# Word2Vec

— Assign vectors to words.

— One approach

     — Co-occurence matrix.

     — SVD ( Singular value decomposition )

     — Word2vec — CBOW ( continuous bag of words model ).

           — Skipgram.

— Word2vec.

$$W_p \cdots W_m \; W_{m+1} \; W_{m+2} \cdots \cdots W_{2m+1}$$

— Suppose we can assign vectors to words.

— Using those assignment we will define

(CBOW)

$$p(\,W_{m+1} \mid W_1 \cdots W_m, W_{m+2} \cdots W_{2m+1}\,)$$

center word        context

(Skipgram)

Opposite.

$$= \quad \frac{e^{u_c^T v_o}}{\displaystyle\sum_I} \qquad e^{v'^T_{w_{m+1}} v_{w_c}}$$
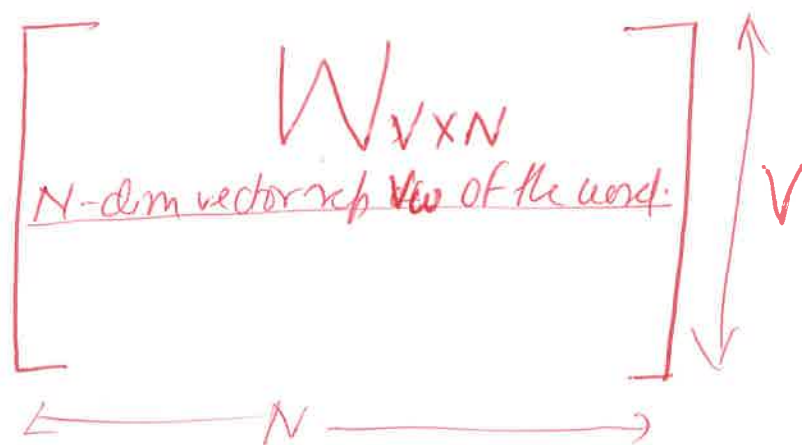
$V'$ — outside rep

$V$ — inside rep.

Vocabulary size — V

Hidden layer size — N. (we want each word to be represented in a vector of size N.)

*Input is a one-hot encoded vector.*

Only one of $\{x_1 \cdots x_V\}$ is 1 rest is 0.

Weights between input layer & hidden layer is a V×N matrix. W.

$$\left[\begin{array}{c} W_{V \times N} \\ \underline{N\text{-dim vector rep } V_W \text{ of the word.}} \end{array}\right] \updownarrow V$$

$$\longleftarrow N \longrightarrow$$

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_k \\ x_V \end{bmatrix} \qquad W^T \cdot X = \begin{bmatrix} h_1 \\ \vdots \\ h_N \end{bmatrix} \qquad h_{N \times 1}$$

When $x_k = 1$ & $x_{k'} = 0$ for $k' \neq k$.

$$h = W^T x \quad \text{~~~~} \text{ is the } \underline{k\text{th row of } W}$$
in the as column

Use Prop. $x_k$

$V_{W_I}$ — vector rep of input word $w_I$

So to write it as a column we hav $V_{W_I}^T$.

$$\boxed{h = V_{W_I}^T}$$

$x_1$

$h_1$

$y_1$

$x_k$ $W_{V\times N}$ $h_i$ $W'_{N\times V}$ $y_j$

$x_V = \{w_{ki}\}$ $h_N = \{w'_{ij}\}$ $y_V$

$V = $ Size of vocabulary.

$$u_j = v'^T_{w_j} h$$

when $v'_{w_j}$ is the $j$th column

of $W'$

(Becomes $j$th row in $W'^T$)

$=$ denotes the score of the $j$th word in the vocabulary.

$$p(w_j \mid w_I) = y_j = \frac{e^{u_j}}{\sum^V_{j'=1} e^{(u_{j'})}}$$

$$= \frac{e^{v'^T_{w_j} v^T_{w_I}}}{\sum^V_{j'=1} e^{v'_{w_{j'}}^T v^T_{w_I}}}$$

$v_w$ — input vector

$v_{w'}$ — output vector

(for one training sample)

Training Objective is to maximize the last equation.

the conditional probability of observing the actual word (W_O) given the input context word $w_I$

denote its index as $j^*$

$$\log p(w_O/w_I) = \log y_{j^*}$$

$$= \log \frac{e^{u_j^*}}{\sum_{j'=1}^{V} e^{u_{j'}}} = u_j^* - \log \sum_{j'=1}^{V} e^{u_{j'}}$$

$$= -E$$

→ So minimize $E = -\log p(w_O/w_I)$.

$$E = -u_j^* + \log \sum_{j'=1}^{V} e^{u_{j'}}$$

$$\frac{\partial E}{\partial u_j} = -1 + \frac{\partial \left[ e^{u_{j'}} \right] / \partial u_j}{\sum_{j'=1}^{V} e^{u_{j'}}} \qquad \text{if } j^* = j$$

$$= 0 \qquad \neq \qquad )) \qquad j^* \neq j$$

$$= -t_j + \frac{e^{u_j}}{\sum_{j'=1}^{V} e^{u_{j'}}} = -t_j + y_j$$

$$= y_j - t_j \quad = e_j$$

Where $t_j = 1$ if $j = j^*$
$= 0$ o'wise

$$\frac{\partial E}{\partial W'_{ij}} = \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial W_{ij}}$$

$$= e_j \cdot \frac{\partial u_j}{\partial W_{ij}}$$

---

$$W_{V \times N} = \{w_{ki}\} = \begin{bmatrix} w_{11} & \cdots & w_{1i} & w_{1N} \\ w_{k1} & \cdots & w_{ki} & w_{kN} \\ & \vdots & & \\ w_{v1} & \cdots & w_{vi} & w_{vN} \end{bmatrix} \quad W' = \begin{bmatrix} w'_{11} & \cdots & w'_{1j} & w'_{1V} \\ w'_{i1} & \cdots & w'_{ij} & w'_{iV} \\ & \vdots & & \\ w'_{N1} & \cdots & w'_{Nj} & w'_{NV} \end{bmatrix}$$

$$u_j = V'_{w_j} h = w'_{1j} h_1 + w'_{2j} h_2 + \cdots w'_{Nj} h_N$$

$$= \sum_{i=1}^{N} w'_{ij} h_i$$

So $\dfrac{\partial u_j}{\partial w'_{ij}} = h_i$ $\quad\Bigg|\quad$ $\dfrac{\partial u_j}{\partial h_i} = w'_{ij}$

So $\dfrac{\partial E}{\partial w'_{ij}} = e_j \cdot h_i$

---

Using stochastic gradient $\quad w'_{ij}{}^{(new)} = w'_{ij}{}^{(old)} - \eta e_j h_i$

or

$$V'_{w_j}{}^{(new)} = V'_{w_j}{}^{(old)} - \eta \cdot e_j \cdot h \qquad \text{for } j=1 \cdots V$$

$$\boxed{h = W^T \cdot X \quad \text{So } h_i = w_{1i} x_1 + w_{2i} x_2 + \cdots w_{vi} x_v \\ h_i = \sum_{k=1}^{V} w_{ki} x_k \quad\Bigg|\quad \frac{\partial h_i}{\partial w_{ki}} = x_k}$$

*4

# Update equation for input → hidden weights ⑤

$$\frac{\partial E}{\partial w_{ki}} = \frac{\partial E}{\partial h_i} \cdot \boxed{\frac{\partial h_i}{\partial w_{ki}}}$$

$x_k$

$$E = -u_{j^*} + \log\left( e^{u_1} + e^{u_2} + \cdots e^{u_v} \right)$$

page 3.

$$E = -u_{j^*} + \log\left( \right)$$

each $u_j$ has $h_i$ in them.

$$\frac{\partial E}{\partial h_i} = \frac{\partial E}{\partial u_1}\frac{\partial u_1}{\partial h_i} + \frac{\partial E}{\partial u_2}\frac{\partial u_2}{\partial h_i} + \cdots$$

$$= \sum_{j=1}^{v} \frac{\partial E}{\partial u_j}\frac{\partial u_j}{\partial h_i} = \sum_{j=1}^{v} e_j \cdot w_{ij} = EH_i$$

page 3    page 4.

So $\quad \dfrac{\partial E}{\partial w_{ui}} = EH_i \cdot x_k$

Gradient step is detailed hand on this.

Multi-word context: Defn if $h = V_{w_I}^T$ is changed to

$$h = \frac{1}{C}\left( V_{w_1} + V_{u_2} + \cdots V_{w_C} \right)^T$$

when $C$ is the # words in the context.

Everything else stays same.