

Procesamiento de Lenguaje Natural

Olivia Gutú y Julio Weissman

Maestría en Ciencia de Datos
Semana 8: Modelos continuos de bolsa de palabras



- Identificar los conceptos claves de la representación de palabras
- Generar *word embedding vectors*
- Preparar el texto para el aprendizaje automático
- Implementar el modelo continuo de bolsa de palabras

- Un *integer i* por palabra: muy simple pero ‘cero’ sentido semántico.
- *one-hot vectors*: vectores e_i en $\mathbb{R}^{|V|}$, simple, no hay orden, pero vectores enormes, sin sentido semántico, la distancia entre cualquier par de vectores es 1.
- *word embedding vectors*: baja dimensión, significado semántico, inferencia. vectores cercanos semánticamente tienen distancia coseno cercana.

Significado de los vectores



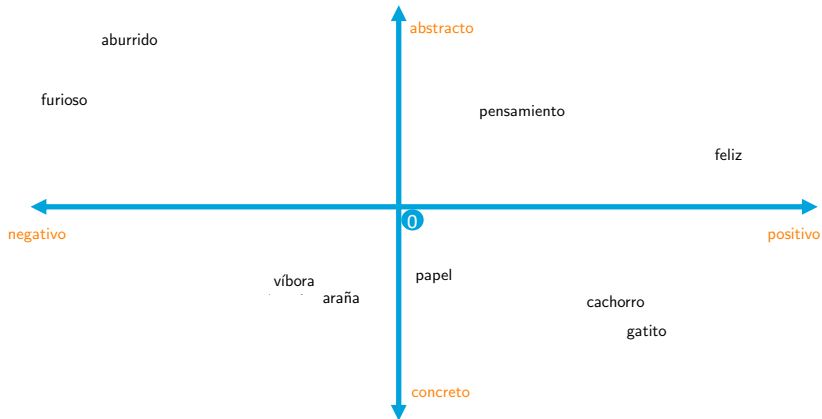
Universidad de Sonora



Significado de los vectores



Universidad de Sonora

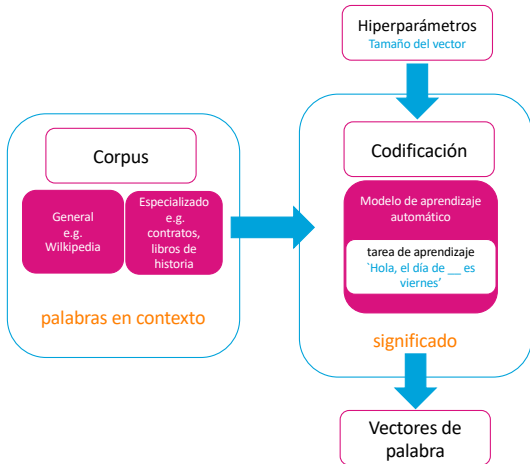


Básicos:

- **word2vec** (Google, 2013)
 - CBOW
 - Skip-Gram
- **GloVe** (Stanford, 2014)
- **fastText** (Facebook, 2016)
 - Soporte de OOV

Avanzados (deep learning)

- **BERT** (Google, 2018)
- **ELMo** (Allen Institute for AI, 2018)
- **GPT-2** (OpenAI, 2018)



Limpieza y tokenización

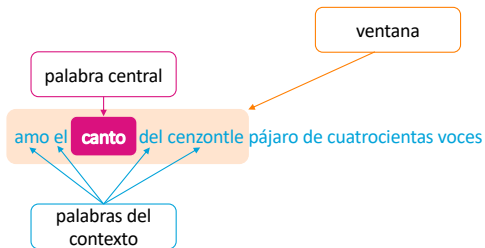
- Uniformizar mayúsculas y minúsculas `canto` = `Canto` = `CANTO`
lowecase/uppercase
- Puntuación , ! . ? ??? !!! \rightarrow . “” «» ’ \rightarrow \emptyset
- Números 1 2 3 9 10 \rightarrow \emptyset 3.1416 \rightarrow < NUMBER >
- Caracteres especiales % \$ § \rightarrow \emptyset
- Tokens especiales 😊 #nezahualcoyotl \rightarrow : *happy* : #nezahualcoyotl

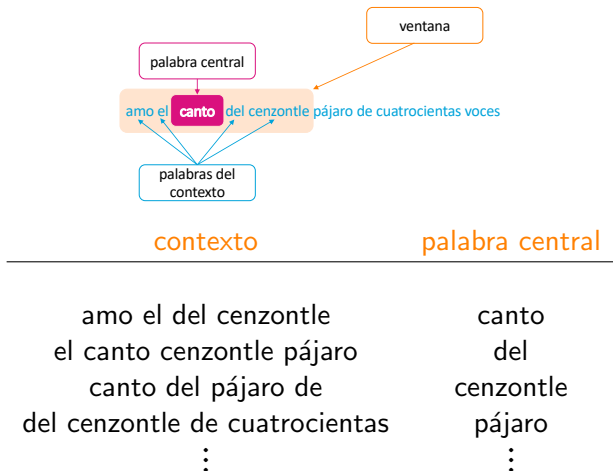
♡ 1000 el canto del cenztotl!!!! #nezahualcoyotl



[‘♡’, ‘el’, ‘canto’, ‘del’, ‘cenztotl’, ‘.’, ‘#nezahualcoyotl’]

- tamaño de ventana = 5 (window size)
- $C = 2$ (context half-size)





Codificación inicial de las palabras centrales

$V = \{\text{amo, canto, cenizonte, cuatrocientas, de, del, el, pájaro, voces}\}$

1	0	0	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	0	0	1	0	0	0	0	0
0	0	0	0	0	1	0	0	0
0	0	0	0	0	0	1	0	0
0	0	0	0	0	0	0	1	0
0	0	0	0	0	0	0	0	1

Codificación inicial del **contexto**: 'promedio' de las palabras involucradas.

	amo el del cenizonte				contexto
Promedio	1	0	0	0	0.25
	0	0	0	0	0
	0	0	0	1	0.25
	0	0	0	0	0
	0	0	0	0	0
	0	0	0	0	0
Promedio	0	0	1	0	0.25
	0	1	0	0	0.25
	0	0	0	0	0
	0	0	0	0	0
	0	0	0	0	0
	0	0	0	0	0

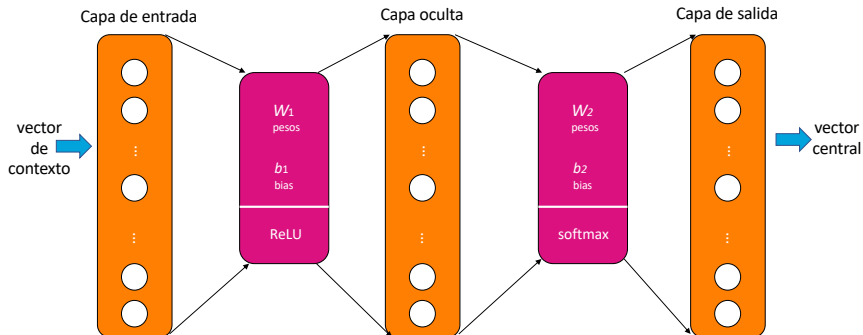
Preparación final

contexto	vector contexto	palabra central	vector palabra central
amo el del cenizontle	[0.25,0,0.25,0,0,0.25,0.25,0,0]	canto	[0,1,0,0,0,0,0,0,0]
⋮	⋮	⋮	⋮

Arquitectura del modelo CBOW



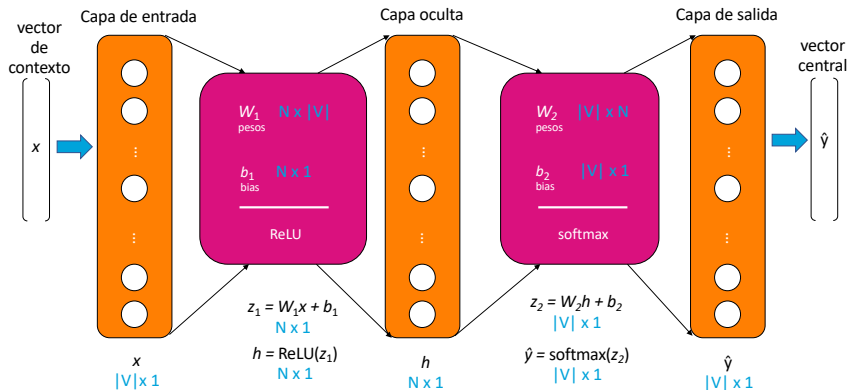
Universidad de Sonora



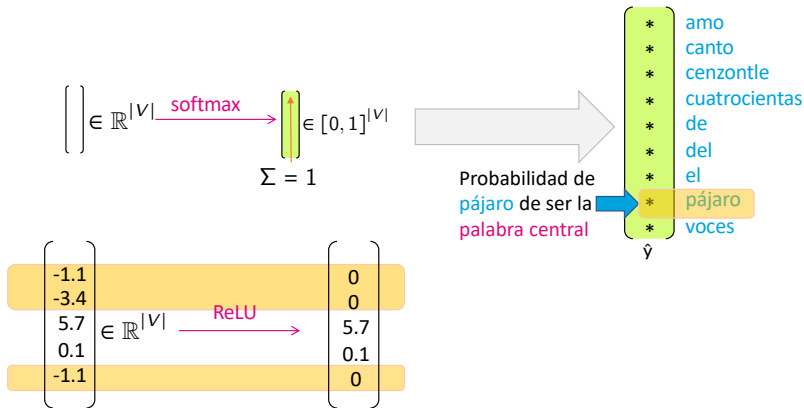
Dimensiones (una entrada)



Universidad de Sonora



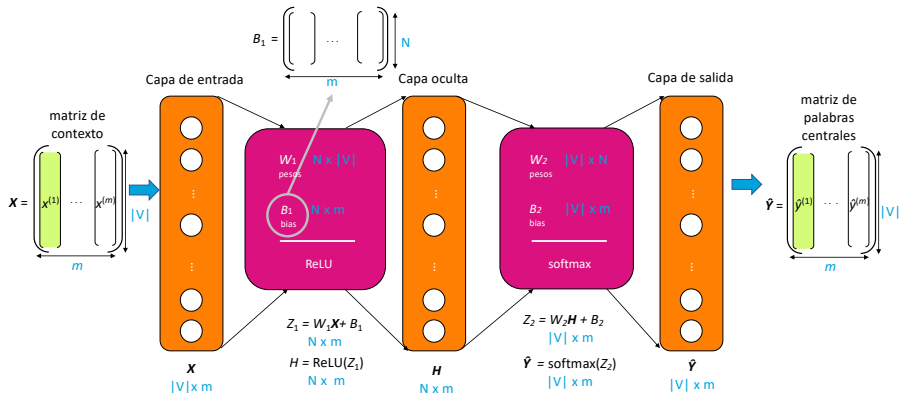
- $\text{ReLU}(z) = \max(0, z)$ entrada por entrada (*rectified linear unit*)
- $\text{softmax}(z)_i = \frac{\exp(z_i)}{\sum_{k=1}^{|V|} \exp(z_k)}$

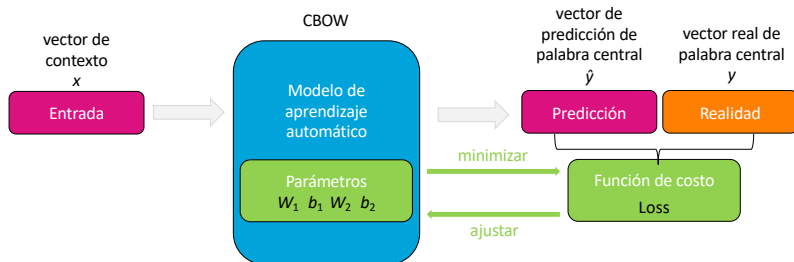


Dimensiones (entrada en *batch*)



Universidad de Sonora





$$J = -(y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + \dots + y_{|V|} \log \hat{y}_{|V|})$$

0 amo
0 canto
0 cenizontle
0 cuatrocientos
1 de
0 del
0 el
0 pájaro
0 voces

y

vector real de
palabra central

0.05
0.05
0.05
0.05
0.4
0.1
0.1
0.1
0.1

\hat{y}

-1.301
-1.301
-1.301
-1.301
-0.397
-1.000
-1.000
-1.000
-1.000

$\log \hat{y}$

vector de
predicción de
palabra central

0
0
0
0
-0.397
0
0
0
0

$y \log \hat{y}$

$$J = - \sum \downarrow = 0.397$$

$$J_{\text{batch}} = \frac{1}{m} \sum_{i=1}^m J_i$$

- Backpropagation

$\text{Grad} J_{\text{batch}}$

- Descenso de gradiente

$$(W_1, W_2, b_1, b_2) := (W_1, W_2, b_1, b_2) - \alpha \text{Grad} J_{\text{batch}}$$

Backpropagation

$$\frac{\partial J_{\text{batch}}}{\partial W_1} = \frac{1}{m} \text{ReLU}(W_2^T(\hat{Y} - Y)) X^T$$

$$\frac{\partial J_{\text{batch}}}{\partial W_2} = \frac{1}{m} (\hat{Y} - Y) H^T$$

$$\frac{\partial J_{\text{batch}}}{\partial b_1} = \frac{1}{m} \text{ReLU}(W_2^T(\hat{Y} - Y)) (1, \dots, 1)_m^T$$

$$\frac{\partial J_{\text{batch}}}{\partial b_2} = \frac{1}{m} (\hat{Y} - Y) (1, \dots, 1)_m^T$$

Descenso de gradiente

$$W_1 := W_1 - \alpha \frac{\partial J_{\text{batch}}}{\partial W_1}$$

$$W_2 := W_2 - \alpha \frac{\partial J_{\text{batch}}}{\partial W_2}$$

$$b_1 := b_1 - \alpha \frac{\partial J_{\text{batch}}}{\partial b_1}$$

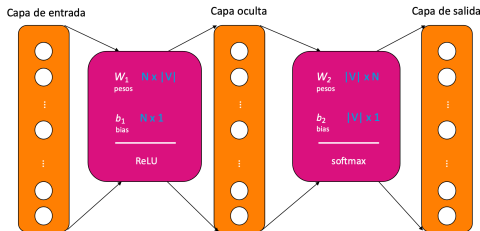
$$b_2 := b_2 - \alpha \frac{\partial J_{\text{batch}}}{\partial b_2}$$

α : tasa de aprendizaje

Extracción de vectores de palabra: opción 1



Universidad de Sonora



$$W_1 = \begin{pmatrix} \begin{bmatrix} w^{(1)} \end{bmatrix} & \dots & \begin{bmatrix} w^{(|V|)} \end{bmatrix} \end{pmatrix} \begin{matrix} \uparrow \\ N \end{matrix}$$

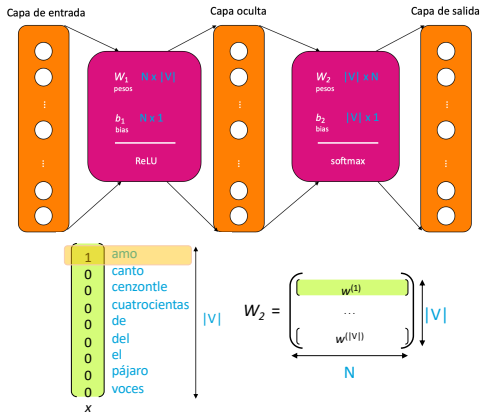
$|V|$



Extracción de vectores de palabra: opción 2



Universidad de Sonora



$$W_3 = \frac{1}{2}(W_1 + W_2^T)$$

$$W_3 = \left(\begin{array}{c|c|c} \boxed{w^{(1)}} & \dots & w^{(|V|)} \end{array} \right) \begin{array}{l} \updownarrow N \\ \hline \leftarrow |V| \rightarrow \end{array}$$

■ Analogías

- Semánticas: 'México' es a 'CDMX' lo que 'Argentina' es a $\langle ? \rangle$
- Sintácticas: 'cantó' es a 'cantar' lo que 'bailó' es a $\langle ? \rangle$

■ Clustering

fuelle: Google news

