

Word Embeddings

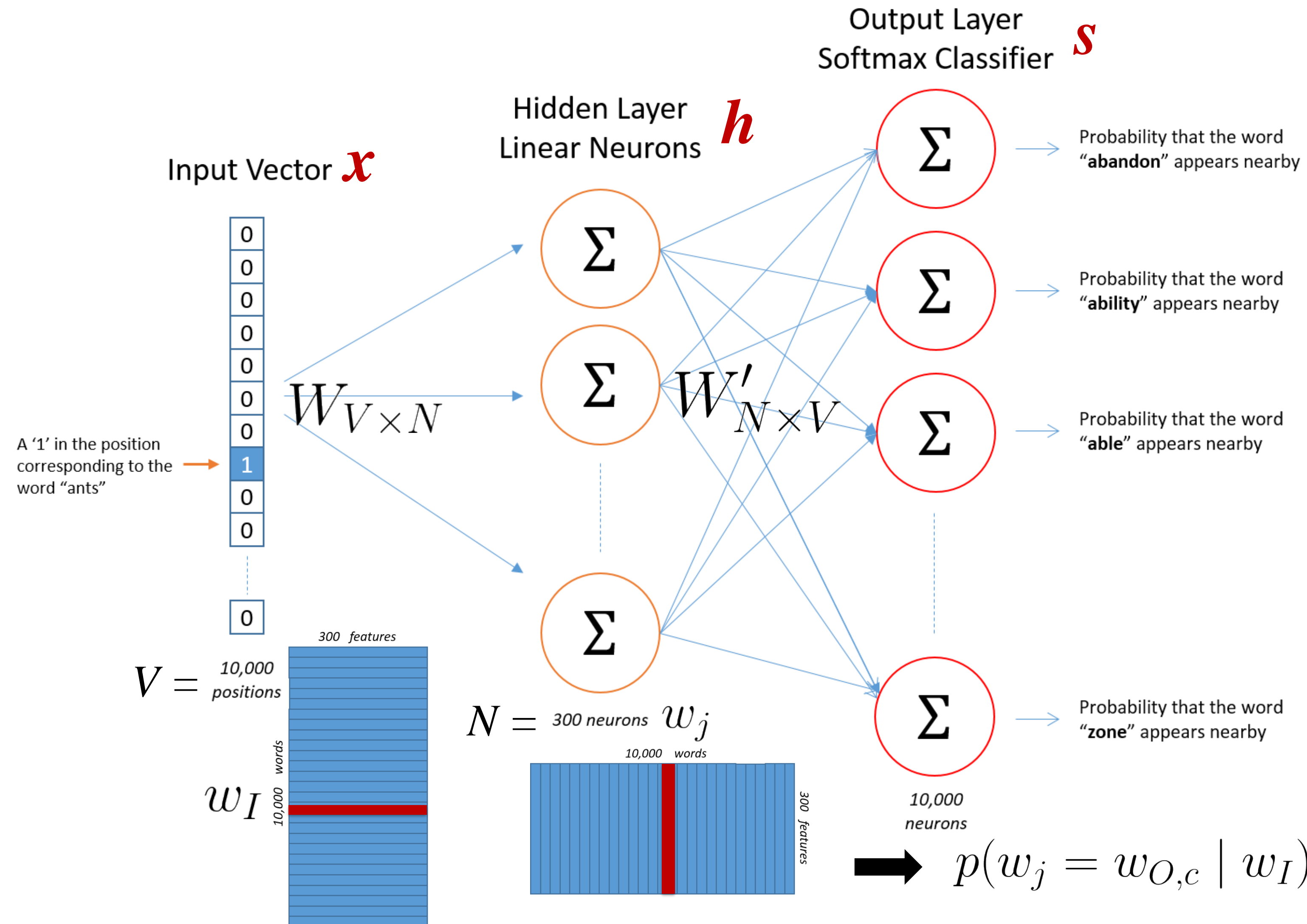
Word2Vec Training



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 Illustration



Loss Function

- Given a target word (w_I)

$$\begin{aligned} C(\theta) &= -\log p(w_{O,1}, w_{O,2}, \dots, w_{O,C} \mid w_I) \\ &= -\log \prod_{c=1}^C \frac{\exp(s_{j_c})}{\sum_{j'=1}^V \exp(s_{j'})} \\ &= -\sum_{c=1}^C s_{j_c} + C \log \sum_{j'=1}^V \exp(s_{j'}) \end{aligned}$$



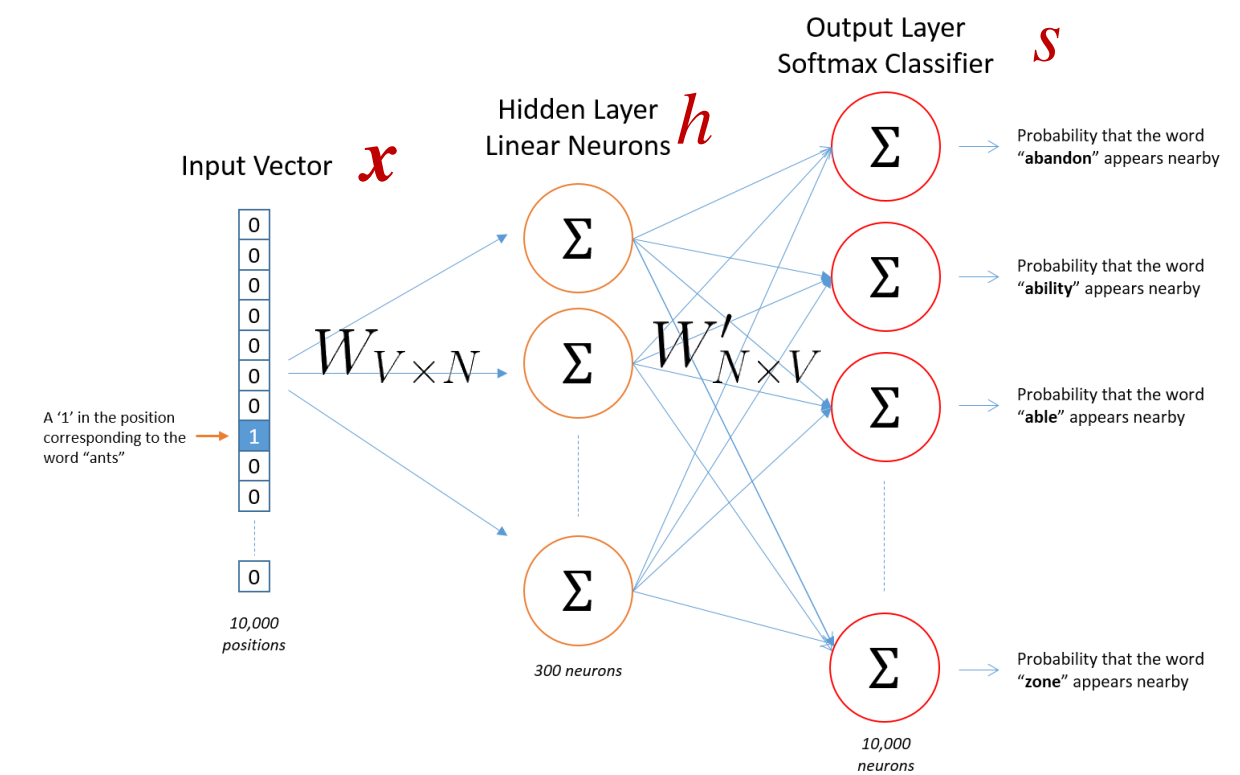
SGD Update for W'

- Given a target word (w_I)

$$\frac{\partial C(\theta)}{\partial w'_{ij}} = \sum_{c=1}^C \frac{\partial C(\theta)}{\partial s_{jc}} \frac{\partial s_{jc}}{\partial w'_{ij}} = \sum_{c=1}^C (y_{jc} - t_{jc}) \cdot h_i$$

$$\frac{\partial C(\theta)}{\partial s_{jc}} = y_{jc} - \underbrace{t_{jc}}_{=1, \text{ when } w_{jc} \text{ is within the context window}} := \underbrace{e_{jc}}_{=0, \text{ otherwise}} \text{ error term}$$

$$s_j = v'_{w_j}^T \cdot h$$



$$w'_{ij}^{(t+1)} = w'_{ij}^{(t)} - \eta \cdot \sum_{c=1}^C (y_{jc} - t_{jc}) \cdot h_i$$



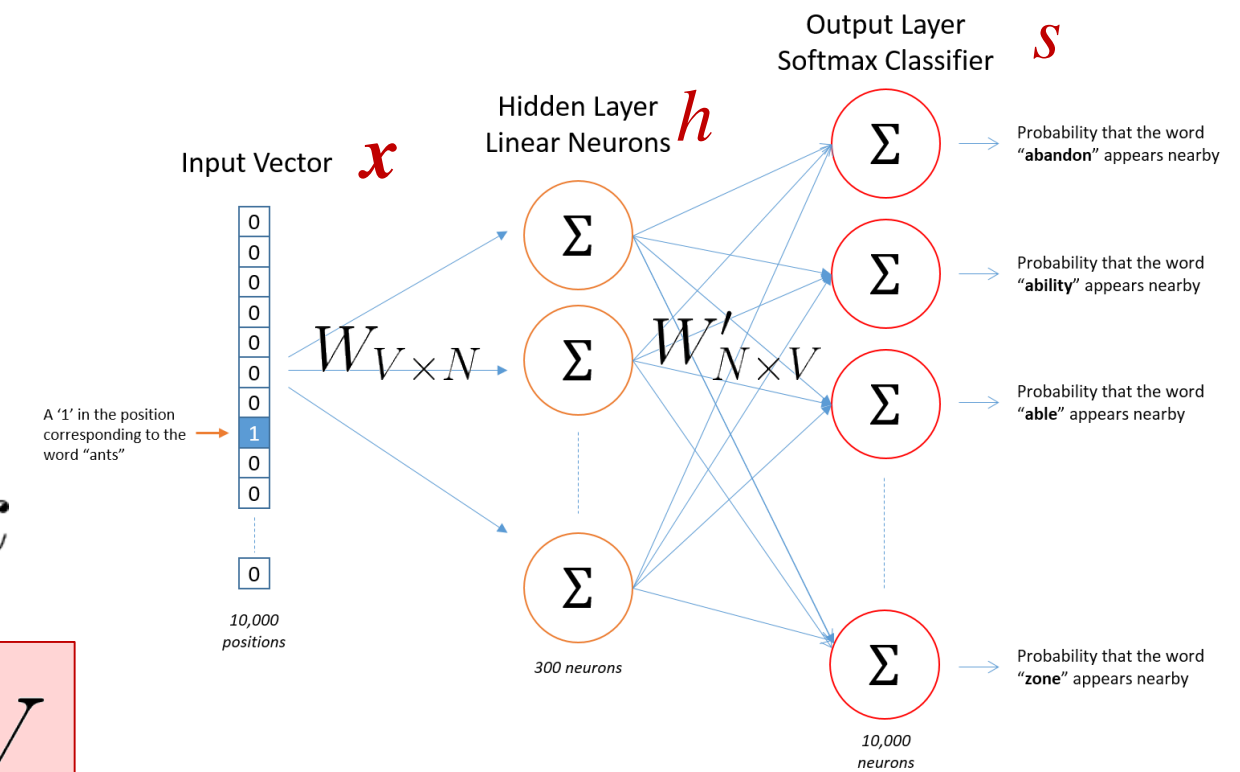
SGD Update for W

$$\frac{\partial C(\theta)}{\partial w_{ki}} = \frac{\partial C(\theta)}{\partial h_i} \frac{\partial h_i}{\partial w_{ki}} = \sum_{j=1}^V \sum_{c=1}^C (y_{jc} - t_{jc}) \cdot w'_{ij} \cdot x_k$$

$$h = x^T W$$

$$\frac{\partial C(\theta)}{\partial h_i} = \sum_{j=1}^V \frac{\partial C(\theta)}{\partial s_j} \frac{\partial s_j}{\partial h_i} = \sum_{j=1}^V \sum_{c=1}^C (y_{jc} - t_{jc}) \cdot w'_{ij}$$

$$s_j = v'_{w_j}{}^T \cdot h$$



$$w_{ij}^{(t+1)} = w_{ij}^{(t)} - \eta \cdot \sum_{j=1}^V \sum_{c=1}^C (y_{jc} - t_{jc}) \cdot w'_{ij} \cdot x_j$$



SGD Update

$$w'_{ij}{}^{(t+1)} = w'_{ij}{}^{(t)} - \eta \cdot \sum_{c=1}^C (y_{jc} - t_{jc}) \cdot h_i$$

$$EI_j = \sum_{c=1}^C (y_{jc} - t_{jc})$$

$$v'_{w_j}{}^{(t+1)} = v'_{w_j}{}^{(t)} - \eta \cdot EI_j \cdot h$$

$$w_{ij}{}^{(t+1)} = w_{ij}{}^{(t)} - \eta \cdot \sum_{j=1}^V \sum_{c=1}^C (y_{jc} - t_{jc}) \cdot w'_{ij} \cdot x_j$$

$$EH_i = \sum_{j=1}^V EI_j \cdot w'_{ij} \cdot x_j$$

$$v_{w_I}{}^{(t+1)} = v_{w_I}{}^{(t)} - \eta \cdot EH^T$$

large vocabularies or large training corpora → expensive computations

limit the number of output vectors that must be updated per training instance
→ hierarchical softmax, sampling

