



Procesamiento de Lenguaje Natural

Modelos de Lenguaje

Mauricio Toledo-Acosta
mauricio.toledo@unison.mx

Departamento de Matemáticas
Universidad de Sonora



Section 1

Introducción



Referencias

- **Chapter 3.** Jurafsky, D., Martin, J. H. (2019). Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition.
- **Chapter II.6.** Eisenstein, J. (2018). Natural language processing. Jacob Eisenstein.



Objetivo

Procesamiento
de Lenguaje
Natural

Introducción

Modelo de
 n -gramas

Generalizaciones

Evaluando
Modelos de
Lenguaje

Aspectos
adicionales

Algunas
aplicaciones

¿Cuál es la siguiente palabra?

En el parque, los niños juegan con ...



Objetivo

Introducción

Generalizaciones

¿Cuál es la siguiente palabra?

En el parque, los niños juegan con ...

¿Cuál es la siguiente palabra?

El agua hierve a 100 grados ...



- *café el en mesa la sobre libro un dejé olvidé y*



- *café el en mesa la sobre libro un dejé olvidé y*
- *dejé un libro sobre la mesa en el café y lo olvidé*



6/37





- 1 **Modelos de Lenguaje Estadísticos.** Predicen la siguiente palabra dadas las palabras que le preceden, esto lo hacen usando conteos de co-ocurrencias. El modelo de n -gramas es el más sencillo.
- 2 **Modelos de Lenguaje Neuronales.** Predicen la siguiente palabra usando embeddings que capturan diversos fenómenos lingüísticos a partir del uso de redes neuronales.
 - Modelos neuronales *clásicos*: Word2Vec, FastText, GloVe, ConceptNet, ...
 - Modelos basados en mecanismos de atención: Bert, LLaMa, GPT, Claude, ModernBert, DeepSeek, ...



¿Cómo se calculan las probabilidades?

Procesamiento
de Lenguaje
Natural

Introducción

Modelo de
 n -gramas

Generalizaciones

Evaluando
Modelos de
Lenguaje

Aspectos
adicionales

Algunas
aplicaciones

Queremos calcular la probabilidad de una secuencia de palabras $W = w_1, w_2, \dots, w_n$. Esto lo hacemos con la probabilidad conjunta:

$$P(W) = P(w_1, w_2, \dots, w_n).$$

Una tarea relacionada es calcular la probabilidad condicional

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

Ambas están relacionadas por la regla de la cadena

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



10/37



La distribución de probabilidad del valor futuro de una variable aleatoria depende únicamente de su valor presente, siendo independiente de la historia de dicha variable.

Podemos hacer las siguientes aproximaciones:

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

$$P(w_n|w_1w_2...w_{n-1}) \approx P(w_n|w_{n-1})$$

$$P(w_n|w_1w_2...w_{n-1}) \approx P(w_n|w_{n-2}w_{n-1})$$

• • •



Al decir n -gramas nos referimos a las secuencias o al modelo predictivo que asigna probabilidades.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 13/37



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 14/37



Ejemplo: probabilidades

Procesamiento de Lenguaje Natural

Introducción

Modelo de n -gramas

Generalizaciones

Evaluando Modelos de Lenguaje

Aspectos adicionales

Algunas aplicaciones

	me	gusta	comer	tacos	tamales	mucho	hoy
me	0.07	0.72	0	0.03	0.01	0	0.17
gusta	0.04	0	0.80	0.02	0.08	0.12	0
comer	0	0.04	0	0.55	0.45	0	0
tacos	0.10	0	0.30	0	0.50	0.10	0
tamales	0	0.14	0.29	0.71	0	0.14	0
mucho	0.22	0.33	0	0	0.06	0	0.39
hoy	0.18	0	0.09	0.09	0	0.64	0



Algunas ventajas

Procesamiento de Lenguaje Natural

Introducción

Modelo de n -gramas

Generalizaciones

Evaluando Modelos de Lenguaje

Aspectos adicionales

Algunas aplicaciones

- Los n -gramas ayudan a capturar la información contextual y la semántica dentro de una secuencia de palabras, proporcionando una comprensión más matizada del lenguaje.
- En tareas de recuperación de información (information retrieval), los n -gramas ayudan a emparejar y clasificar documentos según la relevancia de los patrones de n -gramas.
- Los n -gramas sirven como características poderosas en la clasificación de texto y el análisis de sentimientos, capturando patrones significativos que contribuyen a la caracterización de diferentes clases o sentimientos.





the in	rain Spain	in falls
Spain mainly	falls on	mainly the
on plain		



n -gramas sintácticos

Procesamiento
de Lenguaje
Natural

Introducción

Modelo de
 n -gramas

Generalizaciones

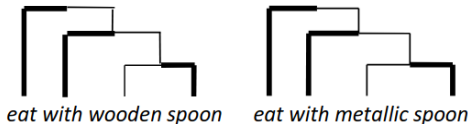
Evaluando
Modelos de
Lenguaje

Aspectos
adicionales

Algunas
aplicaciones

A diferencia de los n -gramas donde las subsecuencias se toman en el orden en el que aparecen en el texto, en los **n -gramas sintácticos** los vecinos se toman siguiendo las relaciones sintácticas de los árboles de dependencia sintáctica.

eat with wooden spoon eat with metallic spoon



¿Cuántos bigramas y bigramas sintácticos en común tienen?



22/37



La **perplejidad** es una métrica para evaluar el rendimiento de un modelo de lenguaje. Mide la incertidumbre del modelo para predecir la próxima palabra en una secuencia. Cuanto menor sea la perplejidad, mayor será la capacidad del modelo para predecir la palabra siguiente.

$$\begin{aligned} \text{Pp}(W) &= \sqrt[N]{\frac{1}{P(w_1 w_2 \cdots w_N)}} \\ &= \sqrt[N]{\prod_i \frac{1}{P(w_i | w_{i-1})}} \end{aligned}$$

W es la secuencia entera de palabras de un conjunto de prueba.



Considerar textos escritos en AAVE (African American Vernacular English):



Perplejidad

Generalizaciones

- Al calcular la perplejidad, el modelo de n -gramas debe construirse sin ningún conocimiento del conjunto de prueba. De otra forma la perplejidad puede ser artificialmente baja.
- La perplejidad de dos modelos lingüísticos sólo es comparable si utilizan vocabularios idénticos.
- Una mejora en la perplejidad (m. intrínseca) no garantiza una mejora del rendimiento de una tarea de PLN como el reconocimiento del habla o la traducción automática (m. extrínseca). Sin embargo, la perplejidad suele estar correlacionada con dichas mejoras.



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻ 26/37



They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions



Si un n -grama aparece un número suficiente de veces, podemos tener una buena estimación de su probabilidad. Sin embargo, es probable que falten algunas secuencias de palabras perfectamente aceptables. Tendremos casos de n -gramas de probabilidad 0 que deberían tener probabilidad distinta de 0.

denied the allegations:	5
denied the speculation:	2
denied the rumors:	1
denied the report:	1

denied the offer: 0
denied the loan: 0



Frame Title

Procesamiento de Lenguaje Natural

Introducción

Modelo de n -gramas

Generalizaciones

Evaluando Modelos de Lenguaje

Aspectos adicionales

Algunas aplicaciones

En los modelos de lenguaje estadísticos, como los n -gramas, a menudo nos encontramos con una secuencia de palabras que nunca apareció en el texto de entrenamiento. Si le asignamos una probabilidad de 0, el modelo se *rompe* porque cualquier oración que contenga esa secuencia tendrá una probabilidad total de 0. El suavizado es la técnica para evitar este problema, asignando una pequeña parte de la probabilidad a eventos no vistos.



Suavizado: Palabras desconocidas

Procesamiento de Lenguaje Natural

Introducción

Modelo de *n*-gramas

Generalizaciones

Evaluando Modelos de Lenguaje

Aspectos adicionales

Algunas aplicaciones

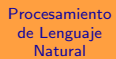
¿Qué pasa con palabras que nunca ha visto en el entrenamiento?

En un sistema de **vocabulario cerrado** el conjunto de prueba no contiene palabras desconocidas. En un sistema de **vocabulario abierto**, tenemos que lidiar con palabras que no hemos visto antes, a las que llamaremos **palabras fuera de vocabulario (OOV)**. Podemos lidiar con estas palabras desconocidas añadiendo una pseudopalabra llamada $\langle \text{UNK} \rangle$.



- Convertir un sistema abierto en uno cerrado:
 - Escoger un vocabulario fijo.
 - Convertir cualquier palabra OOV en $\langle \text{UNK} \rangle$.
 - Estimar las probabilidades para $\langle \text{UNK} \rangle$ de la forma usual, como si fuera una palabra *normal*.
- Crear un vocabulario fijo implícito. Reemplazamos palabras por $\langle \text{UNK} \rangle$ en el entrenamiento basándonos en su frecuencia.





Algunas aplicaciones





36/37

