



SENA PASCUAL-LAVILLA | P.J. MULAS CÁMARA | R. FERNÁNDEZ-CALVILLO CÁCERES | C. MARTÍNEZ CABEZALI  
ME. MOLINA CAÑIZARES | EMILIA DE LOS ÁNGELES IGLESIAS ORTUÑO | MARÍA CONCEPCIÓN ARROYO  
PERLA VANESSA DE LOS SANTOS | GERARDO VÉLEZ VILLAFañE

# Análisis de los datos obtenidos de la red social Twitter para la identificación precoz de la tendencia al suicidio de los usuarios

## Analysis of data obtained from the social network Twitter for the early identification of users' suicidal tendencies

PJ. Mulas Cámara, R. Fernández-Calvillo Cáceres, C. Martínez Cabezali y ME. Molina Cañizares\*

\* Procesado Masivo de Datos, Universidad Rey Juan Carlos, Madrid, España,  
pj.mulas, ra.fernandezcal, c.martinezcab, me.molina. 2018@alumnos.urjc.es

---

### Abstract:

Although not everyone is aware of it, data available on the Internet are very useful and have a great potential to help our society. The digital platform Twitter is a social network where people sometimes express their feelings and emotions. And this paper arises from the idea of doing an analysis of these data through a Machine Learning tool, to find a psychiatric picture of depression, and if it is possible, the associated suicidal tendency. Twitter data extraction tool has been Tweepy, and with the profile data users, it has been made, an excel database that collects the information. Next, with the Machine Learning tool called UMAP, an unsupervised analysis of the database has been carried out, thanks to which it has been possible to differentiate three groups, with a very low inter cluster distance, which suggest that each observation looks a lot like its neighbors. From these three groups, we find one which behavior or use of the platform would be associated with a normal or standard way. The two other two group of meet part of the characteristics associated with depression.

**Keywords:** Twitter, suicide, Machine Learning, UMap, Tweepy, unsupervised, cloud, early detection.

---

### Resumen:

Aunque no todo el mundo sea consciente todos los datos disponibles en la red son útiles y tienen un gran potencial de ayuda a nuestra sociedad. La plataforma digital Twitter es una red social donde en ocasiones las personas expresan sus sentimientos y emociones, y este proyecto surge de la idea de hacer un análisis de estos datos de que se pueda realizar a través de una herramienta de Machine Learning un perfil típico un cuadro psiquiátrico de depresión, y si es posible la tendencia al suicidio asociada.

La herramienta para la extracción de datos de Twitter utilizada ha sido Tweepy, y con los usuarios obtenidos con ésta y las características definidas para ellos, se ha generado una

base de datos en formato excel en la nube One Drive que recoge toda esta información. A continuación, con la herramienta de Machine Learning llamada UMAP, se ha realizado un análisis de forma no supervisada de la base de datos, gracias al cuál se han podido diferenciar tres grupos, con una distancia intercluster muy baja, lo que quiere decir que cada observación se parece mucho a sus vecinos. De estos tres grupos hay uno al que se asociaría una conducta o uso de esta plataforma de una forma normal o estándar, y otros dos de diferente dimensión que cumplen parte de las características asociadas al trastorno de depresión.

**Palabras clave:** Twitter, suicidio, Machine Learning, UMap, Tweepy, no supervisado, nube, detección precoz.

---

### **Article info:**

*Received:* 19/10/2021 / *Received in revised form:* 01/12/2021

*Accepted:* 15/03/2022 / *Published online:* 03/02/2023

DOI: 10.5944/comunitania.24.2

---

## **1. Introducción**

Actualmente vivimos en un mundo globalizado en el plano económico, político, social gracias al avance tecnológico. Esto implica que parte de los aspectos de nuestra vida (por no decir toda ella) queden reflejados en la red.

Si hablamos del plano social la clave está en las redes sociales o cualquier plataforma digital que facilite la conexión entre personas. No somos conscientes de la cantidad de datos que dejamos en la red debido al uso de estas plataformas, como pueden ser nuestras emociones, nuestros recuerdos, nuestros sueños, hobbies o gustos, quiénes son nuestros seres queridos, los sitios que frecuentamos... Y aunque lo que más resuena en nuestras cabezas es una connotación negativa de este hecho, lo cierto es que también tiene un gran potencial positivo que bien usado puede ayudar a sus usuarios.

En concreto Twitter es una red social que se caracteriza por sus textos cortos y concisos a través de los cuáles reportar información de forma pública principalmente. Este proyecto surge de la idea de realizar un análisis del contenido de los usuarios de forma que se pueda realizar a través de una herramienta de Machine Learning un perfil típico un cuadro psiquiátrico de depresión, y si es posible la tendencia al suicidio asociada, y así se pueda realizar una detección precoz de este problema en usuarios concretos.

---

## 2. Metodología

El proceso de elaboración del proyecto se ha realizado en tres etapas que se definirán a lo largo del artículo. En primer lugar, se lleva a cabo una extracción de características, en segundo lugar, se genera una base de datos que es posteriormente almacenada y, por último, se aplican las herramientas de Machine Learning pertinentes para el análisis de los datos.

### 2.1. Tweepy

Antes de comenzar con el proceso de selección de características se hace una mención especial a la herramienta utilizada para la extracción de datos de Twitter. Tweepy es una API (Application Programming Interfaces) que, conectándose a Twitter a través de las credenciales de un usuario, permite obtener recursos de esta red social de forma fácil y rápida mediante Python [1]. El usuario utilizado se ha creado desde cero y se le ha dado de alta como Twitter Developer de nivel Elevate [2] con el fin de aumentar el número de peticiones realizadas por ejecución y así aumentar el volumen de datos obtenidos.

### 2.2. Selección de características

Para comenzar se va a realizar una primera filtración de usuarios, esto se realiza con el fin de focalizar a los sujetos, y así, partir de usuarios que potencialmente puedan tener signos de alguna afección psicológica que los lleve a tener pensamientos suicidas. Esta primera criba no puede ser tampoco muy exhaustiva ya que se busca tener perfiles con alta variabilidad entre ellos para poder agruparlos posteriormente. Teniendo en cuenta lo expuesto anteriormente, el primer paso a realizar es una búsqueda de usuarios que en un periodo corto de tiempo desde la ejecución del programa hayan utilizado la palabra "die" en alguno de sus tuits, para ello se buscan publicaciones con esa palabra y se obtiene la Id de usuario. La búsqueda va a realizarse en inglés ya que es el lenguaje predominante en Twitter.

Una vez obtenida la Id, se obtiene el historial del sujeto del que se extraen las características. Para ello y debido a las limitaciones que expone Twitter para extraer datos, se adquieren 1000 tuits por usuario, lo que ya es un perfil bastante completo de una persona y su actividad. Además, en esta búsqueda ya se realiza la extracción de dos características, el número de followers y amigos.

A continuación, se realizará el estudio del perfil del usuario, para ello se extraerán características que provienen directamente del análisis de los tuits, y se seguirán pautas que se han llevado a cabo en estudios similares al expuesto en este trabajo. Primeramente, se analiza la longitud de cada uno de los tuits y se realizará la media y varianza

con todos los tuits, esto se realiza debido a que los usuarios con depresión tienden a realizar tuits más cortos y menos elaborados debido a su falta de interés generada [3]. En segundo lugar, se realiza un conteo de palabras claves en todos los tuits, estas palabras claves son palabras relacionadas con pensamientos suicidas, emociones de carácter negativo o tendencia a generalizar, presentes en sujetos con depresión o enfermedades similares que pueden buscar el anonimato de Twitter como una vía de desahogo emocional y una forma de canalizar y exponer sus sentimientos, ya que los usuarios piensan que aquí no están tan solos y se sienten apoyados [4]. Estas palabras se han agrupado por similitud contextual, y se han agrupado siguiendo un patrón expuesto en la tabla 1.

**TABLA 1. Tabla de palabras clave**

CATEGORÍA	PALABRAS BUSCADAS
<b>SENTIMIENTOS NEGATIVOS Y PENSAMIENTOS SUICIDAS</b>	
TRISTEZA	Sad, sadness
PENSAMIENTOS SUICIDAS	Dying, kill myself, disappear
DEPRESIÓN	Depressed, depression, depressive
LLORAR	Cry
DORMIR	Sleep, nightmare
SOLEDAD	Alone, solo, loneliness,
CANSANCIO	Tired
<b>GENERALIZACIONES</b>	
NADIE	Nobody, no one, anybody
NADA	Nothing, anything
NUNCA	Never, ever
SIEMPRE	Always, forever
TODO	Everyone, everything, all

Por último, se elaborará un perfil temporal del individuo, es decir, se analizarán parámetros relativos al tiempo con respecto a las publicaciones. Una de las características que se obtienen es la frecuencia de publicación, para ello se observa la diferencia de tiempo entre tuits y se hace una media. Se extraen además el número de tuits por día de la semana ya que perfiles de personas depresivas suelen ser más activas los fines de semana porque suelen disminuir sus relaciones sociales y su actividad de ocio [3,5]. Siguiendo el mismo criterio se han observado las horas a las que publican los usuarios, debido a que los sujetos con algún trastorno psicológico de los mencionados tienden a presentar como síntoma la somnolencia y por tanto suele ser más activo en horario nocturno, entre las 22:00-8:00 [3].

### 2.3. Almacenamiento en la nube

El algoritmo desarrollado selecciona 20 usuarios, de los cuales escoge 1000 tweets para realizar el análisis. De esta manera, creamos una base de datos con únicamente 20 observaciones y las 29 características mencionadas en el apartado anterior. Puesto que se está trabajando con datos que no están etiquetados, el método a utilizar es aprendizaje no supervisado. Con ello, el objetivo es formar grupos de observaciones que poseen características similares, e identificar si estos corresponden a personas con tendencia al suicidio o no.

Para obtener resultados consistentes y evitar el 'underfitting' es necesario añadir más observaciones a nuestro conjunto de datos. Para ello, se debe aumentar la base de datos

añadiendo usuarios de 20 en 20. Puesto que el número de datos con el que se trabaja es bastante elevado: 20000 cada vez que se añade un conjunto nuevo a nuestra base, el tipo de almacenamiento a utilizar debe tener memoria suficiente para cargarlos.

En este caso, se ha optado por el almacenamiento en la nube, concretamente OneDrive. Este centro de archivos otorga al usuario registrado 5120 GB, espacio suficiente para tratar con BigData. Además, permite acceder a los archivos desde cualquier dispositivo y trabajar en tiempo real con aplicaciones como Word o Excel garantizando su seguridad.

Primero, se debe sincronizar la nube de OneDrive con el equipo. Esto permite trabajar los archivos almacenados en la nube utilizando el *path* que tienen en el ordenador. El siguiente paso es incluir un archivo Excel, en el que estén los 20 primeros usuarios, sobre el cual se irán añadiendo los demás. Por último, se crearán subconjuntos de 20 usuarios para incluirlos en el conjunto principal y así obtener la base de datos final.

Finalmente, la base de datos consta de 300 usuarios, con 1000 tweets de cada uno, y 29 características.

#### 2.4. UMAP

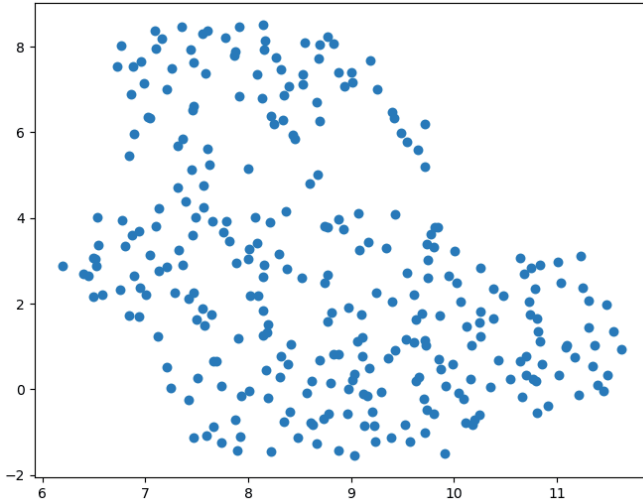
Al no disponer de una etiqueta, se ha optado por emplear aprendizaje no supervisado, de forma que se obtengan grupos de usuarios similares y se pueda analizar que características comparten entre sí, con el fin de diferenciar aquellos con tendencias suicidas claras.

El algoritmo elegido ha sido UMAP (Uniform Manifold Approximation and Projection). Este se basa en una reducción de la dimensionalidad buscando equivalentes topológicos del conjunto de datos. La reducción de la dimensionalidad que permite observar en un espacio 2D la distribución de los usuarios e intuir de forma subjetiva los grupos que se formarán.

A la hora de realizar el clustering propiamente dicho, UMAP se basa en el algoritmo HDBSCAN, que toma los datos con la reducción de dimensionalidad proporcionada por UMAP y los agrupa atendiendo a parámetros que se han establecido:

- Número de vecinos que escanear con respecto a cada muestra, que se han fijado en 50.
- La métrica empleada para el cálculo de distancias, que en este caso la escogida ha sido "Manhattan".
- La mínima distancia entre dos clústeres, 0.5 ha sido la seleccionada.

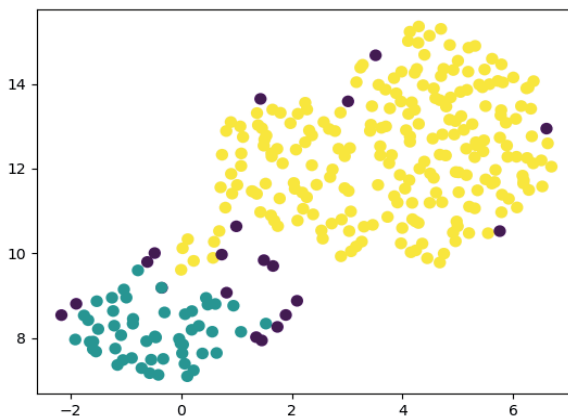
**FIGURA 1.: Representación de las muestras tras la reducción de dimensionalidad de UMAP. Cada punto es un usuario**



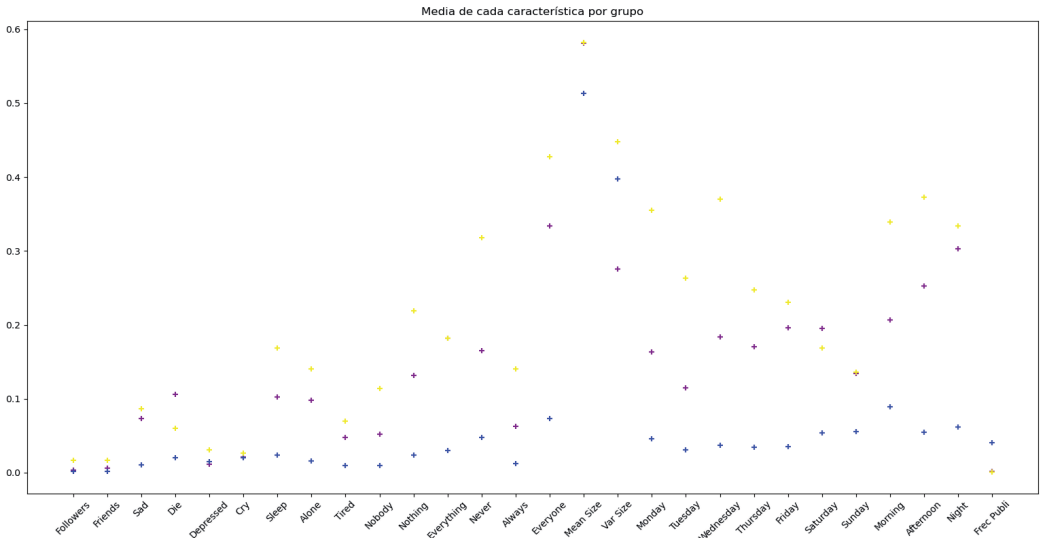
En los apartados posteriores se verán los resultados obtenidos tras el empleo de estos algoritmos.

### 3. Resultados

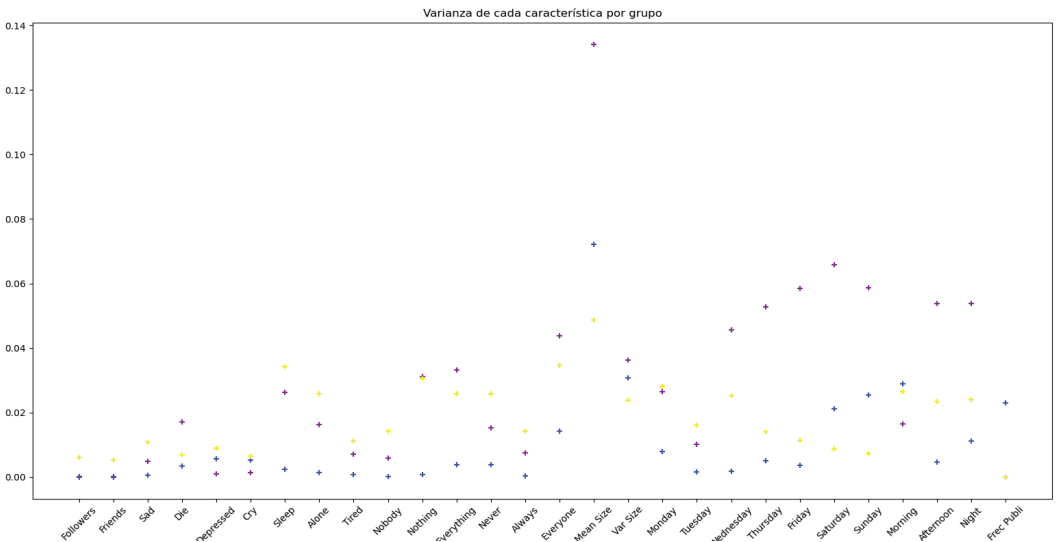
**FIGURA 2. Agrupación de los datos en 2D usando UMAP. Cada punto corresponde a una observación y cada color a un grupo de usuarios**



**FIGURA 3. Gráfico para interpretar las características de cada grupo. Cada punto representa el valor medio de las características en los usuarios de cada grupo. Cada color corresponde a un grupo**



**FIGURA 4. Gráfico para interpretar las características de cada grupo. Cada punto representa la varianza de las características en los usuarios de cada grupo. Cada color corresponde a un grupo**





En la figura 1 se ven representados 3 grupos obtenidos con el método UMAP. Cada conjunto ha sido creado agrupando usuarios que contienen características similares de la base de datos.

Para etiquetar a los usuarios, se ha calculado la media y la varianza de las características de cada grupo, graficadas en la Figura 2 y 3 respectivamente. Con ello, se podría interpretar la cantidad de usuarios que tienen tendencia al suicidio analizando la variabilidad de cada característica.

#### **4. Discusión**

Los grupos obtenidos son muy diferentes en cuanto a número de usuarios. Hay un grupo mayoritario, correspondiente al color amarillo, uno intermedio, de color azul y otro minoritario, de color morado.

Las características más representativas serán aquellas que muestren más variabilidad entre los grupos. Para ello, hay que fijarse en la varianza. Tanto las palabras descriptivas, la longitud de los tweets y los días de publicación, son las que más varían entre los distintos grupos.

El grupo azul es el menos activo en la aplicación de Twitter ya que, según las medias de las características, es el de menor valor en todas. Entre los otros dos grupos, la frecuencia de publicación y la longitud del tweet no varía. El amarillo, publica más los lunes y los miércoles tanto mañana, tarde y noche sin mucha diferencia significativa. Sin embargo, el morado hace más publicaciones los viernes y los sábados, aumentando progresivamente de mañana a noche. Esta diferencia entre ambas es importante, ya que las personas con el trastorno psicológico de depresión suelen publicar más los fines de semana por la noche, ya que tienen poca vida social o sufren insomnio debido a ésta [3]. Por otro lado, el grupo amarillo utiliza todas las palabras seleccionadas, siendo las mayoritarias 'everyone', 'sleep', 'nothing' y 'never'. El grupo morado utiliza más las palabras 'everyone', 'die', 'nothing', 'never', no habiendo mucha diferencia entre ambos.

Por todo lo analizado anteriormente, se puede interpretar que el grupo de usuarios con conducta suicida es el grupo morado, que corresponde al conjunto minoritario, seguido del amarillo, el mayoritario, y del azul, el intermedio.

#### **5. Conclusión**

Con este proyecto se pretende detectar casos de posibles suicidios de forma prematura analizando tuits de usuarios. Tras el análisis de los mismos vemos como sí que hay distinción entre varios grupos y principalmente en uno de ellos se cumplen

la mayoría de las características que hemos expuesto anteriormente. No se puede asegurar que los usuarios pertenecientes a ese grupo sufran de afecciones psicológicas ya que no se dispone de las etiquetas y la generación de los grupos se basa en similitud de características, pero sí que se puede realizar una vigilancia en aquellos usuarios que se introduzcan en el grupo considerado como más vulnerable.

## Referencias

[1] "Getting started — tweepy 4.5.0 documentation". Tweepy Documentation tweepy 4.10.0 documentation. Consultado el 4 de mayo de 2022. ([https://docs.tweepy.org/en/v4.5.0/getting\\_started.html](https://docs.tweepy.org/en/v4.5.0/getting_started.html))

[2] Portal Developer de Twitter. Consultado el 25 de abril de 2022. (<https://developer.twitter.com/en/portal/products>)

[3] "Así se podrá detectar la depresión por Twitter". Saber Vivir. Consultado el 4 de mayo de 2022. ([https://www.sabervivirtv.com/actualidad/detectar-depresion-por-twitter-tuits-depresivos\\_3659](https://www.sabervivirtv.com/actualidad/detectar-depresion-por-twitter-tuits-depresivos_3659))

[4] FitaBarcelona, J. (2019) La depresión se puede detectar por Twitter, La Vanguardia. Consultado el 1 de Mayo de 2022. (<https://www.lavanguardia.com/vivo/psicologia/20190710/463353450225/detectar-depresion-twitter.html>)

[5] Unidad de Coordinación Académica de Ciencias de la Salud y de la Vida (UPF) Upf. edu. Consultado el 4 de mayo de 2022. ([https://www.upf.edu/web/biomed/inici/-/asset\\_publisher/Us0jfwFAevmx/content/id/243188507/maximized](https://www.upf.edu/web/biomed/inici/-/asset_publisher/Us0jfwFAevmx/content/id/243188507/maximized))

## ARTICULOS/ARTICLES

La familia: desde el inicio hasta los últimos cambios en España / The family: from the beginning to the latest changes in Spain Sena Pascual-Lavilla .....	Págs 9-24
Análisis de los datos obtenidos de la red social Twitter para la identificación precoz de la tendencia al suicidio de los usuarios / Analysis of data obtained from the social network Twitter for the early identification of users' suicidal tendencies P.J. Mulas Cámara, R. Fernández-Calvillo Cáceres, C. Martínez Cabezali y ME. Molina Cañizares .....	Págs 25-33
Transformaciones de la familia mexicana y su incidencia en la convivencia y la gestión de los conflictos / Transformations at mexican family and its impact in coexistence and conflict management Emilia de los Ángeles Iglesias Ortuño .....	Págs 35-57
Trabajo social y cuidados en la vejez: un tema emergente para la intervención profesional / Social work and care in the elderly: an emerging topic for professional intervention María Concepción Arroyo y Perla Vanessa de los Santos .....	Págs 59-73
Del desvanecimiento del sujeto moderno al in-surgir. Aportes desde el Trabajo Social Decolonial / From the vanishing of the modern subject to the in-emergence. Contributions from Decolonial Social Work Gerardo Vélez Villafañe .....	Págs 75-92

## RESEÑAS/REVIEWS

Dubet, F. (2022): Tous inégaux, tous singuliers. Paris: Seuil / Dubet, F. (2022): Todo desigual, todo singular. París: Umbral (por Eguzki Urteaga) .....	Págs 93-97
Lewin, K. (1951). La teoría de campo en la ciencia social / Lewin, K. (1951). Field theory in social science (por José Javier Miranda Mayo) .....	Págs 99-102