Minería de textos

Máster Online en Ciencia de Datos





Dr. José Raúl Romero

Profesor Titular de la Universidad de Córdoba y Doctor en Ingeniería Informática por la Universidad de Málaga. Sus líneas actuales de trabajo se centran en la democratización de la ciencia de datos (*Automated ML* y *Explainable Artificial Intelligence*), aprendizaje automático evolutivo y análitica de software (aplicación de aprendizaje y optimización a la mejora del proceso de desarrollo de software).

Miembro del Consejo de Administración de la *European Association for Data Science*, e investigador senior del Instituto de Investigación Andaluz de *Data Science and Computational Intelligence*.

Director del **Máster Online en Ciencia de Datos** de la Universidad de Córdoba.



UNIVERSIDAD Ð CÓRDOBA

Agrupamiento de textos

Introducción al agrupamiento de documentos





Introducción

- El agrupamiento (*clustering*) determina cómo se deben reunir documentos en grupos que compartan características
- Es habitual hacer uso de los algoritmos de agrupamiento más habituales:
 - Modelos jerárquico. Se basan en el concepto de que objetos similares estarán más cercanos a objetos relacionados en el espacio del vector que los no relacionados.
 Permite la visualización en un diagrama de árbol (dendrogram) mostrando una jerarquía exhaustiva de clústeres
 - Modelos basados en centroide. Construyen grupos que tienen un elemento central representativo para cada clúster, de forma que exhibe las características distintivas que lo distingue del resto de grupos. k-means y k-medoids son ejemplos representativos



Introducción

- Modelos basados en distribución. Aplican distribuciones de probabilidad de forma que objetos con distribuciones similares estén en el mismo grupo. Estos modelos permiten capturar dependencias y correlaciones entre características y atributos, aunque son propensos al over-fitting
- Modelos basados en densidad. Generan clústeres a partir de puntos de datos que se agrupan juntos en áreas de alta densidad, lo que podría ocurrir de forma aleatoria en zonas dispersas del espacio del vector. Estas zonas dispersas se tratan como ruido o sirven de frontera.
- El agrupamiento puede tener similitudes con la obtención de n-gramas



n-gramas Vs. clustering de documentos

- El agrupamiento puede verse con similitudes con la obtención de n-gramas
- La ventana deslizante se utiliza para determinar si 2 o más palabras son similares
 - Cuenta las apariciones conjuntas de un número de palabras y se mueve por el texto
 - El conteo se refleja en una tabla o matriz que muestra la cercanía en la relación entre palabras
 - Trata todo el texto como una única unidad, sin diferenciar inicios, rupturas, paradas, etc. →
 solo importa la secuencia de palabras → modelos secuenciales
- El método de la **bolsa de palabras** (*bag of words*) utiliza el vector de palabras
 - También realiza un conteo por filas (documento) como bloques de texto
 - En este caso, la secuencia no es determinante
 - Una bolsa de palabras es la metáfora de un espacio físico alojando palabras relacionadas
 - Muy criticado por no considerar aspectos clave relacionados con el significado de las palabras



n-gramas Vs. clustering de documentos

Se puede considerar que los métodos de obtención de n-gramas realizan un agrupamiento de palabras (word clustering) frente a métodos requeridos para poder agrupar documentos (document clustering)



Clustering de palabras

- ¿Por qué deberíamos agrupar (clustering) si ya hemos creado una matriz TDM? Esta matriz ya está agrupando variables en forma de tabla
- El análisis factorial es el método para generación de n-gramas más utilizado cuando el texto está dispuesto en forma tabular
 - Los factores ofrecen una vista numérica, simplificando los datos al conjunto de relaciones más simple
 - El análisis factorial es un método de reducción de datos





Clustering de palabras

- Los métodos de clustering ofrecen patrones de similitud en los datos más sutiles
 - Los métodos de clustering trabajan sobre las filas de la matriz, es decir, los documentos
 - Es importante escoger adecuadamente el método porque tienden a tener estructuras "preferidas", esto es, tienden a estructuras en los datos que otros métodos podrían describir mejor
 - En texto es complicado definir la "similitud", por lo que distintos métodos pueden llegar a distintas conclusiones sobre el mismo documento
- El **objetivo subyacente** a los métodos de analítica de textos es extraer información de forma concisa y condensada
 - Todo proceso de reducción implica pérdidas de información respecto al original → en texto, no podemos cuantificar cuánta información perdemos ni qué información necesitamos mantener
 - El problema reside en la transformación de texto a variables → no hay métricas de pérdida
 - Esto hace difícil decidir qué método de clustering de palabras es mejor que otro

Agrupamiento de textos

Clustering de documentos





Clustering de documentos

- Puede interesarnos categorizar el documento completo en base a la información aportada por los términos que lo contienen
- Recordemos que un documento puede venir representado por algo muy breve (tweet) o algo muy extenso (libro completo)



Ejemplo práctico

- Vamos a realizar clustering sobre una serie de tweets de Elon Musk
- Lo primero que haremos, una vez hayamos cargado los tweets, será preprocesarlos siguiendo el pipeline explicado en los temas anteriores

```
index
                                                                                                                                                       Cada tweet se corresponde con un
   14 Prevideas for paying ~$10B devicost incl. Kickstarter & collecting underpants, which turned out to be um... less lucrative than expected
                                                                                                                                                       documento y, por lo tanto, va a ser
      Tesla Semi truck unveil & test ride tentatively scheduled for Oct 26th in Hawthorne. Worth seeing this beast in person. It's unreal
                                                                                                                                                       preprocesado de manera independiente.
      China, Russia, soon all countries w strong computer science. Competition for Al superiority at national level most likely cause of WW3 imo.
      Putting together SpaceX rocket landing blooper reel. We messed up a lot before it finally worked, but there's some epic explosion footage?
      To be clear, a Hyperloop passenger version wouldn't have intense light strobe effect (just for testing), nor uncomfortable ad
                                                                                                                          def preprocess tweet(tweet text):
      Will run the SpaceX pusher sled later this week and see what it can do
                                                                                                                             # pasamos todo el texto a miniscula
      Btw, high accel only needed because tube is short. For passenger transport, this can be spread over 20+ miles, so no sp 31
                                                                                                                             tweet text = tweet text.lower()
                                                                                                                             # tokenizamos
      Might be possible to go supersonic in our test Hyperloop tube, even though it's only 0.8 miles long. Very high accel/decel
                                                                                                                             words = [word tokenize(sentence) for sentence in sent tokenize(tweet text)]
      Max recovered booster velocity was Mach 7.9 (Bulgarian Sat). Energy is velocity squared, so this is a bigger difference the
                                                                                                                             words = [x for xs in words for x in xs] # flatten the list
      Max velocity: Mach 6.9 Max altitude: 247 km Highest so far, but velocity matters much more
                                                                                                                             # eliminamos las stop words y los signos de puntuación
      Touchdown: Vertical Velocity (m/s): -1.47 Lateral Velocity (m/s): -0.15 Tilt (deg): 0.40 ? Lateral position: 0.7m from target
                                                                                                                             stop words = nltk.corpus.stopwords.words("english")
 124 Watching eclipse with sunglasses on through the Model S glass roof. Wow!
                                                                                                                             stop words.extend(list(punctuation))
 129 Pics of SpaceX spaceSuit developed for NASA commercial crew program coming out next week. Undergoing ocean landi
                                                                                                                             words = [word for word in words if word not in stop words]
  140 Wigs me out too much. Deleting.
                                                                                                                             # realizamos el steamming
 152 Would like to express our appreciation to Microsoft for use of their Azure cloud computing platform. This required massive
                                                                                                                             words = [PorterStemmer().stem(word) for word in words]
                                                                                                                             return " ".join(words)
```



Ejemplo práctico

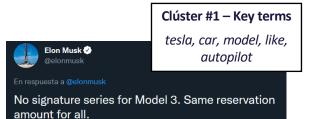
- Utilizaremos el TF-IDF de los términos de los tweets como entrada del algoritmo de clustering
- El objetivo es que cada clúster quede caracterizado por una serie de términos o n-gramas

	capitol	capitol build	captain	captain hercul	car
nobodi like regul everyth car plane food drug etc 's danger public regul ai	0.0	0.0	0.0	0.0	0.102583
openai first ever defeat world 's best player competit esport vastli complex tradit board game like chess go	0.0	0.0	0.0	0.0	0.000000
could n't believ incred inspir creativ	0.0	0.0	0.0	0.0	0.000000
project loveday winner	0.0	0.0	0.0	0.0	0.000000
want happen fast pleas let local feder elect repres know make big differ hear	0.0	0.0	0.0	0.0	0.000000





Esta figura es un ejemplo reducido con solo 5 tweets y n-gramas. Cabe destacar que se obtiene una matriz muy sparse





Clúster #2 – Key terms rocket, launch, falcon, good, land

Rocket is extra toasty and hit the deck hard (used almost all of the emergency crush core), but otherwise good





