

Introducción a Natural Language Processing

Modelos de Lenguaje

Ronald Cárdenas Acosta

Setiembre, 2016

Outline

- 1 Definición
- 2 Modelado
- 3 Estimación de Probabilidades
- 4 Evaluación

Outline

- 1 Definición
- 2 Modelado
- 3 Estimación de Probabilidades
- 4 Evaluación

Modelos de Lenguaje (*Language Models*)

Objetivo

Asignar una probabilidad a una oración

Aplicaciones

- Traducción: $P(\text{'atrapo al raton gato el'}) < P(\text{'el gato atrapo al raton'})$
- Corrección ortográfica: $P(\text{'el gato atrapo al rtaon'}) < P(\text{'el gato atrapo al raton'})$
- Reconocimiento de voz: $P(\text{'el tubo de agua'}) > P(\text{'el tuvo de agua'})$
- Generación de resúmenes, question-answering, etc.

Modelos de Lenguaje Probabilísticos

- Se representa la oración como una secuencia de palabras.
$$P(S) = P(w_1, w_2, \dots, w_n)$$
- Alternativamente, se puede hallar la probabilidad de la siguiente palabra:
$$P(w_3 | w_1, w_2)$$
- Un modelo que calcula cualquiera de estas probabilidades se llama *Modelo de Lenguaje*.

Outline

- 1 Definición
- 2 Modelado**
- 3 Estimación de Probabilidades
- 4 Evaluación

Recordatorio de Probabilidades

La regla de la cadena:

$$P(A, B, C, D) = P(A).P(B|A).P(C|A, B).P(D|A, B, C)$$

En general:

$$P(x_1, x_2, \dots, x_n) = P(x_1).P(x_2|x_1)...P(x_n|x_1, x_2, \dots, x_{n-1})$$

Probabilidad de una oración: Regla de la cadena

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i | w_1, w_2, \dots, w_{i-1})$$

- Ejemplo:
 $P(\text{hoy esta soleado en Lima}) =$
 $P(\text{hoy}).P(\text{esta}|\text{hoy}).P(\text{soleado}|\text{hoy,esta}).$
 $P(\text{en}|\text{hoy,esta,soleado}).P(\text{Lima}|\text{hoy,esta,soleado,en})$
- ¿Como estimar estas probabilidades? Conteo de ocurrencias en corpus
 $P(\text{Lima}|\text{hoy,esta,soleado,en}) = \frac{\text{cuenta}(\text{hoy,esta,soleado,en,Lima})}{\text{cuenta}(\text{hoy,esta,soleado,en})}$
- Si la oración es muy larga, el número de oraciones posibles se vuelve muy largo, y la cuenta muy baja.

Suposición de Markov

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i | w_{i-k}, \dots, w_{i-1})$$

- Donde k es el número de palabras anteriores a considerar (*grado de suposición de Markov*).
- Para $k = 1$:
 $P(\text{Lima} | \text{hoy, esta, soleado, en}) \approx P(\text{Lima} | \text{en})$
- Para $k = 2$:
 $P(\text{Lima} | \text{hoy, esta, soleado, en}) \approx P(\text{Lima} | \text{soleado, en})$

Modelos N-grama

- **Modelo Unigrama:** Cuando $k = 1$

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

- **Modelo Bigrama:** Cuando $k = 2$

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

- **Modelo Trigrama:** Cuando $k = 3$

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i | w_{i-2}, w_{i-1})$$

Outline

- 1 Definición
- 2 Modelado
- 3 Estimación de Probabilidades**
- 4 Evaluación

Estimación de Probabilidades: Bigramas

Recibe el nombre de *Maximum Likelihood Estimate* (estimación de máxima probabilidad)

$$P(w_i | w_{i-1}) = \frac{\text{cuenta}(w_{i-1}, w_i)}{\text{cuenta}(w_{i-1})}$$

Se suele agregar indicadores de inicio y termino de oración. Ejemplo: $\langle \text{START} \rangle$ hoy esta soleado en Lima $\langle \text{STOP} \rangle$. De manera que se debe considerar también:

$P(\text{hoy} | \langle \text{START} \rangle)$: probabilidad que hoy sea la primera palabra

$P(\langle \text{STOP} \rangle | \text{Lima})$: probabilidad que Lima sea la última palabra

Outline

- 1 Definición
- 2 Modelado
- 3 Estimación de Probabilidades
- 4 Evaluación**

Evaluación de modelos de lenguaje

- Evaluación Extrínseca:
evaluar la tarea en la que el modelo es parte (traducción, corrector ortográfico)
- Evaluación Intrínseca:
estimar el grado de ajuste de las probabilidades a la data [Perplexity]

Perplexity

- Cuantifica las ramificaciones de las posibles combinaciones de palabras:
"Factor de ramificación"
- Se calcula usando probabilidades logarítmicas para evitar *underflow*
- Minimizar el Perplexity es equivalente a maximizar la probabilidad
- Mientras menor Perplexity, mejor.

Perplexity

Probabilidad de la data de testeo, normalizado por el número de palabras.

$$PP(test) = P(w_1 w_2 \dots w_N)^{-\frac{1}{N}}$$

Expresado normalmente en prob. logarítmicas ($\log(P(test - data))$):

$$PP(test) = 2^{-L}$$

Donde:

$$L = \frac{1}{N_{test}} \sum_{i=1}^{N_{test}} \log(p(w_i | w_1 \dots w_{i-1}))$$