

SENTIBOX: Marco Híbrido para Detección de Bots en X

Dr. Axel Rodríguez

Facultad de Ingeniería en Electricidad y Computación (FIEC)
Universidad de Panamá: Panamá, Panamá
Universidad Superior de Guadalajara
Guadalajara, Jalisco, México
<https://orcid.org/0000-0002-2485-1382>

Dr. Janzel Villalaz Guerra

Director de Investigación y Postgrado de la Vicerrectoría de Investigación y Postgrado
Universidad de Panamá: Panamá, Panamá
<https://orcid.org/0000-0001-8914-3216>

Mgter. Gabriel Vergara

Facultad de Ingeniería Industrial
Universidad Tecnológica de Panamá: Panamá, Panamá
gabriel.vergara@utp.ac.pa
<https://orcid.org/0000-0002-2247-7876>

Abstract

We propose SENTIBOX, a new approach for detecting bots on Twitter was officially rebranded as X, leveraging various data Types including profile, tweet, and neighbor-based features. By incorporating sentiment and emotion analysis and using an adversarial representation learner, SENTIBOX uniquely improves bot detection. The framework includes a multi-type feature extractor and an adaptive feature selection mechanism. To integrate different types of data, SENTIBOX further designs a feature extractor that computes statistical values to enrich tweet-based features, producing deep profile-based features by leveraging relationships among user properties. Experiments show that SENTIBOX improves accuracy by 1.9% and MCC (Matthews Correlation Coefficient) by 3.5% compared to state-of-the-art baselines in three real-world datasets. Moreover, under different scenarios with various combinations of features, SENTIBOX consistently achieves the best performance.

Keywords: Bot detection, sentiment analysis, emotion analysis, hybrid model, feature selection.

Resumen

Proponemos SENTIBOX, un nuevo marco para la detección de bots en Twitter reubautizada recientemente como la red social X, aprovechando varios tipos de datos, incluyendo características basadas en el perfil, tweets y vecinos. Al incorporar el análisis de sentimientos y emociones utilizando un aprendizaje de representación adversarial, SENTIBOX mejora notablemente la detección de bots. El nuevo enfoque incluye un extractor de características multi-tipo y un mecanismo adaptativo de selección de características. Para integrar diferentes tipos de datos, SENTIBOX además diseña un extractor de características que calcula valores estadísticos para enriquecer las características basadas en tweets, produciendo características profundas basadas en perfiles al aprovechar las relaciones entre las propiedades del usuario. Los experimentos muestran que SENTIBOX mejora la precisión en un 1.9% y MCC (Coeficiente de Correlación de Matthew) en un 3.5% en comparación con los estados del arte más avanzados en tres conjuntos de datos del mundo real. Además, en diferentes escenarios con varias combinaciones de características, SENTIBOX logra consistentemente el mejor rendimiento.

Palabras Clave: Detección de bot, Análisis de sentimientos, Análisis de emociones, Modelo híbrido, Selección de características.

1. INTRODUCCIÓN

Las plataformas de redes sociales se han convertido en fuentes principales de información, pero facilitan la propagación de desinformación mediante bots sociales (Gorodnichenko et al., 2021). Se estima que los bots constituyen el 15% de las cuentas activas de Twitter en 2022, con un tercio destinado a propósitos maliciosos, manipulando discusiones públicas e interfiriendo en procesos electorales (Graham et al., 2020).

La detección de bots se aborda mediante dos enfoques: ingeniería de características, basada en atributos estadísticos del usuario, y aprendizaje profundo, que utiliza redes neuronales para extraer características complejas (Feng et al., 2021). Métodos avanzados como LSTM, CNN y GCN han demostrado efectividad, mientras que el análisis de sentimientos introduce variables emocionales para mejorar la precisión. Los enfoques híbridos que combinan múltiples técnicas están ganando relevancia para una detección más robusta.

2. METODOLOGÍA

En este estudio, abordamos la detección de bots en X utilizando el marco SENTIBOX (un nuevo enfoque para la detección de bots en Twitter basado en el análisis de emociones y sentimientos). El objetivo es clasificar a los usuarios de Twitter como cuentas bot (y =

bot) o humanos ($y = \text{human}$) basándose en características del perfil, tweets recientes y la información de los vecinos.

A. Definición del Problema

La API de Twitter proporciona un objeto de usuario que contiene todo tipo de metadatos que describen al usuario de Twitter referenciado. En este estudio, representamos a un usuario como $U = \{P, T, N\}$, donde P es el conjunto de propiedades del perfil, T es el conjunto de tweets recientes y N es el conjunto de información de los vecinos. La tarea de detección de bots en Twitter es clasificar a un usuario dado U como una cuenta bot ($Y = \text{bot}$) o como humano ($Y = \text{human}$). Logramos esto aprendiendo una función $f(U) \rightarrow \hat{y}$ de modo que \hat{y} se aproxime a la verdad absoluta “ y ” para maximizar la precisión de la predicción.

B. Definición del Problema

La Figura 1 ilustra nuestro marco híbrido propuesto, SENTIBOX, que maneja simultáneamente los tweets de un usuario, su perfil y la información de sus vecinos a través de tres canales. Los tweets se procesan con un Calculador de Estadísticas para generar estadísticas del tweet, un Analizador de Sentimiento y Emoción para obtener puntuaciones de sentimiento y emoción, y un Mezclador de Grupos Diversos que, junto con un Aprendiz de Representación Adversarial (ARL), predice si los tweets agrupados son generados por bots. Para los datos de perfil y vecinos, se analizan todas las propiedades disponibles para generar características basadas en el perfil y los vecinos. El Selector Adaptativo de Características determina el conjunto de características más adecuado, el cual se entrena y evalúa con un Clasificador de Potenciación de Gradiente para la clasificación final.

C. Análisis de Sentimientos y Emociones

Para analizar el sentimiento y las emociones en las publicaciones de usuarios, utilizamos un analizador de sentimientos y emociones. Este enfoque se fundamenta en la premisa de que los bots de Twitter tienen una capacidad limitada para imitar cómo las personas reales expresan sentimientos internos. Nuestra metodología incorpora un analizador adicional de emociones para capturar características emocionales refinadas, además de los análisis convencionales de sentimiento.

El proceso comienza alimentando cada tweet a nuestro analizador de sentimientos y emociones, donde para el análisis de sentimiento extraemos los grados de positividad, negatividad y neutralidad de cada palabra utilizada en la oración. Cada clase de sentimiento recibe una puntuación que representa la proporción de ese sentimiento en toda la oración. Además, calculamos un puntaje compuesto normalizado entre -1 (extremadamente negativo) y +1 (extremadamente positivo) sumando los puntajes de valencia de todas

las palabras según el léxico. Para el análisis de emociones, identificamos cinco categorías emocionales diferentes: feliz, enojado, triste, sorprendido y temeroso, asignando a cada una un puntaje proporcional en cada oración. Este enfoque utiliza modelos preentrenados como VADER (Hutto & Gilbert, 2014) y Text2Emotion (Pennebaker et al., 2001), cada uno con sus ventajas en la detección de sentimientos y emociones en contextos de redes sociales. Una vez realizados los análisis de sentimiento y emociones, obtenemos nueve puntajes para cada tweet, que luego se agregan en el generador de características del tweet y se utilizan para formar grupos diversos.

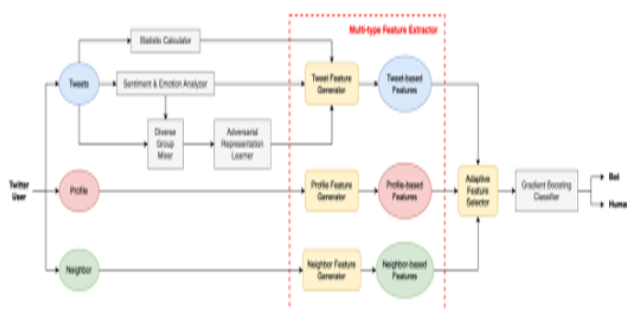


Fig. 1. Visión general de nuestra propuesta de marco híbrido SENTIBOX.

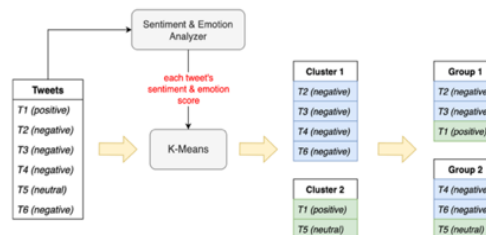


Fig. 1. Visión general de nuestra propuesta de marco híbrido SENTIBOX.

D. Aprendizaje de Representación Adversarial (ARL)

En lugar de entrenar el modelo ARL con cada tweet individualmente, se propone el Diverse Group Mixer (Mezclador de Grupo Diverso), que combina sistemáticamente los tweets para agregar un número adecuado de tweets diversificados, logrando resultados de entrenamiento más estables. Específicamente, se obtiene el puntaje de sentimiento y emoción de cada tweet mediante un analizador de sentimientos y emociones, para posteriormente aplicar el agrupamiento mediante el algoritmo de aprendizaje automático no supervisado K-Means. Para determinar el número óptimo de clusters, se calculan la distancia intra-cluster promedio y la distancia al cluster más cercano promedio para encontrar el puntaje de silueta (Rousseeuw, 1987), para cada valor de k. Se prefiere el puntaje de silueta sobre el Método del Codo (Bholowalia & Kumar, 2014), ya que el punto de “codo” no es muy claro en este conjunto de datos.

Posteriormente, los tweets dentro de cada cluster se distribuyen equitativamente en un cierto número de grupos. Se determina el mejor número de grupos a través de experimentos, cuyos resultados se presentan en la Sección IV-C. De esta manera, cada grupo contiene tweets con un grado similar de diversidad. La Figura 2 muestra un ejemplo de cómo funciona el Mezclador de Grupo Diverso, donde se puede observar que cada grupo se forma siguiendo la distribución de sentimientos y emociones del usuario. Esta técnica de agrupamiento cubre uniformemente los comportamientos generales de un usuario, y se

demuestra que aumenta efectivamente el rendimiento, como se explica en la Sección IV-C.

E. Aprendizaje de Representación Adversarial

Inspirado en los modelos semi-supervisados GAN(Salimans et al., 2016), y GAN-BERT (Croce et al., 2020), extiende las Representaciones de Codificador Bidireccional de Transformers (BERT)(Devlin et al., 2018), utilizando Redes Generativas Adversariales (GAN) (Goodfellow et al., 2014), para la etapa de ajuste fino. Este modelo semi-supervisado mejora la capacidad de generalización cuando faltan datos anotados. La arquitectura de ARL (Figura 3) aprovecha la capacidad de BERT para producir representaciones de alta calidad de textos de entrada y adopta textos generados que imitan los datos reales para ayudar a la red a generalizar sus representaciones para la tarea final. En esta sección, se explica cómo funciona ARL y se describen las funciones de pérdida utilizadas en este modelo.

Primero, se genera una representación vectorial real de los tweets mediante BERT, mientras que el generador (G) crea una representación vectorial falsa que imita la real. El objetivo de G es generar ejemplos falsos similares a los datos reales, de modo que el discriminador (D) no los identifique como falsos. Hay dos tipos de pérdidas asociadas con G: LGfeat, que minimiza la disimilitud entre las representaciones reales y falsas, y LGunsup, que cuenta las veces que D identifica correctamente los ejemplos falsos D, por su parte, clasifica las representaciones generadas por G o BERT como bot o humano, con pérdidas LDsup y LDunsup asociadas a ejemplos reales y falsos, respectivamente. El proceso de entrenamiento minimiza tanto LG como LD, haciendo el modelo más estable y consistente.

F. Extractor de Características Multi-tipo

El Extractor de Características Multi-tipo integra información numérica, binaria y textual de conjuntos de datos de detección de bots, combinando características basadas en tweets, perfiles y vecinos para simplificar el diseño en comparación con otros métodos híbridos. Utiliza la metadata de los usuarios para expandir el conjunto de características y capturar información oculta, generando características como la longitud mínima de los tweets, ratios de menciones y hashtags, y promedios de sentimientos y emociones. También emplea predicciones del Adversarial Representation Learner (ARL) para determinar si un usuario es humano o bot. En cuanto a la información de perfiles, se recopilan y analizan más de 40 características, seleccionando las más relevantes y creando nuevas características como la tasa de crecimiento de seguidores y el tema de la imagen de fondo del perfil. Para la información basada en vecinos, se cuentan los bots que un usuario sigue y los que lo siguen, utilizando una campaña de crowd-sourcing para la anotación de cuentas de bot. Estas características, efectivas para capturar comportamientos adicionales de bots, se incluyen

en el conjunto final de características, proporcionando una cobertura más completa de los comportamientos de bots que las basadas solo en tweets o perfiles.

G. Selector de Características Adaptativo

Para optimizar el rendimiento del clasificador, se aplica el Método Wrapper, que identifica sistemáticamente el mejor conjunto de predictores para un modelo de ML dado. La selección de características busca disminuir el tiempo de computación, aumentar la interpretabilidad del modelo y mejorar la precisión predictiva al reducir predictores redundantes. Inspirado en la selección progresiva, se propone un selector de características adaptativo: comenzando con un subconjunto vacío, se añade en cada iteración la característica que mejora más el rendimiento del clasificador. El subconjunto con la mayor puntuación de evaluación final es seleccionado como el conjunto de características definitivo.

H. Clasificador de Gradient Boosting

Después de seleccionar el conjunto ideal de características, adoptamos Gradient Boosting como nuestro clasificador final debido a su efectividad demostrada en comparación con el bosque aleatorio en la gestión de diversas características. Gradient Boosting ofrece flexibilidad mediante ajuste de hiperparámetros y es robusto, requiriendo poco preprocesamiento de datos y entrenando rápidamente. Este método combina múltiples modelos de aprendizaje débil para construir un modelo predictivo sólido, utilizando árboles de decisión de manera estratégica. Durante el entrenamiento, estandarizamos las características y utilizamos Grid Search en el conjunto de validación para optimizar los hiperparámetros del modelo para demostrar la contribución de cada componente al rendimiento general.

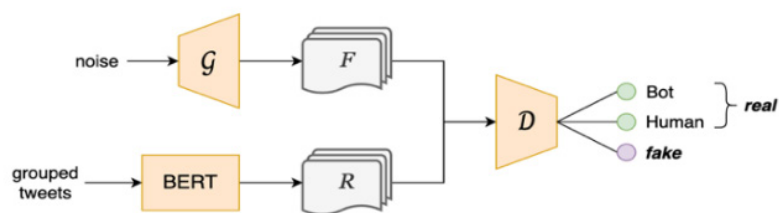


Fig. 3. . La arquitectura de nuestro modelo de aprendizaje de representación adversarial propuesto: dada una distribución aleatoria, un generador G crea un conjunto de ejemplos falsos F. Estos ejemplos falsos luego sirven como entrada para el discriminador D, junto con representaciones vectoriales de ejemplos reales R que son generados por BERT.

TABLA I LISTA DE CARACTERÍSTICAS USADAS EN SENTIBOX

Tweet-based		Profile-based	
Feature Name	Explanation	Feature Name	Explanation
len_min	the shortest tweet's length	joined_days	collect_time - create_time
hashtag	number of hashtags in all tweets	followers_growth_rate	followers_count / joined_days
url_ratio	number of urls in all tweets / tweet_count	username_length	length of username
neu	proportion of words associated with neutral sentiment	digits_in_username	number of digits in username
angry	proportion of words associated with angry emotion	common_color	number of default colors used
cmp_pred	compare the number of labels from ARL	image_theme	profile page's theme setting
cmp_prob	compare the average probability from ARL	domain_count	interest domains within 4 categories

3. EXPERIMENTOS Y RESULTADOS

Esta sección describe en detalle la configuración experimental y demuestra que nuestro modelo es superior a otros modelos de vanguardia utilizando datos recopilados entre 2017 y 2020. Además, realizamos experimentos adicionales para mostrar la efectividad de nuestro mecanismo propuesto con otros modelos.

TABLA II LISTA DE CARACTERÍSTICAS USADAS EN SENTIBO

Dataset	#Accounts		#Tweets	
	Human	Bot	Human	Bot
Twibot-20	5,237	6,589	878,183	1,020,876
Cresci-17	3,474	4,912	2,839,361	3,457,133
PAN-19	3,380	3,380	338,000	338,000

A. Configuración Experimental

1. Conjuntos de Datos: Se utilizan tres conjuntos de datos: TwiBot-20 (Feng, Wan, Wang, Li, et al., 2021b), Cresci-17 (Cresci et al., 2017), y PAN-19 (Pizarro, 2020). Los conjuntos varían en tamaño y período de recolección, haciendo más desafiante el conjunto más reciente.

- TwiBot-20 (Feng, Wan, Wang, Li, et al., 2021b): Representativo de la generación actual de bots y usuarios genuinos de Twitter, incluye información de tweets, perfiles y vecinos.
- Cresci-17 (Pizarro, 2020) : Cubre la mayoría de los bots tradicionales de la era temprana de Twitter, incluye datos de tweets y perfiles.
- PAN-19 (Pizarro, 2020): Conjunto de datos relativamente pequeño, contiene solo información de tweets.

2. Métodos de Referencia: Se comparan varios modelos basados en ingeniería de características y aprendizaje profundo, como SATAR (Feng, Wan, Wang, Li, et al., 2021a) y BotRGCN (Feng, Wan, Wang, & Luo, 2021). Entre otros, se incluyen: Random-forest (Lee et al., 2011), DenStream (Miller et al., 2014), DNA-fingerprinting (Cresci et al., 2016), Botometer (Davis et al., 2016), Contextual-LSTM (Kudugunta & Ferrara, 2018), GCNN (Ali

Alhosseini et al., 2019), BiLSTM (Wei & Nguyen, 2019), Data-selection (Yang et al., 2020), SATAR (Feng, Wan, Wang, Li, et al., 2021a), y BotRGCN (Feng, Wan, Wang, & Luo, 2021).

Hiperparámetros: Se utiliza Grid Search para definir la mejor configuración inicial para los hiperparámetros del clasificador Gradient Boosting, utilizando validación cruzada de cinco pliegues para ajustar el modelo.

3. Métricas de Evaluación: Se utilizan principalmente Accuracy, F1-score y MCC para evaluar la efectividad del modelo, basándose en el usuario como unidad de medida.

B. Bot Detection Performance

La Tabla III resume el desempeño de detección de bots de diferentes métodos en tres conjuntos de datos. Mostramos que el rendimiento de nuestro método propuesto, SENTIBOX, es competitivo en comparación con otros métodos de referencia, demostrando que SENTIBOX es generalmente exitoso en la detección de bots de Twitter. Además, debido a que SENTIBOX supera a los métodos de última generación en el conjunto de datos representativo y comprensivo TwiBot-20, se generaliza a escenarios del mundo real. SENTIBOX, por otro lado, es capaz de ajustarse a las generaciones cambiantes de bots, ya que muestra el mejor desempeño en los tres conjuntos de datos con tiempos de recolección de datos que varían desde 2017 hasta 2020.

Para una comparación más justa, limitamos nuestro método a usar solo ciertos aspectos de la información. Por ejemplo, Data-selection (D'Andrea et al., 2015) maneja únicamente información de perfil de usuario; por lo tanto, solo recuperamos propiedades de usuario del conjunto de datos para comparar nuestro método con los otros métodos en las mismas condiciones. Como se muestra en la Figura 4, SENTIBOX sigue superando a todos los métodos de referencia con tipos limitados de características. Resumimos las razones de su éxito :

1. Tanto Data-selection (Yang et al., 2020), como SENTIBOX usan información de perfil. SENTIBOX supera a Data-selection (Yang et al., 2020), lo que muestra que SENTIBOX utiliza mejor las propiedades del perfil del usuario, sin necesidad de una ingeniería de características tan pesada.
2. Para los métodos basados principalmente en la semántica de tweets, SENTIBOX supera significativamente a DNA-fingerprinting (Cresci et al., 2016), y BiLSTM (Wei & Nguyen, 2019), esto indica que nuestro enfoque puede capturar mejor los patrones de tweets típicos de bots al aprovechar más características basadas en tweets, sugiriendo que los métodos de detección de bots deberían incorporar más aspectos de la información semántica.
3. Métodos como Random-forest (Lee et al., 2011), DenStream (Miller et al., 2014), y Contextual-LSTM (Kudugunta & Ferrara, 2018), consideran conjuntamente la información

de perfil y tweets. Sin embargo, no pueden competir con SENTIBOX debido a la falta de características comprensivas, mientras que SENTIBOX utiliza los metadatos de tantas formas como sea posible.

- Por último, los vecinos de usuario son materiales relativamente nuevos que han aparecido recientemente en los conjuntos de datos de detección de bots, por lo que no hay muchos estudios que cubran esta información. SENTIBOX adopta la misma combinación de información del usuario (perfil y vecino) que GCNN (Ali Alhosseini et al., 2019), pero logra un rendimiento mucho mayor gracias a un mejor diseño.

TABLA III COMPARACIÓN DE RENDIMIENTO PARA MÉTODOS DE DETECCIÓN DE BOTS

		Random-forest	DenStream	DNA-fingerprinting	Botometer	Contextual-LSTM	GCNN	BiLSTM	Data-selection	SATAR	BotRGCN	Ours
TwiBot-20	Acc	0.7456	0.4801	0.4793	0.5584	0.8174	0.6813	0.7126	0.8191	0.8412	0.8462	0.8832
	F1	0.7823	0.6266	0.1072	0.4892	0.7517	0.7318	0.7533	0.8546	0.8642	0.8707	0.8963
	MCC	0.4879	-0.1372	0.0839	0.1558	0.6710	0.3543	0.4193	0.6643	0.6863	0.7021	0.7663
Cresci-17	Acc	0.9750	0.5204	0.4029	0.9597	0.9799	/	0.9670	0.9847	0.9871	/	0.9952
	F1	0.9826	0.4737	0.2923	0.9731	0.9641	/	0.9768	0.9893	0.9910	/	0.9966
	MCC	0.9387	0.1573	0.2255	0.8926	0.9501	/	0.9200	0.9625	0.9685	/	0.9880
PAN-19	Acc	/	/	0.8797	/	/	/	0.9464	/	0.9509	/	0.9621
	F1	/	/	0.8701	/	/	/	0.9448	/	0.9510	/	0.9619
	MCC	/	/	0.7685	/	/	/	0.8948	/	0.9018	/	0.9243

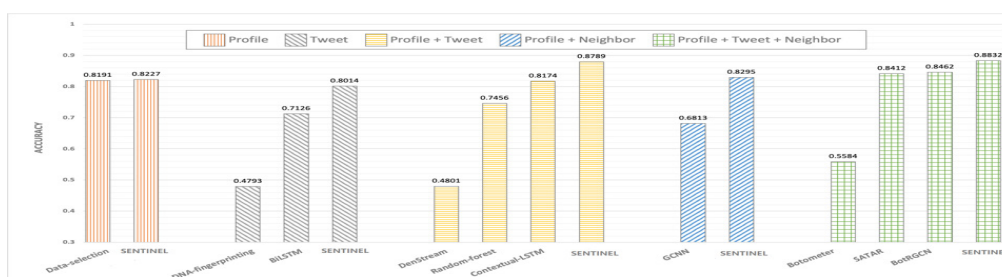


Fig. 4. Comparación de rendimiento para diferentes combinaciones de información de usuario en TwiBot-20. Tenga en cuenta que no comparamos “Vecino” o “Tweet+Vecino” dado que no hay ningún método que utilice tal información.

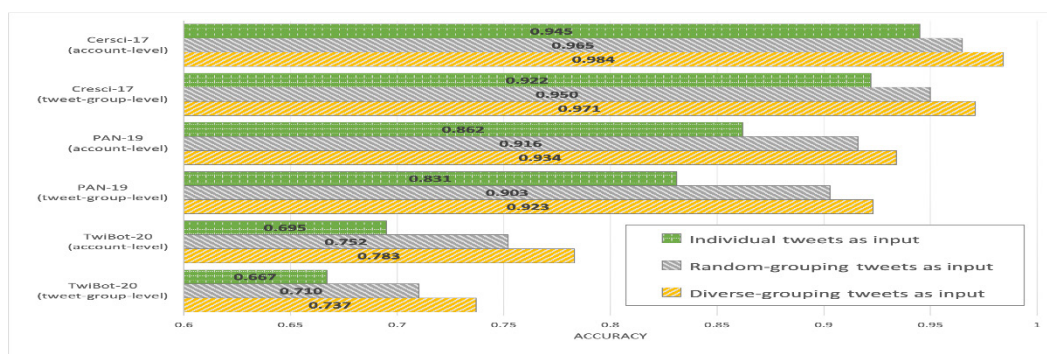


Fig. 5. Comparación de rendimiento para diferentes combinaciones de información de usuario en TwiBot-20. Tenga en cuenta que no comparamos “Vecino” o “Tweet+Vecino” dado que no hay ningún

método que utilice tal información.

4. CONCLLUSIÓN

La detección de bots en redes sociales es un tema cada vez más relevante, y en este artículo presentamos SENTIBOX, un nuevo método híbrido que combina aprendizaje de representación adversarial con análisis de sentimientos y emociones. Nuestra metodología introduce una técnica de agrupación diversa que mejora significativamente la precisión tanto a nivel de grupos de tweets como a nivel de cuentas individuales. Además, hemos refinado el diseño de la pérdida de coincidencia de características para reducir el tiempo de entrenamiento en un 50% y estabilizar el proceso de aprendizaje. SENTIBOX utiliza un extractor de características versátil que integra información de tweets, perfiles y vecinos del usuario, formando así un perfil completo del usuario en cuestión. Nuestro método selecciona automáticamente el conjunto de características más apropiado según las características del conjunto de datos específico, lo cual se refleja en nuestros resultados experimentales. Hemos demostrado que SENTIBOX supera a los métodos actuales en tres conjuntos de datos del mundo real y en diferentes configuraciones de características. Como trabajo futuro planeamos incorporar más información de vecinos mediante enfoques basados en grafos (Zhao et al., 2020) y aplicaciones en sensores inalámbricos (Rodríguez García & Jipsion, 2019) para enriquecer aún más nuestro conjunto de características y mejorar la precisión de detección de bots.

REFERENCIAS

- [1] Gorodnichenko, Y., Pham, T., & Talavera, O. (2021). Social media, sentiment and public opinions: Evidence from #Brexit and #USElection. *European Economic Review*, 136, 103772.
- [2] Graham, M., Avery, E., & Park, S. (2020). The role of social media in local government crisis communications. *Public Relations Review*, 41(3), 386-394.
- [3] Feng, S., Wan, H., Wang, N., & Luo, M. (2021). BotRGCN: Twitter bot detection with relational graph convolutional networks. *Proceedings of the 2021 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*, 236-239.
- [4] Hutto, C. J., & Gilbert, E. (2014). VADER: A parsimonious rule-based model for sentiment analysis of social media text. *Proceedings of the 8th International AAAI Conference on Weblogs and Social Media*, 216-225.
- [5] Pennebaker, J. W., Francis, M. E., & Booth, R. J. (2001). *Linguistic inquiry and word count: LIWC 2001*. Mahwah, NJ: Lawrence Erlbaum Associates.
- [6] Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20, 53-65.
- [7] Bholowalia, P., & Kumar, A. (2014). EBK-means: A clustering technique based on elbow method and K-means in WSN. *International Journal of Computer Applications*, 105(9), 17-24. Claude es IA y puede cometer errores. Por favor, verifica nuevamente las respuestas. Sonnet 4.5 Claude es IA y puede cometer errores. Por favor, verifica nuevamente las respuestas.
- [8] Zhao, K., Kang, J., Jung, J., & Sohn, G. (2020). BotRGCN: Twitter bot detection with relational graph convolutional networks. *Proceedings of the 2020 IEEE/ACM International Conference on Advances*