



词嵌入表示

N-gram Model Recap

N-元模型回顾

- 句子: x_1, \dots, x_m
- 贝叶斯展开:

$$P(x_1, x_2, \dots, x_m) = P(x_1)P(x_2|x_1) \dots P(x_m|x_1, x_2, \dots, x_{m-1})$$

- N-元模型:

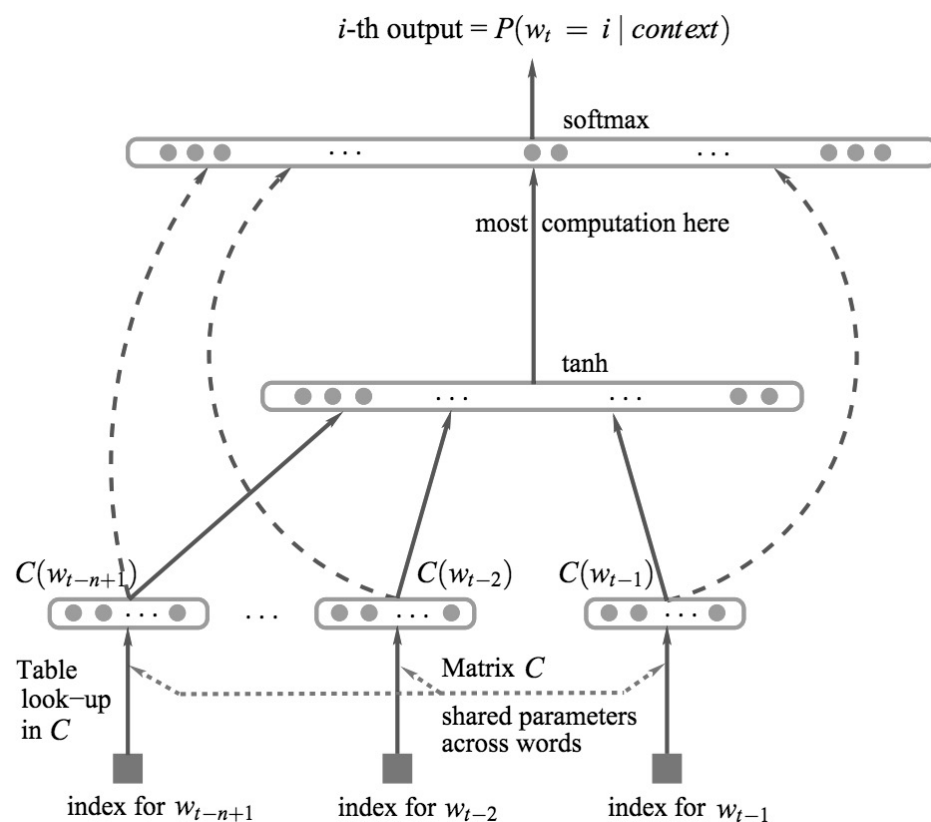
$$P(x_1, x_2, \dots, x_m) = P(x_1)P(x_2|x_1) \dots P(x_m|x_1, x_2, \dots, x_{m-1})$$

- 极大似然估计:

$$P(x_i|x_{i-(n-1)}, \dots, x_{i-1}) = \frac{\text{count}(x_{i-(n-1)}, \dots, x_{i-1}, x_i)}{\text{count}(x_{i-(n-1)}, \dots, x_{i-1})}$$

Neural Language Model

神经语言模型



A Neural Probabilistic Language Model, Yoshua Bengio, etc

Neural Language Model

神经语言模型

- “一只狗在叫”出现一百次 v.s. “一只猫在叫”出现一次
- 对n-gram模型来说, 狗和猫的权重不一致
- 对神经概率语言模型来说:
 - 假定了相似词对于的词向量也相似
 - 概率函数关于词向量是光滑的

RNN Language Model

递归神经网络语言模型

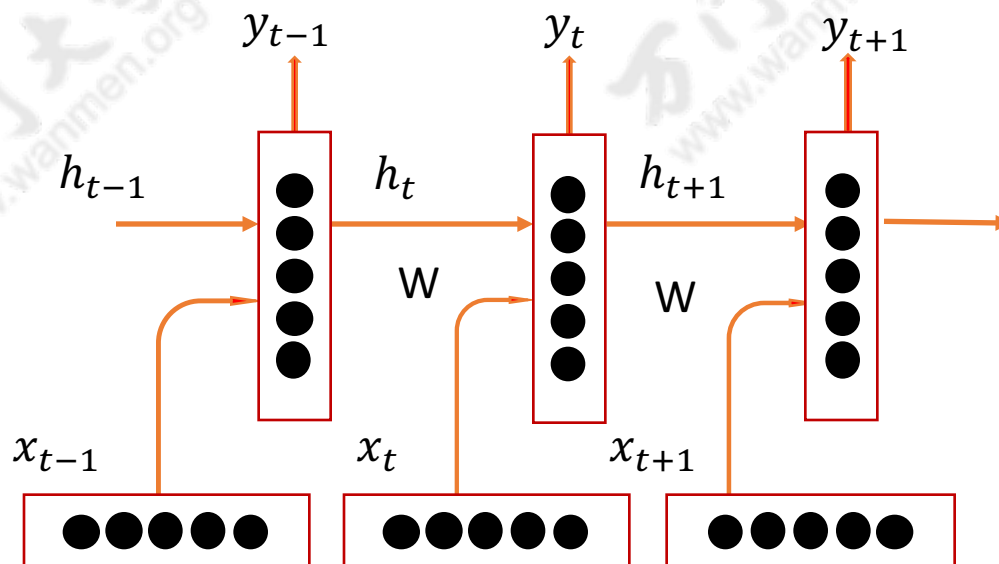
- 句子: x_1, \dots, x_m
- 时间 t 时刻:

$$h_t = \sigma \left(W^{(hh)} h_{t-1} + W^{(hx)} e(x_t) \right)$$

$$\hat{y}_t = \text{softmax}(W^{(s)} h_t)$$

$$\hat{P}(y_{t+1} = v_j | x_t, \dots, x_1) = \hat{y}_{t,i}$$

- 结构



Word Embedding

词嵌入

- 独热表示(One-hot Encoding):

国王	0	1	0	0	0
王后	0	0	0	1	0

- 分布式表示(Distributed Representation):

The advantage of Distributed Representation compared to one-hot encoding is that: It can find the relationship between wordings that are similar in broad category, not only the spelling.

国王	4.3	1.0	0.6	3.2	-0.7
王后	0.25	0.5	-1.2	0.9	1.2

Word Embedding: Skip-Gram & CBOW

词嵌入: Skip-Gram & CBOW

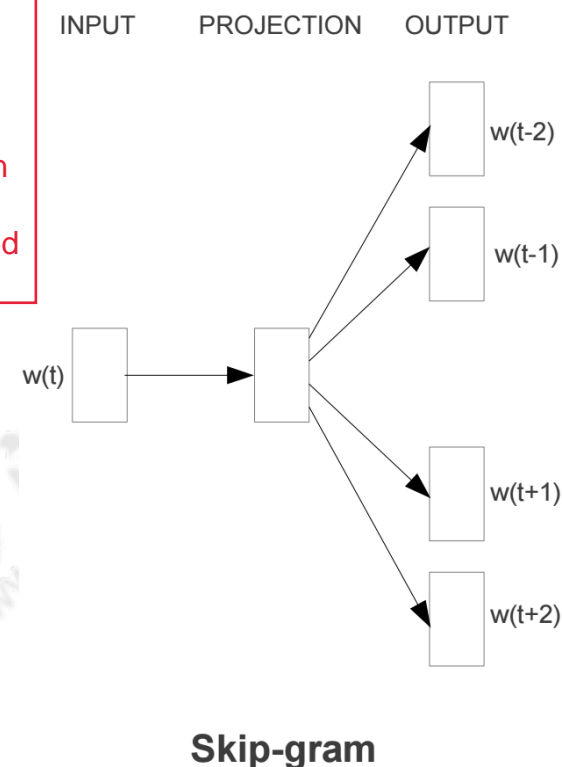
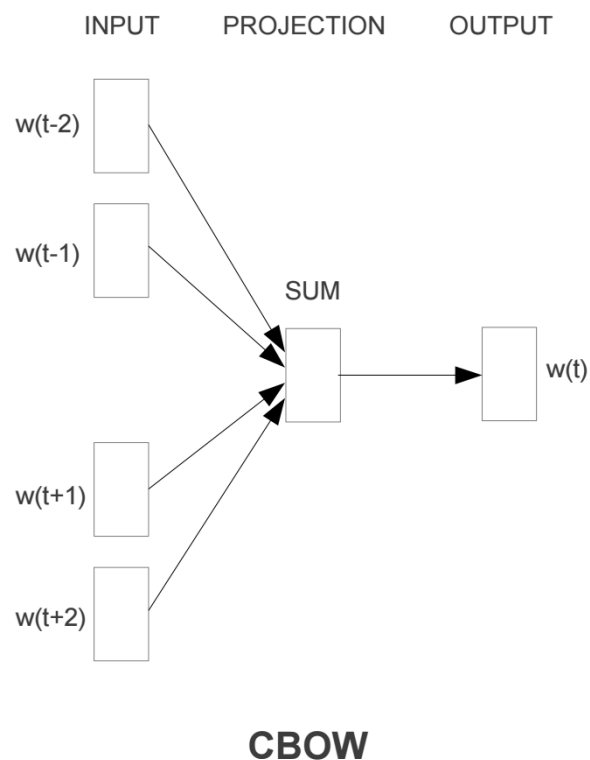
Word Embedding common techniques:

1. Word2Vec; 2. Glove

In Word2Vec technique, there are 2 models of word embedding:
CBOW: Continuous Bag of Words
Skip Gram: reverse of the CBOW.

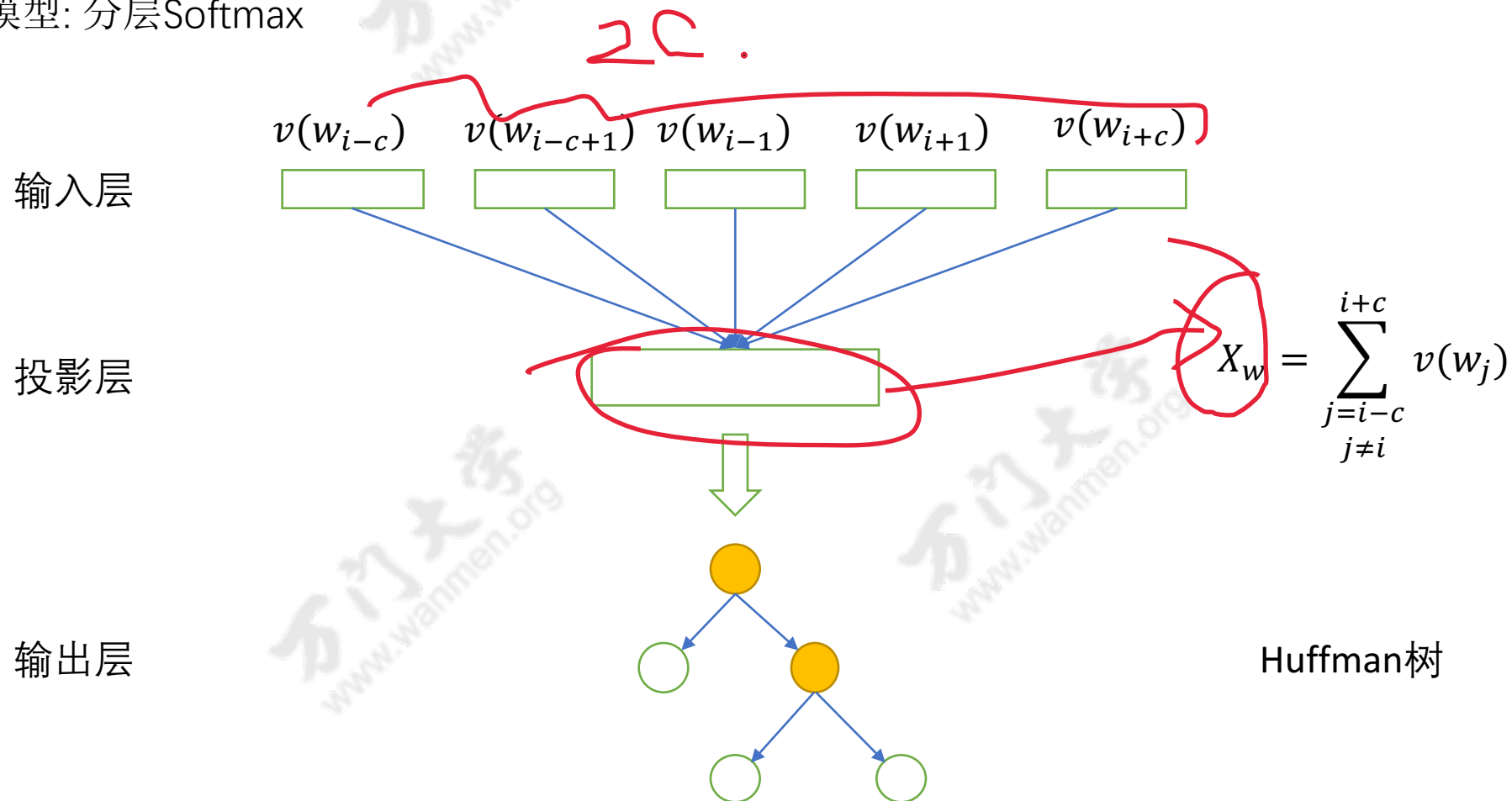
Difference:

In the CBOW model, the distributed representations of context (or surrounding words) are combined to predict the word in the middle. While in the Skip-gram model, the distributed representation of the input word is used to predict the context.



CBOW: Hierarchical Softmax

连续词袋模型: 分层Softmax



CBOW: Hierarchical Softmax

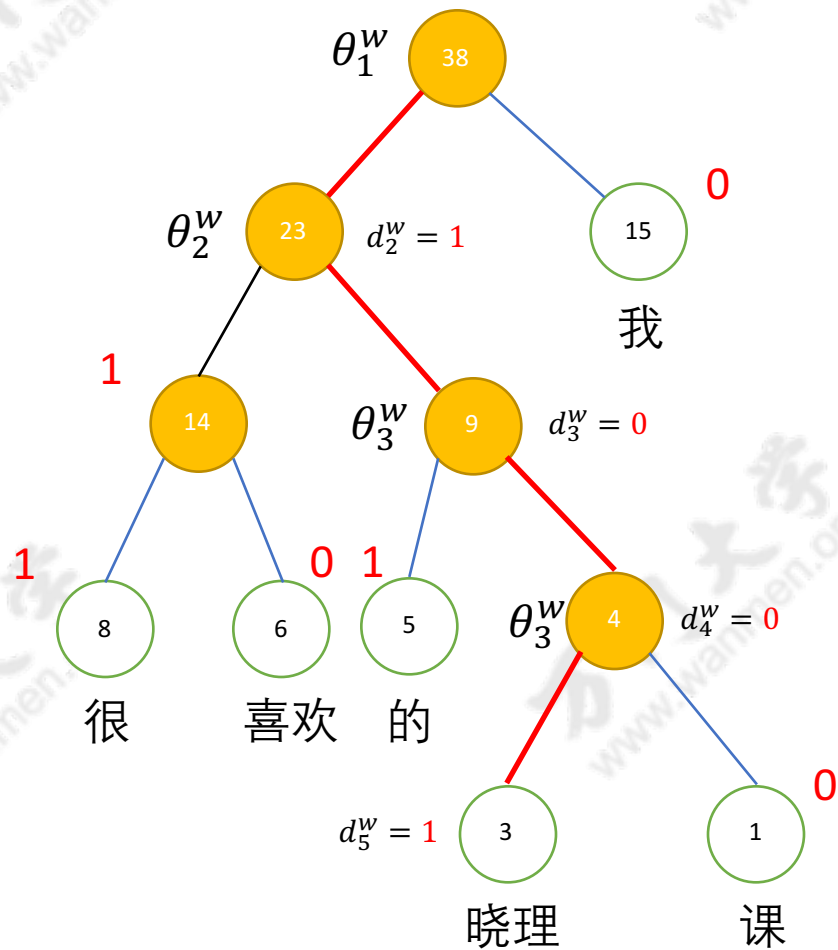
连续词袋模型: 分层Softmax

- p^w : 从根结点出发达到 w 对应的叶子结点的路径
- n^w : 路径 p^w 中包含的结点的个数
- $p_1^w, p_2^w, \dots, p_{n^w}^w$: 路径 p^w 中的 n^w 结点, p_1^w 为根结点, $p_{n^w}^w$ 为词 w 对应的结点
- $d_2^w, d_3^w, \dots, d_{n^w}^w \in \{0,1\}$: 词 w 的Huffman编码, 由 $l^w - 1$ 位编码构成, d_j^w 表示路径 p^w 中第 j 个结点对应的编码
- $\theta_1^w, \theta_2^w, \dots, \theta_{n^w}^w \in R^m$: 路径 p^w 中非叶子结点对应的向量

These parameters
are trained and
calculated

CBOW: Hierarchical Softmax

连续词袋模型: 分层Softmax



- $w = \text{'晓理'}$
- $n^w = 5$
- $d_2^w, d_3^w, d_4^w, d_5^w$ 为 1, 0, 0, 1, '晓理' 的 Huffman 编码
- 每一次分支都视为一次二分类
- $\text{Label}(p_i^w) = 1 - d_i^w, i = 2, 3, \dots, l^w$
- 分为正类概率: $\sigma(X_w^T \theta) = \frac{1}{1 + e^{-X_w^T \theta}}$

CBOW: Hierarchical Softmax

连续词袋模型: 分层Softmax

In total $2c$ words (from $i-c$ to $i+c$)

- 计算 $p(w_i | \underbrace{w_{i-c}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+c}}_{\text{In total } 2c \text{ words (from } i-c \text{ to } i+c)})$
- 第 i 次分支概率: $p(d_i^w | X_w, \theta_{i-1}^w) = \sigma(X_w^T \theta_{i-1}^w)^{1-d_i^w} (1 - \sigma(X_w^T \theta_{i-1}^w))^{d_i^w}$
- $p(w_i | w_{i-c}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+c}) = \prod_{j=2}^{n^w} p(d_i^w | X_w, \theta_{j-1}^w)$
- 例如: $p(\text{晓理} | w_{i-c}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+c}) = \prod_{j=2}^5 p(d_i^w | X_w, \theta_{j-1}^w)$
- 目标函数: $\mathcal{L} = \sum_{w \in \mathcal{C}} \log(p(w | w_{i-c}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+c}))$

CBOW: Hierarchical Softmax

连续词袋模型: 分层Softmax

$$\mathcal{L} = \sum_{w \in C} \log \prod_{j=2}^{n^w} \left\{ \sigma(X_w^T \theta_{j-1}^w)^{1-d_j^w} (1 - \sigma(X_w^T \theta_{j-1}^w))^{d_j^w} \right\}$$

$$\mathcal{L}(w, j) = (1 - d_j^w) \cdot \log \left(\sigma(X_w^T \theta_{j-1}^w) \right) + d_j^w \cdot \log(1 - \sigma(X_w^T \theta_{j-1}^w))$$

$$\frac{\partial \mathcal{L}(w, j)}{\partial \theta_{j-1}^w} = [1 - d_j^w - \sigma(X_w^T \theta_{j-1}^w)] X_w$$

$$\frac{\partial \mathcal{L}(w, j)}{\partial X_w} = [1 - d_j^w - \sigma(X_w^T \theta_{j-1}^w)] \theta_{j-1}^w$$

$$v(\tilde{w}) = v(\tilde{w}) + \lambda \sum_{j=2}^{n^w} \frac{\partial \mathcal{L}(w, j)}{\partial X_w}, \quad \tilde{w} \in \text{Context}(w)$$

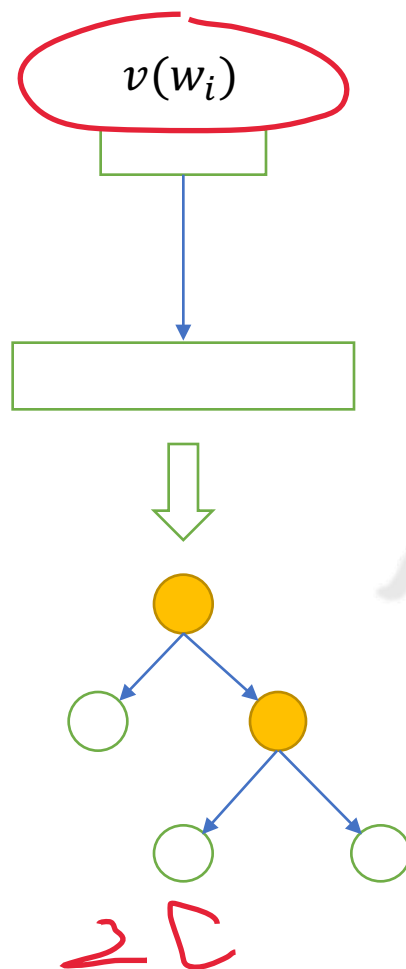
Skip-gram: Hierarchical Softmax

Skip-gram: 分层Softmax

输入层

投影层

输出层



$v(w_i)$

Huffman树

Skip-gram: Hierarchical Softmax

Skip-gram: 分层Softmax

In total $2c$ words (from $i-c$ to $i+c$)

- 计算 $p(w_{i-c}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+c} | w_i)$
- 第 j 次分支概率: $p(d_j^u | v(w), \theta_{j-1}^u) = \sigma(v(w)^T \theta_{j-1}^u)^{1-d_j^u} (1 - \sigma(v(w)^T \theta_{j-1}^u))^{d_j^u}$
- $p(w_{i-c}, \dots, w_{i-1}, w_{i+1}, \dots, w_{i+c} | w_i) = \prod_{j=2}^{n^u} p(d_j^u | v(w), \theta_{j-1}^u)$
- 目标函数: $\mathcal{L} = \sum_{w \in \mathcal{C}} \log \prod_{u \in \text{context}(w)} \prod_{j=2}^{n^u} \left\{ \sigma(v(w)^T \theta_{j-1}^u)^{1-d_j^u} (1 - \sigma(v(w)^T \theta_{j-1}^u))^{d_j^u} \right\}$

Skip-gram: Hierarchical Softmax

Skip-gram: 分层Softmax

$$\mathcal{L} = \sum_{w \in C} \log \prod_{w \in \text{Context}(w)} \prod_{j=2}^{n^u} \left\{ \sigma(v(w)^T \theta_{j-1}^u)^{1-d_j^u} (1 - \sigma(v(w)^T \theta_{j-1}^u))^{d_j^u} \right\}$$

$$\mathcal{L}(w, u, j) = (1 - d_j^u) \cdot \log(\sigma(v(w)^T \theta_{j-1}^u)) + d_j^u \cdot \log(1 - \sigma(v(w)^T \theta_{j-1}^u))$$

$$\frac{\partial \mathcal{L}(w, u, j)}{\partial \theta_{j-1}^u} = [1 - d_j^u - \sigma(v(w)^T \theta_{j-1}^u)] v(w)$$

$$\frac{\partial \mathcal{L}(w, u, j)}{\partial v(w)} = [1 - d_j^u - \sigma(v(w)^T \theta_{j-1}^u)] \theta_{j-1}^u$$

$$v(w) = v(w) + \sum_{u \in \text{Context}(w)} \sum_{j=2}^{l^u} \frac{\partial \mathcal{L}(w, u, j)}{\partial v(w)}, \quad \tilde{w} \in \text{Context}(w)$$

CBOW: Negative Sampling

连续词袋模型: 负采样

• 给定 w 的上下文, 预测 w , w 为正样本, 其余词为负样本

$$L^w(\tilde{w}) = \begin{cases} 1, & \tilde{w} = w \\ 0, & \tilde{w} \neq w \end{cases}$$

$$p(u|\tilde{w}) = [\sigma(v(\tilde{w})^T \theta^u)]^{L^w(u)} \cdot [1 - \sigma(v(\tilde{w})^T \theta^u)]^{1-L^w(u)}$$

$$\text{目标函数: } \mathcal{L} = \sum_{w \in C} \log \prod_{\tilde{w} \in \{w\} \cup \text{NEG}(\tilde{w})} \prod_{j=2}^{n^u} \{[\sigma(v(\tilde{w})^T \theta^u)]^{L^w(u)} [1 - \sigma(v(\tilde{w})^T \theta^u)]^{1-L^w(u)}\}$$

$$\frac{\partial \mathcal{L}(w, \tilde{w}, u)}{\partial \theta^u} = [L^w(u) - \sigma(v(\tilde{w})^T \theta^u)] v(\tilde{w})$$

$$\frac{\partial \mathcal{L}(w, \tilde{w}, u)}{\partial v(\tilde{w})} = [L^w(u) - \sigma(v(\tilde{w})^T \theta^u)] \theta^u$$

$$v(\tilde{w}) = v(\tilde{w}) + \sum_{u \in \{w\} \cup \text{NEG}(w)} \frac{\partial \mathcal{L}(w, u)}{\partial X_w}, \quad \tilde{w} \in \text{Context}(w)$$

Skip-gram: Negative Sampling

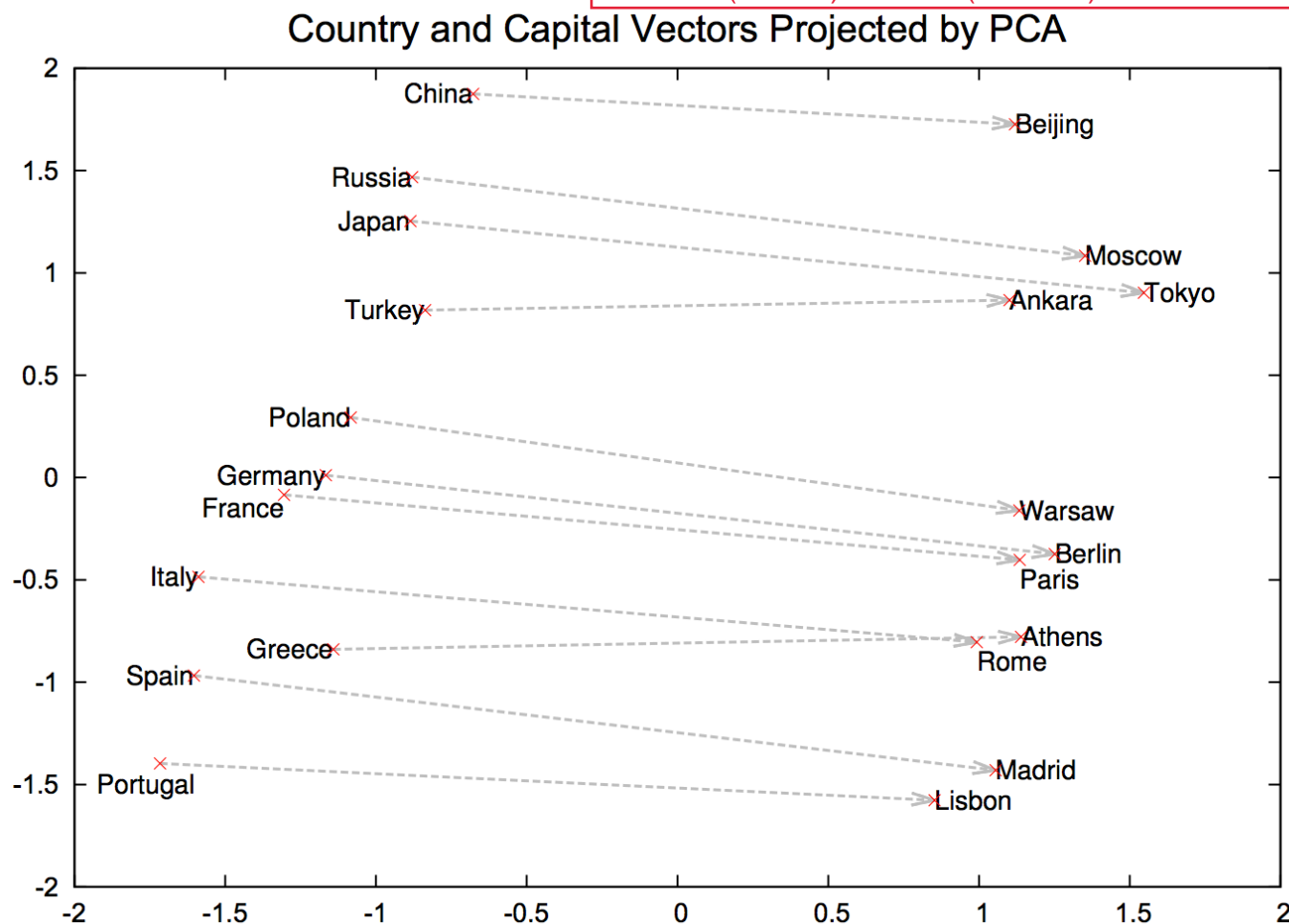
Skip-gram: 负采样

- 给定 w 的上下文, 预测 w , w 为正样本, 其余词为负样本
- $L^w(\tilde{w}) = \begin{cases} 1, & \tilde{w} = w \\ 0, & \tilde{w} \neq w \end{cases}$
- $p(u|\text{Context}(w)) = [\sigma(X_w^T \theta^u)]^{L^w(\tilde{w})} \cdot [1 - \sigma(X_w^T \theta^u)]^{1-L^w(\tilde{w})}$
- 目标函数: $\mathcal{L} = \sum_{w \in C} \log \prod_{u \in \{w\} \cup \text{NEG}(w)} \prod_{j=2}^{n^u} \{[\sigma(X_w^T \theta^u)]^{1-d_j^u} [1 - \sigma(X_w^T \theta^u)]^{d_j^u}\}$
- $\frac{\partial \mathcal{L}(w, u)}{\partial \theta^u} = [L^w(u) - \sigma(X_w^T \theta^u)] X_w$
- $\frac{\partial \mathcal{L}(w, u)}{\partial X_w} = [L^w(u) - \sigma(X_w^T \theta^u)] \theta^u$
- $v(\tilde{w}) = v(\tilde{w}) + \sum_{u \in \{w\} \cup \text{NEG}(w)} \frac{\partial \mathcal{L}(w, u)}{\partial X_w}, \tilde{w} \in \text{Context}(w)$

Word Vectors: Visualization

词向量: 可视化

vector (countries) - vector (capital) is similar no matter what countries or capital.
e.g. Vector (China) - Vector (Beijing) almost equal to Vector (Russia) - Vector (Moscow)



Distributed Representations of Words and Phrases and their Compositionality