

word2Vec

刘绪成

YNU

December 1, 2020

- 1 Continuous Bag-of-Words Model
 - One-word context
 - Multi-word context
- 2 Skip-Gram Model
- 3 Optimizing Computational Efficiency
 - Hierarchical Softmax
 - Negative Sampling

Model of CBOW

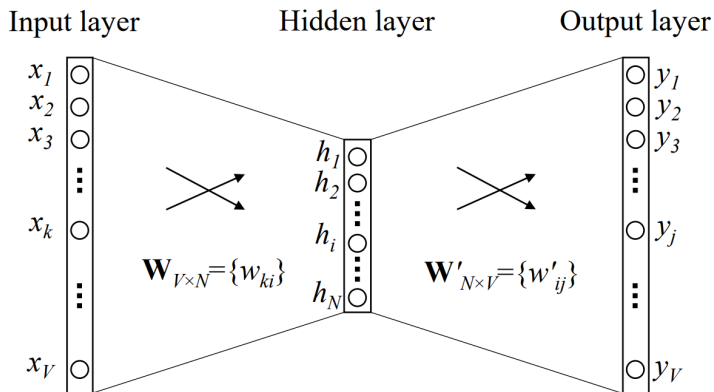


Figure 1: A simple CBOW model with only one word in the context

Derivation

随机初始化 W , 由于输入向量 x 是 One-hot Vector, 隐藏层 h 就是 W 的某一行。

$$\mathbf{h} = \mathbf{W}^T \mathbf{x} = \mathbf{W}_{(k, \cdot)}^T := \mathbf{v}_{w_l}^T \quad (1)$$

从隐藏层到输入层:

$$u_j = \mathbf{v}_{w_j}'^T \mathbf{h} \quad (2)$$

使用 softmax 求输入层第 j 个神经元的值

$$p(w_j | w_l) = y_j = \frac{\exp(u_j)}{\sum_{j'=1}^V \exp(u_{j'})} \quad (3)$$

将公式 (1) 和 (2) 代入

$$p(w_j | w_l) = \frac{\exp(\mathbf{v}_{w_j}'^T \mathbf{v}_{w_l})}{\sum_{j'=1}^V \exp(\mathbf{v}_{w_{j'}}'^T \mathbf{v}_{w_l})} \quad (4)$$

反向传播从输出层到隐藏层：我们的训练目标是最大化公式 (4)，即在实际输出 w_O 的条件概率

$$\begin{aligned}\max p(w_O \mid w_I) &= \max y_{j^*} \\ &= \max \log y_{j^*} \\ &= u_{j^*} - \log \sum_{j'=1}^V \exp(u_{j'}) := -E\end{aligned}\tag{5}$$

$$\frac{\partial E}{\partial u_j} = y_j - t_j := e_j\tag{6}$$

$t_j = 1 (j = j^*)$ 其他情况下 $t_j = 0$

E 对 W^{prime} 第 j 列, 第 i 行的偏微分

$$\frac{\partial E}{\partial w'_{ij}} = \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial w'_{ij}} = e_j \cdot h_i \quad (7)$$

使用随机梯度下降更新

$$w'_{ij}^{(new)} = w'_{ij}^{(old)} - \eta \cdot e_j \cdot h_i \quad (8)$$

从隐藏层-> 输入层:

$$\frac{\partial E}{\partial h_i} = \sum_{j=1}^V \frac{\partial E}{\partial u_j} \cdot \frac{\partial u_j}{\partial h_i} = \sum_{j=1}^V e_j \cdot w_{ij} := \text{EH}_i \quad (9)$$

$$h_i = \sum_{k=1}^V x_k \cdot w_{ki} \quad (10)$$

$$\frac{\partial E}{\partial w_{ki}} = \frac{\partial E}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_{ki}} = \text{EH}_i \cdot x_k \quad (11)$$

$$\frac{\partial E}{\partial \mathbf{W}} = \mathbf{x} \otimes \text{EH} = \mathbf{x} \text{EH}^T \quad (12)$$

输入向量的更新:

$$\mathbf{v}_{w_l}^{(\text{new})} = \mathbf{v}_{w_l}^{(\text{old})} - \eta \text{EH}^T \quad (13)$$

Model of CBOW

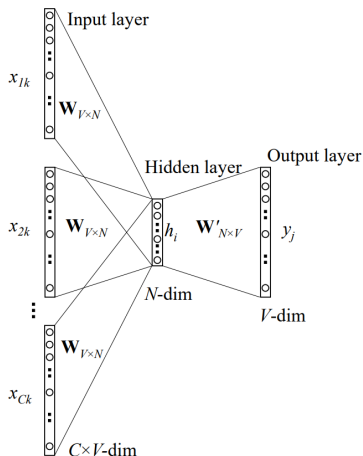


Figure 2: Continuous bag-of-words model

用输入层→隐藏层的 W^T 与平均向量 x 的积作为输入

$$\begin{aligned} h &= \frac{1}{C} W^T (x_1 + x_2 + \dots + x_C) \\ &= \frac{1}{C} (v_{\omega_1} + v_{\omega_2} + \dots + v_{\omega_C}) \end{aligned} \quad (14)$$

$$\begin{aligned} E &= -\log p(\omega_O | \omega_{I,1}, \dots, \omega_{I,C}) \\ &= -u_{j_*} + \log \sum_{j'=1}^V \exp(u_{j'}) \\ &= -V_{\omega_O}^T \cdot h + \log \sum_{j'=1}^V \exp(\omega_{j'}^T \cdot h) \end{aligned} \quad (15)$$

Figure

$$v_{\omega_j}^{(new)} = v_{\omega_j}^{(old)} - \eta \cdot e_j \cdot h \quad \text{for } j = 1, 2, \dots, V \quad (16)$$

$$v_{\omega_{l,c}}^{(new)} = v_{\omega_{l,c}}^{(old)} - \frac{1}{C} \cdot \eta \cdot EH^T \quad \text{for } j = 1, 2, \dots, C \quad (17)$$

Model of Skip-gram

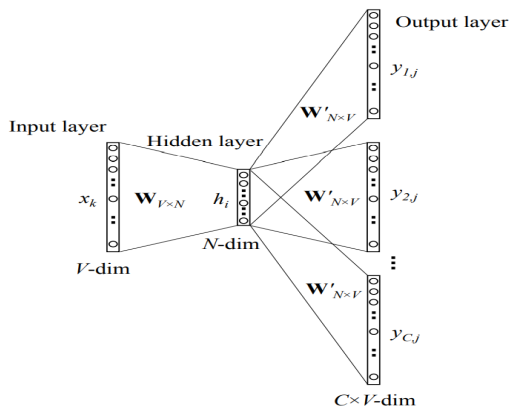


Figure 3: The skip-gram model.

$$\begin{aligned}
 E &= -\log p(w_{O,1}, w_{O,2}, \dots, w_{O,C} \mid w_I) \\
 &= -\log \prod_{c=1}^C \frac{\exp(u_{c,j_c^*})}{\sum_{j'=1}^V \exp(u_{j'})}
 \end{aligned} \tag{18}$$

$$\begin{aligned}
 &= -\sum_{c=1}^C u_{j_c^*} + C \cdot \log \sum_{j'=1}^V \exp(u_{j'}) \\
 \frac{\partial E}{\partial u_{c,j}} &= y_{c,j} - t_{c,j} := e_{c,j}
 \end{aligned} \tag{19}$$

$$El_j = \sum_{c=1}^C e_{c,j} \tag{20}$$

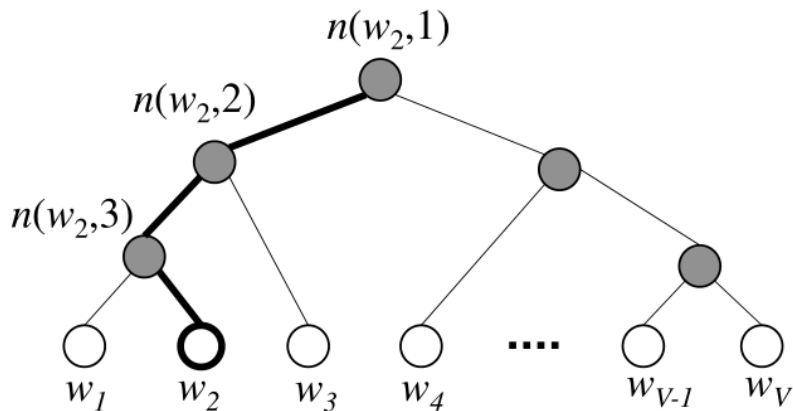
$$\frac{\partial E}{\partial w'_{i,j}} = \sum_{c=1}^C \frac{\partial E}{\partial u_{c,j}} \cdot \frac{\partial u_{c,j}}{w'_{i,j}} = El_j \cdot h_i \tag{21}$$

$$v_{w_j}^{(new)} = v_{w_j}^{(old)} - \eta \cdot El_j \cdot h \quad \text{for } j = 1, 2, 3, \dots, V. \quad (22)$$

$$V_{w_j}^{(new)} = V_{w_j}^{(old)} - \eta \cdot El^T \quad (23)$$

$$EH_i = \sum_{j=1}^V El_j \cdot w'_{ij} \quad (24)$$

Hierarchical Softmax



$$p(w = w_O) = \prod_{j=1}^{L(w)-1} \sigma(\llbracket n(w, j+1) = ch(n(w, j)) \rrbracket) \cdot v_{n(w, j)}^T \cdot h \quad (25)$$

$$\llbracket x \rrbracket = \begin{cases} 1 & \text{if } x \text{ is true} \\ -1 & \text{otherwise} \end{cases} \quad (26)$$

$$p(n, \text{left}) = \sigma(v_n^T \cdot h) \quad (27)$$

$$p(n, \text{right}) = 1 - \sigma(v_n^T \cdot h) = \sigma(-v_n^T \cdot h) \quad (28)$$

$$\begin{aligned} p(w_2 = w_O) &= p(n(w_2, 1), \text{left}) \cdot p(n(w_2, 2), \text{left}) \cdot p(n(w_2, 3), \text{left}) \\ &= \sigma(v_{n(w_2, 1)}^T h) \cdot \sigma(v_{n(w_2, 2)}^T h) \cdot \sigma(v_{n(w_2, 3)}^T h) \end{aligned} \quad (29)$$

$$\sum_{i=1}^V p(w_i = w_O) = 1 \quad (30)$$

$$\llbracket \cdot \rrbracket := \llbracket n(w, j+1) = ch(n(w, j)) \rrbracket \quad (31)$$

$$\mathbf{v}'_j := \mathbf{v}'_{n(w, j)} \quad (32)$$

$$E = -\log p(w = w_O | w_I) = - \sum_{j=1}^{L(w)-1} \log \sigma(\llbracket \cdot \rrbracket \mathbf{v}'_j{}^T \mathbf{h}) \quad (33)$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{v}'_j \mathbf{h}} &= (\sigma(\llbracket \cdot \rrbracket \mathbf{v}'_j{}^T \mathbf{h}) - 1) \llbracket \cdot \rrbracket \\ &= \begin{cases} \sigma(\mathbf{v}'_j{}^T \mathbf{h}) - 1 & (\llbracket \cdot \rrbracket = 1) \\ \sigma(\mathbf{v}'_j{}^T \mathbf{h}) & (\llbracket \cdot \rrbracket = -1) \end{cases} \\ &= \sigma(\mathbf{v}'_j{}^T \mathbf{h}) - t_j \end{aligned} \quad (34)$$

$$\frac{\partial E}{\partial \mathbf{v}'_j} = \frac{\partial E}{\partial \mathbf{v}'_j \mathbf{h}} \cdot \frac{\partial \mathbf{v}'_j \mathbf{h}}{\partial \mathbf{v}'_j} = (\sigma(\mathbf{v}'_j{}^T \mathbf{h}) - t_j) \cdot \mathbf{h} \quad (35)$$

$$\mathbf{v}_j^{(' \text{ new })} = \mathbf{v}_j^{(' \text{ old })} - \eta \left(\sigma \left(\mathbf{v}_j'^T \mathbf{h} \right) - t_j \right) \cdot \mathbf{h} \quad (36)$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{h}} &= \sum_{j=1}^{L(w)-1} \frac{\partial E}{\partial \mathbf{v}_j' \mathbf{h}} \cdot \frac{\partial \mathbf{v}_j' \mathbf{h}}{\partial \mathbf{h}} \\ &= \sum_{j=1}^{L(w)-1} \left(\sigma \left(\mathbf{v}_j'^T \mathbf{h} \right) - t_j \right) \cdot \mathbf{v}_j' := \text{EH} \end{aligned} \quad (37)$$

Negative Sampling

$$E = -\log \sigma(\mathbf{v}_{w_O}^T \mathbf{h}) - \sum_{w_j \in \mathcal{W}_{neg}} \log(-\mathbf{v}_{w_j}'^T \mathbf{h}) \quad (38)$$

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{v}_{w_j}'^T \mathbf{h}} &= \begin{cases} \sigma(\mathbf{v}_{w_j}'^T \mathbf{h}) - 1 & \text{if } w_j = w_O \\ \sigma(\mathbf{v}_{w_j}'^T \mathbf{h}) & \text{if } w_j \in \mathcal{W}_{neg} \end{cases} \\ &= \sigma(\mathbf{v}_{w_j}'^T \mathbf{h}) - t_j \end{aligned} \quad (39)$$

$$\frac{\partial E}{\partial \mathbf{v}_{w_j}'} = \frac{\partial E}{\partial \mathbf{v}_{w_j}'^T \mathbf{h}} \cdot \frac{\partial \mathbf{v}_{w_j}'^T \mathbf{h}}{\partial \mathbf{v}_{w_j}'} = (\sigma(\mathbf{v}_{w_j}'^T \mathbf{h}) - t_j) \mathbf{h} \quad (40)$$

$$\mathbf{v}_{w_j}'^{(new)} = \mathbf{v}_{w_j}'^{(old)} - \eta (\sigma(\mathbf{v}_{w_j}'^T \mathbf{h}) - t_j) \mathbf{h} \quad (41)$$

$$\begin{aligned}\frac{\partial E}{\partial \mathbf{h}} &= \sum_{w_j \in w_O \cup \mathcal{W}_{\text{neg}}} \frac{\partial E}{\partial \mathbf{v}'_{w_j}{}^T \mathbf{h}} \cdot \frac{\partial \mathbf{v}'_{w_j}{}^T \mathbf{h}}{\partial \mathbf{h}} \\ &= \sum_{w_j \in w_O \cup \mathcal{W}_{\text{neg}}} \left(\sigma \left(\mathbf{v}'_{w_j} \mathbf{h} \right) - t_j \right) \mathbf{v}'_{w_j} \\ &:= \text{EH}\end{aligned}\tag{42}$$

References



Xin Rong (2016)

word2vec Parameter Learning Explained

arXiv e-prints arXiv:1411.2738.



Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a)

Efficient estimation of word representations in vector space

arXiv preprint arXiv:1301.3781.



Le, Quoc and Mikolov, Tomas (2014)

Distributed Representations of Sentences and Documents

Proceedings of the 31st International Conference on International Conference on Machine Learning II–1188–II–1196.

Example

python	科比	云南
perl,0.8808220624923706	杜兰特,0.73625	贵州,0.72855
lua,0.856161892414093	布莱恩特,0.72571	云南省,0.70297
smalltalk,0.8457338809967041	奥尼尔,0.72011	昆明,0.69951
javascript,0.8404377102851868	罗德曼,0.71371	西双版纳,0.67357
脚本语言,0.8397466540336609	欧尼尔,0.70804	腾冲,0.66212
僵尸	特朗普	开心
丧尸,0.6313216686248779	川普,0.78675	放暑假,0.53373
狼人,0.5508135557174683	奥巴马,0.77438	华之里,0.53014
吸血,0.5407103300094604	尼克松,0.71462	欢乐,0.51938
幽灵,0.5303795337677002	小布什,0.70919	欢笑,0.51552
奇鸭,0.4994601011276245	希拉里,0.70734	天天,0.51292

Example

伦敦美国英国	king women man
伦敦之于英国，如美国之于	king 之于 man，如 women 之于
纽约,0.70653	narai,0.48565
纽约市,0.63568	queen,0.47201
芝加哥,0.57615	kings,0.46423
美国纽约,0.57061	scotland,0.46213
白天太阳晚上	武松红楼梦水浒传
白天之于晚上，如太阳之于	武松之于水浒传，如红楼梦之于
月亮,0.48442	林黛玉,0.72075
不西沉,0.467524	贾宝玉,0.70228
白昼,0.46403	王熙凤,0.69696
永昼,0.46176	大观园,0.68577

The End