

## Generación de un corpus lingüístico digital en español enfocado a la depresión

César-Jesús Núñez-Prado<sup>1,2</sup>, Claudia Talavera Ortega<sup>1</sup>,  
Liliana Chanona-Hernández<sup>1</sup>, Grigori Sidorov<sup>2</sup>

<sup>1</sup> Instituto Politécnico Nacional,  
Escuela Superior de Ingeniería Mecánica y Eléctrica,  
México

<sup>2</sup> Instituto Politécnico Nacional,  
Centro de Investigación en Computación,  
México

{cesar.jnprado, claudiatalaveraor, lchanona}@gmail.com,  
sidorov@cic.ipn.mx

**Resumen.** La depresión es un desorden mental que afecta a miles de personas alrededor del mundo y que de no ser detectada a tiempo puede llevar a un desenlace mortal como el suicidio. Este tipo de trastorno es considerado como un asesino silencioso debido a que no es de fácil identificación ya que los individuos con esta clase de desorden generalmente intentan ocultarlo, por lo que cualquier tipo de investigación que provea herramientas para su detección temprana siempre será de gran utilidad. Por otra parte, hoy en día, las redes sociales son una fuente ininterrumpida de información que puede ser analizada con la intención de reconocer patrones asociados a la resolución de alguna tarea de investigación, como por ejemplo el análisis de emociones. En esta investigación buscamos crear un corpus (banco de información) que contenga mensajes publicados en la red social *Twitter* y cada mensaje estará asociado con dos posibles clasificaciones: mensaje depresivo o mensaje no depresivo.

**Palabras clave:** Depresión, análisis de emociones, reconocimiento de patrones, corpus.

### Compilation of a Digital Textual Corpus in Spanish Focused on Depression

**Abstract.** Depression is a mental disorder that affects thousands of people around the world and if it is not detected in time, it can lead to a fatal outcome such as suicide. This type of disorder is considered a silent killer because it is not easy to identify, since individuals with this type of disorder generally try to hide it, so any type of research that provides tools for early detection will always be very useful. On the other hand, nowadays, social networks are an uninterrupted source of information that can be analyzed with the intention of recognizing patterns associated with the resolution of some research task, such as the analysis of emotions. In this research we seek to create a corpus (information bank) that

contains messages published on the Twitter social network and each message is associated with two possible classifications: depressive message or non-depressive message.

**Keywords:** Depression, analysis of emotions, recognizing patterns, corpus.

## 1. Introducción

De acuerdo con el psicoanalista Erich Seligmann Fromm<sup>1</sup>, en el momento en el que el ser humano se separó de la naturaleza a través de percibir su propia autoconciencia, se dio paso a la generación de todo el espectro de sentimientos conocidos, entre los que se encuentran; el amor, el odio, la tristeza, la euforia, la envidia, el enojo, la impaciencia, la satisfacción, la culpa, la preocupación, entre algunas. Tales sentimientos repercuten de manera directa en el estado de ánimo y la manera conductual de las personas, es decir; sentimientos positivos serían asociados con conductas y estados de ánimo relacionados con la felicidad, mientras que, por otro lado, los sentimientos negativos son capaces de provocar disociaciones en los individuos.

Uno de los desórdenes mentales más relacionados con los sentimientos de impacto negativo es la depresión, [1] Fromm la define como la incapacidad de sentir alegría o tristeza, es opuesta a la razón debido a que impide el desenvolvimiento humano y por tanto, irracional ya que quien la presenta rehúye a experimentar ese tipo de sentimientos necesarios para que en un futuro pueda alcanzar un crecimiento completo.

La depresión es considerada como un trastorno mental común que afecta tanto a hombres como a mujeres, sin importar si se trata de un adolescente, adulto o adulto mayor y puede desembocar en problemáticas aún más grandes como lo pueden ser las autolesiones, tendencias suicidas o finalmente, cometer suicidio. Lamentablemente, las personas que sufren de este tipo de trastorno experimentan dificultades para poder expresar libremente y de manera cómoda sus sentimientos con otras personas, especialmente con sus familias y generalmente acuden por ayuda cuando el episodio depresivo ya es muy grave.

Según la Encuesta Nacional de Epidemiología Psiquiátrica (ENEP)<sup>2</sup>, entre los años 2001 y 2002 se reportó que un 9.2 % de la población en México había padecido algún tipo de trastorno depresivo durante su vida y desafortunadamente esta cifra continúa incrementándose ya que en el año 2022 se publicó el 2º. Diagnóstico operativo de salud mental<sup>3</sup> y adicciones realizado por los Servicios de Atención Psiquiátrica (SAP) de la Secretaría de Salud en México, en donde se estima que la depresión es uno de los trastornos más frecuentes en la población que cuenta con derechohabiencia con más de 3 millones de casos confirmados.

La derechohabiencia indica que las personas tienen algún servicio de salud público como el IMSS o el ISSSTE. En algunas ocasiones, las personas que presentan tendencias depresivas o suicidas, se apoyan del uso de las redes sociales, ya que encuentran en ellas un medio para poder expresarse con confianza, de una manera libre

<sup>1</sup> Erich Seligmann Fromm fue un psicoanalista, psicólogo social y filósofo alemán.

<sup>2</sup> Véase: <https://www.insp.mx/avisos/sintomas-depresivos-y-atencion-a-la-depresion>

<sup>3</sup> Véase: <https://www.gob.mx/salud/prensa/008-en-mexico-3-6-millones-de-personas-adultas-padecen-depresion>

y con menos prejuicios, dada la posibilidad de desarrollar una identidad digital y con ello poder mantenerse en el anonimato.

Para este proyecto, nosotros utilizaremos la red social *Twitter* debido a que la longitud de los mensajes no sobrepasa los 250 caracteres y ello permite realizar el procesamiento y clasificación de los mensajes en un tiempo menor, y también porque la cuenta de desarrollador de la red social, otorga los permisos necesarios para realizar la descarga de los mensajes publicados.

Consideramos que es de vital importancia el generar recursos lingüísticos enfocados a poder determinar, si algún texto publicado en redes sociales contiene un posible discurso depresivo o suicida, con la intención de poder utilizar a favor de la sociedad, el poder de cómputo que se sigue desarrollando hasta nuestros días y con ello alertar a las personas que sufren de este tipo de trastorno con la finalidad de reducir el índice de depresión y suicidio dentro de la sociedad.

La estructura de la presente investigación mostrará los trabajos relacionados, la aplicación de la metodología empleada para la generación del corpus, los resultados obtenidos, conclusiones y trabajo a futuro.

## 2. Trabajos relacionados

En esta sección se describirán algunas investigaciones a fin al objetivo principal de nuestro trabajo.

En [2] se define a un corpus como un conjunto de textos producidos en condiciones naturales, representativos de una lengua, almacenado en formato electrónico y codificado con la intención de ser analizados científicamente. Refiere que la condición natural se establece, cuando existe la intención real de comunicar una idea y no ser concebidos para ilustrar algún fenómeno lingüístico. Identifica que el término representativo se da cuando se respeta un momento determinado de la historia.

Los corpus actuales contienen miles de millones de formas y el único modo en el que se puede recuperar esta cantidad de información es de manera electrónica. Dentro de la información que proporciona un corpus se puede encontrar a los metadatos, los cuales son información interna que incluye, el nombre de la fuente de información, fecha y lugar de la publicación, editorial, entre algunos.

En [3] crearon un corpus etiquetado para el español con la finalidad de realizar análisis sobre sentimientos. Decidieron realizar la descarga de mensajes publicados en la red social *Twitter* y para considerar que mensajes se procesarían, estos debían contar con por lo menos 5 palabras; para el etiquetado de las palabras utilizaron *Freeling* y excluyeron las palabras de parada con la lista que proporciona la biblioteca *Natural Language Toolkit (NLTK)*. Con dichas especificaciones, generaron un corpus con más de 20 mil mensajes que fueron clasificados de acuerdo a los *hashtags*<sup>4</sup> contenidos en cada mensaje. Las emociones que identificaron fueron la alegría, asco, tristeza, ira, miedo y sorpresa.

Ya combinando la creación de un corpus y el tema de la depresión en [4] desarrollaron un corpus de mensajes de ideación suicida extraídos de la *web* y la *deep web* tanto en español (33 %), como en inglés (67 %). Algunas de las categorías con las

---

<sup>4</sup> Hashtag es el símbolo «#».

que realizaron la clasificación de los textos incluye: depresión, ironía, tristeza, auto-pro-suicida e indefinido y el tamaño final del corpus generado contiene más de 7 mil *tokens*.

### 3. Aplicación de la metodología

En la siguiente sección se mostrará la metodología empleada dentro de la presente investigación.

#### 3.1. Fase 1

En esta fase inicial se generó una cuenta regular de usuario de *Twitter*, en donde la red social solicitó información personal tal como; nombre completo, nombre de usuario (es el nombre que aparece visible para los demás usuarios de la red social), fecha de nacimiento, cuenta de correo electrónico y número de teléfono celular. Una vez que se llenaron todos los campos requeridos, la red social envió un correo electrónico de verificación y con ello se completó con éxito la apertura de la cuenta.

#### 3.2. Fase 2

Para poder acceder a la información disponible en la red social es indispensable contar con una cuenta de desarrollador desde la *API* de *Twitter* (interfaz de programación de aplicaciones de *Twitter*), para solicitarla se debe generar una aplicación y agregar la información de los campos solicitados entre los que se incluyen:

- Pertenece a una institución educativa o a una empresa (para instituciones educativas no hay costo, pero se limita la cantidad de información disponible).
- ¿Cuál es el uso que le darás a la información?
- ¿Qué clase de algoritmos aplicarás?
- ¿Qué resultados esperas obtener?

La apertura de la cuenta de desarrollador no es automática, el personal de *Twitter* se pone en contacto con el cliente en busca de obtener información más detallada con un intercambio de correos electrónicos, que en nuestro caso duró casi 15 días (cabe mencionar que todos los correos son enviados en inglés).

#### 3.3. Fase 3

De manera paralela, mientras se estaba llevando a cabo la fase anterior, se realizó una búsqueda en internet de las palabras más utilizadas y asociadas a la tristeza, depresión y suicidio. Se encontraron 3 listas de las cuales se aplicó la unión de los elementos entre ellas para generar una lista con entradas sin repetir. Sobre cada elemento de la lista final se aplicó la lematización (el cual es el proceso de encontrar las palabras sin flexionar, por ejemplo; el lema de la palabra «abrumado» es «abrumar») para obtener la normalización de las entradas. En la Tabla 1 se presenta una muestra de la lista de palabras depresivas.

**Tabla 1.** Muestra de la lista de palabras depresivas.

Morir	Abrumar	Miseria
Tristeza	Suicidar	Odiar
Terapia	Torturar	Terminar

### 3.4. Fase 4

La cuenta de desarrollador proporciona 4 claves secretas, únicas y renovables para poder realizar la conexión entre diferentes lenguajes de programación y *Twitter*. Nosotros desarrollamos un código fuente en *Python*<sup>5</sup> y apoyándonos de la biblioteca *Tweepy*<sup>6</sup> utilizamos las claves para realizar la conexión y con ello poder acceder a la información disponible en la plataforma, entre la que se encuentra:

- Texto publicado (con dos posibilidades, texto completo o texto parcial).
- Nombre del usuario.
- Fecha y hora de publicación.
- Ubicación de la publicación.
- Cantidad de «me gusta».
- Cantidad de comentarios en la publicación.
- Cantidad de «*retweets*»<sup>7</sup>.

Una ventaja de utilizar *Tweepy* es que contiene un módulo de búsqueda especializada en el que se debe estipular la palabra a buscar (*query*), el periodo de tiempo (ejemplo: mayo 2022) y el número de mensajes a recuperar. Este último punto es muy importante porque la cuenta de desarrollador es académica y únicamente permite la descarga de 800 mensajes por periodos de 30 minutos, si esta cantidad se excede, entonces *Twitter* congela las credenciales y se debe esperar de 2 a 3 días para que las activen de nuevo.

Siendo cuidadosos con la cantidad de mensajes y los periodos de tiempo, nosotros estuvimos descargando mensajes durante 30 días, utilizando como palabras de búsqueda, las palabras de la lista de la Fase 3 y con ello se realizó la descarga de un poco más de 12 mil mensajes.

### 3.5. Fase 5

De igual forma, se llevaron de manera casi paralela la Fase 4 y la Fase 5. Los mensajes descargados de *Twitter* se escribieron en archivos de texto plano y a todos estos mensajes se les hizo un procesamiento automático el cual consistió en lo siguiente:

- Conversión del texto a minúsculas.
- Tokenización (es el proceso mediante el cual se obtienen las unidades mínimas de las oraciones, es decir; se separan las oraciones en palabras).
- Eliminación de los «*retweets*» para procesar mensajes únicos.

<sup>5</sup> Se utilizó la versión 3.9.12.

<sup>6</sup> Se utilizó la versión 4.10.0.

<sup>7</sup> Un «*retweet*» es la acción de publicar en tu propio muro la publicación de otro usuario y se identifica por las letras «RE» al inicio de cada mensaje.

- Eliminación de enlaces a internet.
- Eliminación de imágenes.
- Eliminación de emoticones.
- Eliminación de signos de puntuación.
- Eliminación de caracteres numéricos.
- Lematización.
- Etiquetado de las palabras (part of speech tagging)<sup>8</sup>.
- Eliminación de stop-words<sup>9</sup>.
- Eliminación de nombres propios.
- Eliminación de mensajes repetidos.

Cabe mencionar que la gran diversidad de emoticones disponibles en internet representa un esfuerzo mayúsculo para el etiquetado manual y es por ello que se decidió eliminarlos de los mensajes descargados.

Una vez que se realizó este procedimiento sobre todos los mensajes descargados de *Twitter*, se creó un banco de datos con un total de 1,623 mensajes únicos.

### 3.6. Fase 6

En esta última fase, se pidió el apoyo de 5 psicólogas para realizar la evaluación de los mensajes, tomando como único criterio dos posibilidades, mensaje depresivo o mensaje no depresivo.

El perfil de las sicólogas cumplió los siguientes 3 aspectos:

- Experiencia mínima de 5 años ininterrumpidos.
- Atención a pacientes con depresión o tendencia suicida.
- Uso y conocimiento de redes sociales.

Cada sicóloga debía leer mensaje por mensaje y de acuerdo a su propia experiencia de consulta, debía clasificar cada uno de los mensajes. El periodo de revisión duró 2.5 meses.

Una vez que se contó con los archivos calificados, se inició el proceso de agregar los mensajes al corpus, para ello se consideró la siguiente métrica:

- Si por lo menos 3 sicólogas calificaron como mensaje depresivo al mensaje, se asignaba la clase 1.
- Si por lo menos 3 sicólogas calificaron como mensaje no depresivo al mensaje, se asignaba la clase 0.

En la Tabla 2 se muestran algunos mensajes descargados y etiquetados por los sicólogos. En la columna de la izquierda se visualiza el mensaje y en la columna de la derecha la etiqueta final de cada mensaje.

En la Tabla 3 se visualiza el proceso de clasificación de los mensajes, cada uno de los sicólogos evaluó cada uno de los mensajes y se realizó el conteo de los votos. Cuando un mensaje contenía por lo menos 3 votos de una clase, entonces se le asignó esa clase final al mensaje.

<sup>8</sup> Ejemplo de etiquetado: si aparece la palabra «correr» la etiqueta asociada es «verbo».

<sup>9</sup> *Stop-words* es el conjunto de palabras que incluyen las palabras gramaticales como los artículos definidos, indefinidos, preposiciones, etc.

**Tabla 2.** Muestra de mensajes clasificados

Mensaje	Etiqueta
Sonreír siempre ha sido más fácil que explicar por qué estoy triste, es mejor morir en silencio	1
Sé que estoy sólo incluso cuando me encuentro rodeado de tantas personas.	1
La depresión ha sido mi fiel compañera desde que te marchaste y sé que me acompañará hasta el fin de mi vida	1
Me siento tan triste porque mi América perdió, pero aun así estoy alegre de gritar al mundo que soy lgbt	0

**Tabla 3.** Muestra de las evaluaciones de los psicólogos

Eval 1	Eval 2	Eval 3	Eval 4	Eval 5	Clase
1	1	1	1	1	1
1	1	0	1	0	1
1	1	1	1	1	1
1	0	0	0	0	0

#### 4. Conclusiones y trabajo a futuro

Con el desarrollo de esta investigación se pudo crear un corpus lingüístico digital enfocado a la depresión que cuenta con 1,623 mensajes con dos clasificaciones; 1 si el mensaje es depresivo y 0 si el mensaje no es depresivo y se encuentra almacenado en un archivo de texto plano.

El corpus se encuentra desbalanceado, ya que contiene 985 mensajes clasificados como no depresivos y 638 mensajes clasificados como depresivos. Esta herramienta se pondrá a disposición de investigaciones académicas de manera gratuita vía correo electrónico.

Como trabajo a futuro, estamos pensando en dos posibilidades a corto plazo, por una parte, es importante balancear el corpus, por lo que es necesario seguir descargando más mensajes y solicitar a nuevos psicólogos que nos aporten con su conocimiento y experiencia en la clasificación de los mensajes y, por otra parte, es importante aplicar métodos de inteligencia artificial sobre el corpus generado para saber si es posible que identifiquen un patrón a través de los mensajes.

#### Referencias

1. Moreno-López, S.: La condición humana según Erich Fromm. Pensamiento, Papeles de filosofía, vol. 3, pp 151–171 (2017)
2. Rojo, G. Introducción a la lingüística de corpus en español. New York: Routledge (2021)
3. Sidorov, G., Galicia, S. N., Camacho, V. A.: Construcción de un corpus marcado con emociones para el análisis de sentimientos en Twitter en español. Revista Escritos, BUAP, vol. 1, no. 1 (2016)

4. Zafra, S., Gómez, J. M., Navarro-Colorado, B.: Diseño, compilación y anotación de un corpus para la detección de mensajes suicidas en redes sociales. *Procesamiento de lenguaje natural*, vol. 59, pp. 65–72 (2017)
5. Zucco, C., Calabrese, B., Cannataro, M.: Sentiment analysis and affective computing for depression monitoring. In: *IEEE International Conference on Bioinformatics and Biomedicine*, pp. 1988–1995 (2017) doi: 10.1109/BIBM.2017.8217966
6. Mustafa, R., Ashraf, N., Shabbir, F., Fersund, J., Shahzad B., Gelbukh, A.: A multiclass depression detection in social media based on sentiment analysis. In: *17th International Conference on Information Technology - New Generations*, vol. 1134, pp. 659–662 (2020) doi: 10.1007/978-3-030-43020-7\_89
7. Ameer, I., Arif, M., Sidorov, G., Gómez-Adorno, H., Gelbukh, A.: Mental illness classification on social media texts using deep learning and transfer learning. In: *8th World Conference on Soft Computing (2022)* doi: 10.48550/arXiv.2207.01012
8. Bird, S., Klein, E., Loper, E.: *Natural language processing with Python* (2019)