

# Word Embeddings

## Word2Vec



國立臺灣大學 資訊工程學系  
陳縉儂 助理教授

<http://vivianchen.idv.tw>

# Word2Vec – Skip-Gram Model

- Goal: predict surrounding words within a window of each word
- Objective function: maximize the probability of any context word given the current center word

$$w_1, w_2, \dots, \underbrace{w_{t-m}, \dots, w_{t-1}, \underbrace{w_t}_{w_I}, w_{t+1}, \dots, w_{t+m}}_{w_O}, \dots, w_{T-1}, w_T$$

context window

$$p(w_{O,1}, w_{O,2}, \dots, w_{O,C} \mid w_I) = \prod_{c=1}^C p(w_{O,c} \mid w_I)$$

target word vector

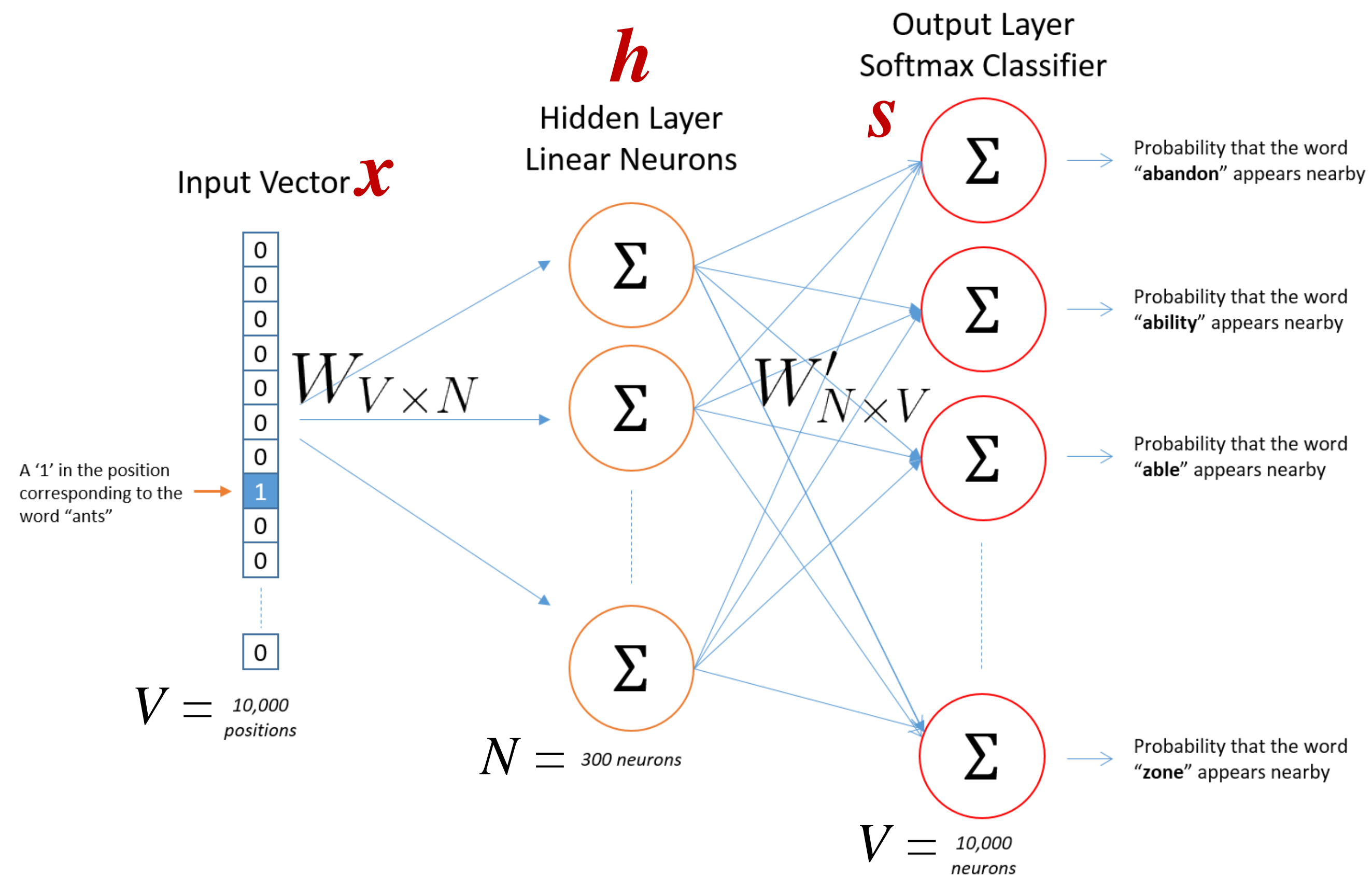
$$C(\theta) = - \sum_{w_I} \sum_{c=1}^C \log p(w_{O,c} \mid w_I) \quad \underbrace{p(w_O \mid w_I)}_{\text{outside target word}} = \frac{\exp(v_{w_O}'^T \underbrace{v_{w_I}}_{\text{target word}})}{\sum_j \exp(v_{w_j}'^T v_{w_I})}$$

Benefit: faster, easily incorporate a new sentence/document or add a word to vocab



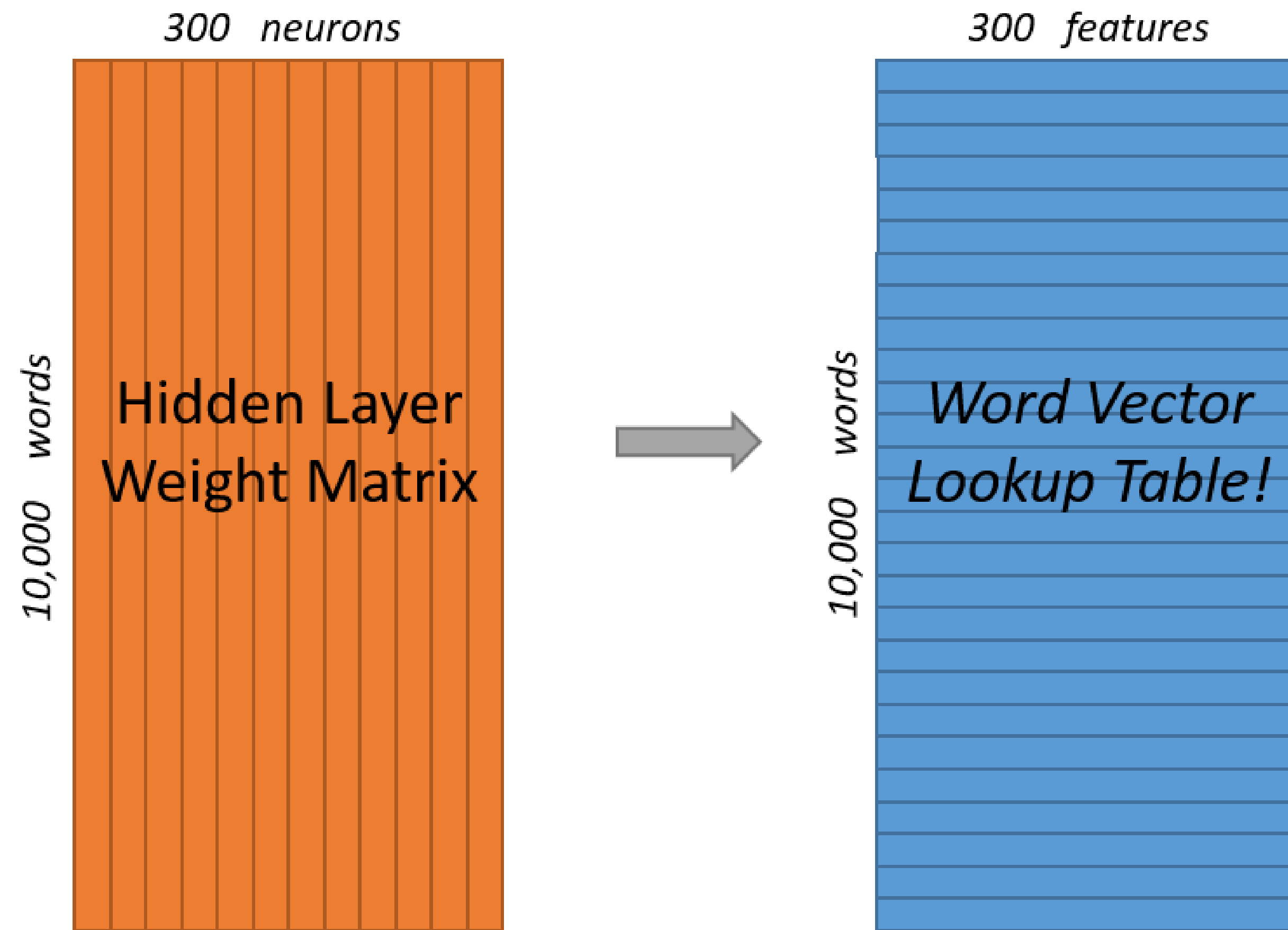
# Word2Vec Skip-Gram Illustration

- Goal: predict surrounding words within a window of each word



# Hidden Layer Matrix $\rightarrow$ Word Embedding Matrix

$$W_{V \times N}$$

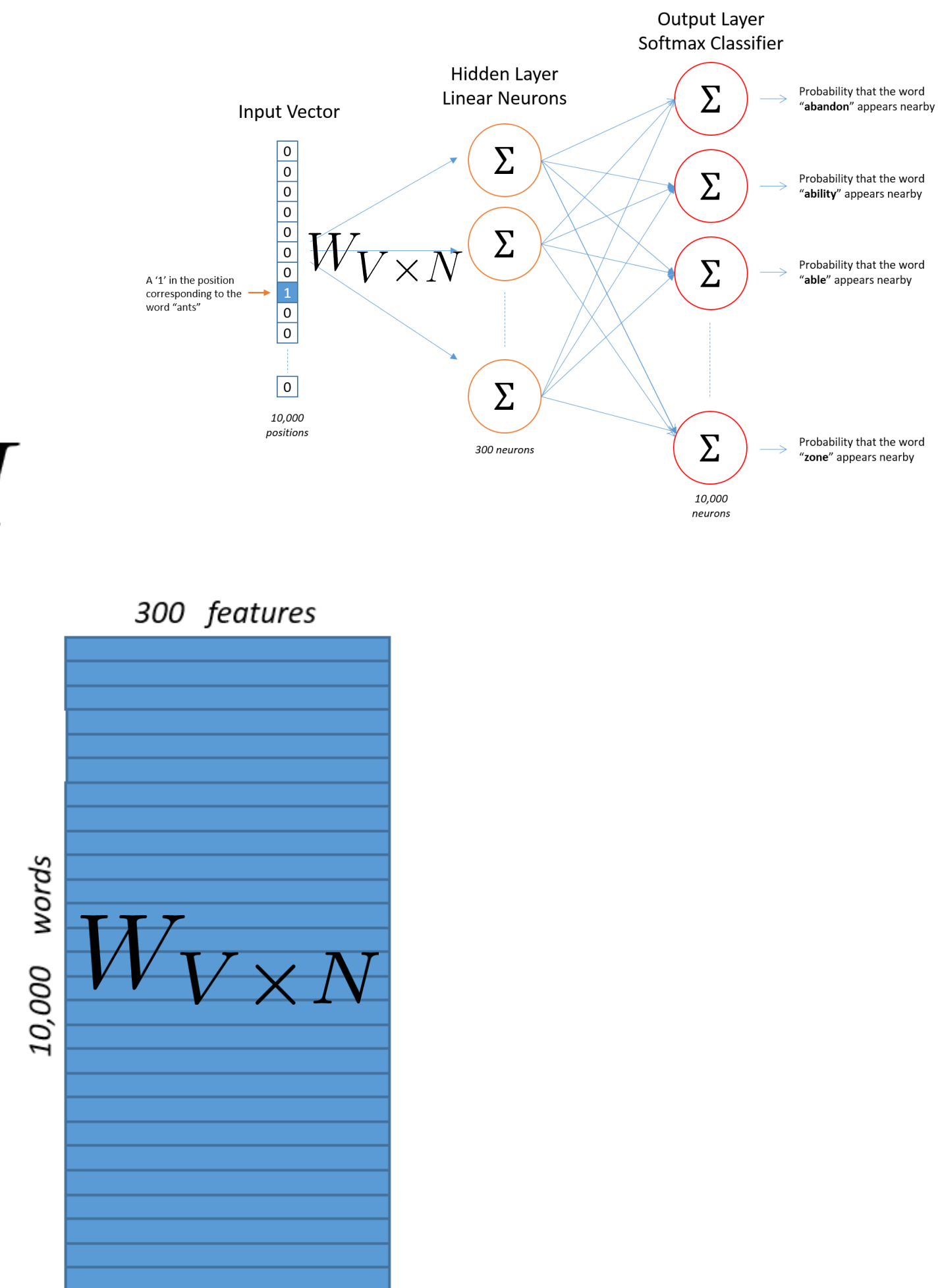


# Weight Matrix Relation

- Hidden layer weight matrix = word vector lookup

$$h = x^T W = W_{(k, \cdot)} := v_{w_I}$$

$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 \end{bmatrix} \times \begin{bmatrix} 17 & 24 & 1 \\ 23 & 5 & 7 \\ 4 & 6 & 13 \\ 10 & 12 & 19 \\ 11 & 18 & 25 \end{bmatrix} = \begin{bmatrix} 10 & 12 & 19 \end{bmatrix}$$



Each vocabulary entry has two vectors: as a **target** word and as a **context** word

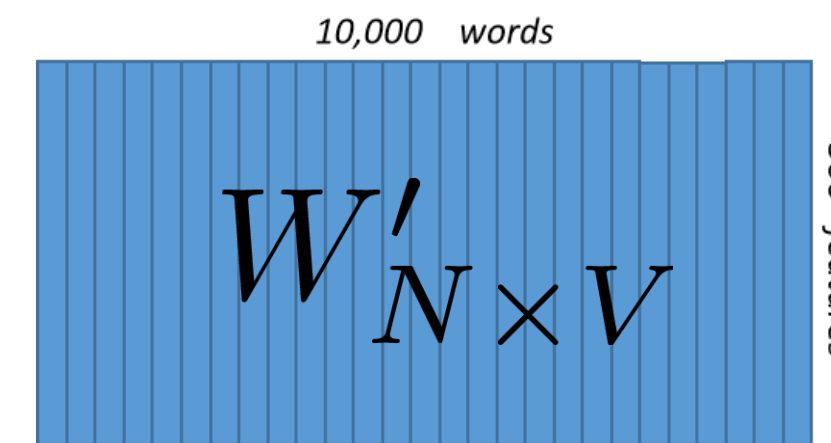




# Weight Matrix Relation

- Output layer weight matrix = weighted sum as final score

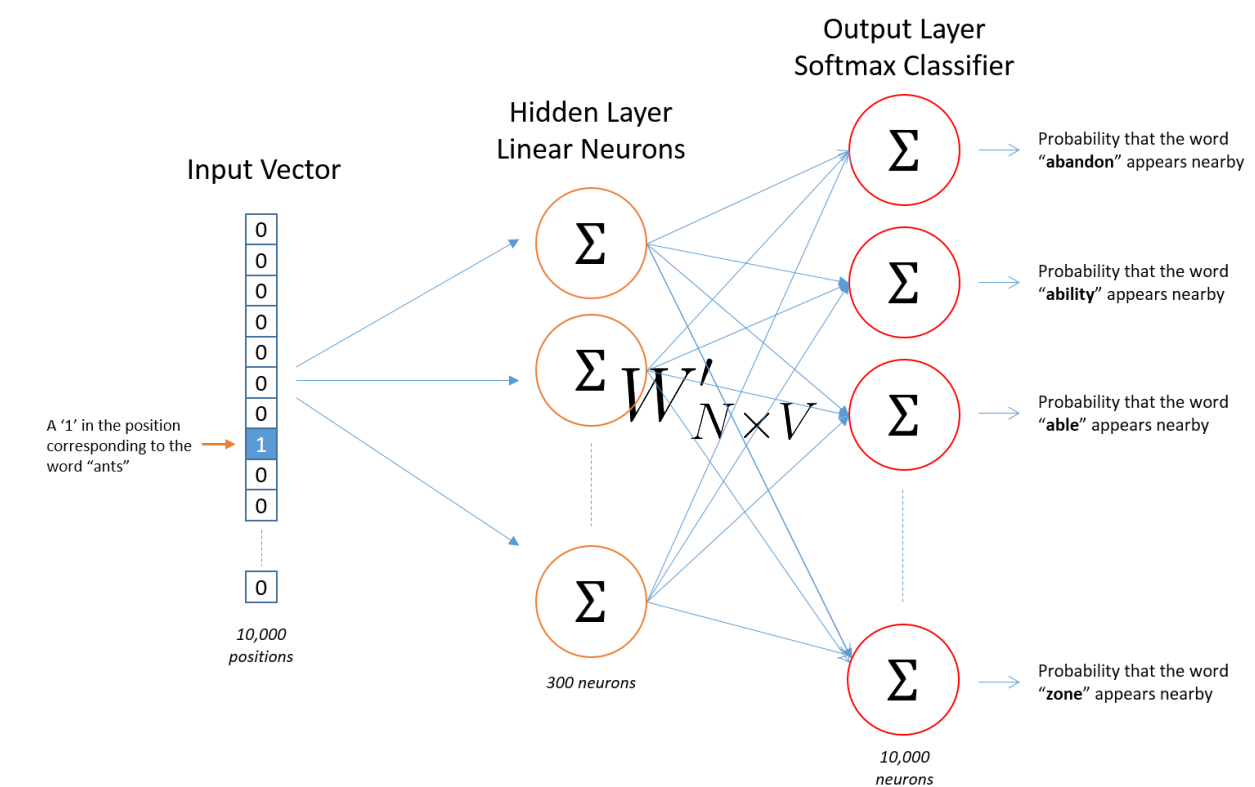
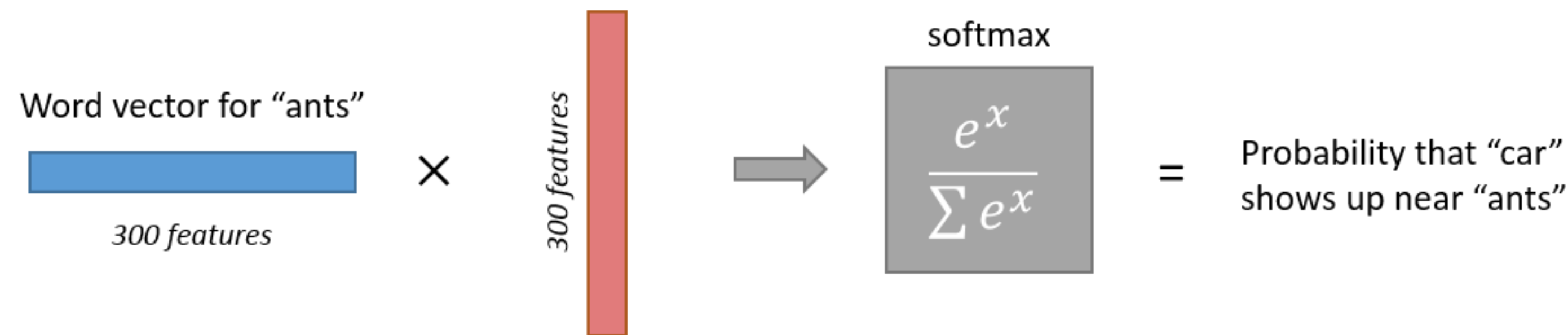
$$s_j = hv'_j w_j$$



$$p(w_j = w_{O,c} \mid w_I) = y_{jc} = \frac{\exp(s_{jc})}{\sum_{j'=1}^V \exp(s_{j'})} \quad \text{softmax}$$

within the context window

Output weights for "car"



Each vocabulary entry has two vectors: as a **target** word and as a **context** word



# Word2Vec Skip-Gram Illustration

