

Lecture 2: Word2Vec

2018年7月4日 星期三 上午9:16

1. Word meaning

Signifier -> Signified idea = denotation

NLTK WordNet

Weakness of Discrete representation:

- Synonyms missing nuance
- Missing new words
- Require human labor
- one-hot coding(Symbolic)
 - Very large dimension
 - No relationship between words(dot product is 0)

? Context words是否真的能够描述center words的语义？普通人说话的上下文一般不会再次解释形容词的含义。
是否可以从词典的释义中提取一定数目的词作为context words?

Distributional similarity

Get a lot of neighbor words to represent its meaning

Distributed Representation

Build a dense vector for each word

2. Word2Vec (Mikolov 2013)

Simple, fast to train, Scalable

E.g. Cat+kitty+dog ~ puppy

通过大量文本以无监督方式学习词向量，在嵌入空间中表征语义信息。

? 多义词的向量仅有一个，语义是否会造模糊

• Two Algorithms

Fake Task, 真正需要的是获得模型参数 (Word Vectors)

- 1.SG (Skip-grams)
 - 给定center word 预测 context
 - $$L = \sum_{w \in C} \log P(\text{content}(w) | w)$$
- 2.CBOW (Continuous Bag of Words)
 - 给定context 预测 center word
 - $$L = \sum_{w \in C} \log P(w | \text{content}(w))$$

• Two Training Methods

Motivation: softmax over all data need $O(N)$ time -> Impractical

- 1.Hierarchical Softmax

• Huffman encoding:	reduce parameters number
• Product along the path:	Reduce compute time to $O(\log(n))$

- 2.Negative Sampling

□ Sample some negative cases rather than whole data	Reduce softmax summation number
---	---------------------------------

3. Softmax :

- Standard map from R^V into probability distribution.
- Using center word C to obtain probability of word O

$$\circ P_i = \frac{e^{u_i}}{\sum_j e^{u_j}}$$

$$\frac{\partial P_i}{\partial u_j} = \begin{cases} P_i(1-P_i) & i=j \\ -P_i P_j & i \neq j \end{cases}$$

In Word2Vec :

$$P(o|c) = \frac{\exp(u_o^\top v_c)}{\sum_{w=1}^V \exp(u_w^\top v_c)}$$

$\langle u, v \rangle$: loose measure of similarity

4. Skipgram

- A bag-of-words model:
 - not care about words' position
- Target :
 - Maximize the product of probability of neighbors
 - Obtain intermediate parameters(word vectors)

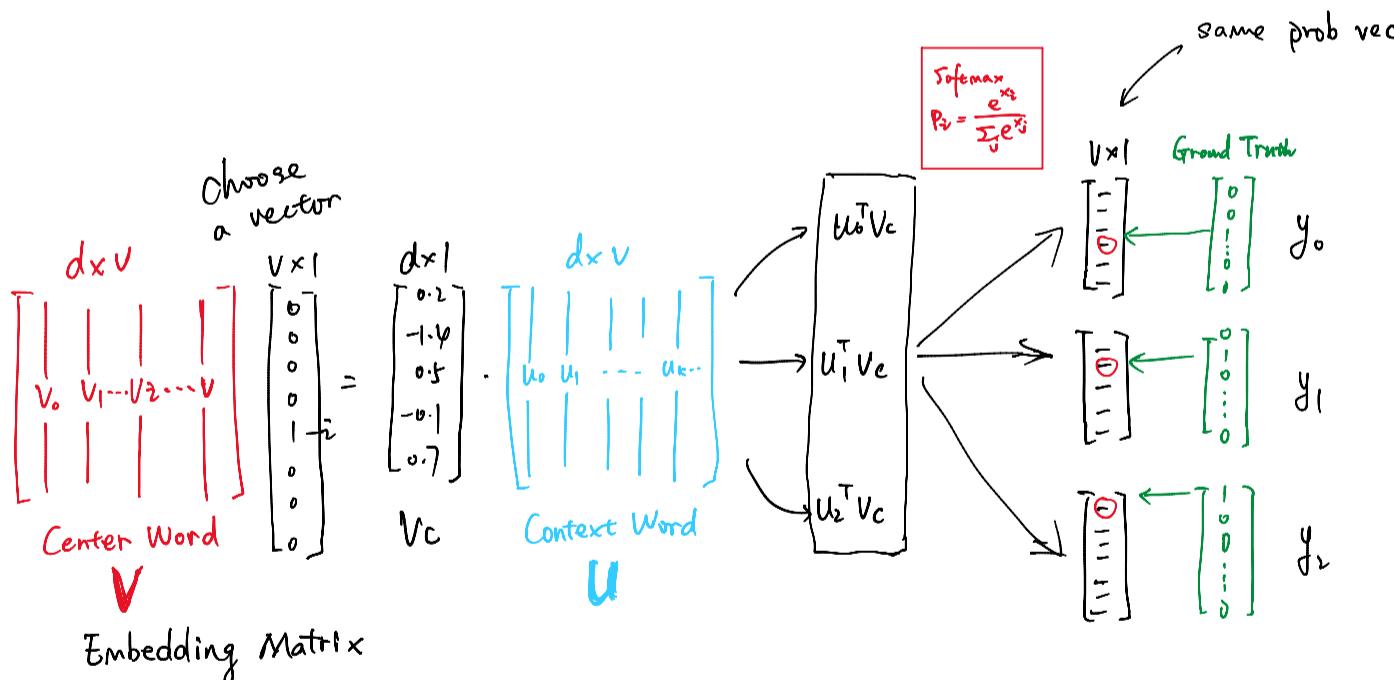
Object function T : text length , m : window size.

$$J(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m} P(w_{t+j}|w_t; \theta)$$

Negative Log Likelihood

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m} \log P(w_{t+j}|w_t; \theta)$$

? 使用对数的原因是求导方便，那为何要加负号？



Parameter:

$$\Theta = [V | U] = [v_0 v_1 \dots v_V | u_0 u_1 \dots u_V]$$

5. Training with Negative Sampling:

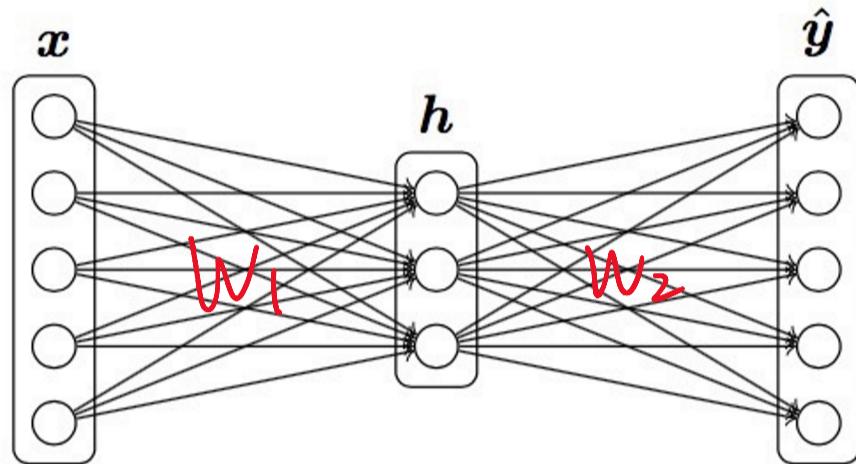
Actually in a 3-layer Neural Network

$$\text{minimize: } J(\theta) = -\frac{1}{T} \sum_t \sum_j \log(p(w_{t+j}|w_t))$$

Global?

$$\frac{\partial}{\partial v_c} \log \frac{e^{u_c^\top v_c}}{\sum_w e^{u_w^\top v_c}} = \frac{\partial}{\partial v_c} \left(u_c^\top v_c - \log \sum_w \exp(u_w^\top v_c) \right)$$

$$\begin{aligned}
 \frac{\partial}{\partial v_c} \log \frac{e^{u^T v_c}}{\sum_w e^{u^T v_c}} &= \frac{\partial}{\partial v_c} \left(u^T v_c - \log \sum_w e^{u^T v_c} \right) \\
 &= u_0 - \frac{\sum_x e^{u_x^T v_c} \cdot u_x}{\sum_w e^{u_w^T v_c}} \\
 &= u_0 - \sum_x \frac{e^{u_x^T v_c} \cdot u_x}{\sum_w e^{u_w^T v_c}} \\
 &= u_0 - \sum_x P(x|c) \cdot u_x \quad u_0 = U \cdot \sigma(U^T V)[:, c] \\
 &= \text{Observed} - \text{Expected}.
 \end{aligned}$$



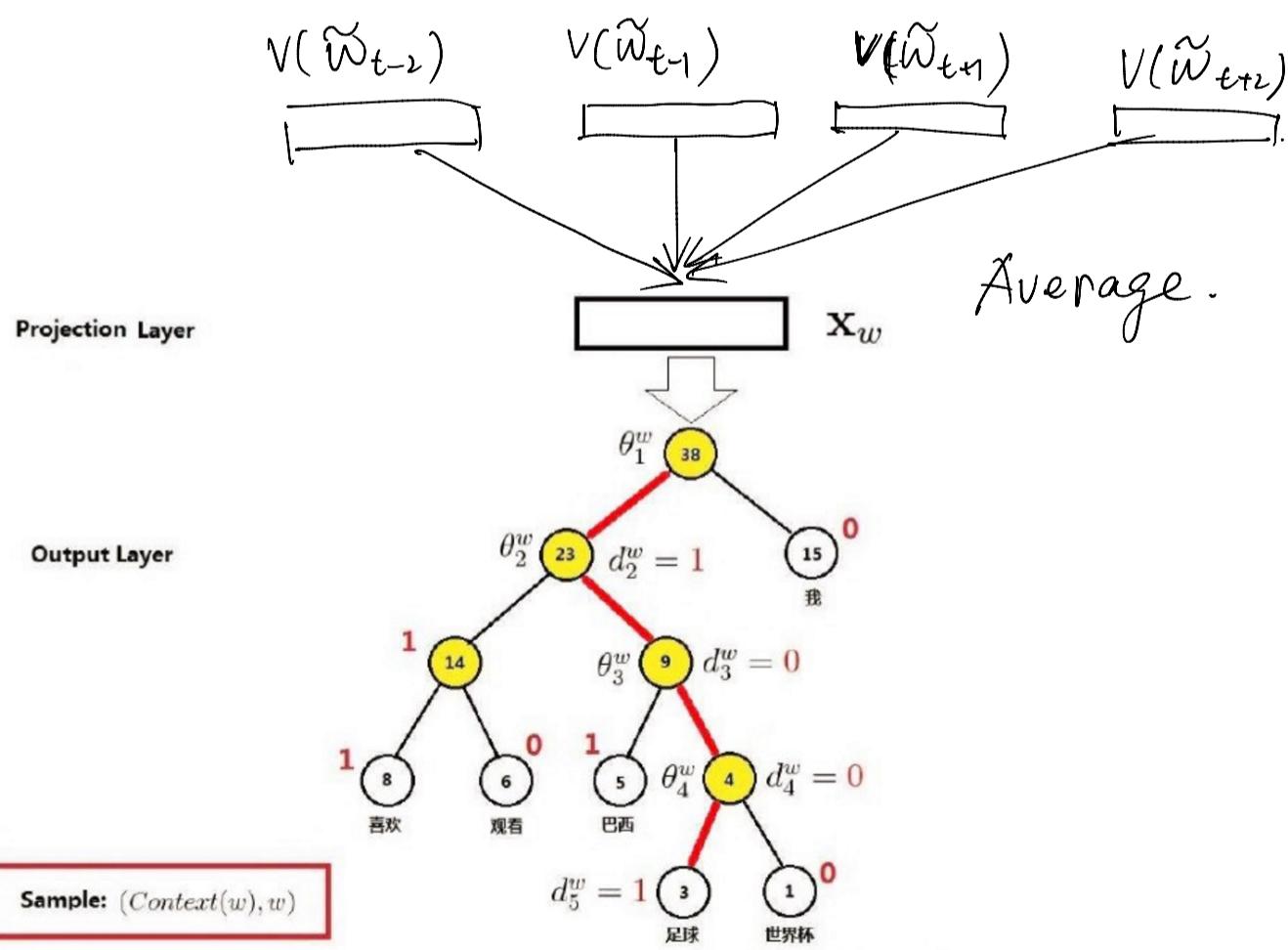
W1: Center Word vectors

W2: Context Word vectors

The final word vectors are W1 concat W2

6.Training with Hierarchical Softmax:

<https://www.cnblogs.com/Determined22/p/5807362.html>



$$\text{P}(d_4^w | x_w, \theta_3^w) = \sigma(X_w^T \theta_3^w)$$

$P(x_4 | x_1, x_2, x_3) = \text{softmax}$

$$L = P(\tilde{x}_w | \text{context}(\tilde{x}_w)) = \prod_{j \sim \text{path}} P(d_j^w | x_w, \theta_{j-1}^w)$$

更新：

$$V(\tilde{w}) := V(w) + \eta \frac{\partial L(w, j)}{\partial x_w}, \quad \tilde{w} \in \text{context}(w)$$

$$\theta_{j-1} = \theta_{j-1} + \eta \frac{\partial L(w, j)}{\partial \theta_{j-1}}$$