

Procesamiento de Lenguaje Natural

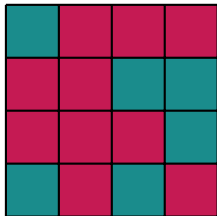
Olivia Gutú y Julio Weissman

Maestría en Ciencia de Datos

Semana 2: Análisis de sentimientos con Naïve Bayes



Corpus de N tuits



$$P(pos) = \frac{N_{pos}}{N} = \frac{6}{16}$$
$$P(neg) = 1 - P(pos) = \frac{10}{16}$$



$$P(feliz) = \frac{3}{16}$$
$$P(pos \cap feliz) = \frac{1}{16}$$

Probabilidades condicionales:

$$P(pos|feliz) = \frac{1}{3}$$



$$\frac{P(pos \cap feliz)}{P(feliz)}$$

$$P(feliz|pos) = \frac{1}{6}$$



$$\frac{P(feliz \cap pos)}{P(pos)}$$

por tanto:

$$P(pos|feliz)P(feliz) = P(feliz|pos)P(pos)$$

luego:

$$P(pos|feliz) = \frac{P(feliz|pos)P(pos)}{P(feliz)}$$

Tuits positivos
me encanta la playa
amo la playa

Tuits negativos
odio la playa
mar molesto

token w	frec(pos)	frec(neg)
me	1	0
encanta	1	0
la	2	1
playa	2	1
amo	1	0
odio	0	1
mar	0	1
molesto	0	1
N_{class}	7	5

token w	$P(w pos)$	$P(w neg)$
me	$\frac{1}{7} \approx 0.14$	0
encanta	$\frac{1}{7} \approx 0.14$	0
la	$\frac{2}{7} \approx 0.28$	$\frac{1}{5} = 0.2$
playa	$\frac{2}{7} \approx 0.28$	$\frac{1}{5} = 0.2$
amo	$\frac{1}{7} \approx 0.14$	0
odio	0	$\frac{1}{5} = 0.2$
mar	0	$\frac{1}{5} = 0.2$
molesto	0	$\frac{1}{5} = 0.2$

Palabras igualmente probables no aportan nada al sentimiento

token w	$P(w pos)$	$P(w neg)$
me	0.134	0.01
encanta	0.134	0.01
la	0.274	0.194
playa	0.274	0.194
amo	0.134	0.01
odio	0.01	0.194
mar	0.01	0.194
molesto	0.01	0.194

Naïve Bayes: radio de probabilidades



Universidad de Sonora

token w	$P(w pos)$	$P(w neg)$	ratio(w)
me	0.134	0.01	13.4
encanta	0.134	0.01	13.4
la	0.274	0.194	1.41
playa	0.274	0.194	1.41
amo	0.134	0.01	13.4
odio	0.01	0.194	0.05
mar	0.01	0.194	0.05
molesto	0.01	0.194	0.05

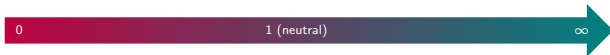
0

1 (neutral)

∞

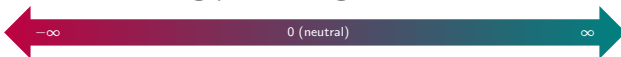
Nuevo tuit: $w_1 w_2 w_3 \cdots w_n$

$$\frac{P(pos|w_1 w_2 w_3 \cdots w_n)}{P(neg|w_1 w_2 w_3 \cdots w_n)} = \frac{P(pos)}{P(neg)} \prod_{i=1}^n \frac{P(w_i|pos)}{P(w_i|neg)}$$



$$\log \frac{P(pos)}{P(neg)} + \sum_{i=1}^n \log \frac{P(w_i|pos)}{P(w_i|neg)}$$

log prior log likelihood



evita errores numéricos

$$\lambda(w) = \log \text{ratio}(w), \quad \text{ratio}(w) = \frac{P(w|\text{pos})}{P(w|\text{neg})}$$

Nuevo tuit: $w_1 w_2 w_3 \cdots w_n$

Si las clases son equilibradas, log prior es igual a cero, en este caso:



si:

$$\sum_{i=1}^n \lambda(w_i) = \log \prod_{i=1}^n \text{ratio}(w_i) > 0$$



si:

$$\sum_{i=1}^n \lambda(w_i) = \log \prod_{i=1}^n \text{ratio}(w_i) \leq 0$$

- Recolectar tuits pre-clasificados (conjunto de entrenamiento)
- Pre-procesar los datos
- Establecer el vocabulario de tipos
- Contar las palabras $\text{frec}(w, \text{class})$, $\text{class} = \{\text{pos}, \text{neg}\}$
- Para cada palabra y cada clase calcular:

$$P(w|\text{class}) = \frac{\text{frec}(w, \text{class}) + 1}{N_{\text{class}} + |V|}$$

- Se calcula $\lambda(w) = \log \frac{P(w|\text{pos})}{P(w|\text{neg})}$
- Se calcula log prior

$$\log \frac{P(\text{pos})}{P(\text{neg})} = \log \frac{\text{núm. tuits positivos}}{\text{núm. tuits negativos}}$$

- convertir todo a minúsculas
- remover signos de puntuación, urls, nombres
- remover palabras vacías (*stop words*)
- aplicar *stemming*
- tokenizar las oraciones

Resultado: [aqui, pasa, nada, kino]

Nuevo tuit $T \rightarrow [w_1, w_2, \dots, w_n]$ (pre-procesamiento)

sentimiento(T) = 👍 si:

$$\log \text{prior} + \sum_{i=1}^n \lambda(w_i) > 0$$

sentimiento(T) = 👎 si:

$$\log \text{prior} + \sum_{i=1}^n \lambda(w_i) \leq 0$$

- Independencia (i.i.d.) ¡falso en PLN!
- Clases relativamente equilibradas