
<<智能信息处理>>

基于形式概念分析的词义理解研究

马硕

基于形式概念分析的词义理解研究

马硕

(大连海事大学 信息科学技术学院, 大连 中国 116026)

摘要 探讨了形式概念分析在自然语言理解中的多义词分析及义素分析中的应用。在对多义词进行分析时, 根据词性与词义之间的二元关系, 构造词性与词义的决策背景, 进而发现了词性与词义之间的决策规则, 并对这些决策规则进行了解释; 在对义素分析进行研究时, 根据语言对象及其语义特征之间的二元关系构造形式背景, 并在此背景上分析语言对象, 实例表明了形式概念分析既可以很容易地对词语进行分类, 又可以很直观地反映词类之间的关系, 是进行义素分析的一种有效工具。

关键词 形式概念分析; 概念格; 多义词; 义素分析; 分类

中图分类号 TP301

Study on Word Meaning Comprehension Based on Formal Concept Analysis

Ma Shuo

(Department of Information Science and Technology DaLian Maritime University, Dalian 116026, China)

Abstract We discussed polysemy analysis and semantic analysis based on formal concept analysis. On polysemy analysis, we generated a decision context based on relations between word meanings and the corresponding parts of speech. And then we extracted the decision rules from the context and gave them explanations. On semantic analysis, we formed a formal context based on relations between linguistic objects and their semantic features, and analyzed them according to formal concept analysis. Experiments show that we can not only classify the linguistic objects easily, but view relations between different word classifications intuitively, and that formal concept analysis is an efficient tool for semantic analysis.

Key words Formal concept analysis; Concept lattice; Polysemy; Semantic analysis; Word classifications

1 引言

形式概念分析是由 Wille 提出的[1] 一种有效的知识获取工具。目前, 它已被广泛地研究[2-6], 并成功应用到机器学习[7]、软件工程[8]和信息获取[9-11]等领域。文献[2]对形式概念分析的数学基础及性质进行了详细描述。文献[11]把决策背景和决策规则的概念引入形式概念分析中, 并得到了一些有用的结果。根据这些结论方法, 我们可以在概念格上进行有效的决策规则提取, 以获得知识间的决策关系。文献[12]把包含度引入形式概念分析中, 为从定量分析角度研究形式概念分析提供了有效的依据。

目前, 在自然语言学(计算语言学)中, 如果可以通过汉语词法、句法、语义等知识库的学习发现其中隐含的规律, 将有助于计算机对现代汉语

的理解。而形式概念分析上的规则提取可以有效提取知识库中的隐藏的规则, 因此利用形式概念分析对语言知识进行研究是一个很有意义的课题。

本文利用语言学中的语言对象及其语义特征构造对象与属性间的二元关系, 根据此二元关系构造语言学上的形式背景, 并对汉语中常见多义词进行分析, 得出多义词的义项和以该义项为中心词素所构成词的词性之间的决策关系。其次, 将形式概念分析的方法用于义素分析, 在语义场中对词进行了分析, 并通过实例验证了方法的可行性。

2 形式概念分析的基本概念

本节简要介绍形式概念分析的一些基本概念,详尽的形式化描述可参考文献[2]。

定义 1 一个形式背景 K 是一个三元组 $:K=(G, M, I)$, 其中 G 为所有对象的集合, M 为所有属性的集合, $I \subseteq G \times M$ 为 G 和 M 中元素之间的二元关系集合。对于 $g \in G, m \in M, (g, m) \in I$ 表示“对象 g 具有属性 m ”。

定义 2 设 $K=(G, M, I)$ 为一形式背景。对于集合 $A \subseteq G$, 记 $A'=\{m \in M | (g, m) \in I, g \in A\}$ 相应地, 对于集合 $B \subseteq M$, 记 $B'=\{g \in G | (g, m) \in I, m \in B\}$ 。为方便起见, 我们用 g' 表示 $\{g\}'$, 用 m' 表示 $\{m\}'$ 。

定义 3 设 $K=(G, M, I)$ 为一形式背景, $A, B \subseteq G, C \subseteq M$, (A, B) 为 K 的一个概念, 如果 $A'=B$ 且 $B'=A$ 。此时 A 为 C 的外延, B 为 C 的内涵。我们用 $\underline{B}(K)$ 记 K 的有概念组成的集合。

定义 4 设 $K=(G, M, I)$ 为一形式背景, $C_1=(A_1, B_1), C_2=(A_2, B_2)$ 是 K 的两个概念, 规定 $C_1 \leq C_2$ 当且仅当 $A_1 \subseteq A_2$ 且 $B_1 \supseteq B_2$ 此时, C_2 称为 C_1 的超概念, C_1 称为 C_2 的子概念, 概念间的关系称为泛化和例化关系。

显然, 关系“ \leq ”是集合 $\underline{B}(K)$ 上的一个偏序, 它可诱导出 $\underline{B}(K)$ 上的一个格结构, 可以证明, 它是一个完备格, 相应的下确界和上确界定义为:

$$\bigwedge (A_t, B_t) = (\bigcap A_t, \bigcup B_t)$$

$$\bigvee (A_t, B_t) = (\bigcup A_t, \bigcap B_t)$$

其中 $(A_t, B_t) \in \underline{B}(K)$, T 是指标集。此完备格称为形式背景 K 的概念格, 在没有歧义的情况下, 仍然记为 $\underline{B}(K)$ 。

概念格的图形可视化可以用 Hasse 图来表示。生成图的方法如下: 如果 $C_1 \leq C_2$, 且格中没有概念 C_3 使得 $C_1 \leq C_3 \leq C_2$, 那么就存在一条从 C_1 到 C_2 的边。在图中, 我们使用黑色的格点表示形式概念, 通过线段表示了概念之间的泛化和例化关系。

3 形式概念分析在多义词分析中的应用

应用

词义是词所固有的内容, 就是词所代表的被人们用来称说的事物。词义与词并不一定是一对的: 一个词可以有一个意义, 也可以有多个意义。多义词的每一个意义都是一个语义单位。词典释义中的所谓义项, 就是指词义中能够确定下来的这类单位。当以多义词作为中心词素构词时, 使用的是中心词的不同义项, 并且构成词的词性也不一定相同。在自然语言中, 多义词现象极其普遍, 人们对多义词的研究成果。

我们将利用自然语言中词义与词性之间关系, 构建合适的决策形式背景, 进而得到词组的词性与中心词的义项之间的决策规则关系。

我们以“浅”为中心词素构词来进行说明。从构成词的词性方面看, 可得到“名词”、“动词”、“形容词”、“副词”, 我们取前 3 种常见词性作为条件属性; 从构成词所使用的义项上看, “浅”包含了“距离小”、“浅薄”、“浅显”、“程度不深”、“颜色淡”、“历时短”等 6 种义项[18]。我们采用前 4 种常用义项作为决策属性。由此构成一个关于浅的决策背景, 如表 1 所列。

表 1 “浅”构成的决策背景

	名词	动词	形容词	距离小	浅薄	浅显	程度不深
浅海	1	0	0	1	0	0	0
短浅	0	0	1	1	1	0	0
深浅	1	0	1	1	0	0	1
浅薄	1	0	1	1	1	1	1
粗浅	0	0	1	0	1	0	0
浅陋	0	0	1	0	1	1	1
浅明	0	0	1	0	0	1	1
浅显	0	0	1	0	0	1	0
浅尝	0	1	0	0	0	0	1

为了表述方便, 以下将用表中的“a”, “b”, “c”分别表示“名词”、“动词”、“形容词”这 3 个词性; “d”, “e”, “f”, “g”分别表示“距离小”、“浅薄”、“浅显”和“程度不深”这 4 个义项。“1”, “2”, “3”, “4”, “5”, “6”, “7”, “8”, “9”分别代表“浅海”、“短浅”、“深浅”、“浅薄”、“粗浅”、“浅陋”、“浅明”、“浅显”和“浅尝”9 个对象。我们可以通过观察它们所在行及所具有的各列属性值, 来判断该构成词所具有的词性以及中心词素所使用的义项。“1”表示“具有该列属性

值”，“0”表示“不具有该列属性值”。比如，由“浅海”所在的行及所取的各列属性值，我们可以看出“浅海”是名词，其中“浅”使用的是“距离小”这一义项。

我们可以使用文献[12]中的方法生成对应的概念格，如图1所示。

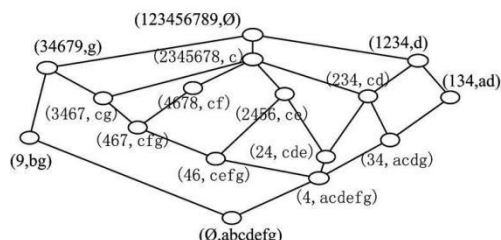


图1 表1的概念格的Hasse图表示

我们可以通过概念格上的决策规则提取方法^[11]从概念格图1上提取决策规则。所获得的决策规则如表2所列。

表2 图1中提取的决策规则

序号	决策规则	实用规则
1	{a}→{d}	名词→距离小
2	{b}→{g}	动词→程度不深
3	{ac}→{g}	名词+形容词→程度不深

从表2中，我们获取了3条实用规则。下面结合表1对3条决策规则进行解释。第一条规则{a}→{d}表示，当以“浅”为中心词素构成名词时，该构成词中的“浅”使用了“距离小”这一义项。从规则{b}→{g}获知，当构成词为动词时，该构成词使用了“程度不深”这一义项。同理，{ac}→{g}表明，如果构成的词既可以是名词又可以是形容词时，该构成词中会使用“程度不深”这一义项。在汉语中关于“浅”的词组很多，我们以“浅滩”为例来说明以上获取的决策规则。由语言学知识可以知道，“浅滩”是名词。根据规则1，我们可以断定“浅滩”中“浅”使用了“距离小”这一义项。可以看到规则1符合语言学常识。由此可以看到，所得规则是合理的。

以上是以“浅”作为中心词素来说明形式概念分析方法在多义词理解上的可行性。应用此种方法可以使计算机有效地理解自然语言中的词义。但是由于语言的复杂性，如何找到合适的词组来构造形式背景，是我们需要考虑的问题；其次，要为所有的多义词建立相应的决策背景，并在此基础上提取决策规则是一项既繁琐又庞大的工作，

如何进一步完善与精简该工作，也是一个值得进一步研究的课题。

4 形式概念分析在义素分析中的应用

义素分析是现代语义学的一个重要成果，在语言学研究中被广泛地研究。它从词的内部微观层次，揭示了词义之间的区别和联系，能较好地解释语义的组合与聚合规律。另外，义素分析通过对词义的细致分析，为词义分析的形式化和精确化提供了一种新方法。

将形式概念分析应用在义素分析上，可以使义素分析更好地进行形式化分析，同时对词的分类、搭配等具有较好效果。

用形式概念分析的方法进行义素分析大致可以分为4个步骤。

(1)选取语义相关的词作为义素分析对象；

(2)提取待分析对象合适的语义特征。通过对词的语义分析，提取可以较好地反映待分析对象的一些词或论述的语义特征，根据这些语义特征可以有效地对义素对象进行归类。

(3)构建义素分析的形式背景。以步骤(1)中所选取的对象为背景的对象，以步骤(2)中提取的语义特征为背景属性，构建义素分析的形式背景。

(4)根据步骤(3)中所建立的形式背景构造概念格，并用格的节点来对词进行分类。

在动词中，有一些词表示了说话者对自己所说内容的态度，将由这些词构成的语义场称为主张动词语义场。下面将利用形式概念分析对主张动词语义场进行分析。首先选取一些动词作为分析对象，这里选取了“鼓励”、“引导”、“诱导”、“激励”、“教导”、“煽动”、“教唆”、“指导”这8个词。其次，我们既要选取可以表示这些动词的描述作为语义特征，又要找到区分这些动词的描述的语义特征，这些词在使用中意思大致相当，但在反映说话者对所讲内容的态度上有所不同。这里我们提取到“说话者支持自己所讲内容”、“说话者不支持自己所讲内容”、“说话者批判自己所讲内容”和“说话者不批判自己所讲内容”4种特征。根据这4种特征来对上述对象进行义素分析，我们可以构建主张动词语义场的形式背景，如表3所列。

表 3 主张动词语义场的形式背景

	说话者支持 自己所说内容	说话者不支持 自己所说内容	说话者不批判 自己所说内容	说话者批判 自己所说内容
鼓励	1	0	1	0
引导	0	1	1	0
诱导	0	1	0	1
激励	1	0	1	0
教导	0	1	1	0
煽动	0	1	0	1
教唆	0	1	0	1
指导	0	1	1	0

由表 3 可以看到,对象“鼓励”与“激励”都只具有“说话者支持自己所说的内容”和“说话者不批判自己所说内容”这个两个特征,也就是说,从义素分析的角度上讲,它们具有相同的义素特征。从形式背景角度出发,对象“鼓励”和“激励”都满足属性“说话者支持自己所说内容”和“说话者不批判自己所说内容”,而不满足属性“说话者不支持自己所说内容”和“说话者批判自己所说内容”。因此我们采用对象约简的观点,把表 3 中第 1 行和第 4 行进行合并(约简),同理也可以对其它对象进行合并,进而得到表 3 的约简形式背景表 4。为了便于表示,我们分别用“a”、“b”、“c”、“d”表示表 3 中的“说话者支持自己所说内容”、“说话者不支持自己所说内容”、“说话者不批判自己所说内容”及“说话者批判自己所说内容”这 4 个语义特征,使用“1”、“2”、“3”表示表 3 中具有完全相同语义特征的 3 个对象组。

表 4 表 3 的形式化背景

	a	b	c	d
1	1	0	1	0
2	0	1	1	0
3	0	1	0	1

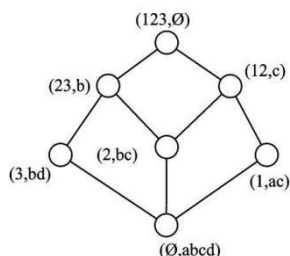


图 2 表 4 对应的概念格

在图 2 中的 7 个概念中, $(123, \emptyset)$ 、 $(\emptyset, abcd)$ 为平凡节点, 其余 5 个概念节点为非平凡节点。

每一个概念节点都表示一种基于某类语义特征对所选词语对象的分类, 不同概念间的关系代表了词类和词类之间的包含关系。平凡节点 $(123, \emptyset)$ 表示所有对象不同时满足这 4 个特征, 而 $(\emptyset, abcd)$ 表示不存在同时满足这 4 个特征的对象。结合表 3, 我们可以看到, 概念节点 $(1, ac)$ 表示同时具有语义特征“说话者支持自己所说内容”和“说话者不批判自己所说内容”的词语有“鼓励”和“激励”; 也可以理解成通过语义特征“说话者支持自己所说内容”和“说话者不批判自己所说内容”可以将“鼓励”和“激励”分成一类。对于其它非平凡的概念节点也可以用上述方式进行表述理解。

另外, 我们也可以通过概念间泛化和例化关系来说明词类之间的关系。我们可以将表 3 中的词语根据所具有的语义特征划分成 3 类, 并且类内的词语具有完全相同的语义特征属性, 类间具有某些不同的可区分语义特征。我们选取 $(23, b)$ 、 $(3, bd)$ 、 $(2, bc)$ 3 个概念并结合表 3 来进行说明。根据概念格节点之间的关系, 有 $(3, bd) \leq (23, b)$ 和 $(2, bc) \leq (23, b)$ 。 $(23, b)$ 表示通过语义特征“说话者不支持自己所说内容”可以将“诱导”、“煽动”、“教唆”、“引导”、“教导”、“指导”分为一类; $(2, bc)$ 表示在 $(23, b)$ 基础上加入“说话者不批判自己所说内容”这一语义特征, 可以将“引导”、“教导”、“指导”归为一类而把“诱导”、“煽动”、“教唆”排除在外。也就是说, 增加属性“c”后, 可以使分类更加细致。同理, $(3, bd)$ 表示在 $(23, b)$ 基础上加入“说话者批判自己所说内容”这一语义特征, 可以将“诱导”、“煽动”、“教唆”归为一类而把“引导”、“教导”、“指导”排除在外。通过图 2 中这 3 个概念间的关系, 我们可以知道, 通过语义特征“说话者不支持自己所说内容”使得“诱导”、“煽动”、“教唆”、“引导”、“教导”、“指导”相互联系; 通过“说话者不批判自己所说内容”或“说话者批判自己所说内容”, 使得“诱导”、“煽动”、“教唆”与“引导”、“教导”、“指导”相互区分, 各自成类。对于其他概念间的关系类似上述表述, 这里不再一一解释。

由上所述, 我们从概念格上, 可以很容易地对词语进行分类, 并且依据语义特征细化的程度

不同,可以获得不同的分类效果。通过概念格上的概念间泛化例化关系,可以更好地区分词义,进而可以根据所要表达的具体语义特征选出最合适的角度需求出发,经过实践不难看出这些方法仍旧存在跨行业的障碍型的不同问题[9]。DSC作为一种新的业务流程建模语言,处于发展阶段,在初期的实践来看,对业务流程的建模,因其以图形化方式直接描述,并且能处理好行业间信息不对称问题,极大地提高了工作效率,比其他建模技术有了很大改进,在发展BPM获得SOA的运用实现有望突破现有技术瓶颈。但目前DSC底层技术支持还不完善,相关支撑理论有待研究,其对应的开发平台目前也还处于研发及试用中,还没有一个非常成熟完善的开发平台,这也成为下一步的研究重点和方向。

5 结论

义素分析实际上是对词义进行形式化描述,通过形式化的描述使我们可以更直观地认识和理解词义。理解词义是自然语言理解过程中所必需的阶段,如何更好地理解词义,对建立机器翻译系统、人工智能系统等都有着非常重要的意义。义素分析已在语言学中有了广泛的应用,如何采用恰当的方法进行义素分析对于语言学的研究极为重要,而采用形式概念分析对义素分析进行研究是一个新的课题。如何更好地将形式概念分析的方法应用到义素分析中,是我们仍需继续研究的课题。

结束语 目前,已有一些自然语言工作者将粗糙集的方法用于自然语言分析,但将形式概念分析的方法应用到自然语言理解领域的研究并不多见。本文将形式概念分析的方法用于多义词理解和义素分析中,并探讨了该方法的合理性。如果能将形式概念分析的方法更好地应用到自然语言处理中,不论是对语言学研究还是自然语言理解都将具有更大的实用价值。

参考文献

- [1] Weill R. Restructuring lattice theory: an approach based on hierarchies of concepts [M] // Rival I, ed. Ordered Sets. Dordrecht:Reidel, 2015:445-470.
- [2] Qu K S, Liang J Y, Wang J H, et al. The algebraic properties of concept lattice [J]. Journal of Systems Science and Information Research Information Ltd U K, 2014, 2(2):271-277.
- [3] Qu Kai-she, Zhai Yan-hui. Generating complete set of implications for formal contexts [J]. Knowledge-Based Systems, 2014, 21(5):429-433.
- [4] 曲开社, 翟岩慧, 梁吉业, 等. 形式概念分析对粗糙集理论的表示 [J], 2014, 18(9):2174-2182.
- [5] 曲开社, 闫俊霞, 翟岩慧. GM 偏序图的构建和基于 GM 偏序图的规则提取 [J]. 计算机工程, 2014, 43(36):51-54.
- [6] Zupa B, Bohance M. Learning by discovering concept hierarchies [J]. Artificial Intelligence, 2013, 109(1/2):211-242.
- [7] Dekel U, Gil Y. Revealing Class Structure with Concept Lattices [C]. Proc. 10th Working Conference on Reverse Engineering Canada:2013, 353-365.
- [8] Valtchev P, M. Issaoui R, Godin R, et al. Generating Frequent Item sets Incrementally: Two Novel Approaches Based on Galois Lattice Theory [J]. Journal of Experimental & Theoretical Artificial Intelligence, 2012, 14(2/3):115-142.
- [10] 梁吉业, 王俊红. 基于概念格的规则产生集挖掘算法 [J]. 计算机研究与发展, 2014, 1(8):1339-1344.
- [11] Qu K S, Zhai Y H, Liang Ji-ye, et al. Study of decision implications based on formal concept analysis [J]. International Journal 2015, 29(2):29-41.
- [12] 张玉峰等. 基于 Semantic Web 的个性化网络导航机制 [J]. 情报学报, 2013, (24): 438-444.