

PROCESAMIENTO DE LA NEGACIÓN EN TWITTER

GEMMA BEL ENGUIX

Y SERGIO LUIS OJEDA TRUEBA

Todos los días se escriben millones de *tuits* en México y en todo el mundo, se dan RT, muchos de ellos se citan u otros usuarios indican que les gustan. Twitter es una red social que prioriza la escritura, con su conocida invitación: “¿Qué está pasando?”. Eso sí, solamente permite mensajes breves, con un máximo de 280 caracteres.

Con este mecanismo se logra una gran red de comunicación instantánea, que permite hacer un seguimiento en tiempo real de la estructura y creatividad del lenguaje, así como sentimientos, opiniones, tendencias sociales o cambios en el comportamiento. Por ello, muchas líneas de investigación en lingüística computacional se centran actualmente en el análisis de las plataformas de redes sociales.

Para estudiar estos mensajes cortos, que no se ajustan a las convenciones retóricas de la lengua escrita, antes hay que comprenderlos bien, ser capaces de identificar automáticamente todos sus componentes. Ello lleva directamente a un tema clave y a la vez básico: la detección de la negación y su alcance.

A continuación, describiremos de forma breve una de las propuestas de investigación elaborada por el Grupo de Ingeniería Lingüística (GIL) del Instituto de Ingeniería para estudiar la negación, en específico la sintáctica, en Twitter, con la colaboración del Instituto de Investigaciones en Matemáticas Aplicadas y en Sistemas.

La negación lingüística

A simple vista, la negación parece un tema trivial o irrelevante para estudiar desde el punto de vista computacional. En cambio, la detección automática de las estructuras negativas es de una importancia capital para la correcta interpretación de los textos por parte de los programas computarizados.

La negación por sí sola tiene la capacidad de cambiar el valor de toda una oración o parte de esta. Por ejemplo, en el enunciado *Quiero leer*, solamente con el adverbio *no* podemos indicar exactamente lo contrario: *No quiero leer*. En cambio, en otras ocasiones, la negación no cambia el sentido completamente, sino de forma parcial. Por ejemplo: *Me ha dicho que le gusta el vino y que no toma licor*.

Si bien es cierto que el uso de la negación lingüística se ha estudiado en campos como los reportes médicos, las reseñas y

en las noticias; en el ámbito de las redes sociales aún no se registran muchos trabajos con dicho enfoque. Investigar este aspecto en una plataforma como Twitter es de sumo interés ya que podría ayudarnos a detectar la postura de los usuarios frente a un tema o tendencia.

Por otra parte, el lenguaje de Twitter tiene ciertas particularidades. Los usuarios continuamente hacen uso de diversos recursos para marcar expresividad en el mensaje emitido. En general, notamos tres usos lingüísticos característicos de Twitter: léxico específico de redes sociales, existencia de errores ortotipográficos y efectos enfático-pragmáticos. Asimismo, el dialecto mexicano fue un reto añadido a la investigación. Por ejemplo, en la variante mexicana, el término negativo *no*, puede encontrarse como *nel*, *Nelson*, *nah*, *ni maiz*, etc. Es importante tener en cuenta estas características de la lengua tecleada al construir los sistemas de detección automática.

Dada la importancia de la identificación y procesamiento de la negación en Twitter, nos fijamos dos objetivos principales:

1. Crear un corpus etiquetado de negación sintáctica en *tuits* escritos en español mexicano, el T-MexNeg Corpus.
2. Diseñar un método capaz de detectar de forma automática la negación sintáctica.

1. El Corpus T-MexNeg

El primer paso para un estudio de lingüística computacional es disponer de un corpus textual. En este caso, primeramente, recopilamos *tuits* mediante la API de Twitter. Después, se extrajo una muestra aleatoria para confeccionar el dataset de estudio. Las fechas de emisión son de septiembre de 2017 a abril de 2019; el origen geográfico, el territorio nacional mexicano.

En total recolectamos 13704 *tuits*, que fueron anotados de forma binaria para separar aquellos que sí contenían la negación sintáctica en su contenido. La clasificación fue realizada en tres ocasiones por un equipo de estudiantes especialmente entrenados. Esto es así para medir el grado de acuerdo entre los etiquetadores y saber si el trabajo se llevó a cabo de forma adecuada.

Una vez completado este proceso, obtuvimos un total de 4892 *tuits* con negación sintáctica, los cuales volvimos a anotar; pero ahora buscando marcar sus principales componentes, siguiendo a Martí *et al.* (2016): *término negativo* (la negación), *event* (lo que se niega a nivel sintáctico) y *scope* (el alcance de la negación dentro de la oración). Además, utilizamos dos etiquetas complementarias: *related negation* (para negaciones dependientes de una negación previa) y *false negation* (términos que propiamente no niegan nada). La Figura 1 ilustra cada uno de estos componentes en tres *tuits* distintos.

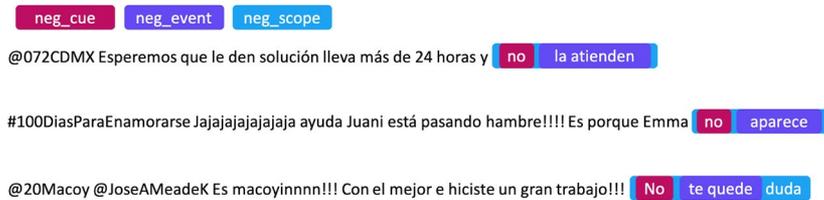


Figura 1. El *scope* (resaltado de azul), incluye el término negativo (rojo) y el *event* (morado). Nótese que el *scope* no necesariamente consta del término negativo más el *event*, sino que puede incluir más palabras

Al igual que la etapa anterior, ésta también se efectuó tres veces, con las mediciones de acuerdo correspondientes. Previo a los procesos de anotado, creamos una guía que nos permitió orientar a los anotadores acerca de cómo y qué etiquetar:

2. Experimentos

Los experimentos se centraron en la detección automática de los componentes de la negación sintáctica: términos negativos, *event*, *scope*, así como en toda la estructura negativa de manera global: términos + *event* + *scope*.

Para ello, usamos dos métodos principales, uno basado en diccionario, otro en Conditional Random Fields (CRF) (Lafferty *et al.* 2001), un modelo estocástico para anotar secuencias de datos que ha dado muy buenos resultados en la identificación

de negación en otros ámbitos, como el biomédico (Agarwal y Yu, 2010). Para comprobar la adecuación y desempeño de nuestro recurso, T-MexNeg, comparamos nuestros experimentos con el otro corpus de negación en español, el SFU ReviewSP-NEG (Kolhakar *et al.* 2019), creado con 400 reseñas de productos extraídas de la web Ciao.es. El estudio de ambos corpus ayuda a ver los términos de negación más frecuentes en cada uno de ellos (Figura 2).

En primer lugar, comparamos los resultados con métodos basados en diccionario y los obtenidos con CRF en ambos corpus. La Figura 3 ilustra los resultados para la detección de los términos negativos.

En la Figura 3 observamos cómo los mejores resultados se obtienen con el método CRF. Así que usamos éste para la detección de todos los componentes de la negación en ambos corpus, con los resultados que se muestran en la Figura 4.

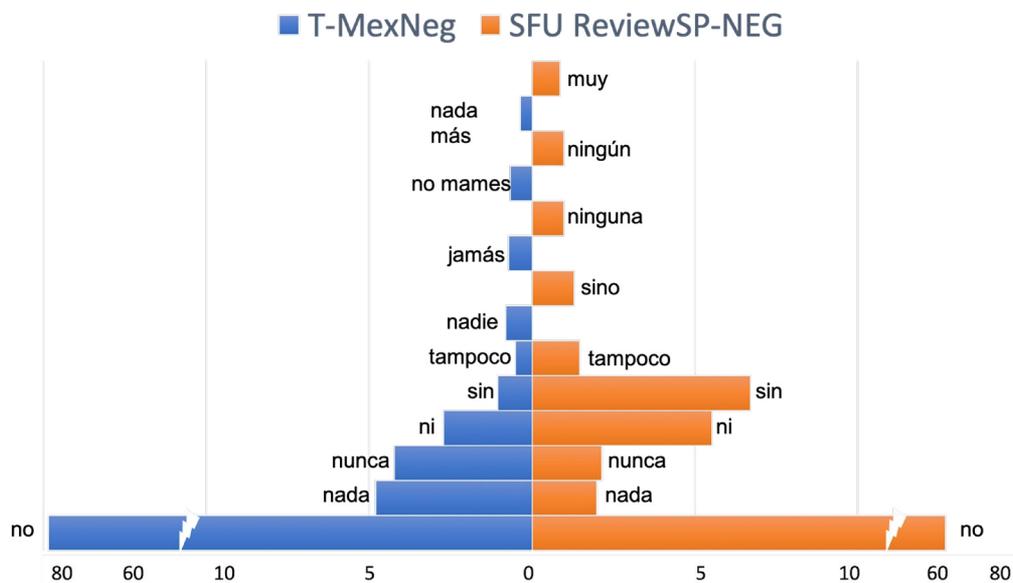


Figura 2. Los términos negativos más frecuentes en el T-MexNeg (azul) y SFU ReviewSP-NEG (naranja)

Detección de Términos Negativos				
	EXP 1: Dictionary		Exp 2: CRF	
	TMN	SFU	TMN	SFU
Precisión	0,89	0,89	0,93	0,79
Recall	0,92	0,83	0,98	0,96
F1	0,91	0,86	0,95	0,88

Figura 3. Comparativa de resultados entre los sistemas basados en diccionario y los basados en aprendizaje automático para la detección del término negativo

Tag	TMN			SFU		
	Precisión	Recall	F1	Precisión	Recall	F1
Neg. Cue	0,93	0,98	0,95	0,79	0,96	0,88
Event	0,86	0,94	0,90	0,51	0,71	0,59
Scope	0,63	0,92	0,75	0,54	0,93	0,68
Global	0,61	0,92	0,73	0,37	0,90	0,50

Figura 4. Aplicación del método CRF a la detección automática de los componentes de la negación en el corpus T-MexNeg y ReviewSP-NEG

Como se puede apreciar, el T-MexNeg parece un corpus muy bien elaborado, en comparación con otros de áreas similares y propósitos cercanos, para la detección de negación. Por otra parte, mientras la detección del término negativo propiamente dicho es muy precisa (F1 de 0.95), la identificación de los otros componentes es más complicada, ya que resulta muy difícil dilucidar los límites de la estructura negativa. En conclusión, este es un campo de investigación abierto que necesita un desarrollo más extenso.

Aplicaciones y uso de corpus

El corpus T-MexNeg es un recurso valioso para lingüistas e ingenieros en computación, con muy diversas aplicaciones. Por una parte, se puede usar desde un punto de vista estrictamente lingüístico para estudiar el lenguaje propio de las redes sociales. Un uso más especializado es la detección automática

de la negación, específicamente en *tuits*. Esto incide en todo tipo de tareas de procesamiento del lenguaje natural, como el resumen automático, el perfilado de autor en redes sociales o los estudios de polaridad y análisis de sentimientos.

Es esta última aplicación la que se está desarrollando actualmente en el GIL, ya que la negación puede cambiar el valor de un enunciado. Esto resulta muy útil para medir el grado de aprobación de un candidato político, de un producto de mercado o de alguna película nueva, por ejemplo. Las herramientas de este tipo tienen una amplia aceptación entre asesores políticos y mercadólogos.

Si se busca saber más del proyecto, los resultados de esta investigación han sido publicados en “Negation Detection on Mexican Spanish Tweets: The T-MexNeg Corpus” de la revista Applied Sciences (Bel-Enguix *et al.* 2021). Además, se puede consultar el corpus en la dirección de GitLab: https://gitlab.com/gil.iingen/negation_twitter_mexican_spanish.

Grupo de Trabajo

Coordinadores:

- Gemma Bel Enguix (II)
- Helena Gómez Adorno (IIMAS)
- Alejandro Pimentel (II, INEGI)

Desarrolladores y equipo de investigación:

- Brian Aguilar Vizuet (Facultad de Ciencias)
- Sergio Luis Ojeda Trueba (II)

Referencias

Agarwal, S. y Yu, H. Biomedical negation scope detection with conditional random fields. *J. Am. Med Inform. Assoc.* 2010, 17, 696-701.

Bel-Enguix, G.; Gómez-Adorno, H.; Pimentel, A.; Ojeda-Trueba, S. L. y Aguilar-Vizuet, B. "Negation Detection on Mexican Spanish Tweets: The T-MexNeg Corpus" In *Applied Sciences*. 2021, 11(9), 3880.

Kolhatkar, V.; Wu, H.; Cavasso, L.; Francis, E.; Shukla, K. y Taboada, M. The SFU opinion and comments corpus: A corpus for the analysis of online news comments. In *Corpus Pragmatics*; Springer: Berlin/Heidelberg, Germany, 2019; pp. 1-36.

Lafferty, J.; McCallum, A. y Pereira, F. C. *Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data*; University of Pennsylvania: Philadelphia, PA, USA, 2001; pp. 282-289.

Martí, M. A.; Taulé, M.; Nofre, M.; Marsó, L.; Martín-Valdivia, M. T. y Jiménez-Zafra, S.M. La negación en español: Análisis y tipología de patrones de negación; *Sociedad Española para el Procesamiento del Lenguaje Natural*: Jaén, Spain, 2016; pp. 41-48.