

王东波 著

# 面向非结构化文本的 知识发现

基于英汉双语平行语料库的  
句法级知识挖掘和抽取研究

中国社会科学出版社

013057352

TP182  
25

王东波 著

# 面向非结构化文本的 知识发现

基于英汉双语平行语料库的  
句法级知识挖掘和抽取研究



TP182/25



中国社会科学出版社

图书在版编目(CIP)数据

面向非结构化文本的知识发现：基于英汉双语平行语料库的句法级知识挖掘和抽取研究 / 王东波著 . —北京：中国社会科学出版社，2013.5

ISBN 978 - 7 - 5161 - 2603 - 5

I. ①面… II. ①王… III. ①知识工程—数据收集—研究 IV. ①TP182

中国版本图书馆 CIP 数据核字(2013)第 100209 号

---

出版人 赵剑英

责任编辑 王琪

责任校对 鲍凤英

责任印制 王超

---

出 版 中国社会科学出版社

社 址 北京鼓楼西大街甲 158 号 (邮编 100720)

网 址 <http://www.csspw.cn>

中文域名：中国社科网 010 - 64070619

发 行 部 010 - 84083685

门 市 部 010 - 84029450

经 销 新华书店及其他书店

---

印刷装订 三河市君旺印装厂

版 次 2013 年 5 月第 1 版

印 次 2013 年 5 月第 1 次印刷

---

开 本 710 × 1000 1/16

印 张 14.75

插 页 2

字 数 227 千字

定 价 39.00 元

---

凡购买中国社会科学出版社图书,如有质量问题请与本社联系调换

电话 : 010 - 64009791

版权所有 侵权必究

## 序

“除了上帝，任何人都必须用数据来说话”这句美国谚语将会在今天这个真正的大数据时代使世人更深刻、更真实、更全面地感受数据的力量和“无微不至”。大数据的迅猛发展和其所具有的战略意义，引起了世界的广泛关注。2012年1月，达沃斯世界经济论坛发布《大数据，大影响》报告，宣称数据已经成为一种新的经济资产类别，就像货币和黄金一样，是21世纪的石油。2012年3月奥巴马政府发布了《大数据研究和发展倡议》书，投资2亿美元启动“大数据研究和发展计划”；2012年11月，国家发展和改革委员会在其发布的高技术服务业研发及产业化专项中，设立了“大数据分析软件开发和服务创新”专项；2013年2月，科技部发布2014年度973计划指南，设立了“大数据计算的基础研究”项目。在以更大规模的数据、更多样化的数据、更实时的数据、数据的价值稀疏性或者不精确性为特征的大数据中，80%的数据是非结构化的。如何从非结构化数据，尤其是非结构化文本数据中挖掘出有价值的知识就成为备受关注和研究的热点。王东波博士从英汉双语平行语料这一非结构化文本入手，通过把情报学的知识、自然语言处理和文本挖掘的技术结合起来的方法，围绕词汇、简单短语结构和复杂短语结构这三个句法层级进行的知识抽取和挖掘正是在这一研究趋势下的积极而有益的尝试。

针对已有面向非结构化文本知识挖掘研究过于依赖模型和统计的现状，王东波博士以其本科阶段所习得的语言学知识为切入点，从英、汉这两种印欧语系和汉藏语系为代表的语言特征入手，面向网络构建了大规模的英汉双语平行语料库，并基于英汉平行语料这一非结构化文本，

## 2 面向非结构化文本的知识发现

挖掘出了词汇句法动态组合的分布规律，以介宾短语结构为例探讨了短语结构的知识抽取，在由复杂短语结构组成的短句基础上构建了类别知识挖掘的模型。本书的主要研究内容如下：

面向网络获取了英汉双语通用和专门平行语料，构建了相应的英汉双语平行语料库。在该部分主要围绕确定抓取网站、制定抓取底表、通过抓取工具获取网页、抽取英汉双语平行语料对、清洗英汉双语平行语料对和对英汉双语平行语料对进行去重处理等问题展开了探讨。

在词汇这一级，结合情报学中的相应方法和知识，挖掘出了英汉双语词汇在句法功能分布复杂度上呈现洛特卡现象的规律。在该部分，基于英汉双语平行语料库、宾州大学英汉树库和清华汉语树库，统计了英语和汉语词汇句法功能的分布，分析了英语和汉语词汇的句法功能分布复杂度情况，计算得出了英语和汉语词汇的平均句法功能分布复杂度值，揭示了词汇句法功能分布复杂度所呈现的洛特卡现象。

在简单短语这一级，通过介宾短语结构，基于条件随机场这一机器学习模型，构建了英汉双语介宾短语结构知识抽取模型，并给出了英汉双语介宾短语结构知识抽取的流程。本书在该部分统计了介宾短语结构的内部和外部句法特征、给出了训练语料的预处理格式、详细说明了自身特征模板和添加特征模板的具体构成内容并与最大熵的性能进行了对比。

在复杂短语这一级，基于已有聚类算法，通过构建词汇和词性特征知识下的类别知识挖掘模型，完成了面向英汉双语专门平行语料的类别知识挖掘的探究。本书通过具体的实验证实了英汉双语词汇特征在类别知识挖掘中的性能，并给出了造成类别知识挖掘性能有差异的原因，同时使用词汇和词性的特征知识，在英汉名词、英汉名词、动词，英汉名词、动词和形容词这三种词汇和词性特征组合的基础上，探究了词性知识在类别知识挖掘上的具体表现。

王东波博士面向英汉双语平行语料这一非结构化数据抽取和挖掘的句法三个层面的知识不仅有利于知识库构建、知识服务、信息检索、信息计量等情报学中的相关研究开展，而且有助于自然语言处理中的歧义消解、知识抽取和机器与辅助机器翻译问题的解决。同时，本书所采用

的把情报学的方法和知识通过自然语言处理的技术和模型融入非结构化数据的知识挖掘和抽取中的策略具有方法论上的意义，在一定程度上具有推广价值。

王东波博士本科主修汉语言文学，硕士为中文信息处理专业，博士主攻情报学，不仅学习了语言学概论、现代汉语、古代汉语和现代汉语语法研究等语言学课程的知识，而且掌握了中文信息处理、计算语言学、语料库语言学和语言统计方法等课程的技术，还系统学习了情报学理论与方法、信息资源管理技术、信息处理与检索技术和信息智能处理与检索等情报学课程的方法，具有跨学科研究的良好基础。同时，读博期间，王东波博士在处理数据、发表论文、承担项目的过程中充分体现了勤奋刻苦、锐意创新、百折不挠的个性。《周易》有言：君子藏器于身，待时而动。我相信，王东波博士在这个大数据时代，依据其“咬合度”极好的跨学科知识、技术和方法，凭借其勤奋、创新和坚持的个性，定能做出更多、更好的科研成果。

是为序。

苏新宁

2013年5月13日

南京大学信息技术开发研究所

# 目 录

<b>第一章 引言</b> .....	(1)
一 课题提出 .....	(1)
二 研究意义 .....	(2)
三 研究方法 .....	(4)
四 研究技术路线 .....	(5)
五 研究创新点 .....	(7)
六 本书结构和所用资源 .....	(8)
<b>第二章 相关研究综述</b> .....	(10)
一 面向网络获取相关英汉双语平行语料的概述 .....	(10)
二 词汇句法功能分布的相关研究概况 .....	(23)
三 短语结构知识抽取的相关研究综述 .....	(33)
四 类别知识挖掘的相关研究 .....	(41)
<b>第三章 面向网络的英汉双语平行语料库自动构建</b> .....	(48)
一 确定获取语料网站和制定抓取词汇底表 .....	(48)
二 网页的抓取和英汉双语平行语料的抽取 .....	(53)
<b>第四章 词汇句法功能分布复杂度呈现规律的知识挖掘</b> .....	(64)
一 词汇句法功能分布复杂度统计数据源简介和句法 结构调整 .....	(64)

2 面向非结构化文本的知识发现	
二 词汇句法功能分布复杂度的获取	(81)
三 词汇句法功能分布复杂度的洛特卡现象揭示	(96)
第五章 基于英汉双语平行语料库的短语结构知识抽取	
——以介宾短语结构为例	(114)
一 英汉双语介宾短语结构句法特征统计分析	(115)
二 有关介宾短语结构知识抽取模型构建的相关介绍 和预处理	(135)
三 介宾短语结构知识抽取模型的确定和英汉双语介宾短语 结构知识的抽取	(147)
第六章 基于英汉双语平行语料库的复杂短语级类别	
知识挖掘	(152)
一 面向英汉双语专门复杂短语平行语料的聚类和词干或词形 算法确定	(152)
二 语料的预处理和相关统计	(157)
三 面向英汉双语复杂短语平行语料的词性选择	(166)
四 基于英汉双语复杂短语平行语料的类别知识挖掘	(171)
结语	(181)
参考文献	(183)
附录 1 宾州大学英语树库的词性标记	(199)
附录 2 宾州大学汉语树库的词性标记	(200)
附录 3 清华大学树库词性标记	(201)
附录 4 汉语自身特征模板	(203)
附录 5 英语自身特征模板	(205)
附录 6 汉语添加特征模板	(207)
附录 7 英语添加特征模板	(210)

目 录 3

附录 8 中国科学院和北京大学核心词性标注集 .....	(213)
附录 9 BNC 统计语料样例 .....	(214)
后记 .....	(217)

(图1)	基于英汉双语平行语料库的句法级知识挖掘和抽取	8
(图2)	南京大学英语单语语料库存储样例	15

## 图 目 录

图 1—1	基于英汉双语平行语料库的句法级知识挖掘和抽取 研究技术流程图	(6)
图 2—1	南京大学英语单语语料库存储样例	(15)
图 3—1	英汉双语平行语料抽样程序	(49)
图 3—2	网页获取词汇与网址链接程序图	(54)
图 3—3	获取网页的配置设置	(55)
图 3—4	GUN Wget 获取含有英汉双语平行语料网页的 进行样例	(56)
图 3—5	英汉双语平行语料对抽取程序图	(57)
图 3—6	英汉双语平行语料去重程序样例	(59)
图 3—7	汉语分词整体系统	(60)
图 3—8	汉语分词预处理系统	(61)
图 3—9	汉语分词还原系统	(61)
图 3—10	汉语分词后处理系统	(62)
图 3—11	英汉双语平行语料库中的语料样例	(63)
图 4—1	多叉树图示	(85)
图 4—2	英汉双语句子树生成流程图	(86)
图 4—3	句法树样例图	(87)
图 4—4	英汉双语词汇句法功能统计流程图	(88)
图 4—5	英语词汇 SFDC 与词汇数目的散点图	(107)
图 4—6	汉语词汇 SFDC 与词汇数目的散点图	(108)
图 4—7	宾州大学汉语词汇 SFDC 与词汇数目的散点图	(110)

图 4—8 宾州英语树库英语词汇句法功能分布复杂度与词汇 数目的散点图	(111)
图 4—9 清华汉语树库汉语词汇句法功能分布复杂度与词汇 数目的散点图	(112)
图 5—1 汉语介宾短语结构内部词性分布状况	(122)
图 5—2 英语介宾短语结构内部词性分布状况	(124)
图 5—3 介宾短语结构内部短语结构分布	(126)
图 5—4 英语介宾短语结构内部短语结构分布	(128)
图 5—5 汉语介宾短语结构长度分布	(130)
图 5—6 英语介宾短语结构长度分布	(131)
图 5—7 汉语介宾短语结构中高频介词的分布状况	(132)
图 5—8 英语介宾短语结构中高频介词的分布状况	(134)
图 5—9 条件随机场的图模型	(136)
图 5—10 最大熵的图形表示	(139)
图 5—11 介宾短语结构训练特征模板样例	(141)
图 5—12 介宾短语结构知识抽取模型构建流程	(146)
图 6—1 基于英汉双语复杂短语平行语料词汇特征聚类流程	(172)
图 6—2 汉语特征的聚类结果	(172)
图 6—3 英语词汇聚类的性能	(173)
图 6—4 基于英汉双语词汇特征的类别挖掘性能	(174)
图 6—5 基于英汉双语词汇和词性的类别知识挖掘流程图	(176)
图 6—6 基于名词词汇和词性的英汉双语复杂短语平行语料 聚类结果	(177)
图 6—7 基于名词、动词词汇和词性的英汉双语复杂短语平行 语料聚类结果	(177)
图 6—8 基于名词、动词和形容词词汇和词性的英汉双语复杂 短语平行语料聚类结果	(177)

## 表 目 录

表 2—1	语料库的发展阶段及典型语料库扫描	(10)
表 2—2	我国目前已建成的主要汉语语料库一览	(12)
表 2—3	我国目前已建成的主要英语语料库一览	(13)
表 2—4	文件头信息	(18)
表 2—5	文件体信息	(19)
表 3—1	英汉双语平行语料分布情况和类别	(50)
表 3—2	基于 BNC 语料库的英语动词词频样例表	(52)
表 3—3	面向网络获取网页的抓取底表样例	(52)
表 3—4	词汇与网址生成的链接样例表	(54)
表 3—5	通用英汉双语平行语料样例	(57)
表 3—6	专门英汉双语平行语料样例	(58)
表 4—1	宾州树库英语语料来源和分布情况	(65)
表 4—2	清华汉语树库功能短语标记分布	(67)
表 4—3	清华汉语树库短语句法结构分布	(67)
表 4—4	宾州英语树库中删除的结构列表	(71)
表 4—5	宾州英语树库中句法结构合并和拆分后的结构列表	(71)
表 4—6	宾州汉语树库中删除的结构列表	(75)
表 4—7	宾州汉语树库中句法结构合并后的结构列表	(75)
表 4—8	清华汉语树库合并、拆分等调整后的句法短语结构分布表	(78)
表 4—9	英汉双语平行语料英语词汇句法功能分布样例	(89)
表 4—10	英汉双语平行语料汉语词汇句法功能分布表样例	(90)

表 4—11	宾州大学英语词汇句法功能分布样例 .....	(91)
表 4—12	宾州大学汉语词汇句法功能分布表样例 .....	(92)
表 4—13	清华大学汉语词汇句法功能分布表样例 .....	(93)
表 4—14	部分英汉双语平行语料英语词汇一句法结构表 和句法功能分布复杂度样例 .....	(95)
表 4—15	部分清华树库汉语词汇一句法结构表和句法功能 分布复杂度样例 .....	(96)
表 4—16	英汉双语平行语料英语词汇和句法功能分布复杂度 对应关系表 .....	(97)
表 4—17	英汉双语平行语料汉语词汇和句法功能分布复杂度 对应关系表 .....	(99)
表 4—18	宾州大学全部树库英语词汇 SFDC 分布情况 .....	(100)
表 4—19	宾州大学三分之二树库英语词汇 SFDC 分布情况 .....	(101)
表 4—20	宾州大学三分之一树库英语词汇 SFDC 分布情况 .....	(101)
表 4—21	宾州大学汉语词汇 SFDC 分布状况 .....	(103)
表 4—22	清华大学汉语树库全部树库数据 .....	(104)
表 4—23	清华大学汉语树库三分之二树库数据 .....	(104)
表 4—24	清华大学汉语树库三分之一树库数据 .....	(105)
表 4—25	英语词汇 SFDC 和词彙总量关系的一元线性 回归结果 .....	(107)
表 4—26	汉语词汇 SFDC 和词彙总量关系的一元线性 回归结果 .....	(109)
表 4—27	宾州大学汉语树库中汉语词汇 SFDC 和词彙总量 关系的一元线性回归结果 .....	(110)
表 4—28	英语词汇句法功能分布复杂度和词彙总量关系的 一元线性回归结果 .....	(111)
表 4—29	汉语词汇句法功能分布复杂度和词彙总量关系的 一元线性回归结果 .....	(113)
表 5—1	汉语介宾短语结构的句法功能分布 .....	(115)
表 5—2	英语介宾短语结构的句法功能分布 .....	(116)

## 8 面向非结构化文本的知识发现

表 5—3 汉语一元右邻接词的词类分布特征	(118)
表 5—4 英语一元右邻接词的词类分布特征	(119)
表 5—5 汉语介宾短语结构右邻接词	(120)
表 5—6 英语介宾短语结构右邻接词	(121)
表 5—7 汉语介宾短语结构高频词性序列	(123)
表 5—8 英语介宾短语结构高频词性序列	(124)
表 5—9 介宾短语结构内部高频的短语结构序列	(127)
表 5—10 英语介宾短语结构内部高频的短语结构序列	(129)
表 5—11 汉语介宾短语结构内部右边界高频词	(132)
表 5—12 英语介宾短语结构内部右边界高频词	(134)
表 5—13 汉英介宾短语结构训练和测试样例	(140)
表 5—14 汉英介宾短语结构预处理结果样例	(142)
表 5—15 基于添加特征模板的汉语介宾短语结构知识抽取 训练语料样例	(144)
表 5—16 基于添加特征模板的英语介宾短语知识抽取训练 语料样例	(145)
表 5—17 基于条件随机场构建的模型性能	(148)
表 5—18 基于最大熵构建的模型性能	(148)
表 5—19 基于添加特征模板的汉语介宾短语结构知识抽取 模型测试性能	(149)
表 5—20 英语添加特征模板的性能	(150)
表 6—1 三种聚类算法在复旦语料上的 Entropy 性能	(155)
表 6—2 三种聚类算法在复旦语料上的 Purity 性能	(155)
表 6—3 Bisecting K-means Clustering 下的词干或词形 还原算法的 Purity	(156)
表 6—4 Bisecting K-means Clustering 下的词干或词形 还原算法的 Entropy	(156)
表 6—5 基于英汉双语复杂短语平行语料类别知识挖掘 语料样例	(157)
表 6—6 形态转换后的英语词汇样例	(161)

## 表 目 录 9

表 6—7 英汉复杂短语平行语料汉语词汇处理结果 .....	(163)
表 6—8 英汉双语复杂短语平行句对处理后的结果样例 .....	(164)
表 6—9 经过词性标注处理后的英汉双语语料 .....	(167)
表 6—10 英汉双语复杂短语平行语料名词类别分布 .....	(169)
表 6—11 英汉双语复杂短语平行语料动词类别分布 .....	(170)
表 6—12 英汉双语复杂短语平行语料形容词类别分布 .....	(170)
表 6—13 三组英汉双语词汇和词性组合的语料规模 .....	(176)

# 第一章

## 引言

### 一 课题提出

课题的研究目标是基于面向网络构建的英汉双语平行语料库，结合自然语言处理和文本挖掘的相关技术和语料资源，使用情报学、机器学习和统计的相关方法，通过比较和使用英汉双语的不同语言特征，从由英汉双语语料组成的非结构化文本中挖掘或抽取出词汇、简单短语结构和复杂短语结构等句法层级上的语言知识。该课题的提出，主要出于以下三个方面的考虑。

(1) 随着自然语言处理和文本挖掘技术的发展，从非结构化文本中挖掘和抽取相应的专门或通用知识以便更好地服务于基础和应用研究日益成为一种趋势。课题从更微观的句法级这一序列层级<sup>①</sup>入手，针对英汉双语语料这一非结构化文本对相应的知识进行挖掘或抽取研究，不仅是在这一大的研究趋势下的一种探究，而且是对整合英汉双语两种语言下的数据预处理、算法和评价方案以及特征在具体知识挖掘和抽取上的一个尝试。

(2) 从英汉双语平行语料库这一非结构化数据集合中挖掘<sup>②</sup>和抽取句法级知识在理论研究和具体应用上具有巨大的需求。构建的英汉双语平行语料库可以直接向相关的研究提供素材，从英汉双语平行语料中挖

① 朱德熙：《语法讲义》，商务印书馆 1983 年版，第 15—28 页。

② Feldman, R., & Sanger, J., *The Text Mining Handbook*, London: Cambridge University Press, 2009, pp. 5—7, p. 225.

## 2 面向非结构化文本的知识发现

掘的词汇级句法知识不仅可以为歧义问题的解决提供相关的理论指导，而且可以对比英汉双语在词汇句法知识上的差异，从而满足解决英汉双语机器翻译、辅助机器翻译的需求。同样地，从英汉双语平行语料库获取的简单短语级和复杂短语级知识更可以直接满足树库构建、信息检索、知识推送服务等的需求。

(3) 已有的英汉双语自然语言处理中的语料和技术与文本挖掘中的聚类算法确保了本研究的可行性。正是有目前标注比较成熟的英语和汉语树库，在进行词汇级知识挖掘的过程中才能完成多角度、各层次的验证，而本课题的研究之所以能得以开展，主要在于自然语言处理技术的发展和逐步成熟，同时各种聚类模型的开发和应用也为研究提供了相应的条件。

综上所述，在文本挖掘这一大的研究趋势下，结合理论和应用的具体需求，在自然语言处理和文本挖掘的相应资源、技术和方法的基础上，本书选取了从英汉双语平行语料中挖掘词汇级、简单短语级和复杂短语级这三个层面的句法知识的研究。

## 二 研究意义

基于英汉双语平行语料库的句法级知识挖掘和抽取对情报学、自然语言处理和相应的研究方法探究具有重要的意义。

### (一) 情报学方面

从具体研究内容出发，结合情报学的相关研究领域，本研究关于英汉双语平行语料库构建、词汇级、简单短语级、复杂短语级的知识挖掘和抽取对知识库构建、知识服务、信息检索、信息计量<sup>①</sup>等研究在某种程度上具有重要的意义。

#### 1. 知识库构建和服务方面

面向网络构建的英汉双语平行语料库，不仅可以为英汉双语知识库的构

<sup>①</sup> 邱均平：《信息计量学》，武汉大学出版社 2007 年版，第 1—35 页。