

# 基于实体网格的语篇表示模型研究\*

邹益民<sup>1,2</sup> 张智雄<sup>1</sup> 曲云鹏<sup>2,3</sup>

<sup>1</sup>(中国科学院国家科学图书馆, 北京 100190)

<sup>2</sup>(中国科学院研究生院, 北京 100190)

<sup>3</sup>(中国国家图书馆, 北京 100081)

**[摘要]** 在自然语言处理中, 人们提出了多种语篇表示模型, 以期实现语篇的自动理解和生成, 而其中能够携带实体分布信息、位置信息、语法信息以及指代信息的基于实体网格的语篇表示模型越来越受到人们的关注。本文主要围绕实体网格的概念以及构造方法、实体网格在连贯性评测上的作用进行分析, 对基于实体网格特征修改和丰富的实体扩展方案以及新的实体网格计算方法进行了深入分析和研究。

**[关键词]** 实体网格; 实体网格扩展; 语篇表示; 语篇质量评价

**[分类号]** TP18

## Discourse Representation Model Based on the Entity Grid

ZOU Yi-min<sup>1,2</sup> ZHANG Zhi-xiong<sup>1</sup> and QU Yun-peng<sup>2,3</sup>

<sup>1</sup>(National Science Library, Chinese Academy of Sciences, Beijing 100190)

<sup>2</sup>(Graduate University of Chinese Academy of Sciences, Beijing 100190)

<sup>3</sup>(National Library of China, Beijing 100081)

**[Abstract]** A variety of discourse representation models were proposed in Natural Language Processing to understand and generate discourse automatically. Entity grid is widely used as it can capture the distributional, syntactic and positional information and co-reference among the entities. This paper mainly focuses on the concept and construction methods of the entity grid, and analyses the effects of the entity grid on the evaluation of coherence. The grid computing and extending methods based on the modification and rich feature of the entity in the grid are also in-depth discussed.

**[Keywords]** Entity Grid; Extending the Entity Grid; Discourse Representation; Evaluation of Discourse Quality

### 1 引言

语篇 (Discourse, Text)<sup>1</sup>通常是指由一系列连续的语段或句子构成的语言整体单位, 能够表达一个特定的意图或一组特定的事件, 其各成分之间, 在形式上是衔接 (Cohension) 的, 在语义上是连贯 (Coherence) 的。随着自动摘要和自动问答等一系列自动语篇生成系统成为人工智能和自然语言处理研究领域的热点, 如何对这些自动生成的语篇质量进行评估也是人们关注的焦点, 语篇连贯性分析作为语篇可读性评测的基础, 对其进行自动评测的研究也显得尤为重要。2005年, 麻省理工学院的 Barzilay 和爱丁堡大学的 Lapata 受向心理论 (Centering Theory) 的启发提出了实体网格 (Entity Grid) 语篇表示模型<sup>[1]</sup>, 用于对语篇的局部连贯性进行评测, 该模型可以通过对原始语篇的计算进行自动的构建。克服了以往评测

---

作者简介: 邹益民(1983-), 男, 河南周口人, 博士生, 主要从事文本挖掘和知识组织的研究; 张智雄(1971-), 男, 北京人, 研究馆员, 博士, 主要从事智能信息处理和信息系统的研究; 曲云鹏(1980-), 男, 黑龙江哈尔滨人, 馆员, 博士生, 主要从事信息系统和长期保存的研究。

\* 本文系国家自然科学基金“基于语言网络的文本主题中心度计算方法研究”(项目编号: 61075047) 以及国家十二五科技支撑计划项目“面向外科技文献信息的知识组织体系建设与应用示范”(项目编号: 2011BAH10B00) 课题五“信息资源自动处理、智能检索与 STKOS 应用服务集成”的研究成果之一。

<sup>1</sup> Discourse 在中文中没有现成的词语与其相对应, 通常被译为“篇章”、“语篇”和“话语”等等, 本文采用语篇来表示 Discourse。

方法中依赖于手工制定的规则、算法在有限的领域内有效并且算法可扩展性和可移植性较差等问题。实体网格已成为目前最流行的连贯性评测模型，被广泛应用于自动摘要、句子排序、作文评分、可读性评测、对话分解和故事生成等任务的连贯性评测中。

## 2 实体网格

### 2.1 实体网格语篇表示模型

在实体网格语篇表示模型中，语篇被映射成一个由实体及其语法角色组成的网格，实体网格是一个二维的数组，网格的行对应语篇中的句子而列对应语篇中的实体，语篇中的每一个实体在网格中都对应于其在给定句子中的语法角色，如图 1 和图 2 所示。在对句子粒度的界定上，由于子句的分割将引入相应的噪声，实体网格将主句作为其分析的基本单元。而篇章实体由通用名词（包括：一般名词和专有名词）和命名实体（例如：人名、地名和机构等）组成，在实体网格中用实体的中心词来代替实体。实体网格用“S”、“O”、“X”和“-”来标识实体在句子中的语法角色，其中“S”对应主语、“O”对应宾语、“X”对应非主语和宾语的其它句法角色、“-”标识相应的实体在给定的句子中不存在，并且实体语法角色具有一定优先级： $S > O > X > -$ ，当同一个实体在给定的句子中具有不同的语法角色时，取优先级最高的角色进行标识<sup>[1, 2]</sup>。

1	[The Justice Department] <sub>S</sub> is conducting an [anti-trust trial] <sub>O</sub> against [Microsoft Corp.] <sub>X</sub> with [evidence] <sub>X</sub> that the company <sub>S</sub> is increasingly attempting to crush competitors <sub>O</sub> .
2	[Microsoft] <sub>O</sub> is accused of trying to forcefully buy into [markets] <sub>X</sub> where [its own products] <sub>S</sub> are not competitive enough to unseat [established brands] <sub>O</sub> .
3	[The case] <sub>S</sub> revolves around [evidence] <sub>O</sub> of [Microsoft] <sub>S</sub> aggressively pressuring [Netscape] <sub>O</sub> into merging [browser software] <sub>O</sub> .
4	[Microsoft] <sub>S</sub> claims [its tactics] <sub>S</sub> are commonplace and good economically.
5	[The government] <sub>S</sub> may file [a civil suit] <sub>O</sub> ruling that [conspiracy] <sub>S</sub> to curb [competition] <sub>O</sub> through [collusion] <sub>X</sub> is [a violation of the Sherman Act] <sub>O</sub> .
6	[Microsoft] <sub>S</sub> continues to show [increased earnings] <sub>O</sub> despite [the trial] <sub>X</sub> .

图 1 语法角色标注后的语篇

	Department	Trial	Microsoft	Evidence	Competitors	Markets	Products	Brands	Case	Netscape	Software	Tactics	Government	Suit	Earnings	
1	S	O	S	X	O	-	-	-	-	-	-	-	-	-	-	1
2	-	-	O	-	-	X	S	O	-	-	-	-	-	-	-	2
3	-	-	S	O	-	-	-	-	S	O	O	-	-	-	-	3
4	-	-	S	-	-	-	-	-	-	-	-	S	-	-	-	4
5	-	-	-	-	-	-	-	-	-	-	-	-	S	O	-	5
6	-	X	S	-	-	-	-	-	-	-	-	-	-	-	O	6

图 2 实体网格

注：此图来源于 Barzilay 和 Lapata 的 Modeling Local Coherence: An Entity-based Approach。

具体的实体网格构建过程如图 3 所示，其中实体共指消解和实体句法角色识别是两个关键的步骤，共指消解对具有共指关系的语篇实体进行识别和聚类，实体网格中采用 Ng 和 Cardie 两位学者提出共指消解算法，该算法利用词汇、语法、语义和位置等多个特征来识别两个实体之间是否构成共指关系。为了在实体网格中填充相应的语法角色，实体网格利用 Collins 提出的统计分析器来识别实体在句子中的作用，而被动语态则通过一个小的模式集合来识别。

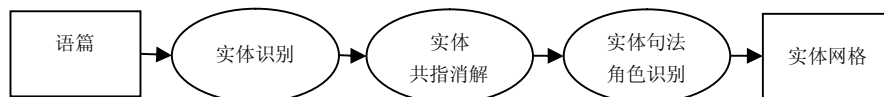


图 3 实体网格的构建过程

实体网格通过图 3 中的构建过程自动将语篇抽象成一组实体转换序列，这种表示模型，能够反映出实体在语篇中的分布信息、位置信息、语法信息以及指代信息，并允许系统通过该模型从给定语料中学习出局部连贯性语篇所具有的特征，而不必基于手工标柱的资源 and 预

先定义的知识库。

## 2.2 实体网络的计算

为了对语篇连贯性进行测量，实体网络模型中一个重要的假设是：“在连贯性的语篇中，实体的分布展现一定的规律性。”，而这种规律性在向心理论，齐普夫定律（Zipf's Law）<sup>[3]</sup>和其它基于实体的语篇理论中也得到了验证，在实体网络中这种分布的规律性，通过网络拓扑结构得到体现。连贯性语篇对应的网络往往包含一些稠密的列（列中有很少的空值）和很多稀疏的列（有很多空值组成），而且在稠密的列中实体往往充当主语或宾语的角色，这种规律性的分布在不连贯的语篇中很少出现。

受向心理论的影响，基于实体网络的分析主要围绕局部实体的转换模式展开，转换模式是一个序列，例如：转换[S,-]在图 1 所示的实体网络中的分布概率为 0.08，实体网络可以看作是一个转换类型的分布，向心理论还认为实体被引入和讨论取决于其在给定语篇中的全局角色，因而，有必要对显著（或重要）实体和非显著（或不重要）实体的转换模式进行了区分，实体网络吸纳了这一思想，通过去除网络中出现频次小于一定阈值的非显著实体来突出显著实体在语篇中的作用<sup>[1,2]</sup>。

在实体网络中，连贯性约束通过实体转换序列隐式地被模式化在网络中，通过特征向量  $\Phi(x_{ij})$  对实体在句子间的转换模式进行编码，其中， $\Phi(x_{ij}) = (p_1(x_{ij}), p_2(x_{ij}), \dots, p_m(x_{ij}))$  表示语篇  $d_i$  的第  $j$  种描述<sup>2</sup>， $m$  表示所有预先定义好的转换模式的数量， $p_1(x_{ij})$  表示转换  $t$  在网络  $x_{ij}$  中的概率。训练集由描述的有序对  $\{x_{ij}, x_{ik}\}$  组成， $x_{ij}$  和  $x_{ik}$  是关于语篇  $d_i$  的不同描述，并且  $x_{ij}$  比  $x_{ik}$  拥有更高的连贯性，不失一般性，假设  $j > k$ ，训练的的目的是找到一个参数向量  $\vec{w}$  使得排序函数  $\Phi(x_{ij})$  最大程度的使得等式  $\vec{w} \cdot (\Phi(x_{ij}) - \Phi(x_{ik})) > 0 \forall j, i, k$  成立。

实体网络通过排序函数  $\Phi(x_{ij})$  和参数向量  $\vec{w}$  对测试集中的任意两个语篇的连贯性进行比较，将连贯性评测转化成了一种学习排序问题，取得了很好的应用效果，目前，它也是最好的句子排序模型的关键组件<sup>[4]</sup>。

## 3 实体网络的扩展

### 3.1 实体网络特征的修改和丰富

实体网络提供了一种新的语篇表示形式，作为一个基本模型，研究者往往根据具体任务的需要对实体网络的特征进行修改和丰富，大致可以分为三种类型：通过新的实体选取方法来进一步挖掘实体在网络中的分布规律、对实体进行语义合并以及在实体网络中融入新的实体特征来满足具体应用的需求。

#### (1) 新的实体选取方法

在标准的实体网络中，通过共指消解，对存在共指关系的实体进行识别和聚类，并利用实体的中心词来代替实体。但 Lapata 等人认为共指消解工具主要是基于在连贯性文本中的训练结果，在有些应用场景下，例如：在对自动摘要的连贯性评测，未必能提高算法的性能<sup>[5]</sup>。在构建实体网络中，其采用了一种简单的方式，将文本中的每个名词对应网络中的一个

<sup>2</sup> 在文本排序任务中，对一个语篇  $d_i$  通过调整其内部的句子顺序，将其分为  $n$  个不同描述的语篇  $d_{i1}, d_{i2}, \dots, d_{in}$ ，并认为原始语篇  $d_i$  的连贯性最好，通过比较不同描述之间的连贯性，来对其进行排序。

单独实体。在这种方式下名词间的精确匹配才会认为是共指关系，命名实体和复合名词被认为是多个实体，例如：名词短语“Former Chilean dictator Augusto Pinochet”被拆分为三个实体：“dictator”、“Augusto”和“Pinochet”，并且认为拆分后的实体和名词短语中心词具有相同的语法角色。

Elsner 等人<sup>[6]</sup>也将语篇中的所有实体都放入实体网格中，认为这种方式能够对那些中心词的修饰语进行建模，例如：“a Bush spokesman”中的“Bush”在 Penn Treebank 算法中其并非为中心词，但同样将作为实体在网格中进行表示，不同的是非中心名词的语法角色用“X”来表示，并通过实验证明了这种改进的方法比标准的实体网格在连贯性判断上提高了大约 4% 的准确性。除了能提高连贯性的判断性能外，这种将非中心词加入实体网格的方式还能够捕获到更多名词的共指信息和实体分布规律。

另外，针对在日语中没有好的共指消解的工具可以使用，以及共指消歧对实体网格性能影响的两面性，Yokono 等人<sup>[7]</sup>只考虑指示代词的情况，这些代词在语篇中有明确的实体与其相对应，而且所指代的实体应出现在其前面的句子中，否则认为该片段是不连贯，并通过引入一个概率值将这种指代关系融入到实体网格中。

## (2) 实体间的语义合并

为了改变网格中实体之间相互独立的假设，Filippova 等人<sup>[8]</sup>通过对语义相关的实体进行合并来扩展实体网格，选取包含命名实体（人名、地名和机构）的新闻文章作为语料，利用 WikiRelate! API 对实体网格中的语义相关实体进行合并。实验结果表明这种实体语义合并的方式并未提高算法整体的性能，但 Filippova 等人坚持认为对语义相关的实体进行聚类合并，可能是未来实体网格的发展方向。并计划将来至于 Wikipedia 和 GermaNet 的信息同现有语料相结合，以及改进实体语义相关度计算方法来提高实体网格扩展方案的整体性能。

Yokono 等人<sup>[7]</sup>也认为实体之间的独立假设不能有效捕捉到实体之间的词汇凝集因素，其利用 Mochizuki 系统来对语义相关的实体进行聚类，形成一条由语义相关实体组成词汇链；实体间的相关性通过它们在不同语篇中的共现关系进行计算，而词汇链间的相关性则是取两条链中实体相似度的最大值。在实体合并的过程中，由于不同实体在给定句子中可能具有不同的语法角色，Yokono 等人提出了两种方式进行解决：1) 取优先级高的语法角色进行代替；2) 在实体网格中保留原始的语法序列；如图 4 所示。在相关日文语料中的实验结果也证明了这种语义合并方式的有效性。

	e1	e2	e3
s1	S	O	-
s2	O	S	O

标准的实体网格

	e1	c1
S1	S	O
S2	O	S

取优先级高的角色

	e1	c1
S1	S	O-
S2	O	SO

对语法角色进行合并

图 4 实体语法角色的选择

这种对语义相关实体进行聚类合并的方法，降低了实体网格的维度，有利于改善实体网格数据稀疏的状况，能够发现实体网格中潜在转换模式，其也将在实体网格中的其它可能应用中发挥重要的作用，例如：利用实体在网格中的分布规律来发现语篇中的重要实体。但实体语义相关度算法以及相关度阈值的选择将极大地影响实体网格的应用效果。

## (3) 新的实体特征的融入

实体网格采用四种语法角色来表示实体在给定句子中的角色，在实施过程中，往往会针对具体的应用场景和不同的语种特点，在实体网格中融入一些新的特征来提高实体网格的性能。Elsner 等人<sup>[6]</sup>认为实体网格没有考虑实体自身的信息，例如：对“Hillary Clinton”和“wheat”而言，它们在下一句子中的同一个语法角色出现的概率相同，但在新闻这类型的文章中人和机构相对重要；而且，包含多个修饰词的实体和以单数出现的实体往往也很重要，其利用命名实体标签对实体的类型进行标识，通过引入基于外部词表的统计信息来标识实体的重要程度。另外，在 Elsner 的其它研究中还通过将实体网格和其提出的新实体语篇模型（Discourse-new Model）<sup>[9]</sup>和代词指代模型（Pronoun Coreference Model）<sup>[9]</sup>相结合来提高整体的性能。

Burstein 等人<sup>[10]</sup>针对作文语篇的特点丰富了实体网格的角色表示，并利用扩展的实体网格对学生的作文质量进行评价。根据作文数据和新闻文本相比噪声大的情况，引入 GUMS 用于标识作文的技术质量，GUMS 包括：语法、词的用法、结构错误和写作风格；引入（\*\_TT）用于揭示同一个实体在语篇中的不同描述形式，例如：转换模式[SS]表示在相邻的句子中使用了同一个实体，但并不能说明在描述该实体时是否使用了相同的词；另外，“外壳名词”（Shell nouns）也可能对语篇的连贯性造成影响，引入（\*\_Shellnouns）对外壳名词的使用情况进行标识；并将以上特征融入实体网格中对来自不同类型作者的 800 篇作文进行了评分，并对结果做了比较和分析。

另外，在一些非英语的应用环境中，研究者们针对各自的语言特点对实体网格进行扩展，Yokono 等人<sup>[7]</sup>将实体网格应用到日语中，将句子间的衔接性关系分为三个组，如表 1 所示，同时把组信息融入到网格中的转换模式中，用于对实体网格的计算；其还增加两种实体语法角色“H”和“R”，其中，“H”用于标识话题的转换，如果一个文章中的话题经常变化，则认为语篇的连贯性较差；“R”用于标识谓语有小品词修饰的实体，认为此类实体比那些谓语没有修饰的实体更加重要。另外，Cheung 等人<sup>[11]</sup>针对德语的特点利用实体在句子拓扑位置上所充当的成分（LK、VC 和 VF 等）来取代标准实体网格中的语言学维度。

表 1 日语中的三组句子衔接关系

Groups	Types of conjunction	Explanation
Group 1	copulative, adversative	the relation connecting two matters in a logical way
Group 2	additive, contrastive, transitive	the relation connecting two separate matters
Group 3	parallel, supplementary, consecutive	the relation connecting two sentences of a matter

注：此表来源于 Yokono 和 Okumura 的 Incorporating Cohesive Devices into Entity Grid Model in Evaluating Local Coherence of Japanese Text

### 3.2 新的实体网格计算方法

标准的实体网格采用机器学习的方法来确定特征向量  $\vec{w}$ ，通过排序函数  $\Phi(x_{ij})$  和参数向量  $\vec{w}$  来比较不同语篇之间的连贯性，但这种方式不容易将其它特征融入到算法当中，本节对其它的几种实体网格计算方法进行分析。

#### (1) 基于联合概率分布的网格计算方法

Lapata 等<sup>[5]</sup>利用句子和实体之间的联合概率分布，来反映实体在语篇句子之间的分布情况，语篇连贯性计算公式如下：

$$P_{coherence}(T) = P(e_1 \dots e_m; s_1 \dots s_n) \approx \prod_{j=1}^m p(e_j; s_1 \dots s_n)$$

其中  $T$  表示语篇， $s_1 \cdots s_n$  表示语篇中的句子， $e_1 \cdots e_m$  表示语篇实体。假设语篇中的实体是相互独立的，将公式转化为实体  $e_j$  在句子中模式转换序列分布概率的乘积形式，并且有：

$$P(e_j; s_1 \cdots s_n) = P(r_{1,j} \cdots r_{n,j}) = \prod_{i=1}^n P(r_{i,j} | r_{1,j} \cdots r_{(i-1),j}) \approx \prod_{i=1}^n P(r_{i,j} | r_{(i-h),j} \cdots r_{(i-1),j})$$

其中  $r_{ij}$  表示  $e_j$  在句子  $i$  中的语法角色， $P(r_{i,j} | r_{1,j} \cdots r_{(i-1),j})$  利用标准马尔可夫模型的独立性假设进行获取， $h$  是依赖的历史模式长度。并利用语篇中句子的数量  $n$  和实体数量  $m$  进行归一化处理：

$$P_{coherence}(T) \approx \frac{1}{m \cdot n} \sum_{j=1}^m \sum_{i=1}^n \log P(r_{ij} | r_{(i-h),j} \cdots r_{(i-1),j})$$

$P(r_{i,j} | r_{1,j} \cdots r_{(i-1),j})$  可以由连贯性文本组成的训练语料中得到，拥有高概率值的语篇会被认为比低概率值的语篇拥有更高的连贯性，这种计算方法容易和基于实体网络特征扩展的方法相融合，此外：Elsner<sup>[12]</sup>、Yokono<sup>[7]</sup>和 Filippova<sup>[8]</sup>等人也根据具体的应用场景提出了相应的融合方法。

这种实体网络计算方法被认为是与实体无关的，在网络中拥有相同列拓扑结构的实体会被分配为相同的概率，而在实际的语篇中，每个句子中包含很少的名词，在这些很少的名词中又只有更少的名词适合充当主语和宾语，换句话说，名词之间通过相互竞争来争取其在句子中的语法角色，一旦一个名词成为主语，其他名词被选择充当句子（一个句子可能包含不同的子句）主语的的概率迅速下降，为此 Elsner 提出了非约束的实体网格（Relaxed Entity Grid）<sup>[13]</sup>把实体相关性引入到网格当中。

## （2）利用统计分析的网格计算方法

在统计学中，肯德尔相关系数（Kendall's  $\tau$ ）<sup>[14]</sup>是一个用来测量两个随机变量相关性的统计值。一个肯德尔检验是一个无参数假设检验，它使用计算而得的相关系数去检验两个随机变量的统计依赖性。肯德尔相关系数的取值范围在-1 到 1 之间，当  $\tau$  为 1 时，表示两个随机变量拥有一致的等级相关性；当  $\tau$  为-1 时，表示两个随机变量拥有完全相反的等级相关性；当  $\tau$  为 0 时，表示两个随机变量是相互独立的。Filippova 等人<sup>[8]</sup>利用肯德尔相关系数对语篇的连贯性进行了计算，首先对训练集中实体网络的模式转换进行排序，记为  $g_1$ 。对测试语篇中的模式转换进行排序，记为  $g_2$ 。通过公式： $\tau = 1 - 4 \frac{t}{N(N-1)}$  来计算语篇的连贯性，其中  $N$  表示总的模式转换个数， $t$  表示由  $g_2$  到正确排列  $g_1$  在相邻句子模式间要经过的交换的总数， $\tau$  越小表示语篇的连贯性越差。Filippova 还利用肯德尔相关系数对转换模式的长度为 3 的情况进行了测试，但与长度为 2 的情况相比算法的性能稍有降低。

Feng 等人<sup>[15]</sup>分别利用肯德尔相关系数、平均连续性（Average continuity）以及编辑距离（Edit distance）来对实体网络中的转换模式进行计算，并同标准的实体网络的计算结果进行了对比分析。同时，其认为标准的实体网络以及基于实体网络的扩展往往需要更多的训练集和更加精巧的特征扩展方案，因而提出了多重排序模型（multiple-rank model）的实体网络扩展方案，允许模型在更细的粒度上对连贯性进行评测，在学习过程中不仅对原始语篇及其派生语篇组成的语篇对进行学习，也对由不同的派生语篇组成的语篇对进行学习。

Pitler 等人<sup>[16]</sup>则是直接通过统计分析的方法对语篇可读性和实体网格中 16 种特征<sup>3</sup>的相关性进行了测试, 在这些转换中用于表示两个相邻句子中使用了同样主语的[SS]转换, 因其保持了句子的焦点, 被向心理论认为是一种非常连贯的转换类型, 但测试结果表明其在可读性判断上却起到消极作用, 另外, 其消极作用的还有[SO]、[SX]、[XX]和[--], 其中和可读性最相关的转换特征是[S-], 考虑这些特征对可读性的影响, 利用线性回归分析同其它可读性指标相结合, 对语篇可读性进行评价。

## 4 总结

本文主要对实体网格的概念、构造方法以及如何利用网格进行连贯性评测进行分析, 并对基于实体网格特征修改和丰富以及新的实体网格计算方法进行了深入的分析和研究。在实体网格的应用上, 除在文中论述的在自动摘要、文本排序、作文评分以及可读性评测等任务中得到广泛应用外, 其还在对话分解 (Chat Disentanglement or Threading)<sup>[17]</sup>和故事生成 (Story Generation)<sup>[18]</sup>等的应用中取得了良好的效果。笔者还认为可以充分利用和挖掘实体在网格的分布规律进行相关的应用, 例如: 可以利用主题在过渡区域相对于主题内部的文本连贯性差的特点进行语篇主题分割; 利用语篇中重要的实体在网格中的分布往往具有“聚集”的特点进行关键词抽取和语篇主题识别等。随着实体网格相关研究的不断深入, 其将拥有更加广阔的应用领域和价值。

## 参考文献

- [1] R. Barzilay, M. Lapata. Modeling local coherence: An entity-based approach[C],//Proceedings of the 43rd Annual Meeting of the ACL, Ann Arbor, 2005: 141-148.
- [2] R. Barzilay, M. Lapata. Modeling local coherence: An entity-based approach[J]. Computational Linguistics, 2008(34):1-34.
- [3] Zipf's Law. [OL]. [2012-07-23]. [http://en.wikipedia.org/wiki/Zipf's\\_law](http://en.wikipedia.org/wiki/Zipf's_law).
- [4] R. Soricut, D. Marcu. Discourse generation using utility-trained coherence models[C],//Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions, Sydney, Australia 2006:803-810.
- [5] M. Lapata, R. Barzilay. Automatic evaluation of text coherence: Models and representations[C],//Proceedings of the 19th International Joint Conference on Artificial Intelligence, Edinburgh, Scotland, UK, 2005:
- [6] M. Elsner, E. Charniak. Extending the entity grid with entity-specific features[C],//Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics, Portland, Oregon, 2011:125-129.
- [7] H. Yokono, M. Okumura. Incorporating Cohesive Devices into Entity Grid Model in Evaluating Local Coherence of Japanese Text[J]. Computational Linguistics and Intelligent Text Processing, 2010:303-314.
- [8] K. Filippova, M. Strube. Extending the entity-grid coherence model to semantically related entities[C],//Proceedings of the 11th European Workshop on Natural Language Generation, Schloss Dagstuhl, Germany, 2007:139-142.
- [9] M. Elsner, E. Charniak. A generative discourse-new model for text coherence[R]. Technical Report CS-07, 2007.
- [10] J. Burstein, J. Tetreault, S. Andreyev. Using entity-based features to model coherence in student essays[C],//Proceedings of the 2010 Annual Conference of the North American Chapter of the ACL, Los Angeles, California, 2010:681-684.
- [11] J.C.K. Cheung, G. Penn. Entity-based local coherence modelling using topological

---

<sup>3</sup> 实体网格中的十六种特征分别为: “SS”、“SO”、“SX”、“S-”、“OS”、“SO”、“OX”、“O-”、“XS”、“XO”、“XX”、“X-”、“-S”、“-O”、“-X”和“-”

- fields[C],//Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, Uppsala, Sweden, 2010:186-195.
- [12] M. Elsner, E. Charniak. Coreference-inspired coherence modeling[C],// Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics, Columbus, Ohio, USA, 2008:41-44.
- [13] M. Elsner, J. Austerweil, E. Charniak. A unified local and global model for discourse coherence[C],//Proceedings of NAACL HLT 2007, Rochester, NY, 2007:436–443.
- [14] Kendall tau rank correlation coefficient. [OL]. [2012-07-23].  
<http://blog.csdn.net/wsywl/article/details/5889419>.
- [15] Extending the Entity-based Coherence Model with Multiple Ranks. [OL]. [2012-07-23].  
<http://aclweb.org/anthology-new/E/E12/E12-1032.pdf>.
- [16] E. Pitler, A. Nenkova. Revisiting readability: A unified framework for predicting text quality[C],//Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing, Honolulu, USA, 2008:186-195.
- [17] M. Elsner, E. Charniak. Disentangling chat with local coherence models[C],//Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics, Portland, Oregon, 2011:1179-1189.
- [18] N.D. McIntyre. Learning to tell tales automatic story generation from Corpora[D]. Edinburgh UK: University of Edinburgh, 2011.