

Word2vec 原理解析

1. 前言

近年来，以深度学习为代表的人工智能技术发展如火如荼，其应用领域也从最开始的计算机视觉扩展到了包括自然语言处理、推荐、语音识别等在内的众多研究方向，并且在大多数领域取得了令人惊艳的效果，以笔者熟悉的自然语言处理（NLP）方向为例，纵观近几年在ACL、EMNLP、CLONG等顶级学术会议上发表的学术论文，随处可见Deep Learning、Neural Network/NN、RNN、Convolutional Neural Networks/CNN、LSTM、Autoencoder、Embedding等字样，传统的自然语言处理的任务比如文本分类、情感分析、命名实体识别、分词、关系抽取、机器翻译等任务，目前取得State of Art的方案也基本被基于深度学习的算法占据。文本表示是NLP中的基础，Word2Vec作为一种基于神经网络思想的文本表示学习方法，对在深度学习盛行的今天想要学习NLP的人来说是必要的入门知识。

word2vec是谷歌开源的一款用于词向量计算的工具，同时也是一套生成词向量的算法方案，其他生成词向量的方案有LSA、PLSA、LDA等等。Word2Vec算法的背后是一个浅层神经网络，其网络深度仅为3层，所以严格说Word2Vec并非深度学习范畴，但其生成的词向量在很多任务中都可以作为深度学习算法的输入，因此在一定程度上可以说Word2Vec技术是深度学习在NLP领域能够更好应用的基础。

本文主要涉及一下知识点：

* 1) 相关背景知识 * 2) Word2Vec中两种词向量模型CBOW和Skip-gram * 3) Word2Vec中两种学习框架Hierarchical Softmax和Negative Sampling的原理分析 * 4) Word2Vec相关拓展与应用

###2. 相关背景知识

####2.1 词向量

计算机不能直接处理文本、图像、声音等内容，需要将其转化为数字特征后才能处理，词向量就是文本中的单词转后的计算机能够处理的数字化特征。词向量的表示方式主要有两种，一种是one-hot representation，就是利用一个高维稀疏（所谓稀疏是指向量中为0的维度特别多）的向量表示一个单词，向量的维度为词典中词的个数，向量中只有一个维度的值为1，其余维度值为0，为1的维度为词在词典中的索引位置，常见的有词袋（BOW，bag of words）模型；另一种是distributed representation，该方法将单词表示成一个低维稠密的向量，常见的方法有LSA、PLSA、LDA、Word2vec、Glove等等。下面以单词powerful为例说明两种不同的表示方式：

$$S=(w_1,w_2,\dots,w_T)$$

w_i 表示 S 中的第 i 的单词，则文本 S 出现的概率为：

![[enter image description here]](<http://images.gitbook.cn/72096170-bd20-11e7-8439-65f318c92b02>)

上面公式利用了贝叶斯公式进行了分解，上式中各个概率值是语言模型的参数，对于一段新给出的文本，只需将分解后的参数进行累乘即可得出新文本出现的概率。

语言模型参数的更一般表示形式为： $p(w/\text{context}(w))$ ，其中 w 表示某个单词， $\text{context}(w)$ 表示单词 w 的上下文，含义为 L 上下文为 $\text{context}(w)$ 时,词 w 出现的概率。

学习语言模型参数 $p(w/\text{context}(w))$ ，目前解决方案有很多，例如N-gram模型、N-pos模型和神经概率语言模型（NPLM）以及本文提到的Word2Vec中的Hierarchical Softmax模型和Negative Sampling模型（其实本质上HS和Negative Sampling是对传统NPLM语言模型的一种简化，其中基于HS的实现是利用Huffman Tree替代了原有的全连接网络）。

这里特别说明一点：词向量是神经网络统计语言模型的副产品。

###3. CBOW模型

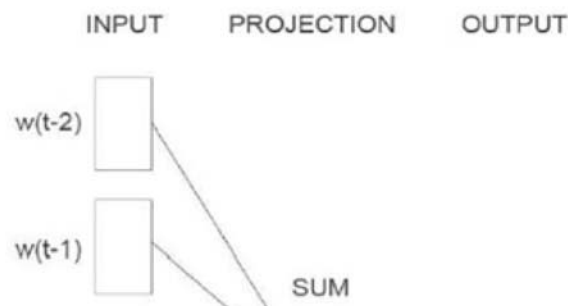
下面介绍连续词袋模型（CBOW model, Continuous Bag of Words Model）。

####3.1 模型结构

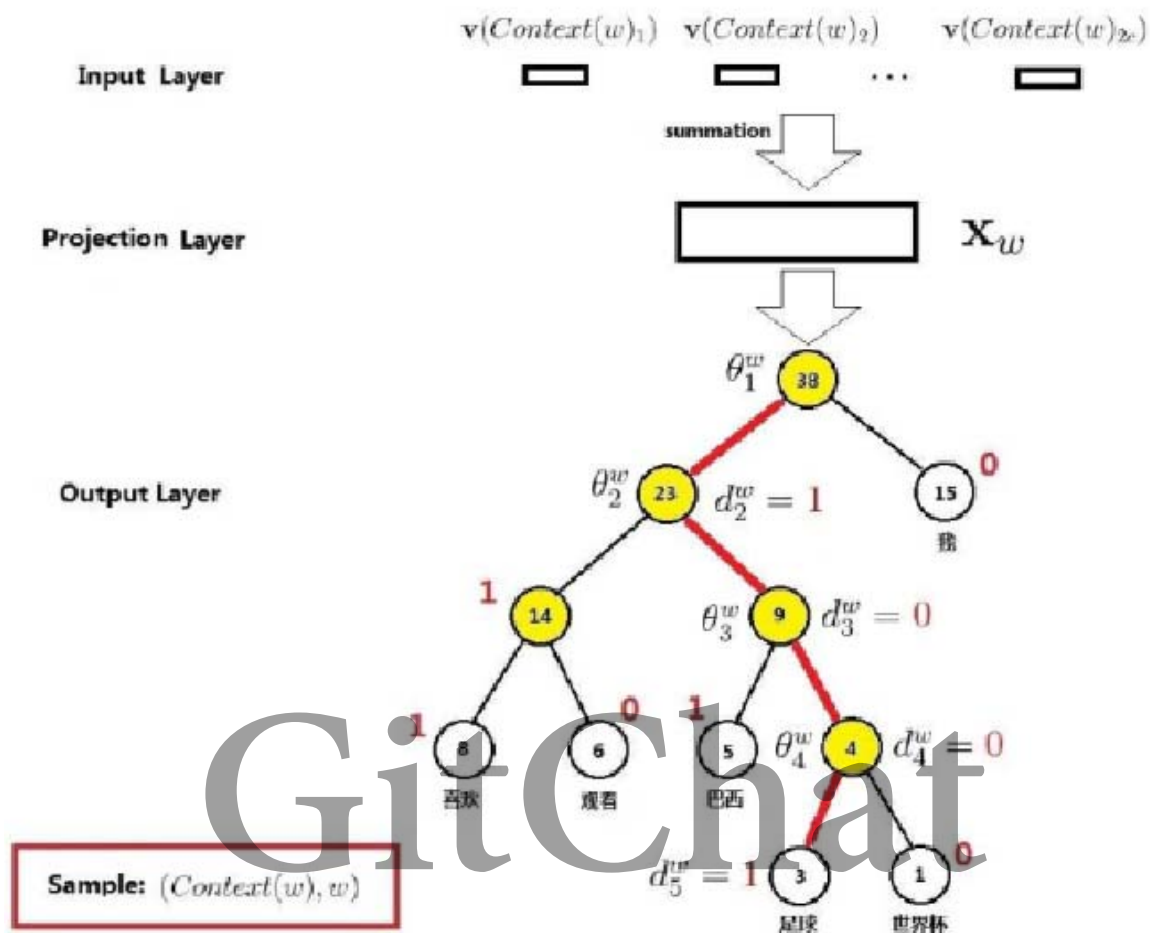
模型包含三层，输入层，映射层和输出层。该网络的含义是已知词 $w(t)$ 的上下文 $w(t-2)$ ， $w(t-1)$ ， $w(t+1)$ ， $w(t+2)$ 的前提下预测词 $w(t)$ 出现的概率，即： $p(w/\text{context}(w))$ 。目标函数为：

$$\mathcal{L} = \sum_{w \in C} \log p(w | \text{Context}(w))$$

CBOW 模型中 $p(w/\text{context}(w))$ 的表示框架有两种，分别为：Hierarchical Softmax和Negative Sampling，下面分别对两种框架下的模型参数求解进行介绍。



基于Hierarchical Softmax框架的CBOW模型将输出层利用一棵Huffman Tree进行表示，该Huffman Tree使用语料中的词根据词频进行构建，树中每个叶子节点代表一个单词，约定（词频较大的）左孩子节点编码为1，（词频较小的）右孩子节点编码为0。同时利用节点的层级关系，将从根节点到单词t的路径中每个节点当做一次二分类（基于Logistic Regression算法），将整个路径产生的概率作为 $p(w/context(w))$ 。



网络包含以下几部分：

- 1) 输入层：包含 $Context(w)$ 中 $2c$ 个词向量 $v(Context(w)_1), v(Context(w)_2), \dots, v(Context(w)_{2c})$, 每个词向量的维度为 m 。
- 2) 投影层：将输入层的 $2c$ 个向量求和累加，即：

$$x_w = \sum_{i=1}^{2c} v(Context(w)_i) \in \mathbb{R}^m.$$

对节点进行分类使用的是逻辑回归算法，因此一个节点被分为正类的概率是

$$\sigma(\mathbf{x}_w^\top \theta) = \frac{1}{1 + e^{-\mathbf{x}_w^\top \theta}},$$

被分为负类的概率为

$$1 - \sigma(\mathbf{x}_w^\top \theta),$$

其中 θ 为待定参数，每个非叶节点对应一个，设为 θ_i^w ，对于从根节点出发到达“足球”叶节点共经历了4次二分类，每次分类结果的概率为：

1. 第 1 次: $p(d_2^w | \mathbf{x}_w, \theta_1^w) = 1 - \sigma(\mathbf{x}_w^\top \theta_1^w);$

2. 第 2 次: $p(d_3^w | \mathbf{x}_w, \theta_2^w) = \sigma(\mathbf{x}_w^\top \theta_2^w);$

3. 第 3 次: $p(d_4^w | \mathbf{x}_w, \theta_3^w) = \sigma(\mathbf{x}_w^\top \theta_3^w);$

4. 第 4 次: $p(d_5^w | \mathbf{x}_w, \theta_4^w) = 1 - \sigma(\mathbf{x}_w^\top \theta_4^w),$

$p(\text{足球} | \text{Context}(\text{足球}))$ 跟这四个概率值的关系有：

$$p(\text{足球} | \text{Context}(\text{足球})) = \prod_{j=2}^5 p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w).$$

综上所述：对于训练语料的词典 D 中任意词 w ，Huffman树上必然存在唯一一条从根节点到词 w 对应节点的路径 p^w ，路径 p^w 上存在 l^w-1 个分支，将每个分支看做一次二分类，每一次二分类产生一个概率，这些概率的乘机即为条件概率 $p(w | \text{context}(w))$ ，其一般形式如下：

$$P(w | \text{Context}(w)) = \prod_{j=2}^{l^w} p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w)$$

其中:

$$\begin{aligned}
\mathcal{L} &= \sum_{w \in \mathcal{C}} \log \prod_{j=2}^{l^w} \{ [\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{1-d_j^w} \cdot [1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)]^{d_j^w} \} \\
&= \sum_{w \in \mathcal{C}} \sum_{j=2}^{l^w} \{ (1 - d_j^w) \cdot \log[\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] + d_j^w \cdot \log[1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] \},
\end{aligned}$$

接下来对 \mathcal{L} 求极大值，需要对 \mathcal{L} 中的变量 $(\mathbf{x}_w$ 和 $\theta_{(j-1)}^w$)进行求解梯度，为方便梯度推导令：

$$\mathcal{L}(w, j) = (1 - d_j^w) \cdot \log[\sigma(\mathbf{x}_w^\top \theta_{j-1}^w)] + d_j^w \cdot \log[1 - \sigma(\mathbf{x}_w^\top \theta_{j-1}^w)].$$

首先计算 $\mathcal{L}(w, j)$ 关于 $\theta_{(j-1)}^w$ 的梯度：

$$\begin{aligned}
\frac{\partial \mathcal{L}(w, j)}{\partial \vartheta_{j-1}^w} &= \frac{\partial}{\partial \vartheta_{j-1}^w} \{ (1 - d_j^w) \cdot \log[\sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w)] + d_j^w \cdot \log[1 - \sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w)] \} \\
&= (1 - d_j^w) [1 - \sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w)] \mathbf{x}_w - d_j^w \sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w) \mathbf{x}_w \\
&= \{ (1 - d_j^w) [1 - \sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w)] \mathbf{x}_w - d_j^w \sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w) \mathbf{x}_w \} \\
&= (1 - d_j^w - \sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w)) \mathbf{x}_w
\end{aligned}$$

其中 η 为更新速率。

同样的， $\mathcal{L}(w, j)$ 关于 \mathbf{x}_w 的梯度为：

$$\vartheta_{j-1}^w := \vartheta_{j-1}^w + \eta [1 - d_j^w - \sigma(\mathbf{x}_w^\top \vartheta_{j-1}^w) \mathbf{x}_w].$$

因此， \mathbf{x}_w 的更新公式为：

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

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

表示w的标签，即正样本标签为1，负样本标签为0，对比HS模型中的非叶子节点的分类。

对于给定的词w和上下文Context(w)，我们的目的是最大化下式：

$$g(w) = \prod_{u \in \{w\} \cup NEG(w)} p(u | Context(w)),$$

其中：

$$p(\mu | Context(\omega)) = \begin{cases} \sigma(x_w^T \vartheta_{j-1}^w), & L^\omega(\mu) = 1 \\ 1 - \sigma(x_w^T \vartheta_{j-1}^w), & L^\omega(\mu) = 0 \end{cases}$$

或者整体表达为：

$$p(\mu | Context(\omega)) = [\sigma(x_w^T \vartheta_{j-1}^w)]^{L^\omega(\mu)} \cdot [1 - \sigma(x_w^T \vartheta_{j-1}^w)]^{1-L^\omega(\mu)}.$$

对于一个给定的语料库C，模型的目标函数为：

$$G = \prod_{w \in C} g(w)$$

为了方便计算，对其取对数操作，即为L

$$\begin{aligned} \mathcal{L} &= \log G = \log \prod_{w \in C} g(w) = \sum_{w \in C} \log g(w) \\ &= \sum_{w \in C} \log \prod_{u \in \{w\} \cup NEG(w)} \left\{ [\sigma(\mathbf{x}_w^\top \theta^u)]^{L^w(u)} \cdot [1 - \sigma(\mathbf{x}_w^\top \theta^u)]^{1-L^w(u)} \right\} \\ &= \sum_{w \in C} \sum_{u \in \{w\} \cup NEG(w)} \left\{ L^w(u) \cdot \log [\sigma(\mathbf{x}_w^\top \theta^u)] + [1 - L^w(u)] \cdot \log [1 - \sigma(\mathbf{x}_w^\top \theta^u)] \right\}. \end{aligned}$$

$$\frac{\partial \mathcal{L}(w, u)}{\partial \mathbf{x}_w} = [L^w(u) - \sigma(\mathbf{x}_w^\top \theta^u)] \theta^u.$$

因此， θ^u 和 \mathbf{x}_w 的更新公式为：

$$\theta^u := \theta^u + \eta [L^w(u) - \sigma(\mathbf{x}_w^\top \theta^u)] \mathbf{x}_w.$$

$$\mathbf{v}(\tilde{w}) := \mathbf{v}(\tilde{w}) + \eta \sum_{u \in \{w\} \cup \text{NEG}(w)} \frac{\partial \mathcal{L}(w, u)}{\partial \mathbf{x}_w}, \quad \tilde{w} \in \text{Context}(w).$$

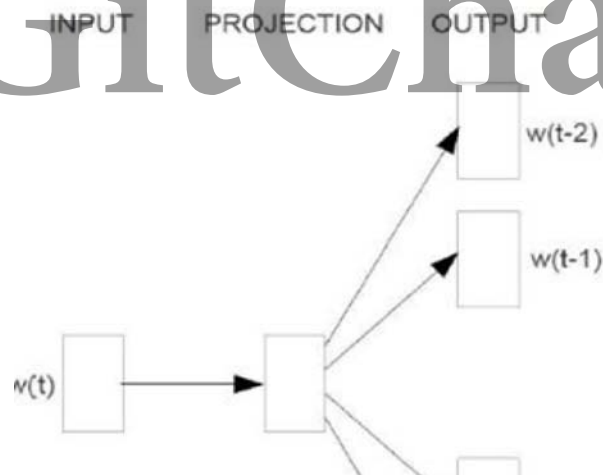
4. Skip-gram模型

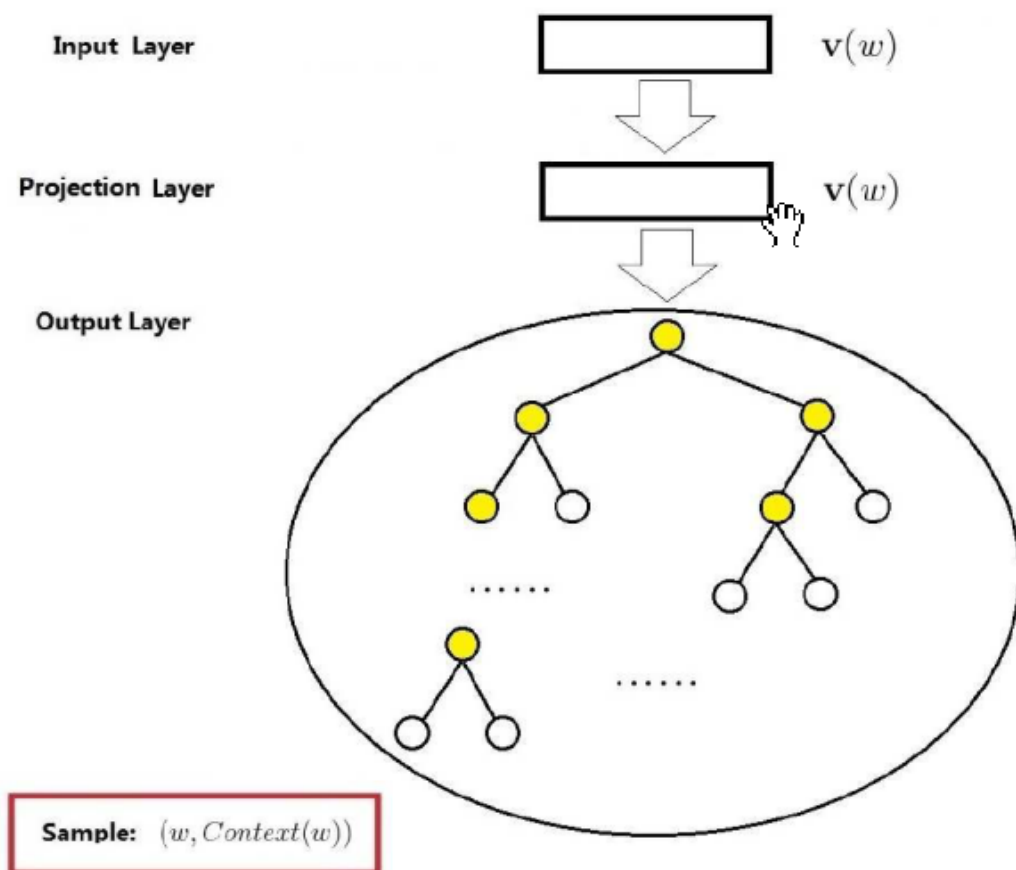
4.1 模型结构

模型同样包含三层，输入层，映射层和输出层。该网络的含义是已知词 $w(t)$ 的前提下预测词 $w(t)$ 的上下文 $w(t-2)$ ， $w(t-1)$ ， $w(t+1)$ ， $w(t+2)$ ，条件概率写为： $p(\text{context}(w)/w)$ 。

目标函数为：

$$L = \sum_{w \in C} \log p(\text{Context}(w)|w).$$





基于Hierarchical Softmax框架的Skip-gram模型同样将输出层利用一棵Huffman Tree进行表示，该Huffman Tree使用语料中的所有词根据词频进行构建，树中每个叶子节点代表一个单词，同时利用节点的层级关系，将从根节点到context(t)中每个单词所在叶子节点的路径中每个节点当做一次二分类（基于Logistic Regression算法），将所有路径产生的概率作为 $p(w/\text{context}(w))$ 。

- 1)输入层：包含1个词向量 $v(w)$ ，词向量的维度为 m 。
- 2)投影层：与输入层相同
- 3)输出层：输出层对应一棵二叉树，它是以语料中快出现过词作为节点，以各词在语料中出现的次数当权值构造出来的Huffman树。这个Huffman树中，叶子节点 N 个，分别对应词典中的词，非叶子节点 $N-1$ 个（上图中黄色结点）。

在skip-gram模型中 $p(w/\text{context}(w))$ 的定义为：

将上面公式带入目标函数

$$L = \sum_{w \in C} \log p(\text{Context}(w)|w)$$

得：

$$\begin{aligned} L &= \sum_{w \in C} \log \prod_{\mu \in \text{Context}(w)} \prod_{j=2}^{l^\mu} \{[\sigma(v(\omega)^T \theta_{j-1}^\mu)]^{1-d_j^\mu} \cdot [1 - \sigma(v(\omega)^T \theta_{j-1}^\mu)]^{d_j^\mu}\} \\ &= \sum_{w \in C} \sum_{\mu \in \text{Context}(w)} \sum_{j=2}^{l^\mu} \{(1 - d_j^\mu) \cdot \log [\sigma(v(\omega)^T \theta_{j-1}^\mu)] + d_j^\mu \log [1 - \sigma(v(\omega)^T \theta_{j-1}^\mu)]\} \end{aligned}$$

下面使用梯度上升对L进行优化，分别对公式中的变量求偏导数：

$$\begin{aligned} \frac{\partial \mathcal{L}(w, u, j)}{\partial \theta_{j-1}^u} &= \frac{\partial}{\partial \theta_{j-1}^u} \{(1 - d_j^u) \cdot \log[\sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)] + d_j^u \cdot \log[1 - \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)]\} \\ &= (1 - d_j^u)[1 - \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)] \mathbf{v}(w) - d_j^u \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u) \mathbf{v}(w) \\ &= \{(1 - d_j^u)[1 - \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)] - d_j^u \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)\} \mathbf{v}(w) \\ &= [1 - d_j^u - \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)] \mathbf{v}(w). \end{aligned}$$

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

因此， $\theta_{(j-1)}^u$ 的更新公式为：

$$\theta_{j-1}^u := \theta_{j-1}^u + \eta [1 - d_j^u - \sigma(\mathbf{v}(w)^\top \theta_{j-1}^u)] \mathbf{v}(w).$$

$\mathbf{v}(w)$ 的更新公式为：

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

NEG(u)表示处理词u时生成的负样本子集，条件概率

$$p(z|w) = \begin{cases} \sigma(v(w)^T \theta^z), & L^u(z) = 1 \\ 1 - \sigma(v(w)^T \theta^z), & L^u(z) = 0 \end{cases}$$

写成完整的表达式

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

为了简化运算，最终目标函数取对数形式：

$$\begin{aligned} \mathcal{L} &= \log G = \log \prod_{w \in \mathcal{C}} \prod_{u \in \text{Context}(w)} g(u) = \sum_{w \in \mathcal{C}} \sum_{u \in \text{Context}(w)} \log g(u) \\ &= \sum_{w \in \mathcal{C}} \sum_{u \in \text{Context}(w)} \log \prod_{z \in \{u\} \cup \text{NEG}(u)} p(z|w) \\ &= \sum_{w \in \mathcal{C}} \sum_{u \in \text{Context}(w)} \sum_{z \in \{u\} \cup \text{NEG}(u)} \log p(z|w) \\ &= \sum_{w \in \mathcal{C}} \sum_{u \in \text{Context}(w)} \sum_{z \in \{u\} \cup \text{NEG}(u)} \log \left\{ [\sigma(v(w)^T \theta^z)]^{L^u(z)} \cdot [1 - \sigma(v(w)^T \theta^z)]^{1-L^u(z)} \right\} \\ &= \sum_{w \in \mathcal{C}} \sum_{u \in \text{Context}(w)} \sum_{z \in \{u\} \cup \text{NEG}(u)} \{ L^u(z) \cdot \log [\sigma(v(w)^T \theta^z)] + [1 - L^u(z)] \cdot \log [1 - \sigma(v(w)^T \theta^z)] \}. \end{aligned}$$

下面对公式中的变量求偏导，进而得出变量的更新公式，具体推导过程参考之前的方法，这里不再赘述。

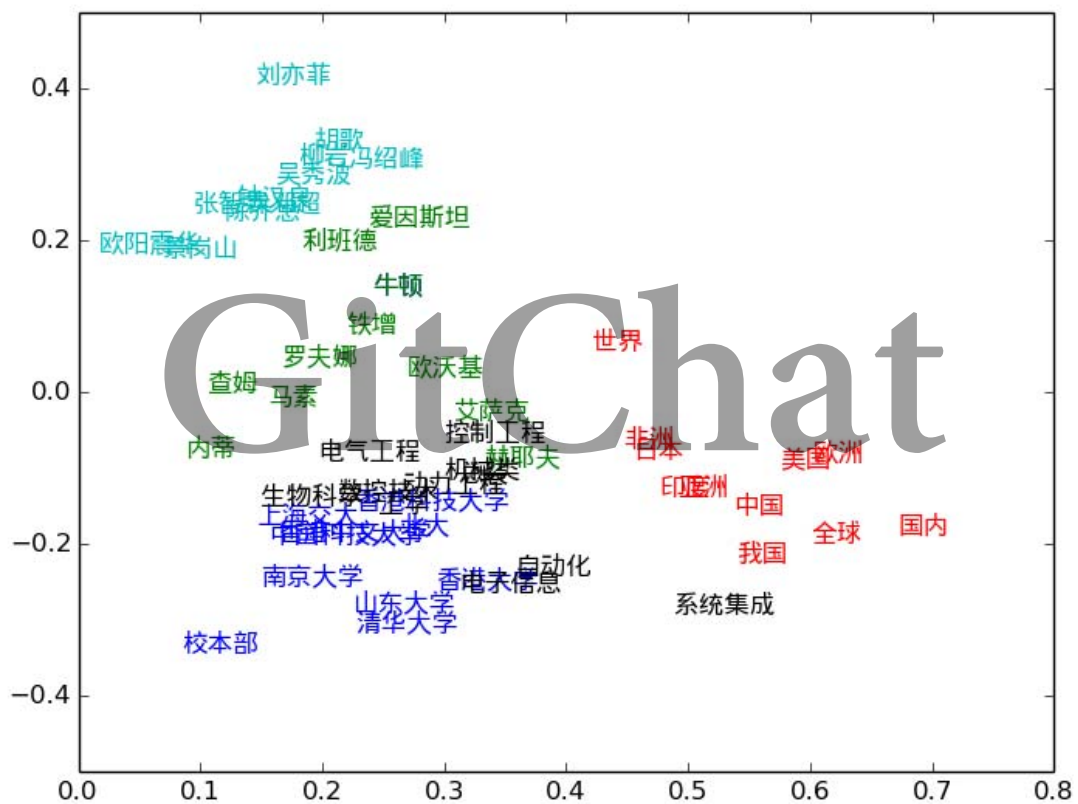
5. 相关应用与拓展

5.1 Word2Vec的应用

1) 近义词/相似文本发现

请输入词语<exit退出>:中国		请输入词语<exit退出>:钓鱼岛		请输入词语<exit退出>:旅游		请输入词语<exit退出>:苹果	
大陆	0.66763467	钓鱼台	0.6219264	观光	0.65619475	三星	0.7224437
中共	0.57856727	钓鱼岛	0.6123347	景点	0.60212	微软	0.7101249
共产党	0.56305367	南海	0.6018163	陆客	0.59477097	Apple	0.66682446
解放军	0.55761635	领土	0.51753837	旅行	0.5677106	iPhone5	0.62071097
台湾	0.5368497	领海	0.4928774	游憩	0.557839	Google	0.597368
反攻	0.5271177	岛屿	0.4853142	赏樱	0.5571045	iPadmini	0.5609188
日本	0.5103535	舰队	0.47854927	游玩	0.52199984	新机	0.559093
王文莹	0.49295437	渔权	0.47229362	观光客	0.51974636	库克	0.5589176
内地	0.48557448	主权	0.46729872	行程	0.51943743	宏达电	0.555409
对岸	0.48428434	东海	0.4613399	参观	0.5077874	产品	0.55437565

下图为使用PCA将词向量进行降维后的效果展示：



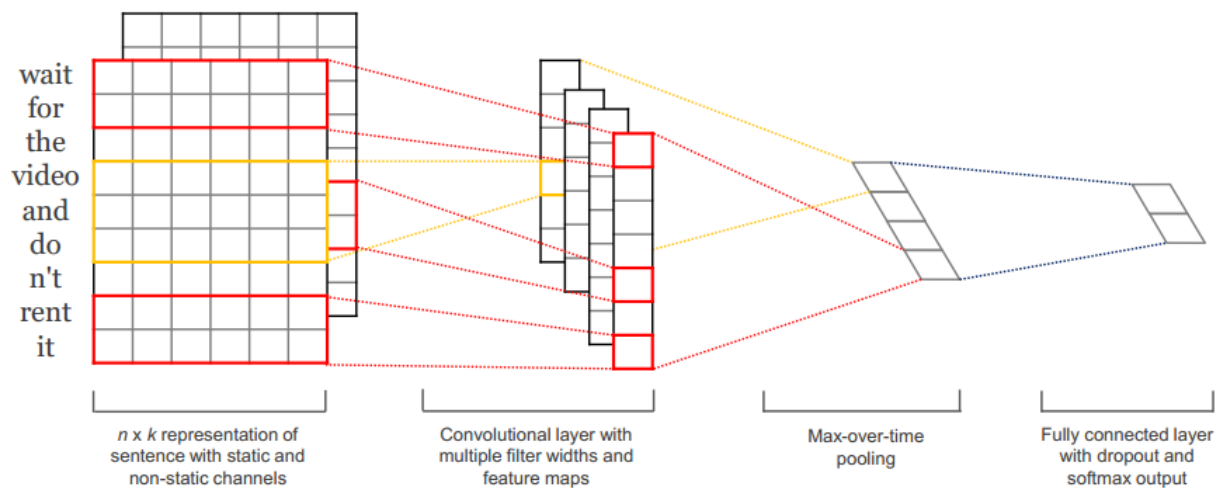
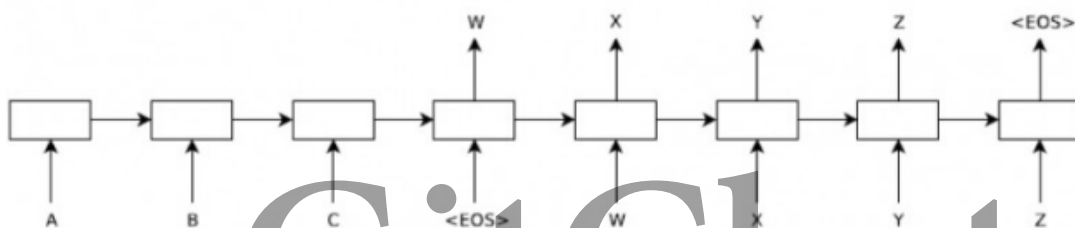
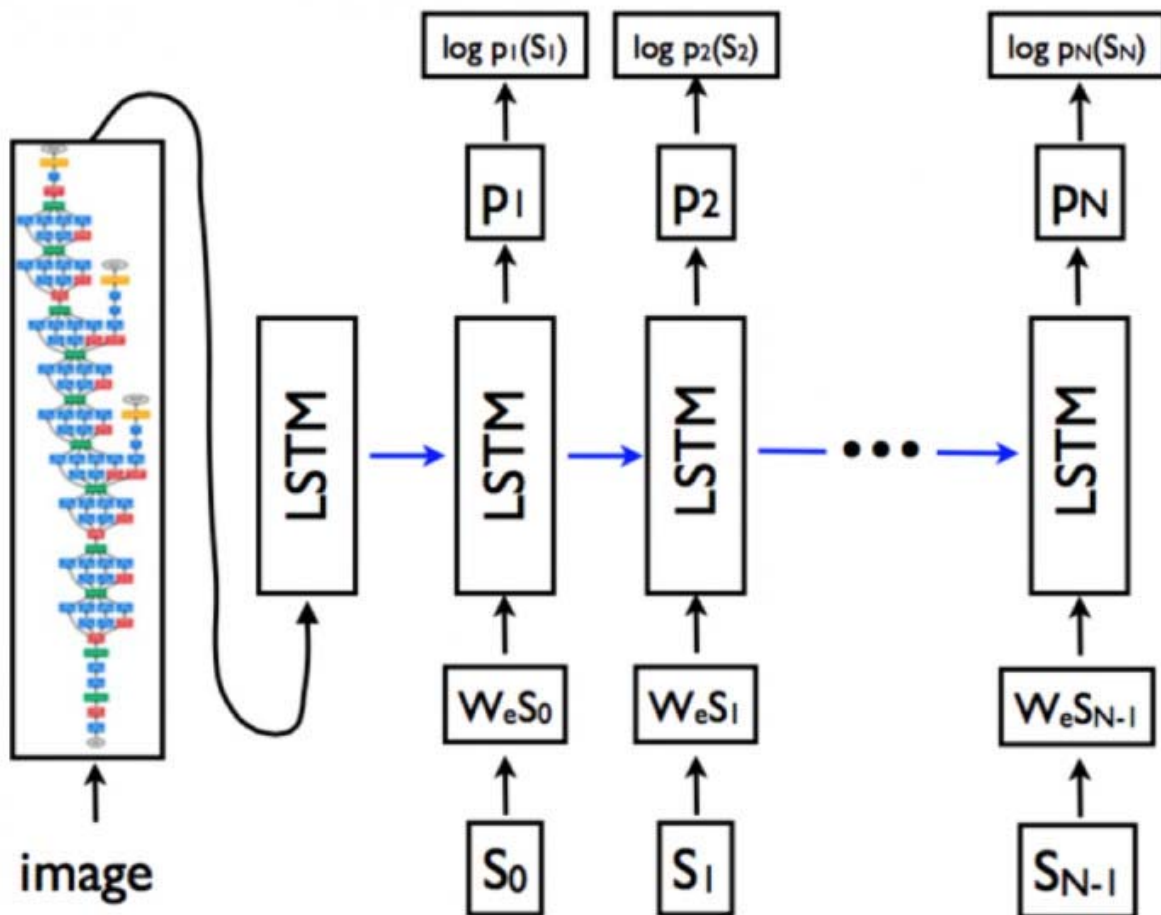


Figure 1: Model architecture with two channels for an example sentence.

同样的，在分词、实体抽取、机器对话、神经网络机器翻译（NMT）等任务中广泛应用的Seq2Seq模型中也有word2vec的应用，类似的词向量也是作为词的特征的进行输入。

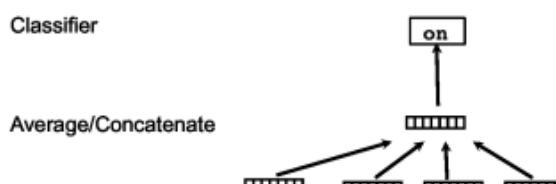


类似的应用还有生成自动图像描述等，模型网络结构如下图：



5.2 其他知识表示

归根到底，word2vec是一种对单词知识的表示方法，word2vec之后出现了大量知识表示方法的研究，例如“段向量”（paragraph vector），提出者将段信息作为一个特殊的词，参与模型的训练，模型结构基本与word2vec相同，具体如下图所示：



类似的知识表示还有graph2vec、thought2vec等等。

6. 参考文献

- [1] Mikolov T, Chen K, Corrado G, et al. Efficient estimation of word representations in vector space[J]. arXiv preprint arXiv:1301.3781, 2013.
- [2] Mikolov T, Sutskever I, Chen K, et al. Distributed representations of words and phrases and their compositionality[C]//Advances in neural information processing systems. 2013: 3111-3119.
- [3] Le Q, Mikolov T. Distributed representations of sentences and documents[J]. 2014, 4:11-1188.
- [4] A. Minnaar, “Part II - Continuous Bag-of-Words Model,” 12 4 2015. [Online]. Available:http://mccormickml.com/assets/word2vec/Alex_Minnaar_Word2Vec_Tutorial_Part_I_The_Skip-Gram_Model.pdf.
- [5] C. McCormick, “Word2Vec Tutorial Part 2 - Negative Sampling,” 11 1 2017. [Online]. Available: <http://mccormickml.com/2017/01/11/word2vec-tutorial-part-2-negativesampling/>.
- [6] C. McCormick, “Word2Vec Tutorial - The Skip-Gram Model,” 19 4 2016. [Online]. Available: <http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-grammodel/>.
- [7] peghoty, “word2vec 中的数学,” 2014, pp.19-20.

GitChat