

Procesamiento del lenguaje natural

Introducción y modelos de lenguaje



INTRODUCCIÓN

Procesamiento del lenguaje natural

- El procesamiento del lenguaje natural (PLN) es un área de intersección entre la **Inteligencia Artificial** y la **Lingüística**, que estudia la comunicación entre personas y máquinas por medio del lenguaje natural.
 - Sistemas capaces de interpretar o generar lenguaje natural
- El lenguaje natural es una forma de comunicación *imprecisa* y *ambigua* que se apoya en el conocimiento *compartido* por los que se comunican

Procesamiento del lenguaje natural

● Tareas de análisis de lenguaje natural

- Análisis morfo-léxico, sintáctico, semántico y pragmático

● Aplicaciones

- Recuperación de información
- Clustering y clasificación de textos
- Análisis de sentimientos
- Extracción de entidades y relaciones
- Traducción automática
- Generación automática de resúmenes
- Sistemas de diálogo
- Generación de lenguaje natural

Análisis del lenguaje natural

- Gran dificultad debido a que el lenguaje es algo vivo: en continua expansión, que va modificándose...
 - El lenguaje se modifica tanto en vocabulario como en sintaxis
- Existencia de jergas locales, profesionales, por franjas de edad...
- Ambigüedad
 - **Léxica**: *Entró en el banco. Se sentó en el banco.*
 - Polisemia: ambigüedad de las palabras
 - **Sintáctica**: *Juan vio a María con unos prismáticos.*
 - A veces es imposible de solucionar
 - **Semántica**: *Los niños compraron el libro de Peter Pan.*
 - **Referencial**: *El jamón está en el armario. Sácalo. Ciérralo.*
- Requiere mucho conocimiento: objetivos del hablante, hipótesis, contexto... No es mera transmisión de palabras...

Niveles de análisis del lenguaje I

- Tokenization:

- Frases -> Palabras

- “kids made good snacks.”=[“kids”, “made”, “good”, “snacks”, “.”]

- Análisis morfológico:

- Palabra -> raíz + sufijos/prefijos + feature

- kids=kid+s+<plural>

- made=make+<past tense>

- Análisis léxico:

- Palabras -> Etiquetas léxicas (POS tags)

- [“kids/NN”, “made/V”, “good/ADJ”, “snacks/NN”, “./PUNCT”]

Niveles de análisis del lenguaje II

- **Análisis sintáctico (parsing):**
 - Palabras + POS tags -> Estructura de la frase
 - 2 tipos:
 - Análisis de constituyentes.
 - Análisis de dependencias.
- **Análisis semántico:**
 - Palabras + POS tags -> Significado de la palabra (desambiguación)
 - Significado de palabras + estructura de frase: significado de frase (normalmente, independientemente del contexto)
- **Resolución de referencias:**
 - Expresiones de referencia + contexto: referente semántico
- **Análisis Pragmático:**
 - Significado de frase + contexto: significado más profundo

Ejemplo de herramienta de PLN

- Demo de herramienta de PLN
 - Freeling:
 - <http://nlp.lsi.upc.edu/freeling/demo/demo.php>
 - Se pueden configurar distintas opciones de análisis y distintas configuraciones de salida, en distintos idiomas

Etapas históricas

- Podemos distinguir 3 grandes etapas históricas
 - PLN simbólico (1950s – 1990s): **Técnicas lingüísticas formales**
 - Se basan en el desarrollo de reglas estructurales que se aplican en las distintas fases del análisis.
 - Encaje de patrones, gramáticas, sistemas de reglas, ontologías, etc.
 - PLN estadístico (1990s-2010s): **Técnicas empiricistas o probabilísticas**
 - Se basan en el estudio de una serie de características de tipo probabilístico asociadas a las distintas fases del análisis del lenguaje.
 - Estas características son extraídas de un conjunto de textos de referencia (corpus)
 - Aprendizaje supervisado, semi-supervisado y no supervisado
 - PLN y redes neuronales (2010s-presente): **Deep Learning**
 - Arquitecturas específicas para trabajar con lenguaje natural
 - Requieren grandes corpus de documentos y poder computacional
 - Grandes avances en los últimos años en áreas como la traducción automática o los sistemas conversacionales

Procesamiento del lenguaje natural

- El problema de los **métodos basados en técnicas lingüísticas formales** es la dificultad de codificar manualmente todo el conocimiento lingüístico necesario (diccionarios, gramáticas, etc.)
 - Esto obliga a trabajar con un lenguaje reducido
 - Este tipo de métodos no suelen contemplar la capacidad de aprendizaje
- En cambio, los **métodos basados en técnicas probabilísticas** aprenden a partir de datos prácticos (corpus de documentos)
- Los sistemas que utilizan métodos necesitan una **fase de entrenamiento** en la que se les debe proporcionar un número suficiente de ejemplos
 - Corpus anotado
- El uso de **redes neuronales profundas** permite encontrar patrones estadísticos complejos en los corpus de documentos

MODELOS DE LENGUAJE

Modelos de lenguaje probabilísticos

- Un **modelo probabilístico del lenguaje** define una distribución de probabilidad sobre el conjunto de elementos a partir de los valores observados en un **corpus de documentos**
 - Según cual sea el objeto de análisis los elementos pueden ser fonemas, letras, sílabas, o palabras
 - Las frecuencias de aparición de cada uno de los elementos son las que se observen en el corpus
- Son realmente útiles hoy día en multitud de tareas de PLN
 - Texto predictivo: $P(\text{Qué tal } \mathbf{estás}) > P(\text{Qué tal } \mathbf{has comido})$
 - Traducción automática: $P(\text{Voy de visita a su } \mathbf{casa}) > P(\text{Voy de visita a su } \mathbf{hogar})$
 - Corrección ortográfica: $P(\text{Tenemos } \mathbf{calor}) > P(\text{Tenemos } \mathbf{color})$
 - Reconocimiento del habla: $P(\text{Se hizo daño } \mathbf{a sí mismo}) > P(\text{Se hizo daño } \mathbf{así mismo})$

Modelos de lenguaje probabilísticos

- Vamos a considerar que nuestros elementos son palabras, pero lo que veamos aplica igual para fonemas, sílabas, etc.
- Un **modelo probabilístico del lenguaje** permite entre otras cosas
 - Calcular la probabilidad de encontrar una frase o secuencia de palabras determinada
 - $P(\text{Yo, quiero, comer, macarrones, con, tomate})$
 - Calcular la probabilidad de la siguiente palabra
 - $P(\text{tomate} \mid \text{Yo, quiero, comer, macarrones, con})$
- Como el texto es secuencial, podemos pensar que la probabilidad de una palabra depende de todas las anteriores
 - Para ello necesitamos refrescar ciertas nociones de probabilidad

Teoría de probabilidad aplicada a PLN

- Probabilidad condicionada: La probabilidad de B habiendo observado A
 - $P(B|A) = P(A, B) / P(A)$
- La probabilidad de que una palabra sea (por ejemplo) ‘perro’, sabiendo que la primera palabra es ‘el’, es la fracción de veces que ‘el’ aparece seguido de ‘perro’ en nuestro corpus
 - $P(\text{perro} | \text{el}) = \text{numVeces}(\text{el}, \text{perro}) / \text{numVeces}(\text{el})$
- La probabilidad condicionada se reescribe como $P(A, B) = P(A) P(B | A)$
 - Si añadimos elementos: $P(A, B, C, D) = P(A)P(B|A) P(C|A, B)P(D|A, B, C)$
 - La probabilidad de cada palabra depende de todas las anteriores
- Esto se generaliza mediante la regla de la cadena

$$P(w_1, \dots, w_n) = P(w_1) P(w_2 | w_1) P(w_3 | w_1, w_2) \dots P(w_n | w_1, \dots, w_{n-1})$$

Relajando la hipótesis de la regla de la cadena

- Según la regla de la cadena, la probabilidad de una palabra depende de todas las anteriores
 - Sin embargo, esta hipótesis es impracticable porque no hay un corpus tan grande para asignar probabilidad a las posibles combinaciones de palabras
 - Siempre habrá alguna combinación que no esté presente en el corpus
 - Se puede relajar la hipótesis de la regla de la cadena, haciendo así factibles los cálculos y obteniendo excelentes resultados
- En lugar de considerar que la probabilidad de un elemento depende de todos los anteriores, supone que solamente los $n-1$ elementos anteriores tienen efecto sobre las probabilidades del siguiente elemento i -ésimo.
 - Hipótesis de Markov: $P(w_i | w_1, \dots, w_{i-1}) \approx P(w_i | w_{i-(n-1)}, \dots, w_{i-1})$
- Esta hipótesis no tiene en cuenta que en lenguaje puede haber dependencias de “larga distancia”, sin embargo, funciona muy bien para algunas tareas

Modelos probabilístico de los n-gramas

- El modelo n-grama es uno de los modelos estadísticos del lenguaje más simples pero más útiles.
 - Se puede aplicar a fonemas, letras, sílabas, palabras... aquí nos centraremos en palabras
- El modelo n-grama utiliza la hipótesis de Markov para indicar que la dependencia es con los n-1 anteriores. Por ejemplo:
 - Modelo bigrama ($n = 2$): $P(w_n | w_{n-1})$
 - Modelo trigramas ($n = 3$): $P(w_n | w_{n-1}, w_{n-2})$

Ejemplo de estimación de bigramas

- Supongamos este corpus
 - <s> Yo quiero patatas </s>
 - <s> Yo quiero patatas con carne </s>
 - <s> No quiero carne </s>
 - <s> Quiero carne con patatas </s>
 - <s> Hoy quiero carne </s>
 - <s> Quiero dormir </s>

Estimación óptima por
máxima verosimilitud

$$P(w_n | w_{n-1}) = \frac{\text{frec}(w_{n-1}, w_n)}{\text{frec}(w_{n-1})}$$

<s> indica “inicio de oración”

</s> indica “final de oración”

- Estos son las probabilidades de algunos bigramas del corpus
 - $P(\text{yo} | \text{<s>}) = \text{frec}(\text{<s>}, \text{yo}) / \text{frec}(\text{<s>}) = 2/6 = 0,333$
 - $P(\text{no} | \text{<s>}) = \text{frec}(\text{<s>}, \text{no}) / \text{frec}(\text{<s>}) = 1/6 = 0,167$
 - $P(\text{quiero} | \text{yo}) = \text{frec}(\text{yo}, \text{quiero}) / \text{frec}(\text{yo}) = 2/2 = 1$
 - $P(\text{patatas} | \text{quiero}) = \text{frec}(\text{quiero}, \text{patatas}) / \text{frec}(\text{quiero}) = 2/6 = 0,333$
 - $P(\text{carne} | \text{quiero}) = \text{frec}(\text{quiero}, \text{carne}) / \text{frec}(\text{quiero}) = 3/6 = 0,5$
 - $P(\text{dormir} | \text{quiero}) = \text{frec}(\text{quiero}, \text{carne}) / \text{frec}(\text{quiero}) = 1/6 = 0,167$
- Según este modelo, la continuación más segura de <s>Yo quiero...
 - **$P(\text{carne} | \text{quiero}) > P(\text{patatas} | \text{quiero}) > P(\text{dormir} | \text{quiero})$**
 - La respuesta cambia si consideramos trigramas $\rightarrow P(\text{patatas} | \text{yo}, \text{quiero}) = 1$

Estimación de la probabilidad de una frase

- Probabilidades de bigramas obtenidos de un corpus supuesto
 - $P(\text{yo} \mid \langle s \rangle) = 0,25$ $P(\text{quiero} \mid \langle s \rangle) = 0,75$
 - $P(\text{quiero} \mid \text{yo}) = 0,5$ $P(\text{tengo} \mid \text{yo}) = 0,2$ $P(\text{soy} \mid \text{yo}) = 0,3$
 - $P(\text{ser} \mid \text{quiero}) = 0,8$ $P(\text{tomar} \mid \text{quiero}) = 0,2$
 - $P(\text{café} \mid \text{tomar}) = 0,6$ $P(\text{leche} \mid \text{tomar}) = 0,3$ $P(\text{distancia} \mid \text{tomar}) = 0,1$
 - $P(\text{artista} \mid \text{ser}) = 0,9$ $P(\text{informático} \mid \text{ser}) = 0,1$

- $P(\text{Yo quiero tomar café}) =$
 - $P(\text{yo} \mid \langle s \rangle) P(\text{quiero} \mid \text{yo}) P(\text{tomar} \mid \text{quiero}) P(\text{café} \mid \text{tomar}) = 0,25 * 0,5 * 0,2 * 0,6 = 0,015$
 - La frase entera puede no estar en el corpus. Lo normal es que no esté.
 - ¡Si un bigrama no está la probabilidad de la frase es cero!

- La frase que se genera siguiendo la opción más probable es: “*Quiero ser artista*”
 - $P(\text{quiero} \mid \langle s \rangle) P(\text{ser} \mid \text{quiero}) P(\text{artista} \mid \text{ser}) = 0,75 * 0,8 * 0,9 = 0,54$

Estimación de probabilidad en casos raros

- Como hemos visto cualquier n-grama que no esté en el corpus recibe probabilidad 0, el alisado de Laplace alivia los problemas de estimación de probabilidades en casos raros
- El **alisado de Laplace** calcula cualquier probabilidad condicional considerando que ha habido unas observaciones adicionales virtuales de todos y cada uno de los n-gramas posibles
- Siendo AB un bigrama (observado o no), el valor de la probabilidad condicional $P(B|A)$ alisado según Laplace

$$P(B|A) = \frac{\text{frec}(A, B) + t}{\text{frec}(A) + t * m}$$

Donde

- t es el número de observaciones virtuales adicionales
- m es el número de monogramas (palabras) existentes en el corpus
- Las probabilidades totales siguen sumando 1

Alisado por interpolación lineal

- El **alisado por interpolación lineal** es ligeramente más sofisticado
- La interpolación lineal usa la probabilidad incondicional $P(w_2)$ calculada a partir de los datos para hacer que la probabilidad condicional $P(w_2|w_1)$ se parezca a ella, de la siguiente forma

$$P_{\text{Int}}(w_2|w_1) = \alpha P(w_2|w_1) + (1 - \alpha)P(w_2)$$

donde $\alpha \in [0,1]$ regula el peso que se le da a la probabilidad condicional y la probabilidad no condicionada

- El valor de α se puede fijar empíricamente con el fin de ajustar el rendimiento
 - También se puede hacer dependiente del “contexto”
 - Si existen muchos bigramas con la palabra w_1 entonces es mejor un valor alto
 - Si no existen muchos bigramas con la palabra w_1 es mejor un valor bajo

Enlaces

- Enlaces interesantes:
 - **Modelos de lenguaje y n-gramas**
 - <https://heuristic-bhabha-ae33da.netlify.app/modelos-de-lenguaje-y-n-gramas.html>