SPRINT 3 - MEETING 40

"We're competing with sleep, on the margin. And so, it's a very large pool of

Natural (Parte 2)

time." -Reed Hastings. CEO Netflix. En la bitácora anterior, conociste el Procesamiento de Lenguaje Natural o NLP,

El equipo docente subió contenido extra de la meeting 💠

Procesamiento Lenguaje

aprendiste a tokenizar el texto, a filtrar stop words, y a usar lexicons o corpus para pre-entrenar algoritmos. Estos pasos forman parte del preprocesamiento del texto, que es necesario para luego entrenar nuestro algoritmo de NLP.

Como podrás imaginar, NLP tiene sus modelos específicos, especializados en análisis de texto. En esta bitácora conocerás algunos de ellos.

TF-IDF significa "Term Frequency-Inverse Document Frequency", jeso ya nos da información sobre lo que hace!

TF-IDF

Term Frequency significa: cuántas veces un término aparece en un documento. La importancia de esto es que se supone que mientras más veces aparece, más

importante es ese término en ese documento. Por ejemplo, un texto sobre tiburones va a contener muchas más veces la palabra tiburón que un texto que no es sobre tiburones. Mientras mayor frecuencia tenga la palabra en el documento, mayor "puntaje" tendrá. Esto nos sirve, por ejemplo, para detectar sobre qué tópico trata un texto. Ahora bien, hay términos (como los artículos "la", "el", o las preposiciones "a", "por") que tienen alta frecuencia en todos los textos. Si les asignamos un

puntaje alto, no ganaremos información con respecto a cuál es el tópico del documento. Aquí es donde entra en juego la segunda parte del nombre TF-IDF: Inverse Document Frequency. Esto significa que se compara la frecuencia de la palabra en ese documento con la frecuencia en todo el corpus de documentos. A mayor frecuencia de la palabra en todos los documentos, menor puntaje tendrá. Así evitamos puntuar alto a stop words o palabras que no agregan información sobre ese documento específico. Word2Vec Word2Vec es una técnica que, como su nombre lo indica, transforma las

ellas. La idea que se encuentra detrás de este método es que las palabras

próximas tienen mayor relación en su significado que las palabras lejanas. ¿Cómo lo hace? Consideremos la siguiente frase extraída de El Jardín de los Senderos que se Bifurcan, por Jorge Luis Borges: "Así combatieron los héroes, tranquilo el admirable corazón, violenta la

espada, resignados a matar y a morir" Para cada palabra vamos a considerar las palabras cercanas. Por ejemplo, si

tomamos una distancia igual o menor a 2, las palabras cercanas a "admirable" son "tranquilo", "el", "corazón" y "violenta". Podemos tomar un texto entero, medir las distancias que todas las palabras tienen entre sí, y luego entrenar una pequeña red neuronal que prediga la probabilidad de que una palabra se

palabra, que nos indicarán su relación con otras palabras.

encuentre cercana a otra. A partir de esto se obtendrán los vectores para cada

Cada vector estará compuesto por n cantidad de números entre -2 y 2 (cada uno

números, y colorear cada uno (rojo si está más cerca del 2, blanco si está cerca del o, azul si está más cerca del -2).

entradas en inglés de Wikipedia:

"Man"

1.5

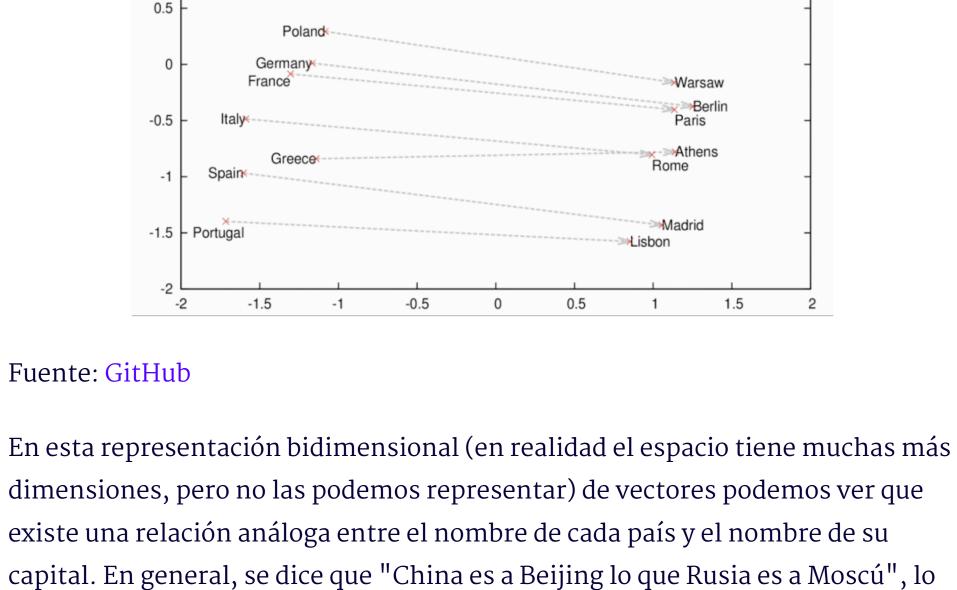
Fuente: GitHub

"Woman" Fuente: GitHub Podemos ver a simple vista ciertas similitudes entre estos vectores,

Country and Capital Vectors Projected by PCA

convertirán en vectores cercanos con direcciones análogas.

Turkey



restamos el vector de la palabra "man" y le sumamos el vector de la palabra "woman", estas serán las predicciones resultantes: model.most_similar(positive=["king", "woman"], negative=["man"]) [('queen', 0.8523603677749634), ('throne', 0.7664333581924438), ('prince', 0.7592144012451172), ('daughter', 0.7473883032798767),

('monarch', 0.721449077129364), ('eldest', 0.7184862494468689), ('widow', 0.7099430561065674)] Fuente: GitHub Esto sucede siempre y cuando le "prohibamos" al algoritmo que la predicción

sea una de las palabras usadas en el cálculo, en cuyo caso será otra vez "king".

('elizabeth', 0.7460219860076904), ('princess', 0.7424570322036743), ('kingdom', 0.7337411642074585),

hacer traducciones.

O four

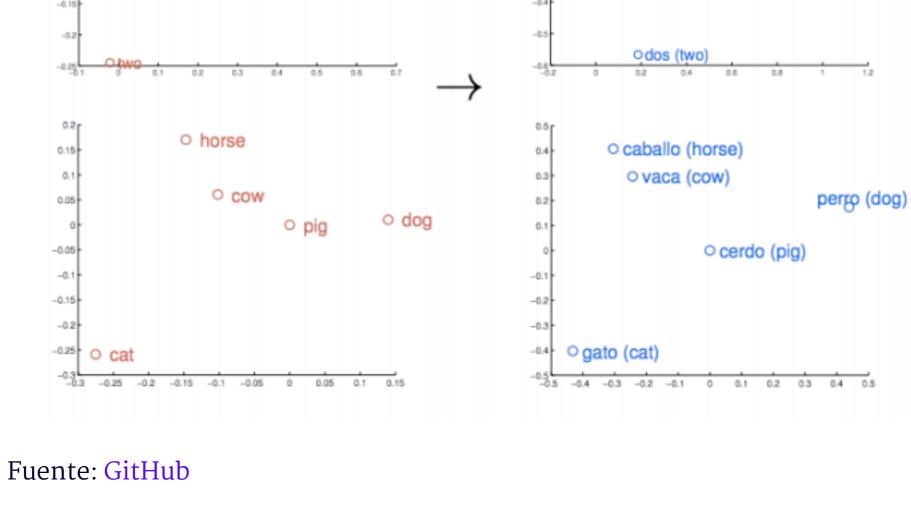
Ofive

womar

king-man+woman

queen Fuente: GitHub Podemos usar esta misma lógica para encontrar similitudes entre palabras de

dos idiomas diferentes, buscando aquellas que ocupan lugares cercanos, para



Sentiment Analysis

My experience The product is Your support team so far has been ok I guess is useless fantastic! NEUTRAL POSITIVE **NEGATIVE** MonkeyLearn

Las relaciones sociales son complejas. Sobre cualquier tema de interés podemos

personas. Esto tiene que advertirte algo importante: el lenguaje está cargado de

Hay muchos casos donde los modelos fracasan en manejar bien determinados

algoritmos que detectan "discurso de odio" tendían a señalar erróneamente

señalar, como mínimo, puntos de vista contradictorios, cuando no verdaderos

conflictos sociales. ¡El lenguaje no es ajeno a estas disputas! Sabes, por

ejemplo, en que existen palabras prejuiciosas para designar a grupos de

subjetividad y de sesgos. Nunca podemos considerarlo como una entidad

plenamente objetiva, aislada de un contexto social. Al analizar nuestros

usos del lenguaje. Por ejemplo, un estudio por Sap et al. (2019) halló que

datasets y armar nuestros modelos, debemos tenerlo en cuenta.

"aprendidos" por algoritmos de procesamiento de lenguaje natural. Por ejemplo, dado que actualmente los hombres son mayoría en el campo de la programación el algoritmo que autocompleta las búsquedas en Google Search tiende a autocompletar las búsquedas con sugerencias que amplifican los prejuicios existentes. Un estudio del Pew Research Center mostró que las búsquedas de imágenes de profesiones comunes sobre-representaban a los hombres (es decir, mostraban menos mujeres de las que realmente hay en esos trabajos). Por ejemplo, la

que se cargan a las páginas web y que el algoritmo de Google levanta sin filtrar por diversidad? ¿O con lo que clickean los usuarios que hacen las búsquedas y el algoritmo aprende subiendo la visibilidad de algún grupo social en desmedro de otros? ¿O alguna otra opción, o más de una a la vez? ¿Cómo podría ser una posible solución? Es importante incluir chequeos contra los sesgos en el diseño de los algoritmos, ya que muchas veces son invisibles para el ojo no entrenado y pasan desapercibidos hasta que es demasiado tarde y ya contribuyó a introducir más desigualdad en el mundo.

Por último, no podemos dejar de mencionar los grandes avances en

Learning como GPT-2 y GPT-3. No entraremos en detalles sobre el

Procesamiento del Lenguaje Natural que representan arquitecturas de Deep

funcionamiento de estos modelos - los detalles técnicos exceden los contenidos

de la carrera - pero sí podemos plantear algunas preguntas asociadas a su uso. ¿Está bien que exista - al menos sin regulación - herramientas capaces de generar noticias y comentarios falsos?¿Qué ocurrirá cuando podamos escribir un texto sobre algún tema dando una simple instrucción a una computadora? ¿Qué impactos tendrá sobre la manera de comunicarnos?¿Y sobre el mercado laboral? Incluso, hay demostraciones donde dando una simple instrucción, estos modelos pueden generar bloques de código - por ahora sencillo - y éste es ejecutado correctamente. ¿Creés que puede ser usado con fines poco éticos? Empieza a trabajar en el notebook

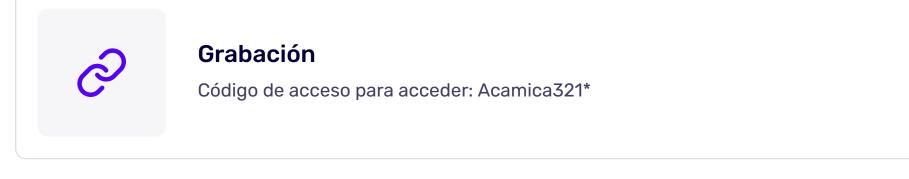
DS_Bitácora_39_y_40_NLP

¡Prepárate para el próximo encuentro!

Profundiza Te invitamos a conocer más sobre el tema de esta bitácora.

Contenido extra de la meeting

Procesamiento Lenguaje Natural (Part...



☐ Grabación Presentación

Lista de recursos

⊘ DS_Bitácora_39_y_40_N Profundiza

Grabación

Contenido extra

palabras en vectores numéricos y así permite operar matemáticamente sobre

es una dimensión o feature). Podemos visualizarlo como una cadena de Para facilitar las cosas, olvidaremos por ahora los números y nos enfocaremos

en los colores. Estos son los vectores para las palabras en inglés "king" (rey), "man" (hombre) y "woman" (mujer), luego de entrenar un Word2Vec sobre las

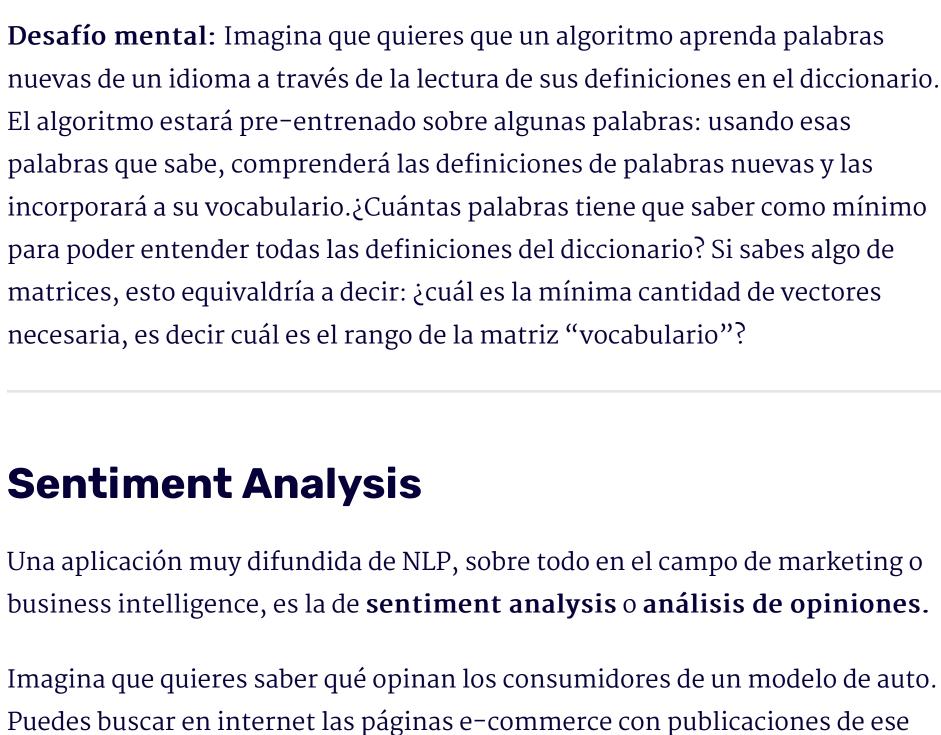
especialmente entre "man" y "woman". Esto significa que ocupan lugares cercanos en el espacio vectorial. Palabras con relación sintáctica o semántica se

cual significa que existe una relación similar entre ambos pares. Como dijimos, podemos hacer operaciones matemáticas con estos vectores. Volviendo al ejemplo anterior, si, tomamos el vector de la palabra "king", le

Aquí vemos exactamente el vector que resulta de hacer la operación numérica king - man + woman, y su cercanía al vector queen (no siempre Word2Vec da resultados tan prolijos, aunque por ahora no veremos las posibles excepciones). king − man + woman ~= queen

Ouno (one)

ocinco (five) oone Otres (three) o three



modelo, grupos en redes sociales y foros de aficionados a los autos y **scrapear**

comentarios en positivos o negativos. Incluso podrías identificar sensaciones

tratar de distinguir qué aspectos del auto están asociados a opiniones positivas

o negativas. Por ejemplo, si la palabra "motor" está asociada a oraciones con

opinión positiva y la palabra "consumo" de combustible está asociada a

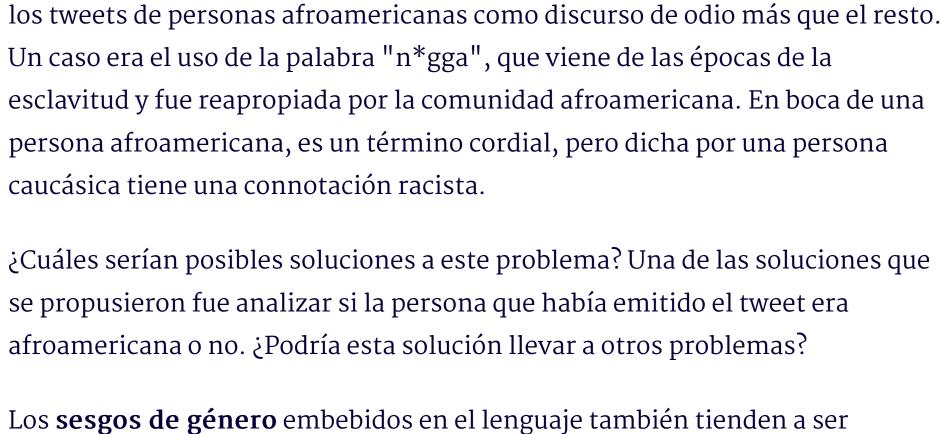
como alegría, enojo o desilusión. Con un poco más de complejidad, podrías

Luego, aplicando un sencillo sentiment analysis, podrías clasificar los

los comentarios de los/as usuarios/as sobre el auto.

oraciones que denotan una opinión negativa.

Sesgos en NLP



búsqueda de CEOs mostraban solamente 10% de mujeres a pesar de que las mujeres representan un 28%. ¿Tendrá que ver esto con sesgos en las imágenes

Para el encuentro, trae ya hecha tal parte, y completen durante la clase esto...