

组合原则和自然语言虚化成分

邹崇理^{1, 2}

(1.四川师范大学 逻辑与信息研究所,成都 610066;2.中国社会科学院 哲学所,北京 100732)

摘要:计算机人工智能时代最重要的任务之一是自然语言的信息处理,逻辑语义学则是其基础理论,而组合原则又是逻辑语义学的基本原则,表现为部分决定整体的函项思想。自然语言的虚化成分是非自然语言复合表达式中对整体意义不起作用的那些部分,自然语言违反组合原则的情况表现为句法和语义的不对应,意味着决定整体意义的“部分”这个概念应该受到限制,组合原则的经典表述在非自然语言的某些场合受到挑战。就自然语言的某些语义领域而言,限制性的组合原则概念是关于组合原则具体精准的表述。

关键词:逻辑语义学;组合原则;自然语言;虚化成分

中图分类号:B815.3;O141 **文献标志码:**A **文章编号:**1000-5315(2017)01-0005-05

逻辑学是人文社会科学和自然科学共同的基础学科。1974年,联合国教科文组织规定的七大基础学科依次为数学、逻辑学、天文学和天体物理学、地球科学和空间科学、物理学、化学、生命科学。由此可见,逻辑学在人类整个知识结构中的基础地位。20世纪中叶以来,计算机科学技术的迅猛发展,导致席卷全球的信息革命,而自然语言是信息的重要载体之一,信息革命离不开自然语言的计算机处理。逻辑作为人类知识结构的基础,同样在非自然语言的信息处理领域内发挥巨大作用。

自然语言的计算机信息处理过程是:首先,把需要研究的语言学问题用数学或逻辑的形式严密而规整地表示出来;其次,把这种严密而规整的表述表示成算法,建立各种自然语言处理系统;第三,对自然语言处理系统进行评测,不断改进质量和性能。逻辑语义学关于自然语言的研究主要对第一步骤发生作用,逻辑语义学是非自然语言信息处理的基础理论学科,前者为后者提供了诸多重要的操作工具和指导思想。

从逻辑语义学角度对自然语言进行形式化研究,组合原则是其重要的方法论,那么怎样理解组合原则?组合原则是什么?

一 什么是组合原则

组合原则是逻辑语义学的基本原则。组合原则直观表述为:如果表达式E依据某个句法规则由部分 E_1 和 E_2 所构成,则E的语义 $M(E)$ 是依据某个语义规则把 E_1 的语义 $M(E_1)$ 和 E_2 的语义 $M(E_2)$ 合并起来而获得的。举例来说,表达式“伟大祖国”的语义是由形容词“伟大”的语义限制名词“祖国”的语义而获得。表达式“戴眼镜的女孩”的语义是以由定语从句“某某戴眼镜”的语义和名词“女孩”的语义合并而成。

严格讲,组合原则意味:一个复合表达式的语义是由部分的语义贴合这些部分的句法运算的意义所构成的函项。组合原则的数学定义深刻揭示了这样的特征^{[1]526}。

令 $A = \langle A, F \rangle$ 和 $B = \langle B, G \rangle$ 都是代数,映射 $h: A \rightarrow B$ 是同态的,当且仅当,存在一个映射 $h': F \rightarrow G$ 使

收稿日期:2016-02-26

基金项目:国家社科基金重大招标项目“自然语言信息处理的逻辑语义学研究”(10&ZD073)。

作者简介:邹崇理(1953—),男,四川成都人,四川师范大学特聘教授、逻辑与信息研究所学术委员会主席,中国社会科学院博士生导师,中国逻辑学会会长,主要研究自然语言逻辑。

得对所有 $f \in F$ 和所有 $a_1, \dots, a_n \in A$ 都有:

$$h(f(a_1, \dots, a_n)) = h'(f)(h(a_1), \dots, h(a_n))$$

在自然语言领域, A 是句法代数, B 是语义代数, h 就是从句法生成到语义组合的意义指派。 A 是句法表达式的集合, B 是语义值的集合, F 是句法算子的集合, G 是语义算子的集合。 f 是 F 中的某个算子, a_1, \dots, a_n 是 A 中的 n 个表达式, $h'(f)$ 是 G 中对应 f 的语义算子, $h(a_1), \dots, h(a_n)$ 是 B 中对应 a_1, \dots, a_n 的 n 个语义值。复合表达式 $f(a_1, \dots, a_n)$ 是 f 对 a_1, \dots, a_n 进行句法生成的结果, 其语义 $h(f(a_1, \dots, a_n))$ 就是语义算子 $h'(f)$ 对 n 个部分语义 $h(a_1), \dots, h(a_n)$ 进行运算的结果, 是语义算子贴合部分语义进行运算的函项。

例子解读: 令句法代数 A 的论域 $A = \{\text{张三}, \text{李四}, \text{散步}, \text{学习}, \text{张三散步}\}$, 语义代数 B 的论域 $B = \{a, b, \{a, b\}, 1, 0\}$ 。再令 $f(\text{张三}, \text{散步}) = \text{张三散步}$, $h(\text{张三}) = a, h(\text{散步}) = \{a, b\}, h'(f) = g$ 。对此进行语义指派得: $h(f(\text{张三}, \text{散步})) = h'(f)(h(\text{张三}), h(\text{散步})) = g(a, \{a, b\}) = 1$, 当且仅当 $a \in \{a, b\}$ 。

可以看出组合原则具有两个重要特征: (1) 复合表达式语义组合“ $h'(f)(h(a_1), \dots, h(a_n))$ ”的根源依据是复合表达式的句法生成“ $f(a_1, \dots, a_n)$ ”, 这就是句法和语义对应的思想; (2) 复合表达式的语义不仅依靠其部分的语义“ $h(a_1), \dots, h(a_n)$ ”, 还取决于合并这些部分的句法生成的意义“ $h'(f)$ ”。

组合原则是现代逻辑的基石, 在构造逻辑系统中起到方法论的作用。组合原则要求逻辑系统中每个句法(语形)形成规则必须对应一个语义解释规则。命题逻辑严格遵循了意义的组合原则, 令 $//$ 为意义指派函项 h , 则有:

- Syn1. 原子公式 $p_1, p_2, \in \text{Form}$;
- Syn2. 若 $\varphi \in \text{Form}$, 则 $(\neg\varphi) \in \text{Form}$;
- Syn3. 若 $\varphi, \psi \in \text{Form}$, 则 $(\varphi \rightarrow \psi) \in \text{Form}$ 。
- Sem1. $\| p_i \| \in \{0, 1\}$;
- Sem2. $\| (\neg\varphi) \| = 1$ 当且仅当 $\| \varphi \| = 0$;
- Sem3. $\| (\varphi \rightarrow \psi) \| = 1$ 当且仅当 $\| \varphi \| = 0$ 或 $\| \psi \| = 1$ 。

句法规则 Syn1 对应语义规则 Sem1, Syn2 对应 Sem2, Syn3 对应 Sem3。显然, 这里复合表达式的语义依据其部分表达式的语义, 复合表达式的所指是部分表达式所指的函项。如 $\| (\varphi \rightarrow \psi) \| = \| \rightarrow(\varphi, \psi) \| = h'(\rightarrow)(\| \varphi \|, \| \psi \|) = 1$, 当且仅当 $\| \varphi \| = 0$ 或 $\| \psi \| = 1$ 。

自 20 世纪 70 年代初开始, 现代逻辑的方法扩展

延伸到自然语言的研究领域, 形成了以蒙太格语法 (Montague Grammar)^{[2][247-270]} 为首的逻辑语义学群体, 组合原则自然也成为逻辑语义学的灵魂。

蒙太格语法是强调组合原则的逻辑语义学理论。在其构造的三个英语部分语句系统那里, 句法和语义处处严格对应。以 PTQ 系统为例, 17 条句法规则对应 17 条语义翻译规则^{[2][247-270]}。每条翻译规则体现出: 复合表达式的翻译是其部分表达式翻译的函项。句子、动词短语和名词短语三类合取复合表达式的句法规则及其翻译规则如下:

- Syn1. 若 $\varphi, \psi \in P_I$, 则 $F_8(\varphi, \psi) = \varphi \text{ and } \psi \in P_I$;
- Syn2. 若 $\delta, \gamma \in P_{IV}, F_8(\delta, \gamma) = \delta \text{ and } \gamma \in P_{IV}$;
- Syn3. 若 $\alpha, \beta \in P_T, F_9(\alpha, \beta) = \alpha \text{ or } \beta \in P_T$ 。

- Tra1. 若 φ, ψ 分别翻译成 φ', ψ' , 则 $\varphi \text{ and } \psi$ 翻译成 $[\varphi' \wedge \psi']$;
- Tra2. 若 δ, γ 别翻译成 δ', γ' , 则 $\delta \text{ and } \gamma$ 翻译成 $\lambda x[\delta'(x) \wedge \gamma'(x)]$;
- Tra3. 若 α, β 分别翻译成 α', β' , 则 $\alpha \text{ or } \beta$ 翻译成 $\lambda P[\alpha'(P) \vee \beta'(P)]$ 。

翻译起意义指派函项的作用。令 T 是翻译函项, 拿 Tra2. 来说, $T(\delta \text{ and } \gamma) = T(\text{and}(\delta, \gamma)) = h'(\text{and})(T(\delta), T(\gamma)) = \lambda x[\delta'(x) \wedge \gamma'(x)]$ 。复合表达式的翻译依赖部分表达式的翻译。

组合原则的作用还体现在更多的领域内。

在计算机科学那里, 连接许多通信处理器的大网络技术发展很快, 人们特别关注超大系统的行为。在有关研究中, 组合原则起到很大的作用: 牵涉整个系统行为的证明应该是各个处理器的证明的函项。这方面的介绍参见文献。

组合原则在形式翻译领域作用更大。为了考察逻辑系统之间的关系, 比较表达力的大小以及获得相对的协调性, 人们往往设立符合组合原则的翻译程序。最著名的例子是 Gödel 把直觉主义逻辑转换成模态逻辑的翻译。在直觉主义逻辑那里, 联接词具有一种构造性解释, 如 $\varphi \rightarrow \psi$ 被解释成: 给定 φ 的证明, 据此构造 ψ 的证明。令 Tr 为翻译函项, 翻译程序定义为:

- a. $\text{Tr}(p) = \Box p$ 对原子命题 p
- b. $\text{Tr}(\varphi \vee \psi) = \text{Tr}(\varphi) \vee \text{Tr}(\psi)$
- c. $\text{Tr}(\varphi \wedge \psi) = \text{Tr}(\varphi) \wedge \text{Tr}(\psi)$
- d. $\text{Tr}(\varphi \rightarrow \psi) = \Box p(\text{Tr}(\varphi) \rightarrow \text{Tr}(\psi))$

复合表达式的翻译, 依据部分表达式的翻译来确定。Gödel 的翻译是一种组合翻译, 逻辑系统之间大量的组合翻译可以参见 Epstein 的著述。

左右范畴 A 和 C 构成一个所谓“省略槽”的复合范畴 $[A\{B\}C]$, 即得: $[A\{B\}C] \rightarrow (A \cdot C)$, 这就是新的推演工具。三元复合范畴 $[A\{B\}C]$ 是删去虚化成分的起点, 据此揭示包含虚化成分的表达式的语义特征, 如“迅速地跑步”, “美丽的女孩”和“玩得高兴”, 其中的“地”、“的”和“得”所属范畴就是起间隔虚化作用的 B。三元复合范畴的语义解释如下:

$$v([A\{B\}C]) = \{x \mid \exists yz[Sxyg(B)z \& y \in \parallel A \parallel \& z \in \parallel C \parallel]\}$$

按照上述定义: $Sxyg(B)z$ 意味: x 是 $y, g(B)$ 和 z 毗连的结果, 具有语义所指的符号串 y 和 z 分别属于 A 和 C, 而 $g(B)$ 指起虚化作用的符号串, B 是 A 和 C 之间的虚化范畴(对应的语义所指为空逻辑式)。于是有:

$$\text{限制 0: } \forall B \forall x[x \sim g(B) \Rightarrow x \in v(B)]$$

这里 $x \sim g(B)$ 的直观理解是: x 是 $g(B)$ (起虚化作用的符号串)。限制 0 表明起虚化作用的符号串是 $\parallel B \parallel$ 中的元素。

于是, 我们提出基于 $[A\{B\}C]$ 的范畴逻辑系统^{[5]370-381}。其公理是:

- 公理 0: $A \rightarrow A$
- 公理 1: $A \cdot B \leftrightarrow B \cdot A$
- 公理 2: $[A\{B\}C] \rightarrow (A \cdot C)$
- 公理 3: $D \cdot [A\{B\}C] \rightarrow [(D \cdot A)\{B\}C]$
- 公理 4: $[A\{B\}C] \cdot D \rightarrow [(A \cdot D)\{B\}C]$
- 公理 5: $D \cdot [A\{B\}C] \rightarrow [A\{B\}(D \cdot C)]$
- 公理 6: $[A\{B\}C] \cdot D \rightarrow [(A\{B\}(C \cdot D))$
- 公理 7: $[A\{B\}C] \cdot [D\{B\}E] \rightarrow [(A \cdot D)\{B\}(C \cdot E)]$

系统的规则有(Lambek 演算的 5 条推演规则):

$$\frac{A \cdot B \rightarrow C}{A \rightarrow C/B} \quad \frac{A \cdot B \rightarrow C}{B \rightarrow A \setminus C}$$

$$\frac{A \rightarrow C/B}{A \cdot B \rightarrow C} \quad \frac{B \rightarrow A \setminus C}{A \cdot B \rightarrow C}$$

$$\frac{A \rightarrow B \quad B \rightarrow C}{A \rightarrow C}$$

此外, 系统还有两条独特的推演规则:

$$\text{规则 6: } \frac{A \rightarrow B}{[A\{D\}C] \rightarrow [B\{D\}C]}$$

$$\text{规则 7: } \frac{A \rightarrow B}{[D\{C\}A] \rightarrow [D\{C\}B]}$$

对构成其他复合范畴的算子, 传承 Lambek 演算 L 系统的语义解释如下:

$$v(A \cdot B) = \{x \mid \exists y \exists z[Rxyz \& y \in v(A) \& z \in v(B)]\}$$

$$v(C/B) = \{y \mid \forall x \forall z[Rxyz \& z \in v(B)] \Rightarrow x \in v(C)\}$$

$$v(A \setminus C) = \{z \mid \forall x \forall v[Rxvz \& v \in v(A)] \Rightarrow x \in v(C)\}$$

按照惯例给出系统的框架语义, 这是一个由三元可及关系 R 和四元可及关系 S 组成的混合框架。系统的语义特色在于下述框架限制:

$$\text{限制 0: } \forall B \forall x[x \sim g(B) \Rightarrow x \in v(B)]$$

$$\text{限制 1: } \forall xyz[Rxyz \Rightarrow Rxzy]$$

$$\text{限制 2: } \forall xyzu[Sxyzu \Rightarrow Rxyu]$$

$$\text{限制 3: } \forall xyzuvw[Rxyz \& Szuwv \Rightarrow t[Sxtwv \& Rtyu]]$$

$$\text{限制 4: } \forall xyzuvw[Rxyz \& Syuvw \Rightarrow t[Sxtvw \& Rtuz]]$$

$$\text{限制 5: } \forall xyzuvw[Rxyz \& Szuwv \Rightarrow t[Sxugt \& Rtyw]]$$

$$\text{限制 6: } \forall xyzuvw[Rxyz \& Syuvw \Rightarrow t[Sxugt \& Rtwz]]$$

$$\text{限制 7: } \forall xyzuvwst[Rxyz \& Syuvw \& Szsvt \Rightarrow \exists ab[Sxavb \& Raus \& Rbwt]]$$

依据上述提供的框架语义解释及其限制, 可以证明系统的可靠性和完全性。可判定性证明也可按照惯例给出^{[5]370-381}。

上文已强调, 系统的最大特色就是公理 2: $[A\{B\}C] \rightarrow (A \cdot C)$ 。意味从 $A: \alpha, B: \emptyset, C: \gamma$ 推出 $A: \alpha, C: \gamma$ 。从句法角度看, 复合表达式“ABC”的部分表达式是“A”、“B”和“C”。公理 2 的潜在显示为: “ABC”即“ $[A\{B\}C]$ ”的整体语义就是“A·C”的语义, 即“ $\alpha(\gamma)$ ”。而这仅仅取决于部分表达式“A”的语义“ α ”和部分表达式“C”的语义“ γ ”, 复合表达式的语义并非如组合原则所要求的是由所有部分表达式的语义来决定。

在自然语言复合表达式中间的部分表达式是虚化成分的条件下, 这时的组合原则就是受限的, 其表述就是: 复合表达式的语义是由除去作为那个虚化成分的部分的语义以外的其他部分的语义贴合这些部分的句法运算的意义所构成的函项。受限组合原则的定义为:

令 $A = \langle A, F \rangle$ 是句法代数和 $B = \langle B, G \rangle$ 是语义代数, 映射 $h: A \rightarrow B$ 是同态的, 当且仅当, 存在一个映射 $h': F \rightarrow G$, 存在 $f \in F$ 并且存在 $a_1, \dots, a_n \in A$ 满足:

$$h(f(a_1, \dots, a_n)) = h'(f)(h(a_1), h(a_{n-1}), h)$$

$(a_{i+1}) \cdots, h(a_n)) (1 < i < n)$

其中, $h(a_i)$ 是作为虚化成分的部分的语义。

例子解读: 动词短语表达式“飞快地跑步”的句法生成: $f(\text{飞快}, \text{地}, \text{跑步})$ 。令“飞快”的逻辑语义为 α , “跑步”的逻辑语义是 γ , “飞快地跑步”的逻辑语义就是 $h(f(\text{飞快}, \text{地}, \text{跑步})) = h'(f)(\alpha, \gamma) = \alpha(\gamma)$ 。从语义角度看, 这里“地”是没有逻辑语义的, 在复合表达式“飞快地跑步”的整体语义组合中不起作用, 是语义虚化的部分表达式。“美丽的姑娘”和“中式的家具”中的

结构助词“的”也都是自然语言中的虚化成分。

三 结论

由于自然语言的丰富多样性, 句法和语义的对应及意义的组合原则往往表现出异彩纷呈的局面。这给人们留下研究的空间, 探讨作为数学概念的组合原则怎样通过具体生动的自然语言而呈现出多种多样的表现形式, 是逻辑语义学介入自然语言信息处理领域所期待的工作, 是逻辑学作为基础工具学科作用于计算机人工智能科学的价值所在。

参考文献:

- [1] Janssen T, Partee. *Compositionality* [C]// Johan van Benthem et al. (eds.). *Handbook of Logic and Language* [M]. Amsterdam: Elsevier, 2011.
- [2] Montague R. *Formal Philosophy* [M]. New Haven: Yale University Press, 1974.
- [3] 张秋成. 类型逻辑语法研究 [M]. 北京: 中国人民大学出版社, 2007.
- [4] 王欣. 类型逻辑语法与现代汉语“是”和“的” [M]. 北京: 北京语言大学出版社, 2009.
- [5] ZOU Chongli et al. The Categorical Logic of Vacuous Components in Natural Language [C]// Van Ditmarsch et al. (eds.). *Logic, Rationality, and Interaction*, LNAI 6953. Berlin: Springer-Verlag, 2011.

Principle of Compositionality and Vacuous Components in Natural Language

ZOU Chong-li^{1,2}

- (1. Institute of Logic and Information, Sichuan Normal University, Chengdu, Sichuan 610066;
2. Institute of Philosophy, Chinese Academy of Social Sciences, Beijing 100732, China)

Abstract: Principle of compositionality is the most important principle for Logical Semantics, a theoretical foundation of NLP (Natural Language Processing), which is one of the most important tasks in the age of artificial intelligence. Principle of compositionality embodies the idea that the meaning of the whole expression is the function of the meanings of its components. However, as a typical phenomenon of anti-syntax-semantics-correspondence, the vacuous components in complex expressions of natural language make no contribution to the meaning of the whole expression. Therefore, the meanings of its parts should be manipulated under a certain restrictions of principle of compositionality. This paper proposes the restricted principle of compositionality as a more accurate expression of the Principle.

Key words: logical semantics; principle of compositionality; natural language; vacuous components

[责任编辑: 帅 巍]