = list(XE_data_train['Textos'])
39 XE_test = list(XE_data_test['Textos'])
40
41 # microTC - micro Text Classification
42 # Params editados
```

```

43 from microtc.textmodel import TextModel
44
45 textmodelEM = TextModel(docs=None,
46 text='text',num_option='delete',usr_option='delete',
47     url_option='delete',emo_option='group',hashtag_option='delete',
48     ent_option='none',lc=True,del_dup=True,del_punc=True,del_diac=True,
49     token_list=[[2, 1], -1,3,4],token_min_filter=0,token_max_filter=1,
50     select_ent=False,select_suff=False, select_conn=False,
51     weighting='tfidf')
52 textmodelEM.fit(list(XE_train))
53
54 from sklearn.svm import LinearSVC
55 svcEM = LinearSVC()
56 svcEM.fit(textmodelEM.transform(XE_train), yE_train)
57
58 y_predEM = svcEM.predict(textmodelEM.transform(XE_test))
59 print('Termino microTC params')
60
61 # Default
62 textmodelEMD = TextModel()
63 textmodelEMD.fit(list(XE_train))
64
65 svcEMD = LinearSVC()
66 svcEMD.fit(textmodelEMD.transform(XE_train), yE_train)
67
68 y_predEMD = svcEMD.predict(textmodelEMD.transform(XE_test))
69 print('Termino microTC default')
70
71 from sklearn.metrics import precision_score, recall_score, f1_score,
72 accuracy_score
73
74 def metricas(modelo,y_test,predict):
75     print(modelo)
76     print('Precision:', round(precision_score(y_test,predict,average=None)[0],3),
77         'Macro:', round(precision_score(y_test, predict, average='macro'),3),
78         'Micro:', round(precision_score(y_test, predict, average='micro'),3))
79     print('Recall:', round(recall_score(y_test,predict,average=None)[0],3),
80         'Macro:', round(recall_score(y_test, predict, average='macro'),3),
81         'Micro:', round(recall_score(y_test, predict, average='micro'),3))
82     print('F1:      ', round(f1_score(y_test, predict, average=None)[0],3),
83         'Macro:', round(f1_score(y_test, predict, average='macro'),3),
84         'Micro:', round(f1_score(y_test, predict, average='micro'),3))
85     print('Accuracy: ', round(accuracy_score(y_test, predict),3))
86
87 print(metricas('\nmicroTC_Params',yE_test,y_predEM))
88 print(metricas('\nmicroTC_Default',yE_test,y_predEMD))
89
90 from sklearn.metrics import confusion_matrix
91 print('\nmicroTC\n',confusion_matrix(yE_test, y_predEM))
92 print('\nmicroTC_Default\n',confusion_matrix(yE_test, y_predEMD))

```

## Código de Python utilizado para el análisis por tweet utilizando microTC (SAP) y fasttext con SVM. Utilizando limpieza previa.

```
1 import fasttext
2 from InferSent.models import InferSent
3 import torch
4 import torchvision
5 import nltk
6 nltk.download('punkt')
7 from sklearn.model_selection import train_test_split
8 import numpy as np
9 import math
10 import tarfile
11 from microtc.utils import tweet_iterator
12
13 tar = tarfile.open("usuarios_depresion.json.tar.gz")
14 tar.extractall()
15 tar.close()
16 data = [i for i in tweet_iterator('usuarios_depresion.json')]
17
18 ktu,klases,tweets,users = [],[],[],[]
19 for i in data:
20     for j in i['tweets']:
21         ktu.append([j['text'],int(i['klass']),i['user']])
22
23 from random import sample
24
25 for i in ktu:
26     tweets.append(i[0])
27     klases.append(i[1])
28     users.append(i[2])
29
30 import re as reg
31 import unicodedata
32 import nltk
33 nltk.download('stopwords')
34 import re, string
35 from nltk.corpus import stopwords
36 lsw = stopwords.words('spanish')
37
38 text2 = tweets.copy()
39
40 def QuitarAcentos(cadena):
41     formaNFKD = unicodedata.normalize('NFKD', cadena)
42     return u"".join([c for c in formaNFKD if not unicodedata.combining(c)])
43 def RP(text):
44     return re.sub('[%s]' % re.escape(string.punctuation), '', text)
45 def RD(digito):
46     return ''.join(i for i in digito if not i.isdigit())
47 #def deEmojify(inputString):
48 #     return inputString.encode('ascii', 'ignore').decode('ascii')
```

```

49 def reducirString(palabra):
50     pattern=reg.compile(r"(\.){1,}", reg.DOTALL)
51     string=pattern.sub(r"\1", palabra)
52     return string
53
54 susp2, susp3, susp5, TextoL, TextoLimpio = [], [], [], [], []
55 QR, stop1, stem1, Tt, KC, Us = [], [], [], [], [], []
56
57 for i in text2:
58     susp2.append(QuitarAcentos(i).lower())
59 for i in susp2:
60     susp5.append(RP(i))
61 for i in susp5:
62     TextoL.append(RD(i))
63 for i in TextoL:
64     aux = i.split()
65     aux2 = []
66     for j in aux:
67         aux2.append(reducirString(j))
68     QR.append(' '.join(aux2))
69
70 for i in QR:
71     aux = []
72     h = i.split()
73     for r in h:
74         if r not in lsw: #and r not in lsw2):
75             aux.append(r)
76     stop1.append(' '.join(aux))
77
78 for i in range(len(stop1)):
79     if stop1[i] == '':
80         continue
81     else:
82         Tt.append(stop1[i])
83         KC.append(klases[i])
84         Us.append(users[i])
85
86 print('Limpieza realizada')
87
88 import pandas as pd
89 dataE = pd.DataFrame({'Textos':Tt, 'Klases':KC, 'Usuarios':Us})
90
91 # FastText
92 ft = fasttext.load_model('crawl-300d-2M-subword.bin')
93 print('Se cargo el modelo de fasttext')
94
95 embeddingsEF = [ ft.get_sentence_vector(v) for v in list(dataE['Textos'])]
96 print('Se crearon los vectores de fasttext')
97
98 dataE = dataE.assign(fft=embeddingsEF)
99
100 # Train Test

```

```

101 XE = dataE
102 yE = dataE['Klases']
103 XE_data_train, XE_data_test, yE_train, yE_test = train_test_split(XE, yE,
104 test_size=0.3)
105 print('se creo el train y el test')
106
107 print(XE.groupby(['Klases']).count())
108 print(XE_data_train.groupby(['Klases']).count())
109 print(XE_data_test.groupby(['Klases']).count())
110
111 yE_train = list(yE_train)
112 yE_test = list(yE_test)
113
114 XEF_train = list(XE_data_train['fft'])
115 XEF_test = list(XE_data_test['fft'])
116
117 # Support Vector Machine (SVM)
118 from sklearn.svm import LinearSVC
119 svceF = LinearSVC(random_state=0)
120 svceF.fit(XEF_train, yE_train)
121 # In[34]:
122 y_predEFO = svceF.predict(XEF_test)
123 print('Termino SVM')
124
125 # microTC - micro Text Classification
126 X_trainEM = list(XE_data_train['Textos'])
127 X_testEM = list(XE_data_test['Textos'])
128
129 from microtc.textmodel import TextModel
130 textmodelEM = TextModel()
131 textmodelEM.fit(list(X_trainEM))
132
133 from sklearn.svm import LinearSVC
134 svcEM = LinearSVC()
135 svcEM.fit(textmodelEM.transform(X_trainEM), yE_train)
136
137 y_predEM = svcEM.predict(textmodelEM.transform(X_testEM))
138 print('Termino microTC')
139
140 from sklearn.metrics import precision_score, recall_score, f1_score,
141 accuracy_score
142
143 def metricas(modelo, y_test, predict):
144     print(modelo)
145     print('Precision:', round(precision_score(y_test, predict, average=None) [0], 3),
146           'Macro:', round(precision_score(y_test, predict, average='macro'), 3),
147           'Micro:', round(precision_score(y_test, predict, average='micro'), 3))
148     print('Recall:', round(recall_score(y_test, predict, average=None) [0], 3),
149           'Macro:', round(recall_score(y_test, predict, average='macro'), 3),
150           'Micro:', round(recall_score(y_test, predict, average='micro'), 3))
151     print('F1:      ', round(f1_score(y_test, predict, average=None) [0], 3),
152           'Macro:', round(f1_score(y_test, predict, average='macro'), 3),

```

```

153     'Micro:', round(f1_score(y_test, predict, average='micro'), 3))
154     print('Accuracy: ', round(accuracy_score(y_test, predict), 3))
155
156 print(metricas('\nFastText_SVM', yE_test, y_predEF0))
157 print(metricas('\nmicroTC', yE_test, y_predEM))
158
159 from sklearn.metrics import confusion_matrix
160
161     print('\nFastText_SVM\n', confusion_matrix(yE_test, y_predEF0))
162     print('\nmicroTC\n', confusion_matrix(yE_test, y_predEM))

```

## Código de Python utilizado para el análisis por tweet utilizando microTC (SEP) y fasttext con SVM. Utilizando limpieza previa.

```

1 import fasttext
2 from InferSent.models import InferSent
3 import nltk
4 nltk.download('punkt')
5 from sklearn.model_selection import train_test_split
6 import numpy as np
7 import math
8 import tarfile
9 from microtc.utils import tweet_iterator
10
11 tar = tarfile.open("usuarios_depresion.json.tar.gz")
12 tar.extractall()
13 tar.close()
14 data = [i for i in tweet_iterator('usuarios_depresion.json')]
15
16 ktu, klases, tweets, users = [], [], [], []
17 for i in data:
18     for j in i['tweets']:
19         ktu.append([j['text'], int(i['klass']), i['user']])
20
21 from random import sample
22
23 for i in ktu:
24     tweets.append(i[0])
25     klases.append(i[1])
26     users.append(i[2])
27
28 import re as reg
29 import unicodedata
30 import nltk
31 nltk.download('stopwords')
32 import re, string
33 from nltk.corpus import stopwords
34 lsw = stopwords.words('spanish')
35

```

```

36 text2 = tweets.copy()
37
38 def QuitarAcentos(cadena):
39     formaNFKD = unicodedata.normalize('NFKD', cadena)
40     return u"".join([c for c in formaNFKD if not unicodedata.combining(c)])
41
42 def RP(text):
43     return re.sub('[%s]' % re.escape(string.punctuation), '', text)
44 def RD(digito):
45     return ''.join(i for i in digito if not i.isdigit())
46 def reducirString(palabra):
47     pattern=reg.compile(r"(\.|\1{1,})",reg.DOTALL)
48     string=pattern.sub(r"\1",palabra)
49     return string
50
51 susp2,susp3,susp5,TextoL,TextoLimpio = [],[],[],[],[]
52 QR,stop1,stem1,Tt,KC,Us = [],[],[],[],[],[]
53
54 for i in text2:
55     susp2.append(QuitarAcentos(i).lower())
56 for i in susp2:
57     susp5.append(RP(i))
58 for i in susp5:
59     TextoL.append(RD(i))
60 for i in TextoL:
61     aux = i.split()
62     aux2 = []
63     for j in aux:
64         aux2.append(reducirString(j))
65     QR.append(' '.join(aux2))
66
67 for i in QR:
68     aux = []
69     h = i.split()
70     for r in h:
71         if r not in lsw: #and r not in lsw2):
72             aux.append(r)
73     stop1.append(' '.join(aux))
74
75 for i in range(len(stop1)):
76     if stop1[i] == '':
77         continue
78     else:
79         Tt.append(stop1[i])
80         KC.append(klases[i])
81         Us.append(users[i])
82
83 print('Limpieza realizada')
84
85 import pandas as pd
86 dataE = pd.DataFrame({'Textos':Tt, 'Klases':KC, 'Usuarios':Us})
87

```



```

88 # FastText
89 ft = fasttext.load_model('crawl-300d-2M-subword.bin')
90 print('Se cargo el modelo de fasttext')
91
92 embeddingsEF = [ ft.get_sentence_vector(v) for v in list(dataE['Textos'])]
93 print('Se crearon los vectores de fasttext')
94
95 dataE = dataE.assign(fft=embeddingsEF)
96
97 # Train Test
98 XE = dataE
99 yE = dataE['Klases']
100 XE_data_train, XE_data_test, yE_train, yE_test = train_test_split(XE, yE,
101 test_size=0.3)
102 print('se creo el train y el test')
103
104 print(XE.groupby(['Klases']).count())
105 print(XE_data_train.groupby(['Klases']).count())
106 print(XE_data_test.groupby(['Klases']).count())
107
108 yE_train = list(yE_train)
109 yE_test = list(yE_test)
110
111 XEF_train = list(XE_data_train['fft'])
112 XEF_test = list(XE_data_test['fft'])
113
114 # Support Vector Machine (SVM)
115 from sklearn.svm import LinearSVC
116 svcEF = LinearSVC(random_state=0)
117 svcEF.fit(XEF_train, yE_train)
118
119 y_predEF0 = svcEF.predict(XEF_test)
120 print('Termino SVM')
121
122 # microTC - micro Text Classification
123 X_trainEM = list(XE_data_train['Textos'])
124 X_testEM = list(XE_data_test['Textos'])
125
126 from microtc.textmodel import TextModel
127 textmodeleM = TextModel(docs=None,
128 text='text', num_option='delete', usr_option='delete',
129 url_option='delete', emo_option='group', hashtag_option='delete',
130 ent_option='none', lc=True, del_dup=True, del_punc=True, del_diac=True,
131 token_list=[[2, 1], -1, 3, 4], token_min_filter=0, token_max_filter=1,
132 select_ent=False, select_suff=False, select_conn=False,
133 weighting='tfidf')
134 textmodeleM.fit(list(X_trainEM))
135
136 from sklearn.svm import LinearSVC
137 svcEM = LinearSVC()
138 svcEM.fit(textmodeleM.transform(X_trainEM), yE_train)
139

```

```

140 y_predEM = svcEM.predict(textmodelEM.transform(X_testEM))
141 print('Termino microTC')
142
143 from sklearn.metrics import precision_score, recall_score, f1_score,
144 accuracy_score
145
146 def metricas(modelo,y_test,predict):
147     print(modelo)
148 print('Precision:',round(precision_score(y_test,predict,average=None)[0],3),
149     'Macro:',round(precision_score(y_test, predict, average='macro'),3),
150     'Micro:',round(precision_score(y_test, predict, average='micro'),3))
151 print('Recall:',round(recall_score(y_test,predict,average=None)[0],3),
152     'Macro:',round(recall_score(y_test, predict, average='macro'),3),
153     'Micro:',round(recall_score(y_test, predict, average='micro'),3))
154 print('F1:      ',round(f1_score(y_test, predict, average=None)[0],3),
155     'Macro:',round(f1_score(y_test, predict, average='macro'),3),
156     'Micro:',round(f1_score(y_test, predict, average='micro'),3))
157 print('Accuracy: ',round(accuracy_score(y_test, predict),3))
158
159 print(metricas('\nFastText_SVM',yE_test,y_predEF0))
160 print(metricas('\nmicroTC',yE_test,y_predEM))
161
162 from sklearn.metrics import confusion_matrix
163
164 print('\nFastText_SVM\n',confusion_matrix(yE_test, y_predEF0))
165 print('\nmicroTC\n',confusion_matrix(yE_test, y_predEM))

```

## Código de Python utilizado para el análisis por usuarios utilizando microTC (SEP) y Tfidf. Utilizando limpieza previa.

```

1 import tarfile
2 tar = tarfile.open("usuarios_depresion.json.tar.gz")
3 tar.extractall()
4 tar.close()
5
6 from microtc.utils import tweet_iterator
7 data = [i for i in tweet_iterator('usuarios_depresion.json')]
8
9 # Extraer los tweets
10 data[0]['tweets'][4099]['text']
11
12 # Extraer lo que se requiere, tres listas de dimension 317 en donde todos
13 los tweets de cada usuario se juntaron para hacer solo un tweet unico por
14 usuario.
15 clases,users,t,tweets = [],[],[],[]
16 for i in data:
17     t = []
18     clases.append(int(i['klass']))
19     users.append(i['user'])

```

```

20     for j in i['tweets']:
21         t.append(j['text'])
22     tweets.append(' '.join(t))
23
24 print('Data Original')
25 print(len(klases))
26 print(len(tweets))
27 print(len(users))
28
29
30 # Limpieza
31 import re as reg
32 import unicodedata
33 import nltk
34 import re, string
35
36 textos = tweets.copy()
37
38 def QuitarAcentos(cadena):
39     formaNFKD = unicodedata.normalize('NFKD', cadena)
40     return u"".join([c for c in formaNFKD if not unicodedata.combining(c)])
41 def RP(text):
42     return re.sub('[%s]' % re.escape(string.punctuation), '', text)
43 def RD(digito):
44     return ''.join(i for i in digito if not i.isdigit())
45 def reducirString(palabra):
46     pattern=reg.compile(r"(\.|\1{1,})",reg.DOTALL)
47     string=pattern.sub(r"\1",palabra)
48     return string
49
50 l1,l2,l3,l4,l5 = [],[],[],[],[]
51
52 for i in textos:
53     l1.append(QuitarAcentos(i).lower())
54 for i in l1:
55     aux = []
56     h = i.split()
57     for r in h:
58         if (r.startswith('http') or r.startswith('rt') or
59 r.startswith('"http')):
60             continue
61         aux.append(r)
62     l2.append(' '.join(aux))
63 for i in l2:
64     l3.append(RP(i))
65 for i in l3:
66     l4.append(RD(i))
67 for i in l4:
68     aux = i.split()
69     aux2 = []
70     for j in aux:
71         aux2.append(reducirString(j))

```

```

72     l5.append(' '.join(aux2))
73
74 nltk.download('stopwords')
75 from nltk.corpus import stopwords
76 lsw = stopwords.words('spanish')
77
78 TweetsLimpios = []
79
80 for i in l5:
81     aux = []
82     h = i.split()
83     for r in h:
84         if r not in lsw:
85             aux.append(r)
86     TweetsLimpios.append(' '.join(aux))
87
88 print('Data Trabajada')
89 print(len(users))
90 print(len(klases))
91 print(len(TweetsLimpios))
92
93 import numpy as np
94 print('Limpieza realizada')
95 print(np.unique(klases, return_counts=True))
96
97 import pandas as pd
98 dataE = pd.DataFrame({'Textos':TweetsLimpios,
99 'Klases':klases, 'Usuarios':users})
100
101 # Particion 70/30
102 from sklearn.model_selection import train_test_split
103
104 X_train, X_test, y_train, y_test = train_test_split(dataE, dataE['Klases'],
105 test_size=0.30)
106 print('Particion realizada')
107
108 print(dataE.groupby(['Klases']).count())
109 print(X_train.groupby(['Klases']).count())
110 print(X_test.groupby(['Klases']).count())
111
112 y_train = list(y_train)
113 y_test = list(y_test)
114
115 # TfIdf
116 XTM_train = list(X_train['Textos'])
117 XTM_test = list(X_test['Textos'])
118
119 from sklearn.feature_extraction.text import TfidfVectorizer
120 from sklearn.naive_bayes import MultinomialNB
121
122 #TfidfVectorizer - entrenamiento
123 vec = TfidfVectorizer()

```

```

124 X = vec.fit_transform(XTM_train)
125 mnb = MultinomialNB()
126 mnb.fit(X, y_train)
127
128 #Tfidf - prueba
129 XTFIDF = vec.transform(XTM_test)
130 mnb_ = mnb.predict(XTFIDF)
131 print('TfIdf realizado')
132
133
134 # microTC - TextModel
135 # TextModel - entrenamiento
136 from microtc.textmodel import TextModel
137 from sklearn.svm import LinearSVC
138
139 textmodel = TextModel(docs=None,
140 text='text',num_option='delete',usr_option='delete',
141 url_option='delete',emo_option='group',hashtag_option='delete',
142 ent_option='none',lc=True,del_dup=True,del_punc=True,del_diac=True,
143 token_list=[[2, 1], -1,3,4],token_min_filter=0,token_max_filter=1,
144 select_ent=False,select_suff=False, select_conn=False,
145 weighting='tfidf').fit(XTM_train)
146
147 lsvc = LinearSVC().fit(textmodel.transform(XTM_train), y_train)
148
149 #TextModel - prueba
150 lsvc_ = lsvc.predict(textmodel.transform(XTM_test))
151 print('microTC realizado')
152
153
154 # Metricas
155 from sklearn.metrics import f1_score, precision_score, recall_score,
156 accuracy_score
157
158 def metricas(modelo,y_test,predict):
159     print(modelo)
160 print('Precision:',round(precision_score(y_test,predict,average=None)[0],3),
161 'Macro:',round(precision_score(y_test, predict, average='macro'),3),
162 'Micro:',round(precision_score(y_test, predict, average='micro'),3))
163 print('Recall:',round(recall_score(y_test,predict,average=None)[0],3),
164 'Macro:',round(recall_score(y_test, predict, average='macro'),3),
165 'Micro:',round(recall_score(y_test, predict, average='micro'),3))
166 print('F1: ',round(f1_score(y_test, predict, average=None)[0],3),
167 'Macro:',round(f1_score(y_test, predict, average='macro'),3),
168 'Micro:',round(f1_score(y_test, predict, average='micro'),3))
169 print('Accuracy: ',round(accuracy_score(y_test, predict),3))
170
171 print(metricas('\ntfidf',y_test,mnb_))
    print(metricas('\nmicroTC',y_test,lsvc_))

# Confusion Matrix
from sklearn.metrics import confusion_matrix

```

```

print('\ntfidf\n', confusion_matrix(y_test, mnb_))
print('\nmicroTC\n', confusion_matrix(y_test, lsvc_))

```

**Código de Python utilizado para el análisis por usuarios utilizando microTC (SAP) y Fasttext e Infersent con SVM, Vecinos más cercanos y Tfidf. Utilizando limpieza previa.**

```

1 import tarfile
2 tar = tarfile.open("usuarios_depresion.json.tar.gz")
3 tar.extractall()
4 tar.close()
5
6 from microtc.utils import tweet_iterator
7 data = [i for i in tweet_iterator('usuarios_depresion.json')]
8
9 # Extraer los tweets
10 data[0]['tweets'][4099]['text']
11
12 # Extraer lo que se requiere, tres listas de dimension 317 en donde todos
13 los tweets de cada usuario se juntaron para hacer solo un tweet unico por
14 usuario.
15 clases, users, t, tweets = [], [], [], []
16 for i in data:
17     t = []
18     clases.append(int(i['klass']))
19     users.append(i['user'])
20     for j in i['tweets']:
21         t.append(j['text'])
22     tweets.append(' '.join(t))
23
24 print('Data Original')
25 print(len(clases))
26 print(len(tweets))
27 print(len(users))
28
29
30 # Limpieza
31 import re as reg
32 import unicodedata
33 import nltk
34 import re, string
35
36 textos = tweets.copy()
37
38 def QuitarAcentos(cadena):
39     formaNFKD = unicodedata.normalize('NFKD', cadena)
40     return u"".join([c for c in formaNFKD if not unicodedata.combining(c)])
41 def RP(text):

```

```

42     return re.sub('[%s]' % re.escape(string.punctuation), '', text)
43 def RD(digito):
44     return ''.join(i for i in digito if not i.isdigit())
45 def reducirString(palabra):
46     pattern=reg.compile(r"(\.|\{|\}|,)",reg.DOTALL)
47     string=pattern.sub(r"\1",palabra)
48     return string
49
50 l1,l2,l3,l4,l5 = [],[],[],[],[]
51
52 for i in textos:
53     l1.append(QuitarAcentos(i).lower())
54 for i in l1:
55     aux = []
56     h = i.split()
57     for r in h:
58         if (r.startswith('http') or r.startswith('rt') or
59 r.startswith('\"http')):
60             continue
61             aux.append(r)
62     l2.append(' '.join(aux))
63 for i in l2:
64     l3.append(RP(i))
65 for i in l3:
66     l4.append(RD(i))
67 for i in l4:
68     aux = i.split()
69     aux2 = []
70     for j in aux:
71         aux2.append(reducirString(j))
72     l5.append(' '.join(aux2))
73
74 nltk.download('stopwords')
75 from nltk.corpus import stopwords
76 lsw = stopwords.words('spanish')
77
78 TweetsLimpios = []
79
80 for i in l5:
81     aux = []
82     h = i.split()
83     for r in h:
84         if r not in lsw:
85             aux.append(r)
86     TweetsLimpios.append(' '.join(aux))
87
88 print('Data Trabajada')
89 print(len(users))
90 print(len(klases))
91 print(len(TweetsLimpios))
92
93 import numpy as np

```

```

94 print('Limpieza realizada')
95 print(np.unique(klases, return_counts=True))
96
97
98 # Tabla en pandas con la data
99 import pandas as pd
100 dataE = pd.DataFrame({'Textos':TweetsLimpios,
101 'Klases':klases,'Usuarios':users})
102
103
104 # Modelo FastText
105 import fasttext
106 ft = fasttext.load_model('crawl-300d-2M-subword.bin')#, encoding="latin1")
107 print('Se cargo el modelo de fasttext')
108
109 embeddingsEF = [ ft.get_sentence_vector(v) for v in list(dataE['Textos'])]
110 print('Se crearon los vectores de fasttext')
111
112 dataE = dataE.assign(fft=embeddingsEF)
113
114 # Modelo InferSent
115 from InferSent.models import InferSent
116 import torch
117
118 def InferSent_model(model_version=1):
119     MODEL_PATH = 'infersent1.pkl'
120     params_model = {'bsize': 64, 'word_emb_dim': 300, 'enc_lstm_dim': 2048,
121                    'pool_type': 'max', 'dpout_model': 0.0, 'version': model_version}
122     # If inferSent1 -> use GloVe embeddings. If inferSent2 -> use InferSent
123     embeddings.
124     if model_version == 1:
125         W2V_PATH = 'glove.840B.300d.txt'
126     else:
127         W2V_PATH = 'crawl-300d-2M.vec'
128     print("loading model ... {}".format(W2V_PATH))
129     model = InferSent(params_model)
130     model.load_state_dict(torch.load(MODEL_PATH))
131     model.set_w2v_path(W2V_PATH)
132
133     model.build_vocab_k_words(K=3000000)
134     return model
135
136 import nltk
137 nltk.download('punkt')
138 model = InferSent_model(model_version=1)
139 sent_detector = nltk.data.load('tokenizers/punkt/english.pickle')
140 print('Se cargo el modelo de inferSent')
141
142 embeddingsEI = [ model.encode(sent_detector.tokenize(v), bsize=128,
143                             tokenize=False, verbose=True) for v in list(dataE['Textos'])]
144 print('Se crearon los vectores de inferSent')
145

```



```

146 dataE = dataE.assign(infer = [i[0] for i in embeddingsEI] )
147
148 # Particion 70/30
149 from sklearn.model_selection import train_test_split
150
151 X_train, X_test, y_train, y_test = train_test_split(dataE, dataE['Klases'],
152 test_size=0.30)
153 print('Particion realizada')
154
155 print(dataE.groupby(['Klases']).count())
156 print(X_train.groupby(['Klases']).count())
157 print(X_test.groupby(['Klases']).count())
158
159 # Vectores fasttext e infersent
160 y_train = list(y_train)
161 y_test = list(y_test)
162
163 XF_train = list(X_train['fft'])
164 XF_test = list(X_test['fft'])
165
166 XI_train = list(X_train['infer'])
167 XI_test = list(X_test['infer'])
168
169 # SVM - fasttext/infersent
170 # Support Vector Machine (SVM)
171 from sklearn.svm import LinearSVC
172 # entrenamiento
173 svcF = LinearSVC(random_state=0)
174 svcI = LinearSVC(random_state=0)
175 svcF.fit(XF_train, y_train)
176 svcI.fit(XI_train, y_train)
177 # prueba
178 svcF_ = svcF.predict(XF_test)
179 svcI_ = svcI.predict(XI_test)
180 print('SVM de fasttext e infersent realizado')
181
182 # K-Vecinos - fasttext/infersent
183 # K-Vecinos Cercanos
184 from sklearn.neighbors import KNeighborsClassifier
185 # entrenamiento
186 knnF = KNeighborsClassifier(n_neighbors=5)
187 knnI = KNeighborsClassifier(n_neighbors=5)
188 knnF.fit(XF_train, y_train)
189 knnI.fit(XI_train, y_train)
190 # prueba
191 knnF_ = knnF.predict(XF_test)
192 knnI_ = knnI.predict(XI_test)
193 print('KNN de fasttext e infersent realizado')
194
195
196 # TfIdf
197 XTM_train = list(X_train['Textos'])

```

```

198 XTM_test = list(X_test['Textos'])
199
200 from sklearn.feature_extraction.text import TfidfVectorizer
201 from sklearn.naive_bayes import MultinomialNB
202
203 # TfidfVectorizer - entrenamiento
204 vec = TfidfVectorizer()
205 X = vec.fit_transform(XTM_train)
206 mnb = MultinomialNB()
207 mnb.fit(X, y_train)
208
209 # Tfidf - prueba
210 XTFIDF = vec.transform(XTM_test)
211 mnb_ = mnb.predict(XTFIDF)
212 print('TfIdf realizado')
213
214
215 # microTC - TextModel
216 # TextModel - entrenamiento
217 from microtc.textmodel import TextModel
218 from sklearn.svm import LinearSVC
219
220 textmodel = TextModel().fit(XTM_train)
221 lsvc = LinearSVC().fit(textmodel.transform(XTM_train), y_train)
222
223 # TextModel - prueba
224 lsvc_ = lsvc.predict(textmodel.transform(XTM_test))
225 print('microTC realizado')
226
227
228 # Metricas
229 from sklearn.metrics import f1_score, precision_score, recall_score,
230 accuracy_score
231
232 def metricas(modelo, y_test, predict):
233     print(modelo)
234     print('Precision:', round(precision_score(y_test, predict, average=None)[0], 3),
235           'Macro:', round(precision_score(y_test, predict, average='macro'), 3),
236           'Micro:', round(precision_score(y_test, predict, average='micro'), 3))
237     print('Recall:', round(recall_score(y_test, predict, average=None)[0], 3),
238           'Macro:', round(recall_score(y_test, predict, average='macro'), 3),
239           'Micro:', round(recall_score(y_test, predict, average='micro'), 3))
240     print('F1:      ', round(f1_score(y_test, predict, average=None)[0], 3),
241           'Macro:', round(f1_score(y_test, predict, average='macro'), 3),
242           'Micro:', round(f1_score(y_test, predict, average='micro'), 3))
243     print('Accuracy: ', round(accuracy_score(y_test, predict), 3))
244
245 print(metricas('\ntfidf', y_test, mnb_))
246 print(metricas('\nmicroTC', y_test, lsvc_))
247 print(metricas('\nSVC_FastText', y_test, svcF_))
248 print(metricas('\nSVC_InferSent', y_test, svcI_))
249 print(metricas('\nKNN_FastText', y_test, knnF_))

```

```

250 print(metricas('\nKNN_InferSent',y_test,knnI_))
251
252
253 # Confusion Matrix
254 from sklearn.metrics import confusion_matrix
    print('\ntfidf\n',confusion_matrix(y_test, mnb_))
    print('\nmicroTC\n',confusion_matrix(y_test, lsvc_))
    print('\nSVC_FastText\n',confusion_matrix(y_test, svcF_))
    print('\nSVC_InferSent\n',confusion_matrix(y_test, svcI_))
    print('\nKNN_FastText\n',confusion_matrix(y_test, knnF_))
    print('\nKNN_InferSent\n',confusion_matrix(y_test, knnI_))

```

## Código de Python utilizado para el análisis por usuarios utilizando microTC (SAP y SEP). Sin limpieza previa.

```

1 import tarfile
2 tar = tarfile.open("usuarios_depresion.json.tar.gz")
3 tar.extractall()
4 tar.close()
5
6 from microtc.utils import tweet_iterator
7 data = [i for i in tweet_iterator('usuarios_depresion.json')]
8
9 # Extraer lo que se requiere, tres listas de dimension 317 en donde todos
10 los tweets de cada usuario se juntaron para hacer solo un tweet unico por
11 usuario.
12 clases,users,t,tweets = [],[],[],[]
13 for i in data:
14     t = []
15     clases.append(int(i['klass']))
16     users.append(i['user'])
17     for j in i['tweets']:
18         t.append(j['text'])
19     tweets.append(' '.join(t))
20
21 print('Data Original')
22 print(len(clases))
23 print(len(tweets))
24 print(len(users))
25
26 import numpy as np
27 print(np.unique(clases, return_counts=True))
28
29 # Tabla en pandas con la data
30 import pandas as pd
31 dataE = pd.DataFrame({'Textos':tweets, 'Klases':clases,'Usuarios':users})
32
33 # Particion 70/30
34 from sklearn.model_selection import train_test_split

```

```

35
36 X_train, X_test, y_train, y_test = train_test_split(dataE, dataE['Klases'],
37 test_size=0.30)
38 print('Particion realizada')
39
40 print(dataE.groupby(['Klases']).count())
41 print(X_train.groupby(['Klases']).count())
42 print(X_test.groupby(['Klases']).count())
43
44 y_train = list(y_train)
45 y_test = list(y_test)
46 XM_train = list(X_train['Textos'])
47 XM_test = list(X_test['Textos'])
48
49 # microTC - TextModel
50 # Params editados
51 # TextModel - entrenamiento
52 from microtc.textmodel import TextModel
53 from sklearn.svm import LinearSVC
54
55 textmodel = TextModel(docs=None,
56 text='text', num_option='delete', usr_option='delete',
57 url_option='delete', emo_option='group', hashtag_option='delete',
58 ent_option='none', lc=True, del_dup=True, del_punc=True, del_diac=True,
59 token_list=[[2, 1], -1, 3, 4], token_min_filter=0, token_max_filter=1,
60 select_ent=False, select_suff=False, select_conn=False,
61 weighting='tfidf').fit(XM_train)
62
63 lsvc = LinearSVC().fit(textmodel.transform(XM_train), y_train)
64
65 # TextModel - prueba
66 lsvc_ = lsvc.predict(textmodel.transform(XM_test))
67 print('microTC params realizado')
68
69 # Params default
70 # TextModel - entrenamiento
71 textmodelD = TextModel().fit(XM_train)
72
73 lsvcD = LinearSVC().fit(textmodelD.transform(XM_train), y_train)
74
75 # TextModel - prueba
76 lsvcD_ = lsvcD.predict(textmodelD.transform(XM_test))
77 print('microTC default realizado')
78
79
80 # Metricas
81 from sklearn.metrics import f1_score, precision_score, recall_score,
82 accuracy_score
83
84 def metricas(modelo, y_test, predict):
85     print(modelo)
86     print('Precision:', round(precision_score(y_test, predict, average=None)[0], 3),

```

```

87     'Macro:', round(precision_score(y_test, predict, average='macro'), 3),
88     'Micro:', round(precision_score(y_test, predict, average='micro'), 3))
89     print('Recall:', round(recall_score(y_test, predict, average=None)[0], 3),
90           'Macro:', round(recall_score(y_test, predict, average='macro'), 3),
91           'Micro:', round(recall_score(y_test, predict, average='micro'), 3))
92     print('F1:          ', round(f1_score(y_test, predict, average=None)[0], 3),
93           'Macro:', round(f1_score(y_test, predict, average='macro'), 3),
94           'Micro:', round(f1_score(y_test, predict, average='micro'), 3))
95     print('Accuracy: ', round(accuracy_score(y_test, predict), 3))
96
97     print(metrics('\nmicroTC_params', y_test, lsvc_))
98     print(metrics('\nmicroTC_default', y_test, lsvcD_))
99
100    # Confusion Matrix
101    from sklearn.metrics import confusion_matrix
102    print('\nmicroTC_params\n', confusion_matrix(y_test, lsvc_))
103    print('\nmicroTC_default\n', confusion_matrix(y_test, lsvcD_))

```