

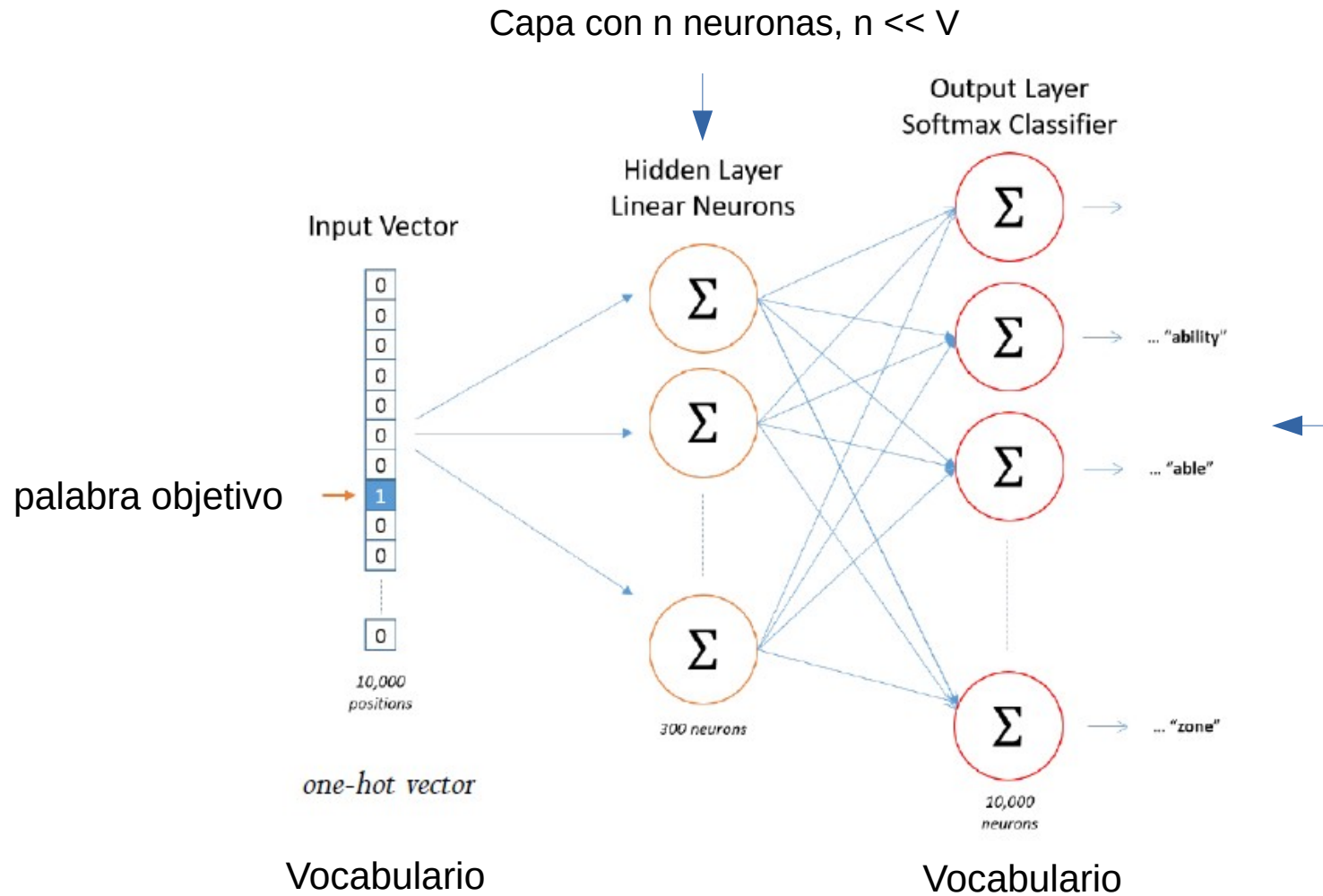


IIC 3800 Tópicos en CC NLP

<https://github.com/marcelomendoza/IIC3800>

- WORD2VEC -

Word vectorization

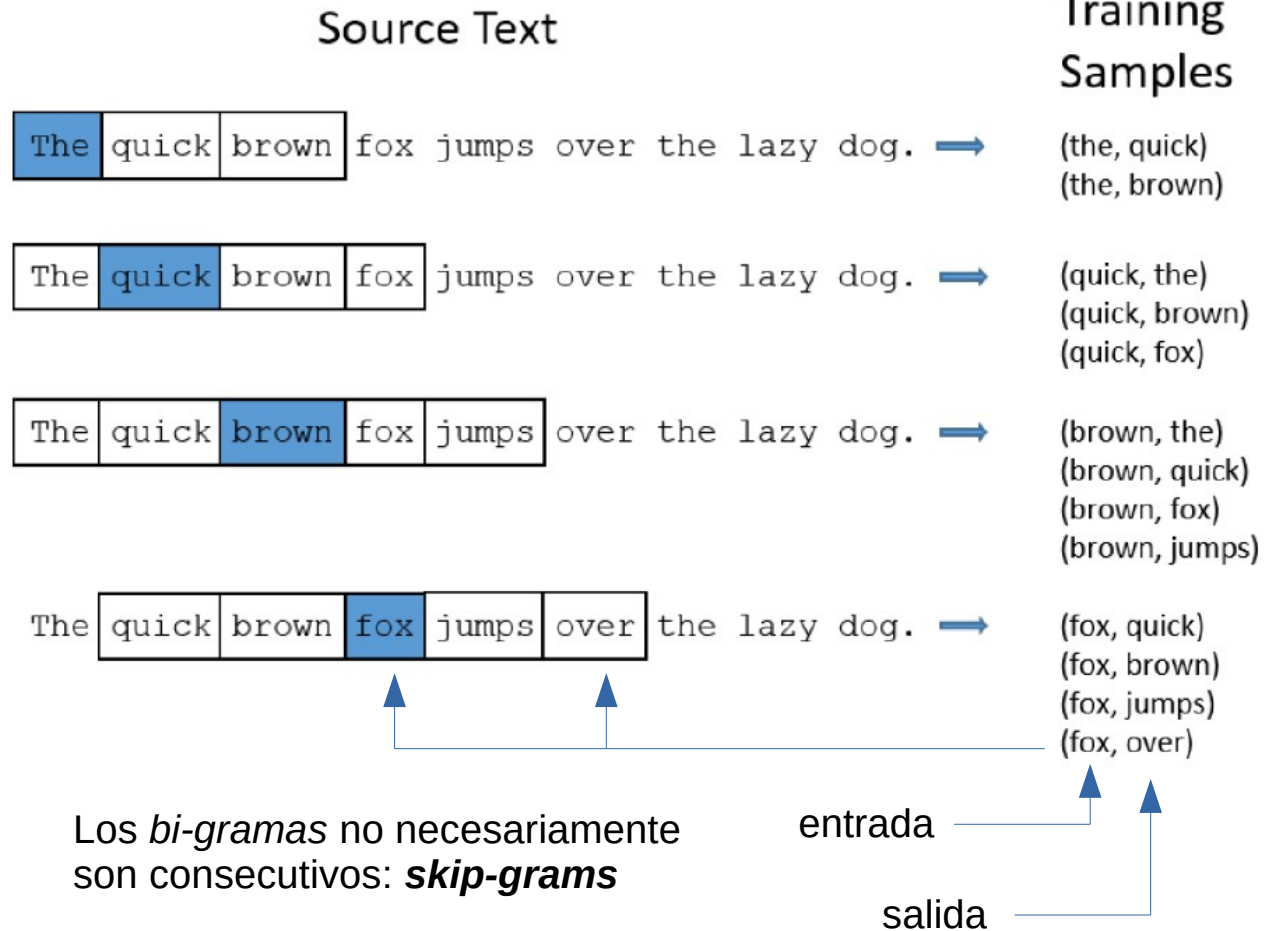


Word vectorization

- Skip-grams (a.k.a. Word2vec): entrenamiento de la red

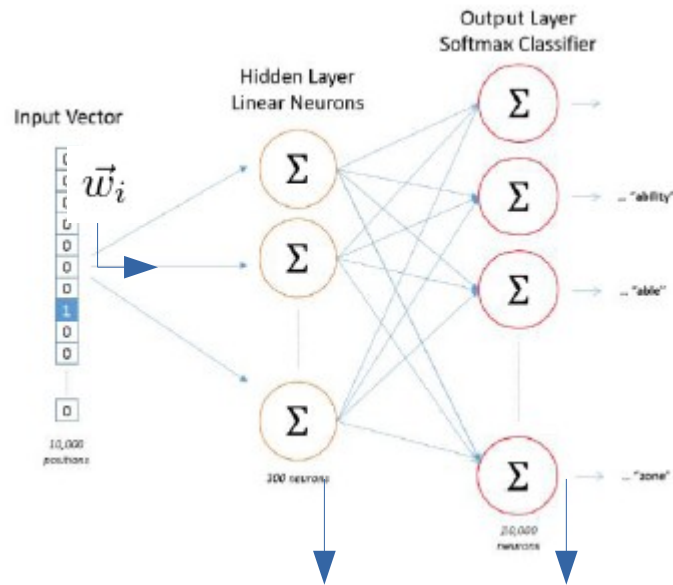
La red se entrena mostrando *bi-gramas*

Una ventana deslizando
recorre el texto →



Word vectorization

- Skip-grams (a.k.a. Word2vec): entrenamiento de la red



Función de pérdida:

$$\mathcal{L}_{SG} = -\frac{1}{|S|} \cdot \sum_{i=1}^{|S|} \sum_{-c \leq j \leq c, j \neq 0} \log(p(w_{i+j}|w_i))$$

$|S|$: # training skip-grams

$-c \leq j \leq c$: tamaño de la ventana de contexto (2c)

W_{in} o W_{out} pueden ser usados como *word embeddings*

Word vectorization

- Skip-grams (a.k.a. Word2vec): Como generar el training set
 - Tratando el desbalance entre skip-grams y pares no observados

Negative sampling:

- ▶ Seleccionamos aleatoriamente k ejemplos negativos (palabras que no están en C). Si no hiciéramos esto, **todas** las palabras que no están en C serían ejemplos negativos ($k = 5$).
- ▶ La probabilidad de seleccionar una palabra como ejemplo negativo es:

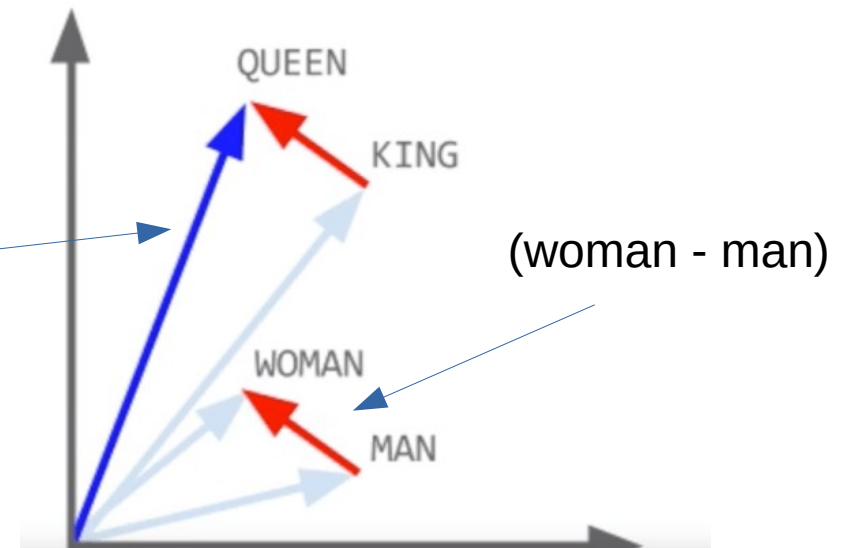
$$P(w_i) = \frac{f(w_i)^\beta}{\sum_{j=0}^n f(w_j)^\beta}$$

donde $0 < \beta < 1$ ($\beta \approx \frac{3}{4}$).

Word vectorization

- Operadores en word2vec: word analogies

king + (woman - man)



$$\arg \max_{b^* \in V} (\text{sim}(b^*, b - a + a^*))$$

Levy & Goldberg, Linguistic Regularities in Sparse and Explicit Word Representations, ACL'14.

Word vectorization

- Operadores en word2vec: `doesnt_match(['king', 'george', 'stephen', 'truck'])`

$$\arg \max_w f(w) = \left\| \sum_{v \in L \setminus w} \vec{v} \right\|, \quad \forall w \in L$$



Word vectorization

- Continuous Bag-of-Words (a.k.a. CBOW)

$$\mathcal{L}_{CBOW} = -\frac{1}{|S|} \cdot \sum_{i=1}^{|S|} \log(p(w_i | w_{i-c}, \dots, w_{i+c}))$$

↓ regularización

$$\mathcal{L} = \mathcal{L}_{CBOW} - \lambda \cdot \sum_V \|\vec{w}_i\|$$

