

文章编号: 1003-0077 (2011) 00-0000-00

基于多维语义关系的谐音双关语识别模型*

徐琳宏¹, 林鸿飞², 祁瑞华¹, 杨亮²

(1.大连外国语大学, 辽宁省 大连市 116044; 2.大连理工大学 辽宁省 大连市 116024)

摘要: 谐音双关语的识别是幽默研究领域的一个重要分支, 并逐渐发展为一个新兴的研究领域。本文提出一种基于四个维度特征集的谐音双关语识别模型, 其中四个维度包括语义透明度、语义相关度、语音扩展性和句式特征集。语义透明度包括词项统计和语句字符长度两个特征, 句式特征集包括人名、大写、时态、词性和位置五个特征。将这四个维度的九个特征加入到二叉判定树中, 使用 K-Means 聚类获取阈值, 完成双关语的识别。本文的实验数据来自于 SemEval2017 任务 7 的语料, 取得了较好的效果, F1 值高于参赛队中的第一名, 实验证明基于四个维度特征的二叉判定树分类方法在谐音双关语识别中是有效的, 且在多个特征中, 语音扩展性和句式特征集的效果比较明显, 这也符合谐音双关语识别中语音作用较大的预测。

关键词: 谐音双关语; 二叉判定树; 语义特征集; 聚类

中图分类号: TP391

文献标识码: A

The Heterographic Pun Identification Model Based on Multi-dimensional Semantic Relationships

XU Lin-hong¹, LIN Hong-fei², QI Rui-hua¹, YANG Liang²

(1. Unit, city, province zip code, China; 2. Unit, city, province zip code, China)

Abstract: The identification of heterographic puns, as an important branch of humor research, has gradually developed into a new research area. This paper presents a heterographic pun identification mechanism based on feature sets in four dimensions, namely, semantic transparency, semantic relevance, phonetic expansibility and syntax feature sets. Semantic transparency feature sets consist of lexical item statistics and character length; syntax feature sets include names, capitalization, tense, part of speech and location. Nine features of the above four dimensions are added to the Binary Decision Tree, with the help of K-Means Cluster, to generate the threshold and complete the identification of a pun. By utilizing the corpus of SemEval2017 Task 7, this paper achieves satisfactory results and its F1 value outcores the top one of all the participating teams, and the experiment proves that the taxonomic approach of the Binary Decision Tree based on four dimensions are effective in identifying heterographic puns. Phonetic expansibility and syntax feature sets, among all, are particularly effective, which is consistent with our presumption that the phonetic feature plays a bigger role in the identification of heterographic puns.

Key words: Heterographic Pun; Binary Decision Tree; Semantic Feature Set; Cluster

1 引言

英文中双关语是一种比较常用的修辞手段, 一个双关句通常采用词义模糊的方式, 使句子产生一个以上的含义^[1]。也就是说作者特意让一个词在句中包含两个或者两个以上的独立含义, 这在许多英文著作中被广泛应用, 例如, 莎士比亚的作品集中经常使用

* 收稿日期: 2017-06-10

定稿日期: 2017-07-18

基金项目: 国家自然科学基金重点项目(61632011); 国家社会科学基金一般项目(15BYY028); 辽宁省自然科学基金(20170540230, 2015020017, 20170540232); 辽宁省优秀人才项目(LJQ2014127)

作者简介: 徐琳宏(1979—), 硕士, 讲师, 主要研究领域为文本情感计算。林鸿飞(1962—), 博士, 教授, 主要研究领域为文本挖掘和信息检索; 祁瑞华(1974—), 博士, 教授, 主要研究领域为自然语言处理; 杨亮(1987—), 博士, 讲师, 主要研究方向为文本情感计算

双关语修辞。此外，双关语通常带有一定的幽默色彩，在人们的日常交际中有较强的交互价值，能更快地增进双方的感情。从类别上来说，双关语一般分为语义双关（Homographic Pun）和谐音双关（Heterographic Pun）。它们在许多领域都有重要的研究价值，如广告语的生成。双关语中人们不仅仅创造幽默，而且会引发听众联想双关语中的目标词。识别出词语的模糊性，从而有益于文本情感识别。

谐音双关是双关语中的一个重要类别，它是指语句中某个词在语音上与目标词相同或相近，用目标词替换后原句产生了新的含义，例如：“Diets are for people who are thick and tired of it all”，其中“thick”为双关键词，它的目标词为“sick”。找到双关键词对应的目标词，才能更好地理解语句的含义，从而体会语句的幽默色彩，可见双关语的理解考验着接受者的智慧。那么，对于机器而言，双关语的理解更是一个难题。

双关语识别和生成是幽默的理解的一个研究分支，在人机交互、问答系统等方面有重要的应用。对幽默的理解，有助于提高人与机器交互的水平，使人与机器之间的沟通变得更加接近人与人之间的沟通。2017年国际语义评测（Semantic Evaluation）发布两个新的任务是关于幽默和双关语的，这也更好地推动了机器对幽默理解力的研究。双关语作为一个重要的幽默来源，是近几年的研究热点，国内外有许多这方面的研究工作。

2 相关工作

谐音双关语是双关语的一种特例，而双关语中很大一部分具有幽默的特性，所以双关语是幽默研究的一个重要分支。三者研究方法上既有各自的特点，也有很多相通的方法和特性。下面从幽默、双关语和谐音双关三个方面介绍目前国内外主要的研究工作。

2.1 幽默的相关研究

很多双关语都有不同程度的幽默感，幽默作为一种语言形式，能够为对话双方营造轻松、愉快的气氛，因而具有独特的交际价值，备受人们青睐。近些年来，自然语言处理中也有很多关于幽默识别方面的研究工作。

英文幽默方面的研究工作比其他语言更多，2009年和2010年Julia M.等人收集幽默文本，结合幽默经典理论，逐句分析幽默语句的特点^[2,3]。2012年，Antonio等人利用模糊性、情感极性和情感场景等特征识别tweet文本的幽默性^[4]。2012年，Yishay Raz等人在tweet文本中，采用模式、词汇、形态和语音等特征识别幽默文本^[5]。2012年和2016年Pascale Fung等人标注了《生活大爆炸》和《宋飞正传》语料集，分别使用CNN、CRF和LSTM分类方法在语音和文本两大类特征基础上识别幽默^[6,7,8]。2014年，Renxian Zhang等人利用GBRT模型识别幽默的tweet文本^[9]。以上是英文幽默的研究工作，其他语言也有一些幽默研究的工作，重要包括：2007年和2009年Paolo Rosso等人使用Ngram完成意大利语幽默研究^[10,11]。2016年，Santiago Castro等人使用KNN和SVM识别西班牙文的幽默。除了以上幽默识别方面的工作，也有一些关于幽默生成的研究。2012年，Igor Labutov等人基于SSTH理论做幽默语句生成的研究，采用人工打分的方式评测生成语句的效果^[12]。2013年，Alessandro Valitutti等人通过词语替换，生成幽默文本，采用人工评估的方式评估幽默等级^[13]。

在幽默计算方面，国内研究者也进行了一定的探讨。2015年张冬瑜等人构建了情感隐喻语料库，这对于幽默的识别提供可以借鉴的方法^[14]。2016年林鸿飞等人回顾了幽默研究的发展历史，详细阐述了幽默计算理论的多种基本理论和应用，对于谐音幽默的处理也给出了相应的讨论^[15]。其他学者对于幽默计算所涉及的语言学理论和方法进行了相应的研究。但目前中文幽默识别还处于起步阶段，经过深入的研究和探索，未来还有很大的发展空间。

2.2 双关语的相关研究

双关语识别作为幽默研究的一个分支，在幽默特性的基础上，还包含一些句式构成和语义模糊等独有的特点。2008年，Hempelmann等人研究双关语目标词的自动识别，及幽默文本的生成^[16]。2005年Richie等在对笑话语言分析的基础上生成双关语^[17]。2009年，Hong等人自动抽取语义模板生成双关语^[18]。2011年，华盛顿大学的Chloe Kiddon等人将原领域的字面意义到目标领域的非字面意义的映射过程定位成一个隐喻

问题，提出委婉语和领域结构两个特征作为语义和文化背景识别双关语^[19]。2015年，Diyi Yang 等人在“pun of the day”的数据集上识别幽默和锚点，将冲突性、模糊性、人际交互和语音风格等特征加入到随机森林和最大衰减法中，完成双关语的识别^[20]。2015年，Tristan Miller 在传统 LESK 算法的基础上研究双关语的词义消歧，提出两个解绑策略^[21]。一是利用词汇的词性和语法特征消歧，二是将 WordNet 中词义聚类的方法。两个方法有效地解决了 LESK 算法词义相近较多的问题。2016年 Tristan Miller 采用一种计算模型检测双关语及特定的语义隔离（the isolation of the intended meanings），并集成 WSD-inspired 系统评估语义双关语^[1]。

2.3 谐音双关语

谐音双关语是双关语的一种类型，其双关词和目标词之间具有发音相同和相似的特性，因而找到发音相似的词对非常重要，这方面的相关研究很早就已经开展。

1986年 Zwicky 等人基于双关语实例，生成 2140 个双关语词对^[22]。1991年 Sobkowiak 在 Zwicky 工作基础上收集了 3850 个双关语和目标词的词对^[23]。1996年，Binsted 等人利用英语发音词典^[24]，制定发音相似的规则生成发音相似的词对，用于双关语生成^[25]。2003年 Hempelmann 等人在 Sobkowiak 研究基础上生成一种简单的语音替换代价表，用于评价生成后的双关语的分值^[26]。2015年，斯坦福大学 Justine T. Kao 等人提出模糊性（ambiguity）和独特性（distinctiveness）两个特征，以语言模型为基础识别幽默双关语^[27]。2016年，Aaron Jaech 等人完成谐音双关语的识别，提出的理解模型分为音素、词汇和文本三个层次，利用语言模型识别谐音双关语，并统计出谐音双关中元音和辅音相互替换的次数^[28]。

英文方面，关于幽默和双关语的研究大致分为三种，一是通过具体双关语实例，建立双关语和幽默的理论，二是利用机器学习方法 SVM、KNN 和深度学习等有监督方法识别双关语，三是完成双关语或幽默文本的生成。目前还很少有研究在无训练集条件下识别谐音双关语。

本文的主要贡献如下：（1）提出谐音双关语识别中的四个维度的特征集（2）在无训练集的条件下，采用二叉判定树识别双关语。文中第二节介绍了幽默和双关语识别方面的研究工作；第三节给出我们的识别模型，将四个维度的特征集加入到二叉判定树中识别谐音双关语；第四节介绍了数据集和相关语言模型；第五节列出了识别模型在 SemEval2017 语料中的实验结果；第六节总结了本文工作，并提出今后工作的设想。

3 谐音双关语的识别模型

谐音双关语识别是找到原句中某些词的同音词或者发音相似的词，放入原句后，语句产生另外一种更易理解的含义。所以理解谐音双关语，首先要通过语音扩展找到合适的目标词集合，然后将目标词逐个放入原句中进行语义匹配，匹配度较高的更可能为双关语。例如，前面给出了例子“Diets are for people who are thick and tired of it all”，双关词“thick”的相似音集合可能包含“sick”，“bick”和“fick”等，但是只有“sick”带回到原句中，与前面“diets”语义相关，与后面“of”形成“be sick of”常用搭配。所以本文从语义透明度、语音扩展性、语义相关度和句式特征集四个维度出发，最后采用二叉判定树识别谐音双关语。

3.1 谐音识别中的四个维度的特征集

3.1.1 语义透明度

所谓语义透明度是指整个句子的语义可以根据合成语句的多个词汇含义来推知的程度^[29,30]，其中包括单词的使用频率，搭配的熟悉程度和典型性等。如 Mok 认为词频越高，其语义透明度则越高^[31]。Pollatsek 以词频和语义透明度为变量，分析失语症患者语言产出的正确率^[32]。在双关语识别中，语义透明度也是一个重要的特征，例如“Psychiatrists like Kentucky Freud Chicken”和“Hotel owners usually have suite dreams”两句中“Freud Chicken”和“suite dreams”的搭配频率较低，更熟悉的搭配是“fried chicken”和“sweet

dreams”。为了解决上述问题，本文使用语句中 Unigram 值和语句长度 (senLength) 两个特征来表征语句的语义透明度。

Unigram 值代表语句中词汇的使用频率，使用频率越低，则语义的透明度越低。本文使用语言模型 (Language Model) 计算语句中词语及搭配的使用频率及熟悉程度。

语句长度值取语句中包含的字符个数，而不是单词的个数，因为单纯使用单词数，不能体现那些用词比较复杂 (单词较长)，但单词数较低的语句的复杂程度。字符个数越大，语句的透明度越低。

3.1.2 语义相关度

语义相关度是计算候选词在原句中的语义流畅程度，即候选词在原句中替换后语义越流畅，则该候选词与句子的相关程度越大，选择相关度最大的作为整个双关句的目标词。候选词与原句的语义相关度是通过 WordNet^[33] 的 Similarity 接口计算，先得到候选词与原句中每个词的相似度，然后选择最大的相似度作为该词的语义相关度。假设 C 为语音扩展词集合， $C = \{c_1, c_2, c_3 \dots c_m\}$ ，则通过公式 (1) 计算 c_i 的语义相关性：

$$\text{Relate}(c_i) = \text{MAX}_{k=1}^n (\text{Sim}(c_i, w_k)) \quad (1)$$

其中 w_k 为原句中第 k 个单词。那么整个语句 s 的语义相关度计算公式如式 (2) 所示：

$$\text{Relate}(s_i) = \text{MAX}_{i=1}^m (\text{Relate}(c_i)), \quad c_i \in C \quad (2)$$

其中 C 为语音扩展后的候选词列表，如果没有语音扩展词，则该句也没有语义相关度。一个语句的语义相关度越大，说明找到的词越可能是目标词，那么该句是双关语的可能性就更大。

3.1.3 语音扩展性

CMU 发音字典 (CMU Pronouncing Dictionary) 是卡内基梅隆大学研发，为语音合成器而设计^[34]。最近发布的版本中包含 134,000 个词条，每个词条的发音采用英文音素集 Arpabet，该音素集由 23 个元音和 31 个辅音组成。例如单词“Hello”发音为“HH AH0 L OW1”，其中数字表示重音。

发音词典中一个词可以包含多个音标，本文在做语音扩展时使用一个词条的所有发音，与任何一个发音匹配都作为该词的语音扩展词。语音扩展词分为两类：发音相同和发音相似。发音相同的词是直接在 CMUdict 查找与该词条任一发音相同的词汇，而寻找发音相似的词汇，需要统计双关语中原词与目标词之间的发音替换规律，我们采用 Aaron Jaech 等人^[28]给出的元音与元音以及辅音与辅音之间的替换矩阵。经过音素替换后，如果在词典中能找到该发音，则作为原词的发音扩展词。

3.1.4 句式特征集

本文采用以下五个特征作为谐音双关语判断的句式特征：

1. 人名：创建人名字典，判断语句中的名词是否为人名，英语双关语中很多存在人名，尤其是“Tom”出现频率较高，用于指代某个人。
2. 全大写：判断语句中是否有单词是全大写的，大写单词一般表示强调或者缩写，表达语句的重要含义，大写的缩写词也容易产生语音扩展。
3. 过去时态：判断一个语句是否包含过去时态。一个句子可能有多个分句，如果有一个分句使用了过去时态，则判定整个语句包含过去时。过去时态一般代表过去发生的动作或状态，而双关语中有一大部分是对过去看见和听见的事情叙述时产生的。
4. 词性特征：在判断人名和大写单词时，选择名词、动词、形容词和副词作为双关词的候选词。
5. 位置特征：一般双关句的双关词都出现在语句的后部，所以完成语音扩展和语义计

算时为语句后半部分的单词分配更高的权重。

前三个为布尔类型，直接用于识别双关语，后两个是为语音扩展或人名和时态服务的。五个特征都是从谐音双关语的句式出发，找到用词、造句以及语法上的一些特征，各种句式特征都可能帮助识别谐音双关语。

3.2 识别算法

本文利用上述四维特征，采用二叉判定树算法识别谐音双关语，使用 SemEval2017 任务 7 的语料，该任务主要是完成双关语的识别。在语料中只提供测试集，没有训练集，通常的有监督分类方法不适用。所以我们采用二叉判定树算法^[35]识别语句是否为双关语。判定树中的每个非叶子结点都包含一个条件，因此对应着一次识别。每个叶子结点对应着种类，因为本文是双关语的二义分类，所以叶子结点分为两类“双关句”和“非双关句”。从上述四个维度出发，每个维度对应一个非叶子结点，识别过程如图 1。

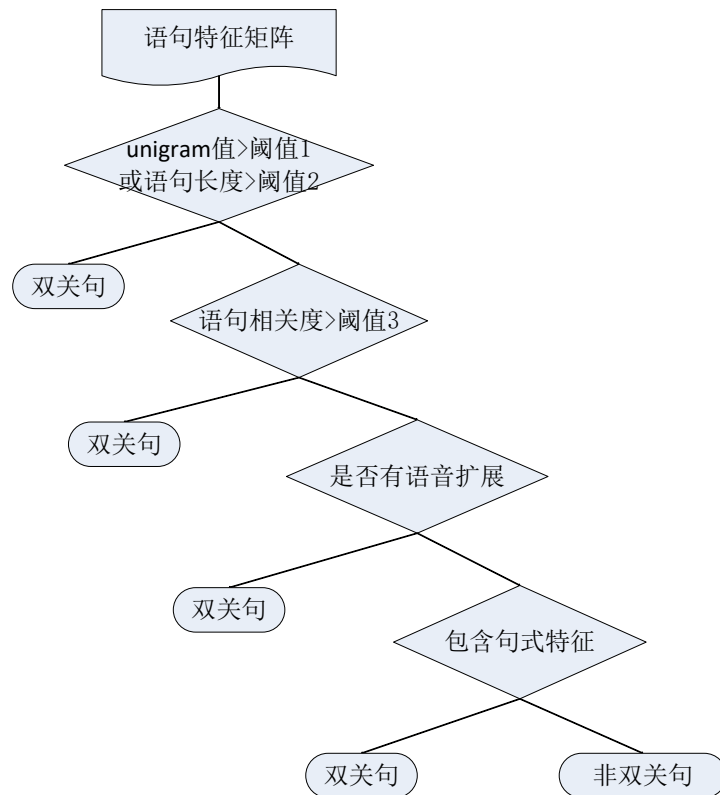


图 1 二叉判定树识别双关语算法

算法的输入是句子 S_i ，首先根据语句 S_i 中各单词计算特征矩阵 M_{ij} ， i 为语句编号， j 代表上述四个维度中的所有特征，然后将特征矩阵 M_{ij} 输入到二叉判定树中，每个非叶子结点根据不同的阈值标准判断当前句是否为双关句，如果是，算法结束。如果不是，则继续到下个结点判断，依次类推。如果所有的特征都不满足，说明不具备成为双关语的任何特征，那么判定为非双关句。从图 1 算法可以看出，语句透明度和语义相关度需要根据阈值识别双关句，较小的阈值会导致召回率较高，错误地将所有的样本分为正例。为了以更合理的方式计算阈值，本文分别将 Unigram 值，语句长度和语义相关度三个特征矩阵，用 K-Means 方法聚类。K-Means 算法是按照最邻近原则把待分类样本点分到各个簇，质心 (centroid) 是指各个类别的中心位置，质心的维数等同于单条数据的维数。针对每个特征矩阵，我们聚类为双关语和非双关语两类，得到左右两个质心。统一选择右质心作为阈值，因为左质心是非双关语的中心，右质心是双关语的中心，这样可以有效防止召回率太高，错误地将所有的样本分为正例，如果选择双关语类别的质心作为阈值。即使一个特征将正例分为负例，后面还有特

征可以矫正分类结果。本文中三类特征都是一维的，得到的质心是一个数值，可以直接作为阈值。

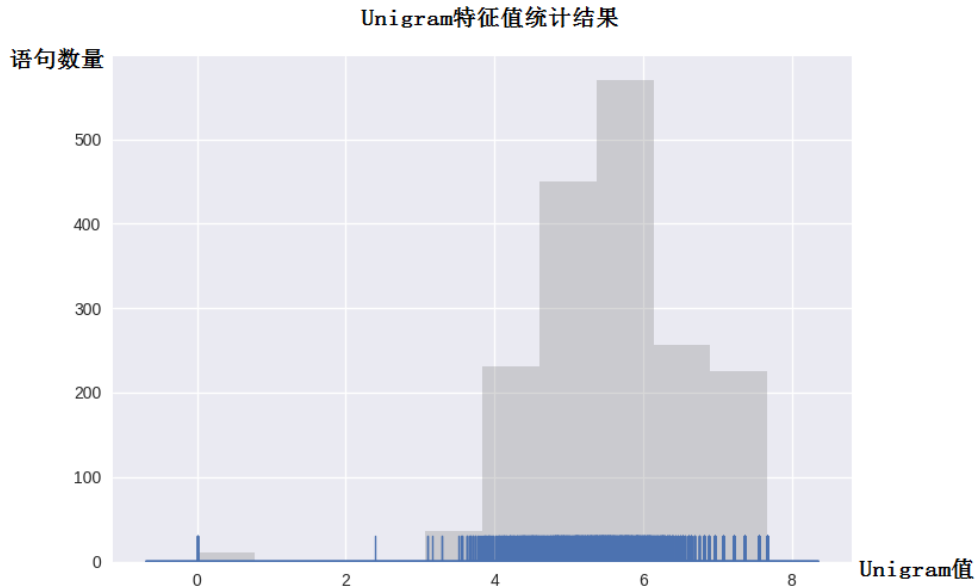


图 2 Unigram 值的语句密度

图 2 以 Unigram 值为例，统计了特征值与同一特征值下语句数量的关系。横坐标为所有可能的 Unigram 取值，纵坐标为语句数量。三个特征矩阵经过 K-Means 聚类后，选择双关语的质心作为该特征的阈值，因此 Unigram 特征、语句长度和语义相关度的阈值分别为：6.5、100 和 0.3。

4 数据集及语言模型

4.1 数据集

本文的实验语料来自 SemEval2017 子任务 7 的谐音双关语(Heterographic pun)^[36]。子任务 7 包含两类语料，语义双关语(Homographic puns)和谐音双关语(Heterographic pun)。评测为两类语料分别制定了三个子任务，子任务一是识别语句是否为双关语，子任务二是识别语句中哪个词为双关词，子任务三是识别双关词对应的目标词及词义。我们的实验选择谐音双关语的语料，完成子任务一，即识别语句是否为双关语。该语料只有测试集，没有训练集，因而不能采用有监督的分类方法完成。数据集采用 XML 格式，每个语句用一个 text 元素表示，为方便任务二和任务三，语句中每个单词用 word 元素表示。无论语句还是单词，都有一个唯一的标识。格式如下：

```
<text id="het_5">
  <word id="het_5_1">Are</word>
  <word id="het_5_2">evil</word>
  <word id="het_5_3">wildebeests</word>
  <word id="het_5_4">bad</word>
  <word id="het_5_5">gnus</word>
  <word id="het_5_6">?</word>
</text>
```

谐音双关语料中，共有 1780 个语句，其中双关语 1270 条，非双关句 510 条。

4.2 评价指标

本文采用准确率、召回率、精确率和 F1 值四项评估双关语的识别效果。如果把正例预测为正例的个数记作 TP，负例预测为正例的个数记作 FP，正例预测为负例的个数记作

FN, 负例预测为负例的个数记作 TN。那么
 准确率=(TP +TN) ÷ (TP+FP+TN+FN)
 召回率= TP ÷ (TP+FN)
 精确率=TP ÷ (TP+FP)
 F1 值=(2 ×精确率×召回率) ÷ (精确率+召回率)

4.3 语言模型

计算语义透明度主要是通过单词的 Unigram 值确定,下面简单介绍计算 Unigram 用到的语言模型 (Language Model)。语言模型是自然语言处理中被广泛应用的统计模型,目前主要采用 n 元语法模型(Ngram Model)。这种模型简单,相关的平滑技术也比较成熟。语言模型主要是构建一个字符串 s 的概率分布 p(s),反应一个字符串 s 出现的频率。假设 s 由单词 $w_1, w_2, w_3 \dots w_n$ 组成,这 P(s)的公式如式 (3)。

$$P(w_1, w_2, w_3, \dots, w_n) = \prod_{i=1}^n p(w_i | w_1 \dots w_{i-1}) \quad (3)$$

上述公式中如果句子太长,则概率会接近 0,由于参数过多,而不能正确训练。所以实际应用中多采用 n 元语法, n=1 时,即一元文法 (Unigram),认为每个词都是孤立存在的。n=2 时,认为第 i 个词 w_i 仅与他前面的一个词相关,成为二元文法 (Bigram),依次类推。下面以二元文法为例,给出 P(s)的计算方法,如公式 (4)。

$$P(w_1, w_2, w_3, \dots, w_n) = \prod_{i=1}^n p(w_i | w_{i-1}) \quad (4)$$

精确率本文集成 KenLM Toolkit^[37]工具包,训练 ngram 语言模型。其中训练语料来自布朗数据集的新闻语料(Newswire Sections Of the Brown Corpus)^[38],语料规模为 6G。

5 实验结果

本文采用的是 SemEva2017 任务 7 的数据集。任务 7 是做双关语的识别,分为两种类型,语义双关和谐音双关,每种类型的语料都分为三个子任务,我们选择的任务是谐音双关语的子任务一,即识别语句是否为双关语。谐音双关语共 1780 条数据,其中双关语 1270 条,非双关语 510 条。本文共选择四个维度的 9 个特征,实验结果如表 1 所示。

表 1 特征迭加的双关语识别效果

特征	准确率	召回率	精确率	F1 值
语句透明度	46.57%	28.80%	88.83%	43.49%
语句透明度+语义相关度	48.88%	33.02%	89.85%	47.22%
语句透明度+语义相关度+语音扩展性	62.87%	57.12%	86.22%	68.72%
语句透明度+语义相关性+语音扩展性+句式特征集	78.48%	82.93%	86.39%	84.62%
SemEval 评测第一名	78.37%	81.90%	87.04%	84.39%

从表 1 中可以看出四个维度的特征是逐个添加到系统中的,单纯的语句复杂度虽然准确率较高,但召回率较低。加入语义相关度后,召回率和 F1 值均提高了 3%左右。再加入语音扩展性后,F1 值大幅度提高了 21%,这是因为本次任务是识别谐音双关语,语句中是否存在语音相同或相似的扩展词是双关语识别的一个重要标志,如果找到语义上关联的语音扩展词,那么这个句子就更可能是谐音双关语。最后加入多个句式特征后,在准确率没有下降的前提下,F1 值又提高 15.5%。这说明谐音双关语在语法和语用上有一定的特点,例如习惯使

用过去式等。四个维度的特征都加入后，分类的 F1 值达到 84.62%，比 SemEval 评测公布结果^[36]的第一名高出 0.3%。说明以上四个维度的特征在谐音双关语识别中有较好的效果，为了定量分析每个特征的作用，我们分别使用每个特征在数据集中分类，结果如表 2。

表 2 分别使用各维度特征的识别效果

特征	准确率	召回率	精确率	F1 值
语义相关度	32.02%	4.88%	98.41%	9.30%
语句透明度	46.57%	28.80%	88.83%	43.49%
语音扩展性	52.58%	39.89%	86.37%	54.57%
句式特征集	69.61%	63.49%	91.29%	74.90%

从表 2 中可以看出，四个维度的精确率都比较高，主要是召回率差别较大，语义相关度的召回率最低，但精确率最好达到 98.41%，召回率较低是因为只有部分谐音双关语具有语义关联性，精确率较高说明具有这个维度特征的语句，基本都是双关语。召回率最高的是句式特征集，达到 63.49%，是语义相关度的 13 倍。该维度中包含多个特征，每个特征都贡献了部分召回率。语句透明度和语音扩展性召回率相差 10%左右，精确率基本一致。从实验结果可以看出四个维度的特征，精确率都比较高，如果单一特征的精确率较低，二叉判定树的方法会把语句错分为双关句，没有修正的机会。而精确率较高，保证每次划分大部分是正确的。单个特征的召回率都比较低，比较适合用二叉判定树的方法，使多个特征结合，每个分支结点都能增加部分召回率，多个特征叠加，使最终的召回率达到合理水平。

从表 2 中可以看出，句式特征集的精确率和召回率都比较高，为细化句式特征集中各个特征的作用，我们只使用句式特征集中的三个特征识别双关语，结果如表 3 所示。该表中给出句式特征集中三个特征逐个用于识别双关语的效果，其中过去时态特征召回率和精确率都较高，所以 F1 值最高，而人名特征具有最高的精确率。可见，时态特征是双关语识别的一个重要特性。

表 3 三个句式特征的识别效果

特征	准确率	召回率	精确率	F1 值
人名特征	42.75%	20.61%	96.32%	33.96%
大写单词特征	30.56%	3.93%	76.92%	7.49%
过去时态特征	67.08%	58.22%	93.08%	71.64%

6 结论及下一步工作

本文提出一种包含四个维度，九个特征的二叉判定树算法，识别谐音双关语。四个维度的特征对双关语的识别都有较好的效果，其中语音扩展性和句式特征集效果最好，大幅度提高了算法的召回率，识别出了更多些谐音双关语，这也符合我们谐音双关语中语音作用较大的猜想。二叉判定树算法将多个特征融合，采用分治的方法，弥补了单个特征召回率较低的问题。分支结点中关键阈值的选择，本文尝试了将特征数据 K-Means 聚类的方法，选取双关语类的质心作为阈值，这样做虽然在单特征下损失了部分召回率，但保证了精确率的水平。而召回率可以通过多特征融合的方法弥补。

语音扩展性虽然在双关语识别中有较高的精确率，但是因为扩展策略的限制，能够找到语音扩展词的语句数量只有一半，下一步工作考虑继续改进语音扩展的策略，使更多的语句找到语音扩展词。另外，还要进一步验证上述特征迁移到其他语料中是否有较好的效果。最后，还考虑继续尝试 SemEval2017 任务 7 中的其他子任务，查找谐音双关

语中双关词及目标词和词义。

参考文献

- [1] Tristan Miller. Towards the automatic detection and identification of English puns[J]. *European Journal of Humour Research*, 2016, 4 (1):59 - 75
- [2] Julia M. Taylor. Ontology-based view of natural language meaning:the case of humor detection [J]. *J Ambient Intell Human Comput*, 2010, 1:221 - 234.
- [3] Julia M. Taylor. Computational Detection of Humor: A Dream or A Nightmare? [J]. *M International Joint Conferences on Web Intelligence and Intelligent Agent Technology*, 2009
- [4] Antonio Reyes, Paolo Rosso, Davide Buscaldi. From humor recognition to irony detection: The figurative language of social media[J]. *Data & Knowledge Engineering*, 2012, 74:1 - 12
- [5] Yishay Raz. Automatic Humor Classification on Twitter [J]. *Proceedings of the NAACL HLT 2012 Student Research Workshop*, 2012:66 - 70
- [6] Dario Bertero, Pascale Fung. Deep Learning of Audio and Language Features for Humor Prediction [J]. *International Conference on Language Resources and Evaluation*, 2016
- [7] Dario Bertero, Pascale Fung. Predicting humor response in dialogues from TV sitcoms [J]. *ICASSP*, 2016
- [8] Dario Bertero and Pascale Fung. A Long Short-Term Memory Framework for Predicting Humor in Dialogues[J]. *Proceedings of NAACL-HLT 2016*: 130 - 135
- [9] Renxian Zhang, Naishi Liu. Recognizing Humor on Twitter[J]. *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, 2014: 889-898
- [10] Antonio Reyes, Davide Buscaldi, and Paolo Rosso. An Analysis of the Impact of Ambiguity on Automatic Humour Recognition[J]. *Springer-Verlag Berlin Heidelberg 2009*: 162 - 169
- [11] Davide Buscaldi and Paolo Rosso. Some Experiments in Humour Recognition Using the Italian Wikiquote Collection[J]. *Springer-Verlag Berlin Heidelberg 2007*:464 - 468
- [12] Igor Labutov, Hod Lipson. Humor as Circuits in Semantic Networks[J]. *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*, 2012:150 - 155
- [13] Alessandro Valitutt, Hannu Toivonen. "Let Everything Turn Well in Your Wife": Generation of Adult Humor Using Lexical Constraints, *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics*, 2013: 243 - 248
- [14] 张冬瑜, 杨亮, 郑朴琪, 徐博, 林鸿飞. 情感隐喻语料库构建与应用 [J]. *中国科学*, 2015, 12(45):1574-1587
- [15] 林鸿飞, 张冬瑜, 杨亮, 徐博. 幽默计算及其应用研究[J]. *山东大学学报*, 2016, 7(51):1-10
- [16] Hempelmann CF, Computational humor: Beyond the pun?[J]. *The Primer of Humor Research. Humor Research*, 2008, (8): 333 - 360.
- [17] Ritchie. Computational mechanisms for pun generation[J]. *Proceedings of the 10th European Workshop on Natural Language Generation*, 2005: 8 - 10
- [18] Hong BA, Ong E. Automatically extracting word relationships as templates for pun generation[J]. in *Proceedings of the 1st Workshop on Computational Approaches to Linguistic Creativity*, 2009:24 - 31
- [19] Chloe Kiddon, Yuriy Brun. That's What She Said: Double Entendre Identification[J]. *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics:shortpapers*, 2011: 89 - 94
- [20] Diyi Yang, Alon Lavie, Chris Dyer, Eduard Hovy. Humor Recognition and Humor Anchor Extraction[J]. *Conference on Empirical Methods in Natural Language Processing*, 2015:2367 - 2376
- [21] Tristan Miller, Iryna Gurevych. Automatic disambiguation of English puns[J]. *Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, 2015:719 - 729
- [22] Arnold Zwicky and Elizabeth Zwicky. Imperfect puns, markedness, and phonological similarity: With fronds like these, who needs anemones. *Folia Linguistica*, 1986, 20(3/4):493 - 503
- [23] Włodzimierz Sobkowiak. Metaphonology of English paronomasic puns, 1991(26)
- [24] Robinson, T. *The British English example pronunciation dictionary* [M], 1996
- [25] Kim Binsted. *Machine humour: An implemented model of puns*. PH.D. University of Edinburgh, 1996
- [26] Hempelmann, C. F. *Paronomasic Puns: Target Recoverability Towards Automatic Generation*. PhD thesis in Purdue University, 2003
- [27] Justine T. Kao, Roger Levy, Noah D. Goodman. A Computational Model of Linguistic Humor in Puns [J]. *Cognitive Science*, 2015:1 - 16

- [28] Aaron Jaech, Rik Koncel-Kedziorski and Mari Ostendorf. Phonological Pun-derstanding[J]. Proceedings of NAACL-HLT 2016:654 - 663
- [29] 王春茂, 彭聘龄. 合成词加工中的词频、词素频率及语义透明度 [J]. 心理学报, 1999, (3):266- 273.
- [30] Cruise DA. Lexical Semantics [M] . Cambridge: Cambridge University Press, 1991
- [31] Mok LW. Word-superiority effect as a function of semantic transparency of Chinese bimorphemic compound words [J] . Language and Cognitive Process, 2009(24) : 1039- 1081
- [32] Pollatsek A. , Hyona J. and Bertram R. The role of semantic transparency in the processing of Finnish compound words [J]. Language and Cognitive Processes, 2005(20) : 261- 290.
- [33] Fellbaum, C. WordNet: An Electronic Lexical Database[M]. Cambridge, MA: MIT Press, 1998
- [34] https://en.wikipedia.org/wiki/CMU_Pronouncing_Dictionary
- [35] https://en.wikipedia.org/wiki/Binary_decision_diagram
- [36] <http://alt.qcri.org/semeval2017/task7>
- [37] Kenneth Heafield, Ivan Pouzyrevsky, Jonathan H. Clark and Philipp Koehn. Scalable Modified Kneser-Ney Language Model Estimation[J]. Annual Meeting of the Association for Computational Linguistics, 2013:690-696
- [38] Henry Kucera and W. Nelson Francis. Computational Analysis of Present-day American English[M]. Brown University Press, 1967

作者联系方式: 徐琳宏 大连市高新园区中海华庭小区 A2-2403 116023 13604116328
qingniao1203@163.com