

基于马尔可夫逻辑网的联合推理开放信息抽取

刘永彬 杨炳儒 李广源 刘英华

(北京科技大学计算机与通信工程学院 北京 100083)

摘 要 在自然语言处理的几个子任务上,传统的方法都是分而治之,例如分词、句法分析、命名实体识别、实体关系识别等。但是,孤立地分析和处理这些子任务会丢失一些彼此之间的内在联系,而这些子任务之间的内在联系往往会对每个子任务有很大帮助。所以,有人提出用联合集成式的模型,从整体上解决这些问题。但是,这些模型都只针对特定领域内的数据进行处理,还未能对开放式的信息进行处理。因此,提出了基于马尔可夫逻辑网的联合推理模型来处理开放式信息抽取(Open IE)。经过大量的实验证明,该模型的执行效率明显高于传统的模型。同时,该模型的适应性更好。

关键词 Open IE, 马尔可夫逻辑网, 联合推理

中图法分类号 TP181 **文献标识码** A

Joint Inference Open Information Extraction Based on Markov Logic Networks

LIU Yong-bin YANG Bing-ru LI Guang-yuan LIU Ying-hua

(School of Computer and Communication Engineering, University of Science & Technology of Beijing, Beijing 100083, China)

Abstract In recent decades, natural language processing has made great progress. Better model of each sub-problem achieves 90% accuracy or better. However, success in integrated, end-to-end natural language understanding remains elusive. The main reasons are that the systems processing sensory input typically have a pipeline architecture; the output of each stage is the input of the next, and errors are cascaded and accumulated through the pipeline of naively chained components, and there is no feedback from later stages to earlier ones. Actually, later stages can help earlier ones to process. Now a number of researchers have paid attention to this problem and proposed some joint approaches. But they do not perform open information extraction (Open IE), which can identify various types of relations without requiring pre-specifications. We proposed a joint inference model, which is based on Markov logic and can perform both traditional relation extraction and open IE. The proposed modeling significantly outperforms the other open IE systems in terms of both precision and recall. The joint inference is efficient and we have demonstrated its efficacy in real-world open IE detection tasks.

Keywords Open information extraction, Markov logic networks, Joint inference

1 引言

最近几年,在自然语言处理的几个问题上,比如词性标注、分词、句法分析、命名实体识别、实体关系识别、语义解析和共指消解等都有很大的进步。现在一些比较好的模型在这些子任务上的精度已经达到了 90% 甚至更好^[1]。然而,近些年来有人认为,孤立地对自然语言处理的几个子任务进行处理是不可取的。比如在实体识别和实体解析两个子任务上,很多实体并不是孤立的,而是有内在联系的,甚至有些实体的识别直接能解决实体的解析问题,同时有些实体的解析也能很好地帮助实体的识别;还有在实体识别和实体关系识别上也是如此。同时,信息抽取的数据常常也是为后期的知识发现而准备的,现在的信息抽取系统和模型主要把精力放在数据的抽取上,而这些抽取的数据之间的一些关系或者规律没

能引起注意,在这样的结果上进行知识发现,必然会丢掉很多有用甚至是非常重要的规律和知识。总之,当今主要的信息抽取技术都忽略了被处理数据之间的相互关系。因此,急需一个集成式的关系模型来应对这个问题,现在已经有部分学者开始重视这个问题,比如 McCallum 在 2009 年提出“联合推理的方法来处理自然语言处理的问题”的设想^[2],还有 Hoifung Poon 和 Pedro Domingos 提出统计关系模型来处理这个问题^[3]。但是,已提出的联合推理的关系模型都是针对特定领域的信息抽取,而没能实现开放式的信息抽取(Open IE)^[4]。因此,本文基于 Domingos 等人提出的马尔可夫逻辑网模型^[5],提出了一种适用开放式的联合推理的信息抽取模型。

本文第 2 节介绍传统的集成式模型的相关工作;第 3 节主要介绍马尔可夫逻辑网的基本内容;第 4 节主要介绍本文

到稿日期:2011-11-14 返修日期:2012-05-24 本文受国家自然科学基金项目(60875029,61175048)资助。

刘永彬(1978-),男,博士生,讲师,主要研究方向为人工智能、机器学习等,E-mail:qingbinliu@163.com;杨炳儒(1943-),男,教授,博士生导师,CCF 高级会员,主要研究方向为人工智能等;李广源(1970-),男,副教授,主要研究方向为知识发现等;刘英华(1976-),女,博士生,主要研究方向为数据挖掘。

提出基于马尔可夫逻辑网的联合推理模型;第5节对本文提出的联合推理方法进行了实验对比。

2 相关工作

许多学者已经注意到上面提到的问题。比如,McCallum等提出的联合概率模型来解决这个问题^[1]。还有,Finkel等人提出的近似贝叶斯模型^[6],该模型是在贝叶斯网络中使用了前馈的方法来推理,但是,该方法的信息传递只允许前馈,不能双向传递相关信息。后来,Wellner等人提出一种集成式的推理应用到信息抽取的共指消解中^[2],该方法是一个无向图模型,该模型允许信息双向反馈,不过它要求严格的逼近条件。最近,有些学者尝试着用 Markov logic 来解决这个问题。比如,Hoifung Poon 和 Pedro Domingos^[3,7]提出的联合推理模型应用于引文的信息抽取,Wanxiang Che 和 Ting Liu^[8]提出的集成模型应用于语义角色识别,Ivan Meza-Ruiz 和 Sebastian Riedel^[9]提出的模型应用于谓语角色识别等。但是,他们的方法都不能适应 OpenIE。在总结了其他学者的模型特点后,提出了基于 Markov logic Networks 的联合推理模型,该模型不仅能够适应 OpenIE,还能实现双向的反馈推理。

3 马尔可夫逻辑网(Markov Logic Networks)

3.1 马尔可夫逻辑网

马尔可夫逻辑网(Markov Logic Networks, MLNs)是由 Richardson 和 Domingos 于 2006 年正式提出^[5]。MLNs 作为统计关系学习的模型之一,结合了一阶谓词逻辑和概率图模型的逻辑结构表达方式,一经提出就显现了其相对于纯逻辑方法或纯概率方法的优势。

一阶逻辑知识库中的每个规则都加上一个权重,就构成了一个马尔可夫逻辑网。原本一个一阶逻辑知识库可以被视为一套硬约束加在可能世界的集合上,即一个世界只要与其中的一条规则冲突,它存在的概率即为零。而马尔可夫逻辑网的基本思想就是软化这些限制:即一个可能世界如果违反了规则,不再是存在性为零,而是这个世界存在的可能性下降,违反的规则数越少,这个可能的世界存在的可能性就越大^[5]。每个规则都有一个相关联的权重,反映出它对可能世界的约束强度:当其它情况一样的前提下,权重越高的规则,满足和不满足此规则的可能世界差别就越大^[5]。

定义 1^[5] 马尔可夫逻辑网 L 是 (F_i, W_i) 对的集合,其中 F_i 代表一阶逻辑规则, W_i 是一个实数;有限的常数集为 $C = \{C_1, C_2, \dots, C_n\}$,马尔可夫逻辑网 $M_{L,C}$ 按如下 1)、2) 来定义:

1) L 中每个谓词的每个可能基元在 $M_{L,C}$ 中有一个二元节点,如果原子为真,节点的值就等于 1,否则为 0。

2) L 中每个规则的每个基本可能在 $M_{L,C}$ 中有一个特征值,当这个规则为真时等于 1,否则等于 0,特征值的权重为 F_i 对应的 W_i 。

从定义 1、1) 和 2) 可以得出,基本马尔可夫逻辑网概率分布如下:

$$P(X=x) = \frac{1}{Z} \exp(\sum_i w_i n_i(x)) = \frac{1}{Z} \prod_i \phi_i(x_{(i)})^{n_i(x)} \quad (1)$$

式中, $n_i(x)$ 是 F_i 在 x 中所有取真值的闭规则的数量,而 $x_{(i)}$ 是 F_i 中为真的原子,又有 $\phi_i(x_{(i)}) = e^{w_i}$ 。

3.2 MLNs 推理及学习

在 MLNs 的最大后验推理中,通过已给出的闭原子 x 寻

找最有可能的闭原子 y 。其中推理任务被定义为式(2)^[5]:

$$\begin{aligned} \arg \max_y P(y|x) &= \arg \max_y \frac{1}{Z_x} \exp(\sum_i w_i n_i(x, y)) \\ &= \arg \max_y \sum_i w_i n_i(x, y) \end{aligned} \quad (2)$$

式中, Z_x 是一个常数, W_i 是规则 F_i 的权重, n_i 是 F_i 在 x 中所有取真值的闭规则的数量。

即使在规模很小的领域,直接计算也是很棘手的,因为马尔可夫逻辑网推理包含了 Sharp-P-complete 复杂度的概率推理,而逻辑推理在有限域也是 NP-complete 复杂度的,所以不能寄予期望。可是,提高推理运算效能的大多数技术都可以运用到马尔可夫逻辑网上,这是因为马尔可夫逻辑网能使用细粒度的知识库,其包含了上下类的独立性,推理比普通的图模型更有效;在逻辑方面,具备了概率语义的马尔可夫逻辑网能够进行高效的、近似推理。原则上,推理可以使用马尔可夫蒙特卡洛算法(MCMC)^[5]近似得出,这种算法还是太慢了。本文使用一种既准确又高效的基于整数线性规划(ILP)的 Cutting Plane Inference(CPI)方法^[10]来进行推理。方法 CPI 是始于边的子集,最终为这个子集解决 MAP 的问题^[9]。

对于 MLNs 的权值学习,本文采用的是 Online Max-Margin 权值学习法^[11]。

4 基于 MLNs 的联合推理模型

4.1 谓词的定义

本文提出了一种针对 Open IE 的联合推理模型。首先,给定抽取出的结果的形式。从 Open 数据集中抽取出的信息的形式是关系三元组 Tuples, $T = (e_i, r_{i,j}, e_j)$ ^[4], $i < j$, 其中 e_i 和 e_j 表示语义实体, $r_{i,j}$ 表示两个实体之间关系的关键词或词组。在本文中, e_i 和 e_j 只表示基本的名词短语,不包括嵌套的名字短语和介词短语等^[4]。如下,三元组的形式示例:

(⟨proper noun⟩, acquired, ⟨proper noun⟩)
 (⟨proper noun⟩, was born in, ⟨proper noun⟩)
 (⟨proper noun⟩, go to, ⟨proper noun⟩)
 (⟨proper noun⟩, come from, ⟨proper noun⟩)
 (⟨proper noun⟩, become, ⟨noun phrase⟩)
 (⟨proper noun⟩, studied at, ⟨proper noun⟩)
 (⟨proper noun⟩, convert, ⟨proper noun⟩)
 (⟨proper noun⟩, derive from, ⟨proper noun⟩)
 (⟨proper noun⟩, was founded by, ⟨proper noun⟩)
 (⟨proper noun⟩, worked in, ⟨proper noun⟩)
 ⋮

在传统的 Open IE 系统中,实体 e_i 和 e_j 首先被识别出来,然后再依据这两个实体来识别它们之间关系的关键词或词组 $r_{i,j}$ 。然而,事实上对于 $r_{i,j}$ 往往是句子或子句中的谓词性质的短语这里包括具有谓语性质的名词短语,更容易被识别,同时这个 $r_{i,j}$ 对识别其它的实体有帮助,传统的方法都忽略了这些帮助信息。所以,本文模型的重点是识别 predicate ($r_{i,j}$),然后再依据这个关键词(谓语短语)来帮助识别相关的两个实体 e_i 和 e_j 。故针对 Open IE 提出了基于 MLNs 的联合推理模型,其过程分为两个阶段:第一阶段,辨别并提取句中的关键词或短语(谓语短语);第二阶段,依据第一阶段提取出的关键词和短语辨别并提取相关语义实体 e_i 和 e_j 。

本文提出的联合推理模型中,针对 Open IE 的特点,定义了 7 个一阶逻辑中的谓词。其中,针对谓语短语的识别,定义

了两个谓词 $isPredicate(p)$ 和 $isRelation(p, t)$ 。 $isPredicate(p)$ 表示在 p 位置的词是一个谓语。 $isRelation(p, t)$ 表示在位置 p 和位置 t 的两个词组成了一个关键短语(谓语短语) r , 例如, $T = (\langle proper\ noun \rangle, graduated\ from, \langle proper\ noun \rangle)$, “ $graduated\ from$ ”就是一个谓语短语 r 。对于实体 e_i 和 e_j 的识别, 定义了 5 个谓词, $isEntity(i, j)$, $hasRelation(p, e_i)$, $preRelation(p, e_i)$, $sucRelation(p, e_i)$ 和 $isTuple(e_i, r_{i,j}, e_j)$ 。谓词 $isEntity(i, j)$ 的含义是从位置 i 到 j 之间的词构成了一个语义实体。谓词 $hasRelation(p, e_i)$ 表示在 i 位置的实体 e_i 是 p 位置的谓语的语义实体。谓词 $preRelation(p, e_i)$ 表示实体 e_i 是 p 位置的谓语的前驱语义实体。谓词 $sucRelation(p, e_i)$ 表示实体 e_i 是 p 位置的谓语的后继的语义实体。谓词 $isTuple(e_i, r_{i,j}, e_j)$ 表示实体 e_i 、关键短语(谓语短语) $r_{i,j}$ 和实体 e_j 构成一个三元组 $T = (e_i, r_{i,j}, e_j)$ 。

针对以上的谓词, 本文还定义了一些可观察到的谓词。谓词 $word(i, w)$ 代表, 标记 i 位置的是一个词 w 。谓词 $pos(i, t)$ 表示在标记 i 位置的词具有词性 t 。谓词 $lemma(i, l)$ 表示在标记 i 位置的词的词条是 l 。谓词 $lt(i, j)$ 表示 i 值小于 j 值。谓词 $Nonsame(e_i, e_j)$ 表示 e_i 和 e_j 不是同一个实体词。同时, 为使模型具有更好的双向驱动能力, 本文还定义了谓词 $Similar(s, i, j, s', i', j')$, 它表示语句 s 和语句 s' 含有相似的语义实体。

4.2 局部规则

如果一个规则的前件中只含有已知的闭原子的成分, 那么该规则就是一个局部规则^[9]。例如,

$$lemma(p, +l_1) \wedge lemma(e, +l_2) \Rightarrow hasRelation(p, e) \quad (3)$$

规则(3)表示如果 p 位置的词条是 l_1 且在 e 位置的词条是 l_2 , 那么在这两个位置的词条可能存在语义关系。其中的“+”号表示其后的元素是 MLNs 中一个规则的实例, 这个实例带有自己的权重^[3]。

对于 $isPredicate(p)$, $isRelation(p, e_i)$ 和 $isEntity(i, j)$, 这些谓词都是在词法和句法的背景下得出的。例如,

$$pos(p, verb) \Rightarrow isPredicate(p) \quad (4)$$

$$lemma(p, +l_1) \wedge pos(t, prep) \Rightarrow isRelation(p, t) \quad (5)$$

规则(4)表达的是如果 p 位置的词性是动词那么它将是一个谓语。规则(5)意味着如果在 p 位置的词条是 l_1 且在 t 位置的词性是介词($p < t$), 那么这两个词可能构成了一个谓语短语。

4.3 全局规则

如果在规则中没有已知闭原子的成分, 那么这个规则就是全局规则。我们使用这种类型的规则是为了确保谓词的各个阶段和结构之间的一致性^[9]。例如,

$$hasRelation(p, e_i) \wedge lt(i, p) \Rightarrow preRelation(p, e_i) \quad (6)$$

$$hasRelation(p, e_j) \wedge lt(p, j) \Rightarrow sucRelation(p, e_j) \quad (7)$$

$$preRelation(p, e_i) \wedge sucRelation(p, e_j) \Rightarrow isTuple(e_i, r_{i,j}, e_j) \quad (8)$$

带权重的规则(6)–(8)表明, 如果 e_i 和 e_j 分别是 p ($i < p < j$) 的语义实体, 那么能得出 $isTuple(e_i, r_{i,j}, e_j)$ 。为了充分利用各个阶段的输出结果, 我们还定义了以下规则集:

$$isTuple(e_i, r_{i,j}, e_j) \Rightarrow preRelation(p, e_i) \quad (9)$$

$$isTuple(e_i, r_{i,j}, e_j) \Rightarrow sucRelation(p, e_j) \quad (10)$$

$$Similar(s, i, j, s', i', j') \wedge isEntity(i, j) \Rightarrow isEntity(i', j') \quad (11)$$

$$Similar(s, p, s', p') \wedge isPredicate(p) \Rightarrow isPredicate(p') \quad (12)$$

规则(9)–(12)的提出是为了模型能够充分地利用各个阶段的输出, 实现信息流的双向流动。同时, 为了确保唯一性, 我们还定义了以下规则。

$$preRelation(p, e_1) \wedge Nonsame(e_1, e_2) \Rightarrow \neg preRelation(p, e_2) \quad (13)$$

$$sucRelation(p, e_1) \wedge Nonsame(e_1, e_2) \Rightarrow \neg sucRelation(p, e_2) \quad (14)$$

MNLs 都有“名字唯一性”的假设, 规则(13)和规则(14)表达的都是实体在规则中的唯一性。

下面将看到所提联合推理模型在 Open IE 中的实验部分。

5 实验结果及分析

5.1 实验数据及评价

本文的实验使用的是 3 个数据集。第一个数据集是 OntoNotes Release 3.0 语料库^[12]。这个语料库用 3 种语言(英语、汉语和阿拉伯语), 包括各种体裁的文字(新闻、电话语音、网络日志、Usenet 新闻组、广播、脱口秀等), 同时这些文字都具有浅层语义结构信息(语法和谓词结构)。本文使用的是 OntoNotes 3.0 英文数据集, 被分成 3 部分: 博客、新闻广播和杂志类。本文使用的第二个数据集是通过爬虫在网上随机收集的。经过去噪预处理, 如去掉网页的导航栏等。同时, 本文还采用了文献[14]中的可视化分析器方法分割网页成数据块。网页中的核心数据块被挑选为实验的数据集, 使用词性标注器解析这些数据块中的所有文本语句^[13]。本文搜集了 3 万多个这样的数据块, 将这个数据集定义成 W3 数据集。第三个数据集特别针对社交网络, 数据集来源于 tweets^[15], 本文随机选取了 1 万多个的 tweets, 将这个数据集定义成 T1。显然, 第二和第三个数据集比第一个数据集的分布更加具有开放性, 同时数据涵盖的类型更加多样。第三个数据集同时具有短语境、口语化、语法混乱等特点。对于以上数据集, 本文使用条件随机场(CRFs)模型标注这些数据集的词性。对于每个词的词条本文使用 WordNet tool 来得到。

为了对比评价所提出的联合推理模型, 本文使用召回率、精准率和 F_1 来评价所提联合推理模型的性能指标。召回率(Recall)和精准率(Precise)是广泛用于信息检索和统计学分类领域的两个度量值, 用来评价结果的质量。其中召回率是检索出的相关文档数和文档库中所有的相关文档数的比率, 衡量的是检索系统的查全率。精准率是检索出的相关文档数与检索出的文档总数的比率, 衡量的是检索系统的查准率。一般来说, 精准率和召回率不会同时很高, 所以在不同的系统需求下所要求的就相应不同, 因此精准率和召回率所要评估的对象是两种不同取向的系统考量, 有的系统需要精准率, 有的系统对召回率有更高的要求, 所以有些时候通过召回率和精准率很难对比评价不同的系统。为此, 本文采用了 F_1 值来对比评价模型的性能质量, 它是召回率和精准率的调和平均, 使用二者的调和平均数来对比评价不同的系统更为科学。

5.2 实验结果

对于第一个 OntoNotes 数据集上的实验结果, 如表 1 所列。其中 Pipeline 是传统的推理模型, 而 Joint 指的是本文提出的联合推理模型, P 是指精准率, R 是指召回率。

表1 OntoNotes 数据集上的实验结果对比

Categories		isPredicate	isEntity	hasRelation	isTuple
Pipeline	P	0.93	0.91	0.87	0.86
	R	0.88	0.85	0.79	0.70
	F ₁	0.90	0.87	0.82	0.77
Joint	P	0.93	0.91	0.89	0.89
	R	0.90	0.87	0.83	0.81
	F ₁	0.91	0.88	0.85	0.84

从表1能够看出,本文提出的联合推理模型在谓语识别、实体识别、语义角色识别和最终的三元组的抽取方面都优于传统的模型;出由于实体组成的复杂性,对于实体的识别要比谓语的识别难度更大些。

表2是对第二个数据集W3进行的实验。从表2能明显发现,本文提出的模型的实验结果优于传统的模型。然而,通过对比这两个实验数据集的结果发现,OntoNotes数据集上的结果要好于W3数据集上的结果,这是因为W3数据集更加具有开放性,是非特定领域数据。

表2 W3数据集上的实验结果对比

Categories		isPredicate	isEntity	hasRelation	isTuple
Pipeline	P	0.83	0.80	0.75	0.73
	R	0.80	0.77	0.70	0.65
	F ₁	0.81	0.78	0.72	0.68
Joint	P	0.89	0.85	0.82	0.80
	R	0.88	0.83	0.76	0.73
	F ₁	0.88	0.83	0.78	0.76

表3是对第三个数据集T1进行的实验。从表3能明显发现,本文提出的模型的实验结果优于传统的模型。但第三个数据集上的实验结果普遍没有第一个和第二个数据集的实验结果好,这是因为社交网络上的数据具有口语化、语法混乱等特点,使得抽取相关语义信息难度加大造成的。通过以上三个数据集的实验结果可得,在Open IE上本文提出的模型要明显好于传统的模型。

表3 T1数据集上实验结果对比

Categories		isPredicate	isEntity	hasRelation	isTuple
Pipeline	P	0.62	0.59	0.53	0.51
	R	0.51	0.53	0.49	0.47
	F ₁	0.56	0.55	0.51	0.49
Joint	P	0.78	0.77	0.73	0.72
	R	0.65	0.63	0.62	0.59
	F ₁	0.71	0.70	0.67	0.65

结束语 本文提出了一个基于MLNs的联合推理模型,该模型应用于Open IE上。它能胜任不同水平的信息抽取任务。最终,本文也通过大量的数据集对比证明了该模型的性能优于传统的模型。

本文提出的联合推理模型有较好的普适性,将来仍然有改进和提升的空间。本文中,对于数据集需要预先处理,尤其是要通过CRFs方法来得到词性的结果,才能进行模型推理。将来我们会把这部分步骤集成到本文提出的模型中,达到完全的、真正的联合推理。同时,本模型对于社交网络上的信息进行抽取时结果还不是很理想,这也是我们下一步主要的工作重点。

参考文献

[1] McCallum A. Joint Inference for Natural Language Processing [C]//Proceedings of the Thirteenth Conference on Computational Natural Language Learning (CoNLL). 2009

[2] Wellner B, McCallum A, Peng Fu-chun, et al. An integrated, con-

ditional model of information extraction and coreference with application to citation matching[C]//Proceedings UAI '04 of the 20th conference on Uncertainty in artificial intelligence. 2004;593-601

[3] Poon H, Domingos P. Joint Inference in Information Extraction [C]//AAAI' 07 Proceedings of the 22nd national conference on Artificial intelligence. 2007;913-918

[4] Banko M, Cafarella M, Soderland S, et al. Open information extraction from the Web[C]//Twentieth International Joint Conference on Artificial Intelligence. 2007;2670-2676

[5] Richardson M, Domingos P. Markov logic networks[J]. Machine Learning, 2006, 62(1/2): 107-136

[6] Finkel J R, Manning C D, Ng A Y. Solving the problem of cascading errors; Approximate bayesian inference for linguistic annotation pipelines[C]//Conference on Empirical Methods on Natural Language Processing (EMNLP). 2006

[7] Singla P, Domingos P. Entity Resolution with Markov Logic[C]//Proceedings of the Sixth International Conference on Data Mining(ICDM' 06). 2006

[8] Che Wan-xiang, Liu Ting. Jointly Modeling WSD and SRL with Markov Logic[C]//Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010). 2010;161-169

[9] Riedel S, Meza-Ruiz I. Collective Semantic Role Labelling with Markov Logic[C]//Proceedings of the 12th Conference on Computational Natural Language Learning. 2008;193-197

[10] Riedel S. Improving the accuracy and efficiency of map inference for markov logic[C]//UAI '08 Proceedings of the Annual Conference on Uncertainty in AI. 2008

[11] Tuyen N, Raymond J H. Mooney, Online Max-Margin Weight Learning for Markov Logic Networks[C]//Proceedings of the Eleventh SIAM International Conference on Data Mining (SDM11). 2011;642-651

[12] [http://www ldc upenn edu/Catalog/CatalogEntry.jsp? catalogId=LDC2009T24](http://www ldc upenn edu/Catalog/CatalogEntry.jsp?catalogId=LDC2009T24)

[13] Zhu Jun, Nie Zai-qing, Liu Xiao-jing, et al. StatSnowball: a Statistical Approach to Extracting Entity Relationships[C]//18th International World Wide Web Conference. 2009;101-110

[14] Zhu J, Nie Z, Wen J-R, et al. Simultaneous record detection and attribute labeling in web data extraction[C]//Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2006

[15] <http://trec.nist.gov/data/tweets/>

[16] Yang Bing-ru, Hou Wei. KAAPRO: An approach of protein secondary structure prediction based on KDD* in the compound pyramid prediction model[J]. Expert Systems With Applications, 2009, 36(1): 9000-9006

[17] Ouyang Chun-ping, Hu Chang-jun, Liu Zhen-yu. Data Grid and GIS Technology for E-Science Application: A Case Study of Gas Network Safety Evaluation[C]//Proceedings of the Fourth International Conference on Complex, Intelligent and Software Intensive Systems (CISIS 2010). IEEE Computer Society, 2010; 520-525

[18] Abd Allah A H. Parametric prediction limits for generalized exponential distribution using record observations [J]. Applied Mathematics & Information Sciences, 2009, 3(2): 135-149

[19] Poon H, Domingos P. Unsupervised Semantic Parsing[C]//Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing (EMNLP '09). 2009

[20] 刘大有, 于鹏, 高滢, 等. 统计关系学习研究进展[J]. 计算机研究与发展, 2008, 45(12)