

# 基于词分布式表征的汉语框架排歧模型

张力文<sup>1</sup>, 王瑞波<sup>1,2</sup>, 李茹<sup>1,3,4</sup>, 张晟<sup>1</sup>

(1.山西大学 计算机与信息技术学院, 山西 太原 030006;

2.山西大学 软件学院, 山西 太原 030006;

3.山西大学 计算智能与中文信息处理教育部重点实验室, 山西 太原 030006;

4.山西省大数据挖掘与智能技术协同创新中心, 山西 太原 030006)

**摘要:** 框架排歧是根据句子中目标词的上下文语境, 从框架库中为该目标词自动选择一个合适的框架。该任务在一定程度上解决了动词中一词多义的现象。本文基于词语及句子的分布式表征, 提出了基于距离和基于词语相似度矩阵的框架排歧模型。与传统方法相比, 本模型有效避免了人工选择特征, 克服了特征空间维度过高、特征之间没有关联性等优点, 使框架排歧的准确率达到 65.71%。并与当前最好的模型, 进行显著性和一致性检验, 进一步验证了词分布式表征对框架排歧任务的有效性。

**关键词:** 汉语框架; 框架排歧; 分布式表征

中图分类号: TP391

文献标识码: A

## Chinese FrameNet Disambiguation Model based on Word Distributed Representation

Zhang Liwen<sup>1</sup>, Wang Ruibo<sup>1,2</sup>, Li Ru<sup>1,3</sup>, Zhang Sheng<sup>1</sup>

(1.School of Computer and Information Technology, Taiyuan, Shanxi 030006, China ;

2. School of Software Shanxi University, Taiyuan, Shanxi 030006, China;

3.Key Laboratory of Computer Intelligence and Chinese Information Processing of Ministry of  
Educatio, Shanxi University, Taiyuan, Shanxi 030006, China;

4.Collaborative Innovation Center Of Big Data Mining and Intelligent Technology in Shanxi,  
Taiyuan, Shanxi 030006, China)<sup>1</sup>

**Abstract:** Frame Disambiguation is based on the context of the target word in the sentence, and automatically selects a suitable framework for the target word from the existing frame library. This task can solve the phenomenon of polysemy in verb. Based on the distributed representation of words and sentences, a framework disambiguation model based on distance and word similarity matrix is proposed. Compared with the traditional methods, the model effectively avoids the artificial selection features, overcomes the shortcomings of the feature space, such as high dimensions and no correlation between features, the accuracy of frame disambiguation is 65.71%. Compare with the current best model, and the validity of the word distributed representation for frame disambiguation is further verified.

**Key words:** Chinese FrameNet; Frame Disambiguation; Word Distributed Representation

### 1 引言

框架排歧任务是框架语义分析的一个子任务<sup>[1]</sup>。具体是指: 从例句库中随机选出一条给

收稿日期:

定稿日期:

**基金项目:** 国家自然科学基金项目 (No.61373082); 国家 863 计划项目 (No.2015AA015407); 国家自然科学基金重点项目 (No.61432011, No.U1435212)

**作者简介:** 张力文(1991-), 男, 硕士研究生, 主要研究方向中文信息处理; 王瑞波(1985-), 男, 博士研究生, 讲师, 主要研究领域为自然语言处理; 李茹(1963-), 女, 博士, 教授, 主要研究领域为中文信息处理和数据库技术; 张晟(1991-), 男, 硕士研究生, 主要研究方向中文信息处理。

定目标词的例句，根据该例句的上下文语境，计算机自动识别出该句所属的框架，排除歧义框架。此任务可以用来消除自然语言中动词的“一词多义”的现象，继而为后续的句子语义分析奠定了重要基础，也为机器翻译、信息检索、自动文摘等应用系统提供语义上的支持。因此，框架排歧任务已经成为框架语义分析中至关重要的一部分。

在框架语义分析中，英文的 FrameNet 作为一种重要的语义资源<sup>[2]</sup>，得到许多研究者的关注。针对汉语，参照英语 FrameNet 构建的汉语框架网<sup>[3]</sup>是一种重要的汉语词语语义分析和理解的资源。它是基于框架语义学理论，以汉语语料为依据，构建的语义知识库。由词元库、框架库和例句库三部分组成。

目前对汉语框架排歧的研究，是将其看作是一个分类问题。利用统计机器学习的方法，人工寻找并选择特征，建立分类器。然而，这样利用特征进行分类的做法主要存在两方面的问题。首先，每种特征的特征标记集合较大，从而导致最终的特征矩阵维度较高且非常稀疏。其次，特征标记之间被认为相互独立，没有任何关联。

针对以上两点不足，本文提出基于距离和基于词相似度矩阵的排歧模型。首先从大规模的无标注语料中，训练词语以及句子的分布式表征，然后应用于上述框架排歧模型中。由于词语和句子分布式表征是低维向量，可以有效的避免特征维度过高，并且这些向量还携带了大量的语义及句法信息，一定程度上解决了特征之间无关联的问题。

本文的组织结构如下：第 2 部分对相关的框架排歧任务进行介绍，并总结了传统方法对该任务的局限；第 3 部分提出了基于距离和基于词矩阵相似度的框架排歧的模型；第 4 部分叙述了实验语料、使用的分布式模型，以及一些实验设置；第 5 部分给出了实验结果及相应的错误分析；第 6 部分与 BASELINE 做了对比，进行了 t 检验和 Kappa 检验，并对错误进行分析；最后，对全文进行了总结，并给出了进一步的研究方向。

## 2 相关工作

针对英文 FrameNet 的框架排歧的研究，一些传统的模型是基于条件随机场、支持向量机、最大熵等分类器建立模型，把框架识别看作多分类问题<sup>[5,6]</sup>。针对汉语框架，李茹<sup>[7]</sup>等提出基于依存分析的条件随机场模型进行汉语框架识别；李国臣<sup>[8]</sup>等研究了基于词元语义特征的汉语框架语义排歧方法。这些传统的方法，不可避免的选择了大量的词和句法特征，使得特征空间维数很大，并且特征之间关联联系较小。党帅兵<sup>[9]</sup>将词语的分布式表征信息，加入到最大熵分类模型中，初步验证了词分布式表征的有效性。本文首次尝试不使用传统的分类模型，直接使用词语的分布式表征，以避免特征选择及降低特征维数，并结合特征之间的关联性，进行框架排歧。

近年来，词语的分布式表示技术受到很多自然语言处理研究者的青睐。分布式语义模型可以从大规模无标注语料中自动学习到句法和语义信息，有针对词语、句子以及文档的分布。Karl Moritz 首次提出将词语的分布式表示应用在词义消歧任务(WSD)<sup>[10]</sup>，词义消歧任务与框架排歧任务的目标相同：根据语境，为目标词选取一个合适的词义。所不同的是这里的词义是框架语义，而不是传统的词的词典义项。再者大部分传统的 WSD 任务只针对名词<sup>[11]</sup>，而框架消歧任务主要是对动词。近年来，有学者将分布式表示技术引入到框架语义分析的任务上。Hermann<sup>[12]</sup>提出基于分布式表征的框架排歧方法，在英文的 FrameNet 语料上取得较好的结果。本文研究的是使用词语或句子的分布式对汉语框架排歧任务的有效性。

## 3 汉语框架排歧模型

### 3.1 汉语框架排歧任务

在一条例句中，给定一个可以激起多个框架的目标词，要求计算机能够基于上下文语境，

从现有的框架库中，自动地为该目标词选择一个适合的框架。形式化表述如下：给定一个句子，记为  $S$ 。 $S$  是由词组成的一个序列，记为  $S=(w_1, w_2, \dots, w_n)$ ， $w_i$  代表组成句子的第  $i$  个词， $1 \leq i \leq n$ 。且目标词可以激起的框架集合记为  $F=\{f_1, f_2, \dots, f_m\}$ 。本文首次使用高斯判别分析来解决汉语框架排歧任务。框架消歧任务认为可以描述为：寻找唯一的  $f$ ，使其满足：

$$f = \arg \max_{f \in F} p(s|f)p(f) \quad (1)$$

本文提出了两种排歧模型：一种假设存在框架向量及其所属的例句向量，基于高斯判别分析，为例句选取合适的框架，具体做法参见 3.2；另一种则假设不存在框架向量和例句向量，利用句子之间的相似度直接选取合适的框架，具体做法参见 3.3。

### 3. 2 基于距离的汉语框架排歧模型

本模型假设存在框架向量，以及框架所属的例句向量。框架代表着一种语义场景，描述某语义场景的例句就属于表示该语义场景的框架；而每条例句又是由词构成。基于以上两点，我们认为：词向量、框架下的例句向量以及框架向量存在于同一空间中。如何表示框架下的例句向量和框架向量，以及如何判别例句的所属框架是本节讨论的内容。

#### 3. 2. 1 框架例句以及框架的分布式表示

框架例句是由词语组成的，可以通过词向量来表示例句向量。另外，也可以直接训练例句向量。因此，本节使用了两种表示方法：一种基于词向量的 MEAN-POOLING，另一种是 Doc2vec。

##### 1. 词向量的 MEAN-POOLING

表 1 构建例句向量的 MEAN-POOLING 方法

MEAN-POOLING
<b>Input:</b> $S = \{w_1, w_2, \dots, w_m\}$
<b>Output:</b> $E(S)$
1. 将例句分词，去掉标点及虚词，得到: $S' = \{w_1, w_2, \dots, w_n\}$
2. 在 Word Embedding 库中，查找例句 $S'$ 中对应词语的词向量: $E(w_1) \in R^d, E(w_2) \in R^d, \dots, E(w_n) \in R^d$
3. $E(S) = (E(w_1) + E(w_2) + \dots + E(w_n)) / n, E(S) \in R^d$

2、Doc2vec。对于构建例句向量，我们使用了 Mikolov 在 Word2vec 的原理上提出了一种将短句变为向量的方法<sup>[13]</sup>。本文使用 PV-DBOW 算法，在中文维基百科上训练例句向量。

由于每个框架下的例句有一定的语义关联性。因而，我们将框架下的例句向量相加后再平均，得到框架向量的表示。

#### 3. 2. 2 框架排歧步骤

首先用词向量把例句  $s$  表示为例句向量，并假设框架下例句的分布服从正态分布，则  $P(f|S)$  服从多维正态分布。这个分布中的两个参数： $\mu$  是均值向量，即上节所提出的框架向量， $\Sigma$  是协方差，在本文中，我们假设是协方差相同。因而有：

$$P(S; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(S - \mu)^T \Sigma^{-1} (S - \mu)\right) \quad (2)$$

所求概率的大小与框架例句向量和框架向量之间的欧式距离呈负相关。因而，将模型简化，我们没有计算概率而是直接通过距离大小来判别例句所属的框架。具体步骤参见下表：

表 2 基于距离的框架排歧步骤

框架排歧步骤
<b>Input:</b> $S = \{w_1, w_2, \dots, w_n\}, w_i \in F, F = \{f_1, f_2, \dots, f_m\}$
<b>Output:</b> $f$
1. 将某词元下的例句分为训练集和测试集;
2. 按上节描述的 mean-pooling 或 Doc2vec 方法, 将所有的例句转换为例句向量, 并用训练集的例句向量构建框架向量;
3. 利用 $Distance(S, f_i) = \sqrt{(E(S) - E(f_i))(E(S) - E(f_i))^T}$ , 计算测试集中的例句向量与每个框架向量欧氏距离;
4. $f = \underset{f_i}{MIN}(\underset{i=1}{\cup} Distance(S, f_i))$ 。

### 3. 3 基于词相似度矩阵的汉语框架排歧模型

本模型假设不存在框架向量, 以及其所属的例句向量。因而对于框架和例句, 没有特定的向量表示。框架代表着一个语义场景, 那么可以描述该语义场景的例句均属于该框架, 因而框架例句之间有一定而定语义联系及相似性。此外, 每条例句都是由词语构成。而词向量包含着的句子的语义和句法信息。基于以上两点, 我们直接使用例句中词语相似度来度量例句的相似度, 最后利用例句间的相似度来判别例句所属的框架。

#### 3. 2. 1 例句间的词相似度矩阵

有两条例句, 训练集例句  $S_1 = \{w_1, w_2, \dots, w_m\}$  测试集例句  $S_2 = \{w'_1, w'_2, \dots, w'_n\}$  在 Word Embedding 库中查找相应的词向量, 构建词余弦相似度矩阵:

$$\cos\langle s_1, s_2 \rangle = \begin{pmatrix} \cos\langle w_1, w'_1 \rangle & \cos\langle w_1, w'_2 \rangle & \dots & \cos\langle w_1, w'_n \rangle \\ \cos\langle w_2, w'_1 \rangle & \cos\langle w_2, w'_2 \rangle & \dots & \cos\langle w_2, w'_n \rangle \\ \dots & \dots & \dots & \dots \\ \cos\langle w_m, w'_1 \rangle & \cos\langle w_m, w'_2 \rangle & \dots & \cos\langle w_m, w'_n \rangle \end{pmatrix} \quad (3)$$

矩阵元素为两个词语的余弦相似度, 即  $\cos\langle w_m, w'_n \rangle = \frac{E(w_m)E(w'_n)^T}{|E(w_m)||E(w'_n)|}$ ,  $E(w)$  为词  $w$  的词向量。矩阵

中行代表测试例句的每一个词与训练例句中的每一个词的相似度, 矩阵中的列代表训练例句中的每个词与测试例句的每个词的余弦相似度。基于上述词相似度矩阵, 本文提出三种计算例句相似度的方法:

表 3 三种例句相似度的测量方法

	相似度矩阵每个值的相加平均 Mean	相似度矩阵的行最大, 列最小值 Min<Max>	相似度矩阵的行最小, 列最大值 Max<Min>
$\text{sim}\langle s_1, s_2 \rangle$	$\frac{\sum_{i=1}^m \sum_{j=1}^n \cos\langle w_i, w'_j \rangle}{m \times n}$	$\underset{i=1}{MIN}(\underset{j=1}{MAX} \cos\langle w_i, w'_j \rangle)$	$\underset{i=1}{MAX}(\underset{j=1}{MIN} \cos\langle w_i, w'_j \rangle)$

### 3.2.2 框架排歧步骤

基于词相似度矩阵的模型，只利用词的相似度，得到框架下例句的相似度，根据例句相似度进而判别其所属框架。

表4 基于词相似度矩阵的框架排歧步骤

框架排歧步骤
<b>Input:</b> $S = \{w_1, w_2, \dots, w_n\}, w_i \in F, F = \{f_1, f_2, \dots, f_m\}$
<b>Output:</b> $f$
1. 将某词元下的例句分为训练集和测试集；
2. 利用上节描述的相似度计算方法，对每一条测试集的例句 $S_i, i \in \{1, 2, \dots, n\}$ 与训练集中每条例句 $S_j, j \in \{1, 2, \dots, m\}$ ，进行相似度计算，得到： $\cos \langle s_i, s_j \rangle$
3. 计算例句与框架的相似度： $sim(s_i, f_k) = (\cos \langle s_i, s_1 \rangle + \cos \langle s_i, s_2 \rangle + \dots + \cos \langle s_i, s_m \rangle) / m$
4. $f = \bigcup_{j=1}^m MAX(sim(s_i, f_j))$

## 4 实验设置

### 4.1 实验语料

本文抽取汉语框架网(Chinese FrameNet)中 88 个可以激起两个以上框架的词元中的 2067 条句子作为语料来构造框架排歧实验。采用中文维基百科<sup>2</sup>做为训练词向量和句向量的语料。Word Embedding 对本语料的覆盖率为 94.58%，具体计算公式如下：

$$\text{覆盖率} = \text{词频向量}^3 * \text{词语出现向量}^4 / \text{语料总词数}$$

### 4.2 评价指标

为了评价本文所提模型的性能，我们采用组块 3×2 交叉验证进行实验。具体做法是，将语料库切分成 4 个大小相同的子集，然后，通过两两组合，形成组块 3×2 份交叉验证实验。组块 3×2 交叉验证在模型估计和选择的优良性能已经被证明，具体可参考 Wang 等的工作<sup>[14]</sup>。

在组块 3×2 交叉验证的条件下，全部目标词的框架分类准确率(Accuracy)的计算公式如下：

$$Accuracy = \frac{\sum_{k=1}^3 \sum_{j=1}^2 \sum_{i=1}^{88} C_{kji}}{\sum_{i=1}^{88} N_{kji}} \quad (4)$$

式中，n 为所选用词元的总数(本文 n=88)， $N_{kji}$  为第 i 个词元的第 j 组块的第 k 折交叉实验中的测试例句总数， $C_{kji}$  为第 i 个词元的第 j 组块的第 k 折交叉实验中的测试例句正确数。

### 4.3 参数设置

在训练词向量时，我们分别采用了 Word2vec 的 CBOW 以及 Skip-Gram 模型，以及 GloVe 型。按之前的经验将各模型的窗口设置为 5，维度设置为 100。分别探究每个模型训练的词向量对框架排歧任务的影响，寻找最优词向量模型。然后设置该模型窗口值及维度，探究不同窗口及维度，对框架排歧模型准确率的影响。

Mu 提出一种词向量后处理的方法<sup>[15]</sup>，该文认为训练好的词向量中都包含着一个公共向

<sup>2</sup> 使用 zhwiki-latest-pages-articles.xml.bz2(2017-1-25)

<sup>3</sup> 语料库中的某词语，在语料库中出现的次数。

<sup>4</sup> 语料库中的词是否在 Word Embedding 库中出现，若该词出现，向量的相应维度为 1，未出现为 0。



量，而且词向量均被相同的方向所支配。这两点会影响词训练出词向量的质量。本文将用这种后处理方法，对 CBOW，Skip-Gram 以及 GloVe 训练出的词向量进行处理。

另外，本文还增加了随机初始化的词向量方法。设置了服从标准正态分布的随机向量。使用该随机向量，是为了与上述词向量模型进行比较。

## 5 实验分析

窗口值为 5，向量维度为 100 时，各类词向量在框架排歧模型中的准确率：

表 5 各类词向量应用于框架排歧模型中的准确率

	Mean-Pooling	Mean	Max<min>	Min<max>	Sentence2vec
CBOW	64.95%	35.65%	36.68%	41.91%	-
后处理的 CBOW	<b>65.71%</b>	35.70%	38.74%	41.88%	-
SGNS	64.07%	35.98%	38.79%	42.31%	-
后处理的 SGNS	64.69%	36.67%	38.53%	42.38%	-
GLOVE	60.20%	40.08%	42.88%	42.83%	-
后处理的 GloVe	62.73%	37.36%	42.87%	42.41%	-
随机向量	42.37%	47.66%	45.29%	46.35%	-
	-	-	-	-	43.97%

接下来，对表 5 的数据进行说明。第一行代表本文实验的所使用的模型方法，第一列代表的是本文实验所用的词向量。其中：“CBOW，SGNS，GLOVE”是分别使用 Continuous Bag-of-Word 模型，Skip-Gram 模型以及 GloVe 模型生成的词向量；“后处理的 CBOW、后处理的 SGNS、后处理的 GLOVE”是使用文献[15]的方法对上述三类词向量进行后处理的词向量；“随机向量”是随机生成均值为 0，方差为 1 的 100 维向量。当采用 Sentence2vec 模型时，输入是句向量，在相应词向量准确率的位置为“-”。

由表 5 可得，用 CBOW 模型训练出的词向量，进行后处理，得到的准确率最高。下面以该词向量作为研究对象。研究当窗口大小不同时，对于框架排歧任务的影响，准确率如下图所示：

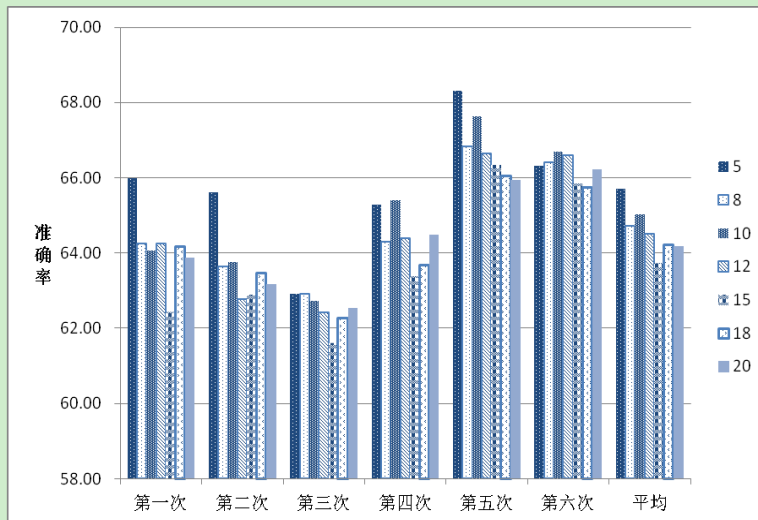


图 1 词向量窗口不同时，框架排歧的准确率

由图 1 可以得出，当窗口设置为 5 时，效果最好。因而，我们选定窗口值为 5。然后比较不同的维度的词向量对于框架排歧任务的影响。具体结果如下：

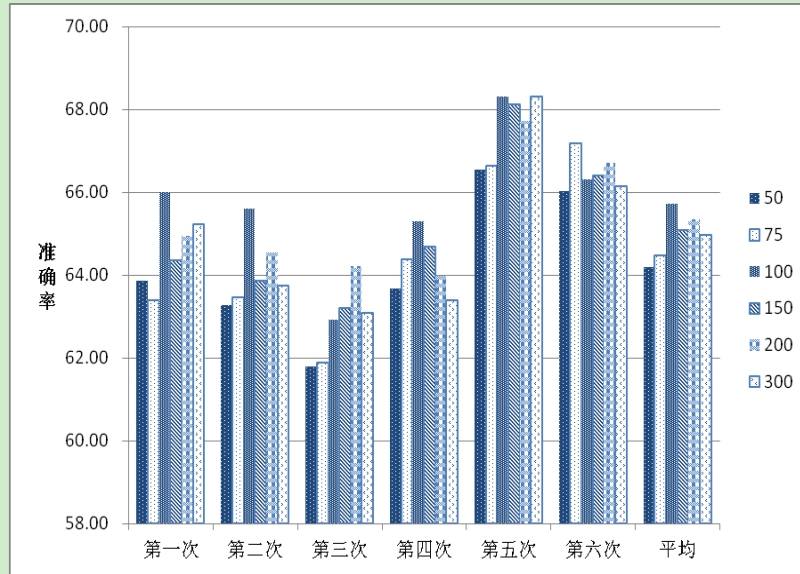


图 2 词向量维度不同时，框架排歧的准确率

综合表 5，图 1，图 2，使用 CBOW 模型，将窗口设置为 5，向量维度设置为 100，并经过后处理的词向量，对框架排歧的效果是最好的。

## 6 错误分析

我们选取高<sup>[4]</sup>的实验做为 BASELINE。由于本文的模型是基于词的分布表征，BASELINE 只选取有关的词特征，不选取句法特征。所选特征为：词包以及窗口为 2 的词特征。即 BASELINE: WordFeature[-2, 2]+BOW。具体结果如下表所示：

表 6 BASELINE 与 MEAN-POOLING 的比较

	第一次		第二次		第三次	
	第一折	第二折	第一折	第二折	第一折	第二折
BASELINE	64.59%	63.78%	67.65%	67.81%	67.74%	67.44%
MEAN-POOLING	62.91%	65.29%	66.32%	68.31%	65.60%	65.99%

### 6.1 T 检验

为了检验我们的方法与 BASELINE 是否有显著性的差异，本文对上述两种方法，全部的 6 次实验的准确率进行了 T 检验。我们的原假设  $H_0: \hat{\mu}_{3 \times 2} = \mu_0$ ，备选假设  $H_1: \hat{\mu}_{3 \times 2} \neq \mu_0$ ，显著性水平  $\alpha = 0.05$ 。其中， $\hat{\mu}_{3 \times 2} = \frac{1}{6} \sum_{i=1}^6 (BASELINE\_ACC_i - MEAN\_POOLING\_ACC_i)$ ， $BASELINE\_ACC_i$  为 BASELINE 第  $i$  次的准确率， $MEAN\_POOLING\_ACC$  为 MEAN-POOLING 第  $i$  次的准确率， $\mu_0 = 0$ 。

$$t_{B3 \times 2cv} = \frac{\hat{\mu}_{3 \times 2} - \mu_0}{\sqrt{\hat{VAR}(\mu_{3 \times 2})}} : t_5 \quad (5)$$

式 (5) 中,  $VAR(\hat{\mu}_{3c2}) = \frac{1}{6} \sum_i \sum_j (\hat{\mu}_i^{(j)} - \hat{\mu}_{3c2})^2$ , 具体推导步骤参见文献 [12],  $\hat{\mu}_i^{(j)}$  为每次实验中 BASELINE 与本文方法得到精确率的差值。经计算, 得到  $t_{B3 \times 2cv} = 0.0204$ 。由于  $t_{B3 \times 2cv} < t_5$ , 概率  $P > 0.05$ 。按  $\alpha = 0.05$ , 接受  $H_0$ , 拒绝  $H_1$ 。使用 MEAN-POOLING 和 BASELINE 进行框架排歧, 实验结果的差异不显著。

## 6.2 Kappa 统计量分析

虽然使用两种方法, 得到实验的结果差异并不显著。但使用 MEAN-POOLING 方法, 实验结果的精确度略低于 BASELINE, 本节对两种方法的实验结果进行 Kappa 统计量分析。对结果一致的部分进行检验, 看一致的部分是否是由偶然因素影响的结果。计算公式如下:

$$K = \frac{P_0 - P_e}{1 - P_e} \quad (6)$$

其中,  $P_0$  为实际一致率,  $P_e$  为理论一致率。具体结果, 参见下表:

表 7 每次实验的 Kappa 值

	第一次	第二次	第三次	第四次	第五次	第六次
kappa	0.51	0.49	0.42	0.43	0.41	0.48

对以上 6 次实验做了 Kappa 统计量分析, 得到结果如上图所示。每次实验的 Kappa 值均在 [0.4, 0.75] 之间。说明两者的一致性可以接受, 表明 MEAN-POOLING 的结果与 BASELINE 的结果中, 一致的部分受偶然因素影响不大。

最后, 我们收集了两类例句: 一、BASELINE 分类正确而 MEAN-POOLING 分类错误, 二、BASELINE 分类错误而 MEAN-POOLING 分类正确。通过分析得到如下结论:

1、对于目标词附近多为代词, 介词以及附近词缺失的例句, BASELINE 往往表现的不好。例如: “分散范围之广, 分布地区之多, 都是其他民族所不能相比的”。目标词是“分散”, 其左边没有词; 以及“她记得离开中国时, 养母带她登过长城。”目标词是“记得”, 左边只有代词“她”。对于那些附近的词缺失或者多为无实际意义的词的例句, MEAN-POOLING 不需要此类特定的特征, 因而可以有效的解决这类句子的分类问题。相反, 目标词附近多为实际的词, BASELINE 要优于 MEAN-POOLING。MEAN-POOLING 的精确率略低, 有可能是在语料库中, 适合于 BASELINE 分类的例句数量占优。

2、若某词元下, 某个框架的训练语料数量偏多。相比于 MEAN-POOLING, BASELINE 受影响较大。例如, 词元“看”有三个框架, 分别为“获知”、“自主感知”、“外观”。其中“获知”框架下的语料例句最少, 而 MEAN-POOLING 正确分类“获知”框架下的例句, 要多于 BASELINE。对于例句较少的框架, BASELINE 由于训练例句较少, 往往表现的不好。

## 7 总结与展望

本文使用 CBOW、Skip-Gram 和 GloVe 等流行的词向量训练模型, 以及使用 Doc2vec 模型, 从无标注的维基百科语料中学习词语及句子的分布式表征, 后直接应用于基于距离和基于词相似度矩阵的框架排歧模型中, 在基于距离的框架排歧模型中, 最高得到了 65.71% 的准确



率。

在表现最好的 MEAN-POOLING 方法中, 使用随机词向量进行实验, 得到的准确率大大低于后处理的 CBOW 词向量。因而, 进一步验证了: 从大量文本中, 通过无监督学习, 得到的词向量, 携带了大量的语法及语义信息, 可以有效地应用在框架排歧任务中。

因为精确率略低于 BASELINE, 随后又使用了 T 检验和 Kappa 检验, 对错误进行分析。可以发现将词的分布式表征应用于框架排歧模型, 与 BASELINE 的实验结果并没有显著的差别, 两者一致的部分受偶然因素影响结果不大, 一致性可以接受。证明了本文方法的有效性。

文章的最后, 通过分析两类例句。我们得到: 词的分布式表征可以有效的避免传统分类模型严重依赖于特征选择的缺陷; 当训练语料不平衡时, 基于词分布式表征的排歧模型受影响较小。本文只是用了词的分布式表征, 并没有考虑句子依存关系等句法特征。下一步的工作过主要是研究如何将依存关系与词的分布式表征有效的结合, 应用于汉语框架排歧模型中。

## 参考文献:

- [1] C. Baker, M. Ellsworth, and K. Erk. SemEval-2007 Task 19: Frame Semantic Structure Extraction[C]//Proceedings of The 4th International Workshop on Semantic Evaluations. Prague, 2007:99-104
- [2] C. J. Fillmore. Frame Semantics[J]. In Linguistics in The Morning Calm, Hanshin Publishing Co.. Seoul, South Korea. 1982: 111-137.
- [3] 王瑞波, 李济洪, 李国臣, 杨耀文, 等. 基于 Dropout 正则化的汉语框架语义角色识别[J] 中文信息学报, 2017, 31(1):147-154
- [4] 李济洪, 高亚慧, 王瑞波, 等. 汉语框架识别中的歧义消解[J]. 中文信息学报, 2011, 25(3):38-44
- [5] Cosmin Adrian Bejan, HaThaway Chris . UTD-SRL: A Pipeline Architecture for Extracting Frame Semantic Structures[C] . //In 45th annual meeting of Association for Computational Linguistics, 2007:460-463.
- [6] Richard Johansson, Nugues Pierre. LTH: Semantic Structure Extraction using Nonprojective Dependency Trees[C] //Proceedings of The 4th International Work on Semantic Evaluations. Prague, 2007:227-230.
- [7] Ru Li, Haijing Liu, Shuanghong Li. Chinese Frame Identification using T-CRF Model[C] //Proceedings of International Conference on Computational Linguistics. Beijing, 2010: 674-682
- [8] 李国臣, 张立凡, 李茹, 等. 基于词元语义特征的汉语框架排歧研究[J]. 中文信息学报, 2013, 27(4):44-51
- [9] 党帅兵, 李国臣, 王瑞波, 等. 基于词分布表征的汉语框架排歧研究[J]. 中北大学学报, 2015, 36(3): 328-332, 337.
- [10] Das D, Ganchev K, Weston J, KM. Hermann Semantic frame identification with distributed word representations[C]. Meeting of The Association for Computational Linguistics, 2016 [C]. Meeting of The Association for Computational Linguistics, 2016
- [11] Navigli, R. 2009. Word Sense Disambiguation: A Survey[J]. ACM Computing Survey. 41, 2(Feb. 2009), 1-69.
- [12] Karl Moritz Hermann, Dipanjan Das, Jason Weston Kuzman Ganchev. Semantic Frame Identification with Distributed Word Representations[C]//Meeting of The Association for Computational Linguistics . Baltimore, USA. 2014:1448-1458
- [13] QVLe, T Mikolov Distributed Representations of Sentences and Documents[J] Computer Science, 2014, 4:1188-1196
- [14] Yu W, Ruibo W, Huichen J, et al. Blocked 3x2 cross-validated T-Test for comparing supervised classification learning algorithms[J]. Neural computation, 2014, 26(1): 208-23
- [15] Mu J, Bhat S, Viswanath P. All-but-the-Top: Simple and Effective Postprocessing for Word Representations

作者联系方式：张力文, 山西省太原市坞城路 92 号(山西大学), 邮编: 030006,

电话: 15634986915 E-MAIL: [505646231@qq.com](mailto:505646231@qq.com)

通讯作者：王瑞波, 山西大学, E-MAIL: [wangruibo@sxu.edu.cn](mailto:wangruibo@sxu.edu.cn)