

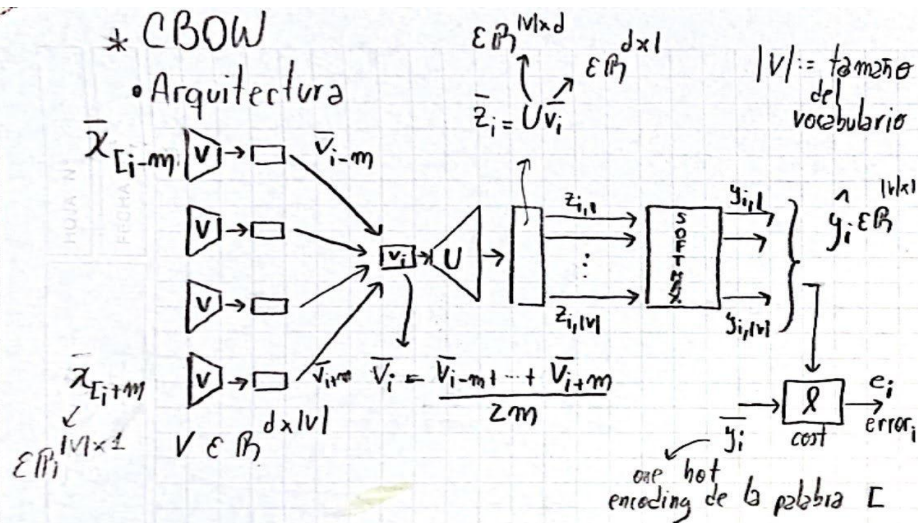
Procesamiento del Lenguaje Natural  
Facultad de Ingeniería  
Universidad de Buenos Aires

Ezequiel Esposito y Pablo Casas  
([eesposito@optiwe.com](mailto:eesposito@optiwe.com), [pcasas.biz@gmail.com](mailto:pcasas.biz@gmail.com) )



# \* CBOW

## • Arquitectura



## • Forward

$$\bar{v}_{t-m} = V \bar{x}_{t-m} \dots \bar{v}_{t+m} = V \bar{x}_{t+m}$$

$$\bar{v}_i = (\bar{v}_{t-m} + \dots + \bar{v}_{t+m}) / 2m$$

$$\bar{z}_i = U \bar{v}_i \quad \bar{z}_i \in \mathbb{R}^{|V| \times 1}$$

$$\hat{y}_i = \text{SOFTMAX}(\bar{z}_i) = \begin{bmatrix} e^{z_{i,1}} / \sum_{h=1}^{|V|} e^{z_{i,h}} \\ \vdots \\ e^{z_{i,|V|}} / \sum_{h=1}^{|V|} e^{z_{i,h}} \end{bmatrix}$$

$\bar{y}_i := \text{one hot encoding de la palabra central } \bar{x}_i$

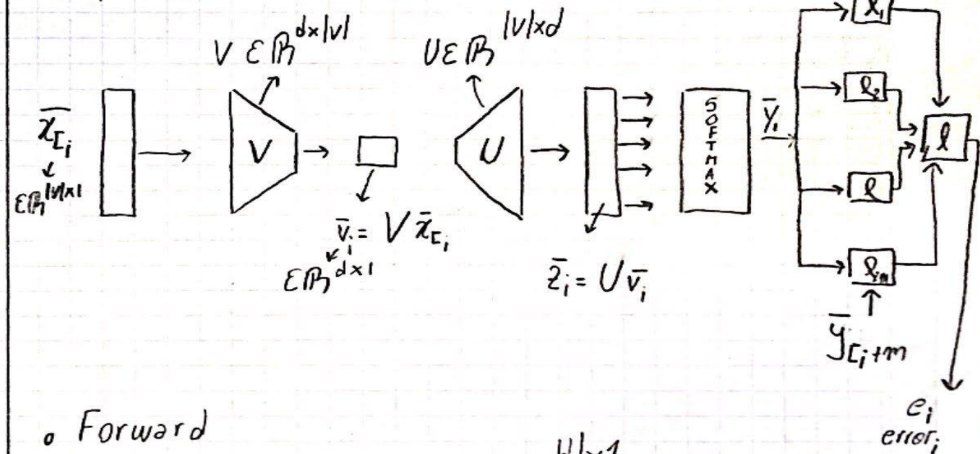
## • Cost

$$J = \prod_{i=1}^n P(\hat{y}_i = \bar{y}_i | \bar{x}_{t-m}, \bar{x}_{t+m})$$

multinomial vectorizada

# \* SKIP-GRAM

## • Arquitectura



## • Forward

$$\bar{v}_i = V \bar{x}_{t_i} \quad ; \quad \bar{v}_i \in \mathbb{R}^{d \times 1}$$

$$\bar{z}_i = U \bar{v}_i \quad ; \quad \bar{z}_i \in \mathbb{R}^{|V| \times 1}$$

$$\hat{y}_i = \text{SOFTMAX}(\bar{z}_i)$$

$\bar{y}_{t-m} = \text{one hot encoding de la palabra } x_{t-m}$

$\bar{y}_{t+m} = \text{one hot encoding de la palabra } x_{t+m}$

## • Cost

$$J = \prod_{i=1}^n \left[ P(\hat{y}_i = \bar{y}_{t-m} | \bar{x}_{t_i}) * \dots * P(\hat{y}_i = \bar{y}_{t+m} | \bar{x}_{t_i}) \right]$$

\* LBOW

• Backward

\* SHIP-GRAM

• Backward

- Recordemos sobre como calcular el cost en una clasificación de  $n$  clases
- Los labels pueden tomar valores entre  $\{0, \dots, L-1\}$  con  $L$  la cantidad de clases. Otra manera de representarlo es con vectores one-hot-encoding

$$y_i = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ para la clase } 0, \quad y_i = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \text{ para la clase } 1, \quad \dots, \quad y_i = \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix} \text{ para la clase } L-1$$

- La salida de nuestra red neuronal es el SOFTMAX:  $\hat{y}_i = \left[ \frac{e^{z_{i,0}}}{\sum_{n=0}^{L-1} e^{z_{i,n}}} \right]$
- Donde cada componente  $j$  de  $\hat{y}_i$  representa la probabilidad de que la salida de nuestro clasificador pertenezca a la clase  $j$ . Entonces podemos escribir de manera genérica

$$P(\hat{y}_i = \bar{y}_i \mid \bar{X}_i = x_i) = \sum_{j=0}^{L-1} y_{i,j} \frac{e^{z_{i,j}}}{\sum_{n=0}^{L-1} e^{z_{i,n}}} \quad \text{donde } y_{i,j} \text{ solo vale } 1 \text{ para un } j \text{ en particular}$$

Para cada muestra  
solo un término de la sumatoria es distinto a cero



\* E-BOW

$$J = \prod_{i=1}^m P(\hat{y}_i = \bar{y}_i | x_{L_i-m}, \dots, x_{L_i+m})$$

$$\arg \min_{\theta} -\frac{1}{n} \log \prod_{i=1}^m P(\hat{y}_i = \bar{y}_i | x_{L_i-m}, \dots, x_{L_i+m})$$

prop del log  $\rightarrow -\frac{1}{n} \sum_{i=1}^m \log [P(\hat{y}_i = \bar{y}_i | x_{L_i-m}, \dots, x_{L_i+m})]$

sustituya  $\rightarrow -\frac{1}{n} \sum_{i=1}^m \log \left[ \sum_{j=0}^{|V|} y_{ij} e^{z_{ij}} / \sum_{k=0}^{|V|} e^{z_{ik}} \right]$

solo el termino  $j=c$  es  $\neq 0$   $\rightarrow -\frac{1}{n} \sum_{i=1}^m \log \left[ \frac{y_{ic}}{e^{z_{ic}}} / \sum_{k=0}^{|V|} e^{z_{ik}} \right]$

prop logs  $\rightarrow -\frac{1}{n} \sum_{i=1}^m \left[ z_{ic} - \log \left( \sum_{k=0}^{|V|} e^{z_{ik}} \right) \right]$

$\Rightarrow$  Como  $z_i = U \bar{v}_i \rightarrow z_{ic} = \bar{\mu}_{ci}^T \bar{v}_i$  embedding de palabra  $\bar{v}_i$

$\Rightarrow \arg \min_{\theta} \sum_{i=1}^m \left( \bar{\mu}_{ci}^T \bar{v}_i - \log \sum_{k=0}^{|V|} e^{z_{ik}} \right)$

$z_{ik} = \bar{\mu}_{ki}^T \bar{v}_i$

\* SKIP-GRAM

(A)  $J = \prod_{i=1}^m \left[ P(\hat{y}_i = y_{L_i-m} | \bar{x}_{L_i}) * \dots * P(\hat{y}_i = y_{L_i+m} | \bar{x}_{L_i}) \right]$

$\arg \min_{\theta} -\frac{1}{n} \log \prod_{i=1}^m \left[ \prod_{\substack{j=0 \\ j \neq m}}^{2m} P(\hat{y}_i = y_{L_i-m+j} | \bar{x}_{L_i}) \right]$   $\rightarrow$  escribi más corto la productoria

(A)  $\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^m \log \left[ \prod_{\substack{j=0 \\ j \neq m}}^{2m} P(\hat{y}_i = y_{L_i-m+j} | \bar{x}_{L_i}) \right]$   $\rightarrow$  prop del log

$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^m \sum_{\substack{j=0 \\ j \neq m}}^{2m} \log \left( e^{z_{ij, L_i-m+j}} / \sum_{k=0}^{|V|} e^{z_{ik}} \right)$   $\rightarrow$  sustituya (solo el termino  $\neq 0$ )

$\arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^m \sum_{\substack{j=0 \\ j \neq m}}^{2m} z_{ij, L_i-m+j} - \log \sum_{k=0}^{|V|} e^{z_{ik}}$   $\rightarrow$  prop log

$\Rightarrow$  Como  $z_{ij, L_i-m+j} = \bar{\mu}_{L_i-m+j}^T \bar{v}_{L_i}$

$\Rightarrow \arg \min_{\theta} -\frac{1}{n} \sum_{i=1}^m \sum_{\substack{j=0 \\ j \neq m}}^{2m} \left[ \bar{\mu}_{L_i-m+j}^T \bar{v}_{L_i} - \log \sum_{k=0}^{|V|} e^{z_{ik}} \right]$

$z_{ik} = \bar{\mu}_{ki}^T \bar{v}_i$