

现代汉语名词歧义度研究

现代汉语名词歧义度研究

李安

2012年10月24日星期三

歧义度的提出、研究目标和框架

- 词义消歧实现
- 义项关系和义项的语义距离

结论1：歧义度与义项距离的关系

结论2：三种义项关系及歧义度成因

结论3：人机词典的不同与现汉的不足

词义消歧

- 用计算机为使用中的多义词标注一个确定的义项。
- 词义消歧在自然语言处理中有重要作用。

词义消歧研究概述

- 人工编写规则
- 统计方法
- 词典方法

词义消歧中存在的问题

- 研究集中于人工智能角度关注，评测多以平均正确率和召回率为标准，对词义的个性研究不足。
- 不同多义词消歧正确率不同。体现的更多的不是方法问题，而是词义问题。

歧义度

- 歧义度是多义词体现在计算机词义消歧中的难度和可能达到效果的歧义性的度量。
- (1)歧义度是多义词义项间区别性形式特征有无、多少的数字描述；(2)歧义度是多义词各个义项在分布上差异大小的数字描述；(3)歧义度是多义词在计算机词义消歧中困难程度的数字描写。

歧义度计算

$$P = \left(1 - \frac{M}{N}\right) \% = \left(1 - \frac{\sum m_i}{N}\right) \% \quad (2.8)$$

歧义度的研究目标和框架

描写歧义度

- 实现词义消歧
- 标注一定规模语料
- 计算歧义度

解释多义词间歧义度差异的原因

- 义项关系和义项的语义距离

计算歧义度

- 采用统计与规则相结合的方式实现。
- 应用到了大规模语料库；搭配库；义类词典。



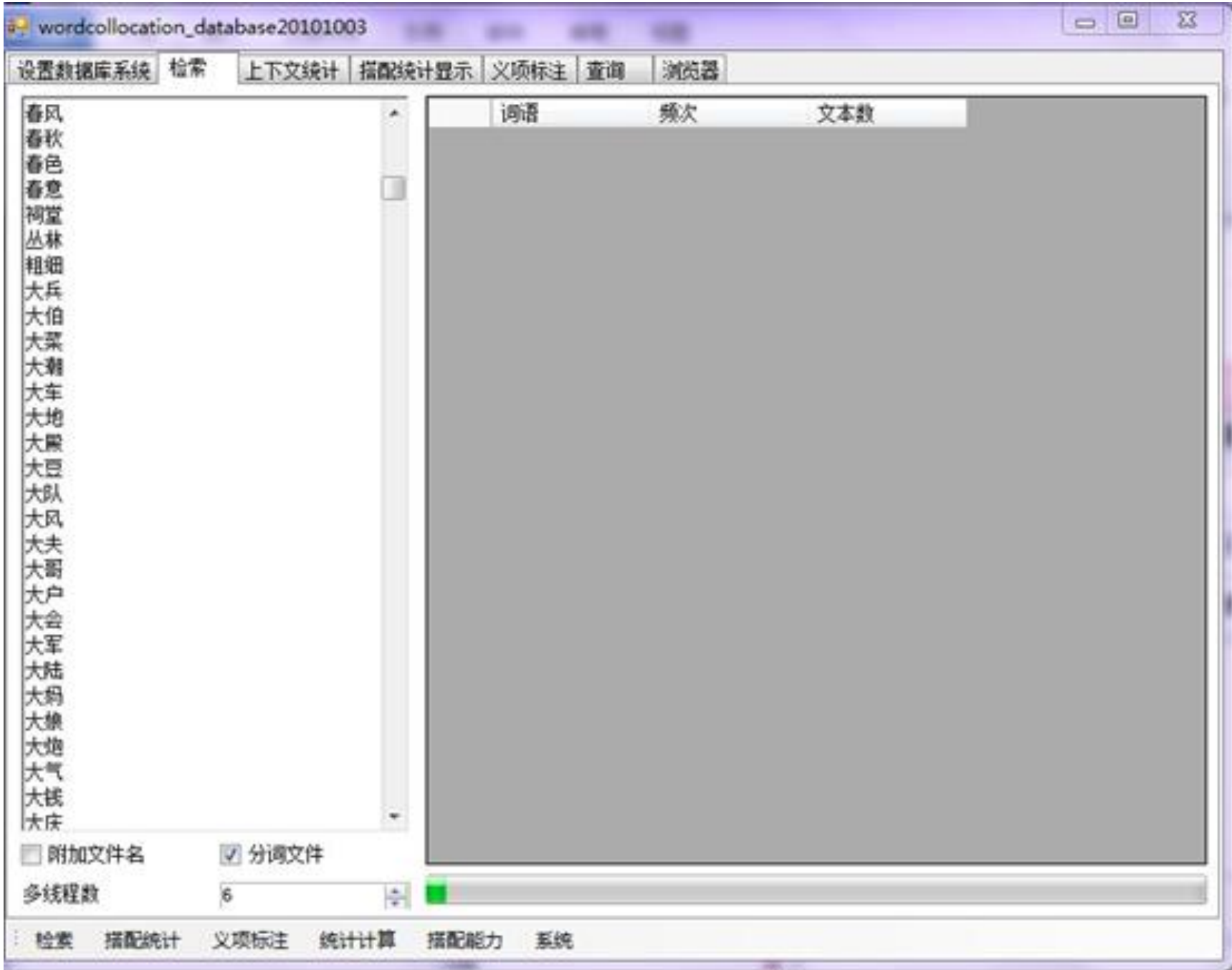
设置数据库系统 检索 上下文统计 搭配统计显示 义项标注 查询 浏览器

数据库服务器	LN-PC	<input checked="" type="checkbox"/> 不使用用户名和密码	连接到Sql Server数据库
使用用户名密码	sa	●●●●●●●●●●●●●●●●	
选择一个数据库	毕业论文	刷新 停止 状态	打开ACCESS数据库
词表 n	▼	词列表	
验证语料 n	▼	检索语料存放表	
统计 n	▼	上下文统计结果表	
	▼	文本类型表	
wuji	▼	义类词典	
	▼	总字频表	
	▼	规则分析表	
词典校对版	▼	现代汉语词典	
名词	▼	名词规则	
动词	▼	动词规则	
形容词	▼	形容词规则	
词性	▼	词性规则	
<pre>create table 词表 (id autoincrement not null,word char (80),pinci int, pinlv single,wenbenshu single)</pre>			执行命令
<pre>create table 检索 (id autoincrement not null,qian ntext,keys</pre>			

原始语料文件夹 统计结果存放目录 保存当前命令

检索 搭配统计 义项标注 统计计算 搭配能力 系统






```
select * from 档案表1
```

结果 消息

id	qian	keys	hou	name
1	就算/只的/得打/的/ / 给/ / 弄/ / 在/ / 外面/ / / 我们/ /	阿婆/ /	就/ / 笨体/ / 了/ / / 口里/ / 说/ / / 不/ / 不/ / 卖/ / /	(三毛《78装模》)
2	嘛/ / 问/ / 然/ /	阿婆/ /	/ / 其/ / 给/ / 儿/ / NUM 生/ / 哥/ / 票/ / 好/ / /	(三毛《78装模》)
3	(/ 由/ / 信/ / 内/ / 取出/ /) /	阿婆/ /	/ / 是/ / 胡/ / APER 伯伯/ / 的/ / - / NUM 张/ / 相片/ / /	(三毛《78装模》)
4	爱好/ / 由来/ / 苦/ / 难/ / 难/ / / - / NUM 字/ / 千/ / NUM 改/ / /	阿婆/ /	还是/ / 初/ / 拜/ / 女/ / NUM / / 头/ / 未/ / 被/ / 成/ / 不/ / 详/ /	(张爱玲《表姨阿姨及其他》)
5	再/ / 临/ / 香/ / 港/ / LOC 表/ / 的/ / 妈/ / 亲/ / 会/ / 不/ / 得/ / 我/ / / 千/ /	阿婆/ /	一/ / 阵/ / 飞/ / 到/ / 香/ / 港/ / LOC / /	(三毛《赴拉萨途中见闻》)
6	巫/ / 婆/ / TIM 给/ / 于/ / 来/ / 了/ / / 前/ / 三/ / 天/ / TIM / / 婆/ /	阿婆/ /	/ / 欲/ / 找/ / / 掉/ / 掉/ / / 夏/ / NUM 兄/ / NUM 堂/ / 妹/ / /	(三毛《亲爱的婆婆大人》)
7	/ / 好/ / 了/ / / 不/ / 要/ / 再/ / 拜/ / 拜/ / 拜/ / 拜/ / 拜/ / 拜/ / 拜/ /	阿婆/ /	/ / 欲/ / 找/ / 掉/ / 掉/ / 每/ / 家/ / 都/ / 有/ / 分/ / 分/ / 一/ / 天/ / TIM / /	(三毛《亲爱的婆婆大人》)
8	黄/ / 昏/ / TIM 的/ / 心/ / 时/ / 时/ / / 父/ / 亲/ / 妈/ / 亲/ / 和/ / 我/ / 带/ /	阿婆/ /	/ / 我/ / 也/ / 要/ / / 拜/ / 拜/ / 拜/ / 拜/ / 拜/ / 拜/ / 拜/ /	(三毛《塑料儿童》)
9	/ /	阿婆/ /	/ / 你/ / 为/ / 什/ / 么/ / 说/ / 我/ / 们/ / 是/ / 塑/ / 料/ / 儿/ / 童/ / 的/ / 心/ / ? /	(三毛《塑料儿童》)
10		阿婆/ /	/ / 我/ / 看/ / 我/ / 还/ / 是/ / 进/ / 去/ / 吧/ / ! /	(三毛《塑料儿童》)
11	这个/ / 生/ / 死/ / 之/ / 心/ / 交/ / NUM 的/ / 女/ / 友/ / / 不/ / 愿/ / 自/ / 己/ /	阿婆/ /	/ / 的/ / 心/ / / 这/ / 种/ / 情/ / 形/ / 在/ / 没/ / 有/ / 亲/ / 属/ / 称/ / 呼/ /	(三毛《回娘家》)
12	在/ / 玛/ / 丽/ / 莎/ / PER 的/ / 家/ / 里/ / / 最/ / 是/ / 自/ / 由/ / /	阿婆/ /	换/ / / 高/ / 心/ / 那/ / 边/ / / 午/ / 饭/ / 的/ / 香/ / 味/ / 早/ / 已/ / 黄/ /	(三毛《回娘家》)
13	用/ / 餐/ / 的/ / 时/ / 候/ / / 我/ / 无/ / 意/ / 间/ / 讲/ / 起/ / 去/ / 妹/ /	阿婆/ /	锅/ / 看/ / 她/ / /	(三毛《美国去妹》)
14	有/ / - / NUM 回/ / 三/ / 毛/ / PER 出/ / 了/ / 新/ / 书/ / / 拿/ /	阿婆/ /	说/ / 了/ / - / NUM 声/ / : / 你/ / 呢/ / ? /	(三毛《第一：我家老二-三小姐》)
15	一/ / 千/ / 零/ / - / NUM 夜/ / 的/ / 阿/ / 婆/ / (/ 十/ / 九/ / NUM 岁/ /) /	阿婆/ /	三/ / 毛/ / PER 在/ / 家/ / 里/ / 是/ / 什/ / 么/ / 样/ / 子/ / /	(三毛《一千零一夜的阿婆》)
16	我/ / 的/ /	阿婆/ /	就/ / 是/ / 一/ / 个/ / NUM 最/ / 普/ / 通/ / 的/ / 阿/ / 婆/ / / 跟/ / 天/ / 下/ /	(三毛《一千零一夜的阿婆》)
17	我/ / 的/ / 阿/ / 婆/ / 就/ / 是/ / 一/ / 个/ / NUM 最/ / 普/ / 通/ / 的/ /	阿婆/ /	/ / 跟/ / 天/ / 下/ / 的/ / 阿/ / 婆/ / 都/ / 差/ / 不/ / 多/ / /	(三毛《一千零一夜的阿婆》)
18	我/ / 的/ / 阿/ / 婆/ / 就/ / 是/ / 一/ / 个/ / NUM 最/ / 普/ / 通/ / 的/ / 阿/ / 婆/ /	阿婆/ /	都/ / 差/ / 不/ / 多/ / /	(三毛《一千零一夜的阿婆》)
19	在/ / 沙/ / 发/ / 上/ / / 那/ / 个/ / 被/ / 称/ / 为/ /	阿婆/ /	的/ / E / X C / H X H / X O / X / / 拿/ / 出/ / 四/ / NUM 个/ /	(三毛《古国出堡》)
20	一/ / 千/ / 零/ / - / NUM 夜/ / 的/ /	阿婆/ /	(/ 十/ / 九/ / NUM 岁/ /) / 黄/ / 昏/ / PER 我/ / 都/ / 没/ / 被/ / 人/ /	(三毛《一千零一夜的阿婆》)
21	说/ / 起/ /	阿婆/ /	的/ / 讲/ / 话/ / / 真/ / 是/ / 一/ / 绝/ / / 她/ / 讲/ / 起/ / 故/ / 事/ /	(三毛《一千零一夜的阿婆》)



规则编号	副+	数+	量+	代+	字句+	名+	形+	动+	+限定对象	施动+
6328						1002C0				

id	词	词	句	句	句	总	频	累	前	后	平	标	表	原	互
1511	的	u	48	44.8		69	5.15	5.15	33	36	4.55	2.84	18838	膀	
1511	着	u	40	37.3		52	3.88	9.04	40	12	2.5	2.05	18838	膀	
1512	在	p	26	24.2		27	2.01	11.0	5	22	4.96	2.53	18838	膀	
1513	光	n	25	23.3		25	1.86	14.7	25	0	1.44	0.50	18838	膀	
1513	晃	v	25	23.3		25	1.86	12.9	25	0	1.08	0.27	18838	膀	
1514	了	u	16	14.9		19	1.42	16.2	10	9	4.31	3.09	18838	膀	
1515	一	NUM	17	15.8		18	1.34	18.9	7	11	3.88	1.90	18838	膀	
1515	了	y	12	11.2		18	1.34	17.5	6	12	4.72	2.73	18838	膀	

15130 等于 搭配 筛选 直接筛选 返回 撤销 上一条 下一条

膀子 n, 鸟类等的翅膀。:6108 all句子 显示词 增加规则

id	qian	ke	hou
18838	工人/n 们/k 跳/v 进/v 2 槽/N _E 米/q 深/a 的/u 沟/n 内/f 甩开/v	膀子/n	大干/v 。/w
18838	“/w 太/d 若/v 驾/v 了/y 的/u 大/a 村长/n 挺/d 忙/a 上/v 去/v 吧/y 甚/v _E 就/d 我/r 上/上/f 有/v 一/NUM 座/g 顶/a 森森/z 的/u 大山/n 压/g 着/u 我/r 品/v 使/v 我	膀子/n	象/N _E 降雨/v 似/v _E 地/u “/w 叭/o 叭/o 起来/v 。/w 打/v 了/u
18838	我/r 上/上/f 有/v 一/NUM 座/g 顶/a 森森/z 的/u 大山/n 压/g 着/u 我/r 品/v 使/v 我	膀子/n	/w 心情/n 舒畅/a 地/u 做/v 买卖/v 。 /w
18838	我/r 甩开/v	膀子/n	/w 装/v 着/u 车箱 /n 上面/f 直到/v 车顶 个/q “/w 漫头/n 顶



如何解释多义词歧义度的差别

- 词汇语义理论一直分为分析主义和功能主义两个派别。
- 歧义度是对义项功能差异的描写，我们希望从词义内涵本身寻找原因，把两种主义结合起来。
- 结构主义语言学认为语言要素都存在于系统当中，组合、聚合是这个系统的主要关系。
- 描写义项的关系，就需要在系统中找到其相对位置。

义类词典

- 义类词典是词义系统的体现。
- 义类的相对位置关系体现了意义异同。
- 以义类词典为描写义项关系的工具。

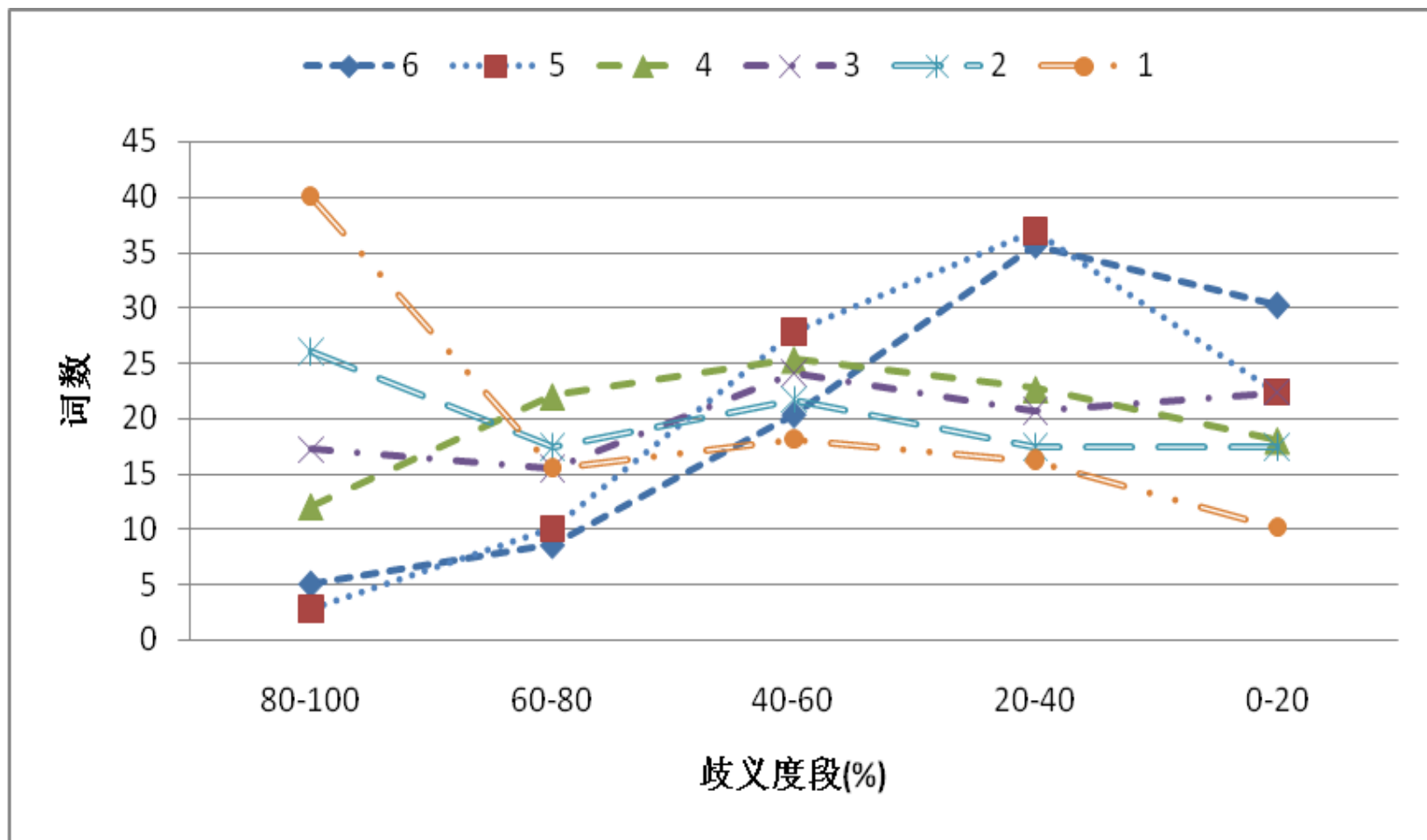
义项关系

- 指多义词义项间的概念关系，具体地用义项在义类词典中所对应的义类间对立关系表示，如“便衣”：
- **【便衣】** ①平常人的服装(区别于军警制服)。②身着便衣执行任务的军人、警察等。
- “便衣” ①为：“具体物-生活用品-服饰-衣服-便装”类
- “便衣” ②为：“生物-人-职业-军人-侦察兵”类，这两个义类的对立就构成了“便衣”的义项间语义关系。

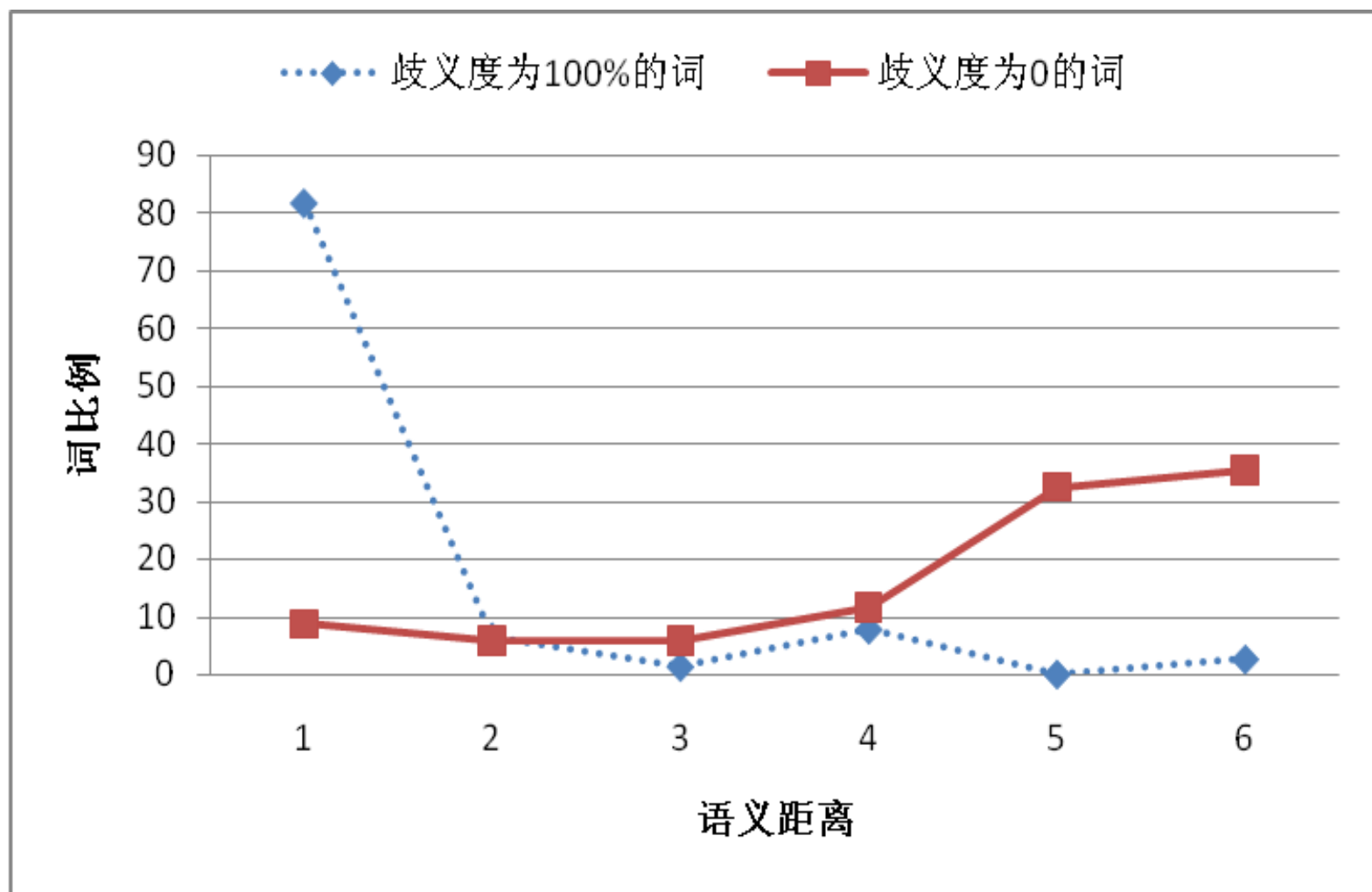
语义距离

- 语义距离为义项在义类词典中相对位置的远近，是义项间语义关系的数字化表示。分1-6级
- 第一级：两个义项属于同一个五级类，如“大雨”：
 - **【大雨】** ①指24小时内雨量达25—50毫米的雨。②指下得很大的雨。
- 第六级：一级类开始不同，如“心脏”：
 - **【心脏】** ①心①。②比喻中心首都。北京是祖国的～。

不同语义距离词的歧义度分布差异明显

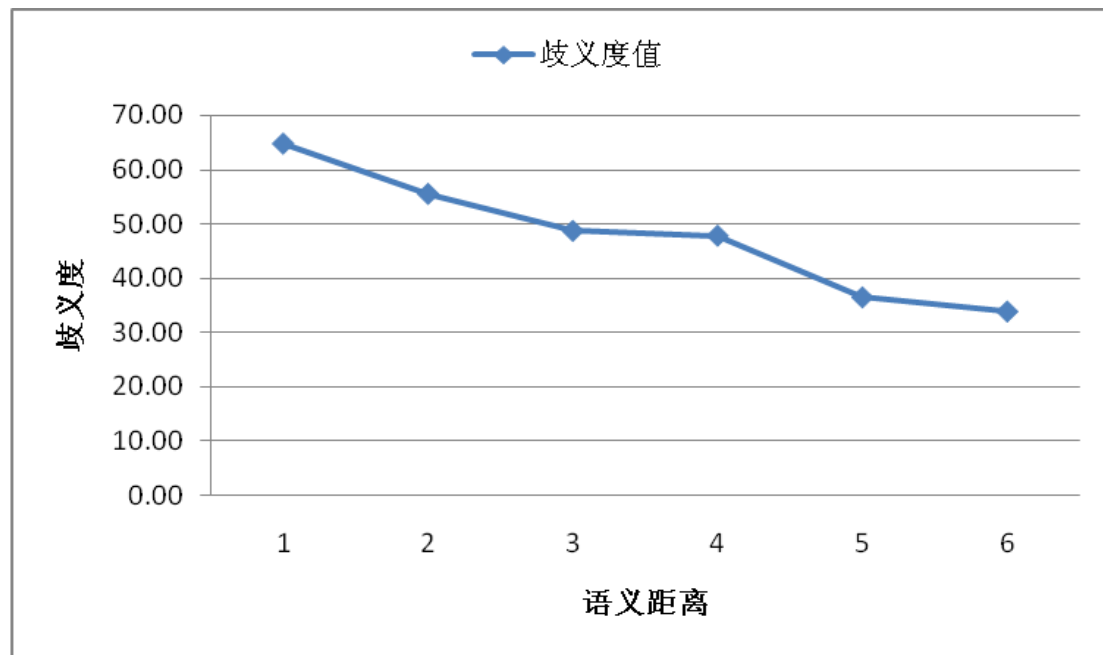


两种极端歧义度词在语义距离的分布迥异



语义距离与歧义度间存在高度负相关关系

语义距离	词数	平均歧义度(%)	歧义度均差(%)
1	265	64.76	18.95
2	69	55.54	9.73
3	58	48.75	2.94
4	150	47.85	2.04
5	219	36.55	-9.26
6	314	33.95	-11.86



使用EXCEL统计得到语义距离与歧义度均值间的相关系数为 $r=-0.98$

三种重要义项关系

同义关系

- 意义基本相同，用法部分或完全相同的关系。近义词指意义相近，用法不完全相同的关系。
- **【暴力】** ①强制的力量；武力。②特指国家的强制力量。
- **【鞭炮】** ①大小爆竹的统称。②专指成串的小爆竹。

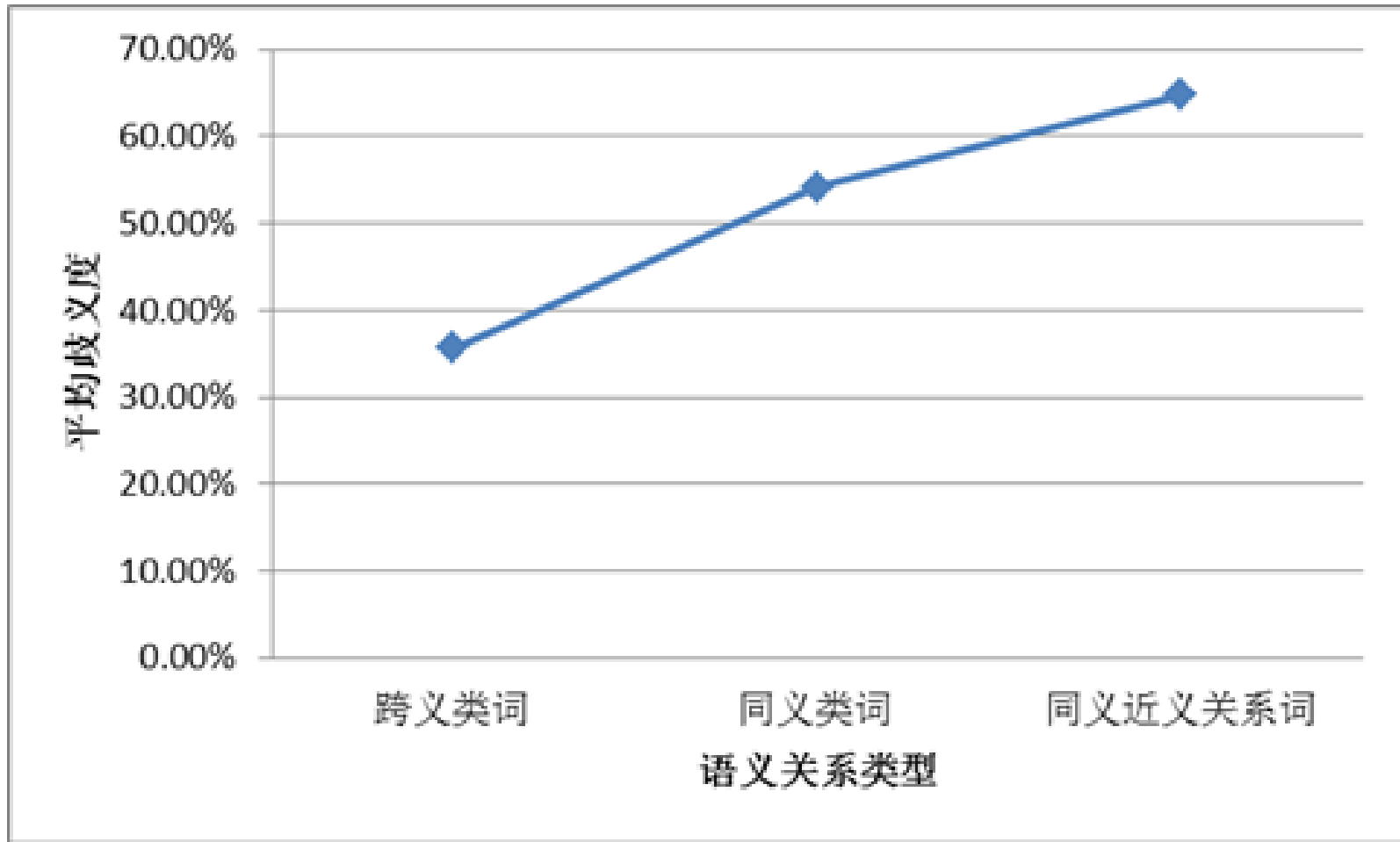
同类关系

- 各个义项属于同一基本层次范畴；语义距离为2或3。这类词“属”相同，但是“种差”较大。
- **【单间】** ①饭馆、旅馆内供单人或一起来的几个人用的小房间。②只有一间的屋子。

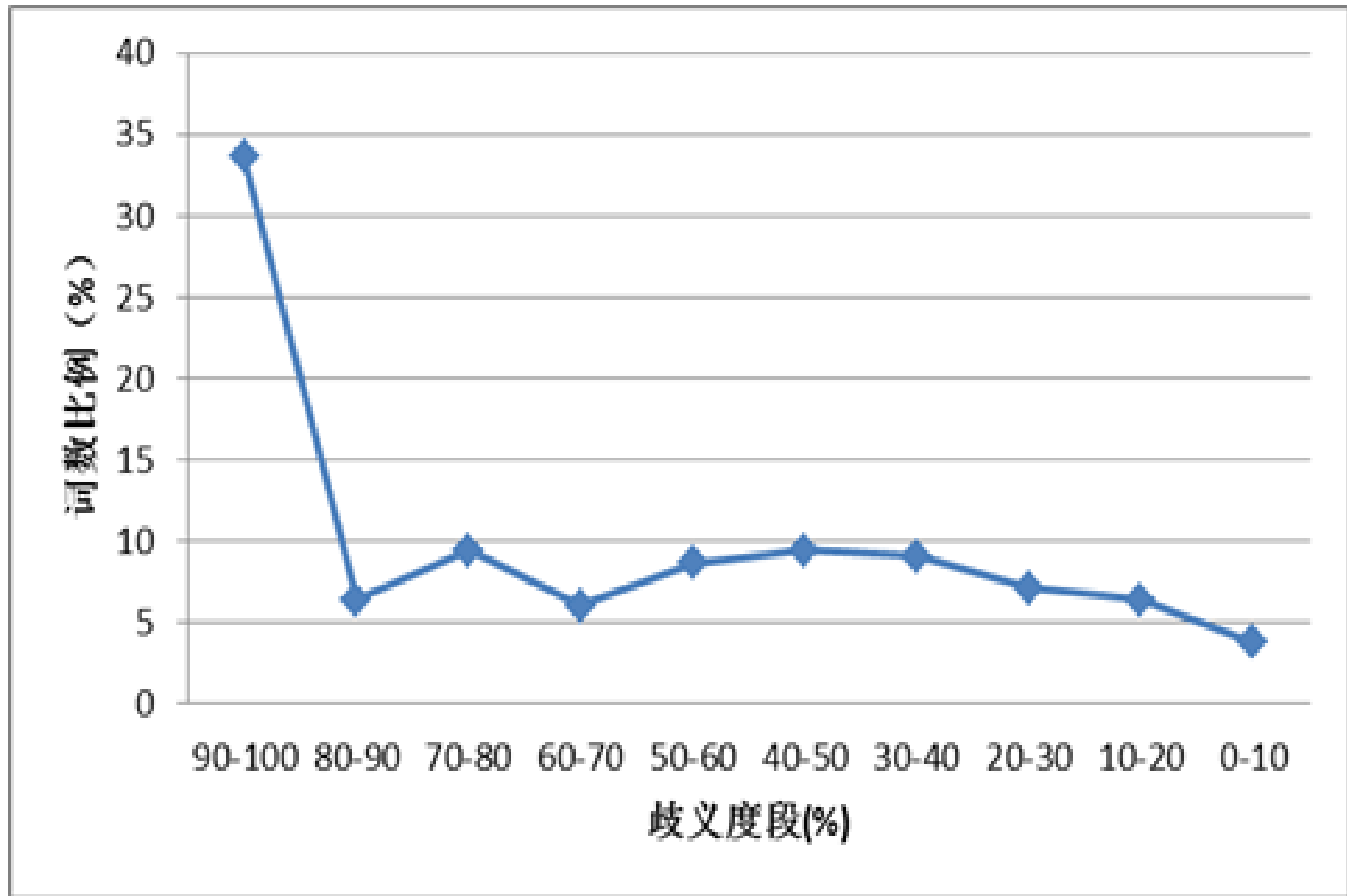
跨类关系

- 其义项不属于同一个基本层次范畴；义项间区别较大；从在词义系统的关系看，语义距离普遍为4、5或6，一个词的义项不同属于一个三级类或二级类甚至一级类。
- **【品质】** ①物品的质量江西瓷~优良。②行为、作风上所表现的思想、认识、品性等本质道德~。

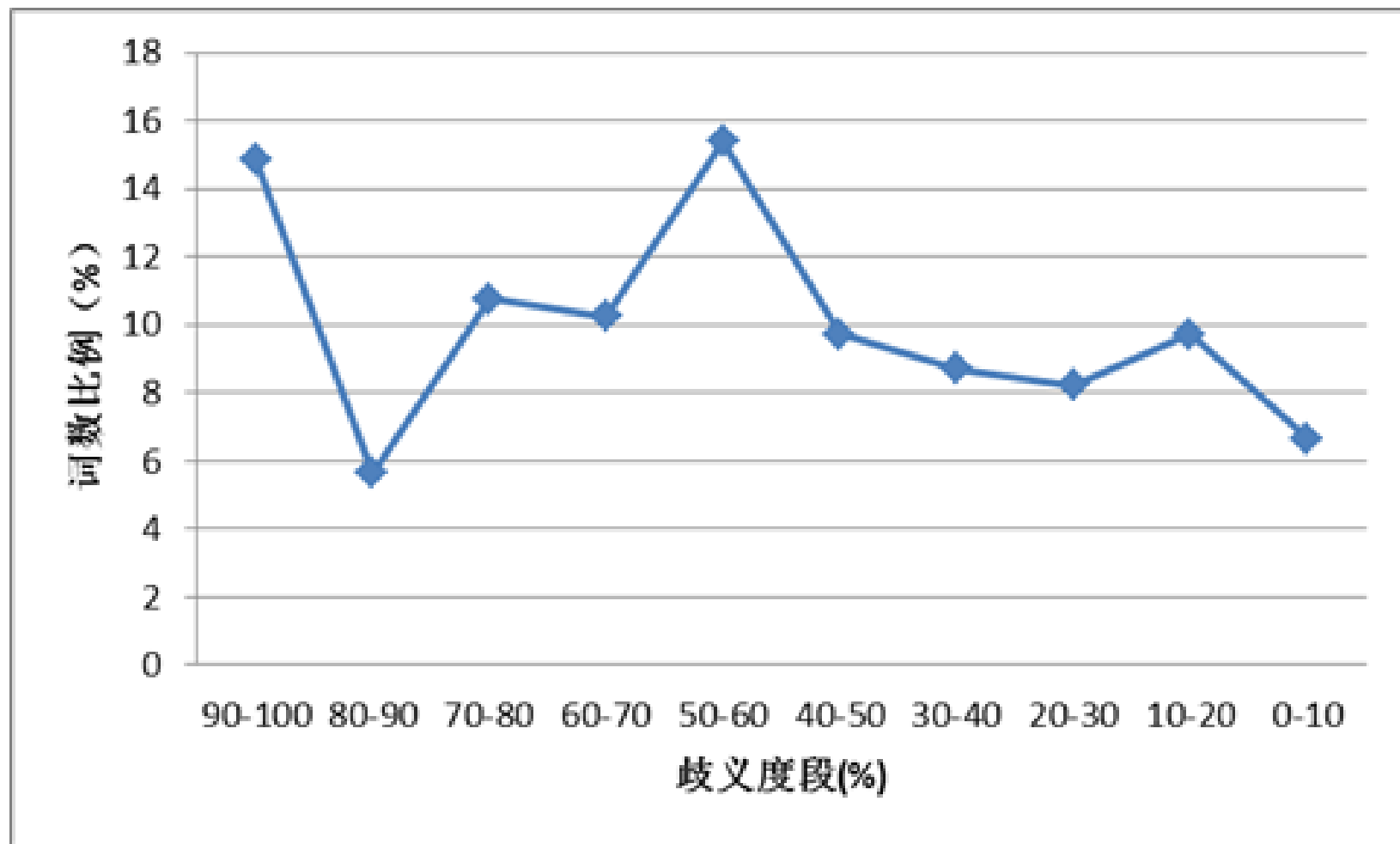
三种词歧义度差异明显



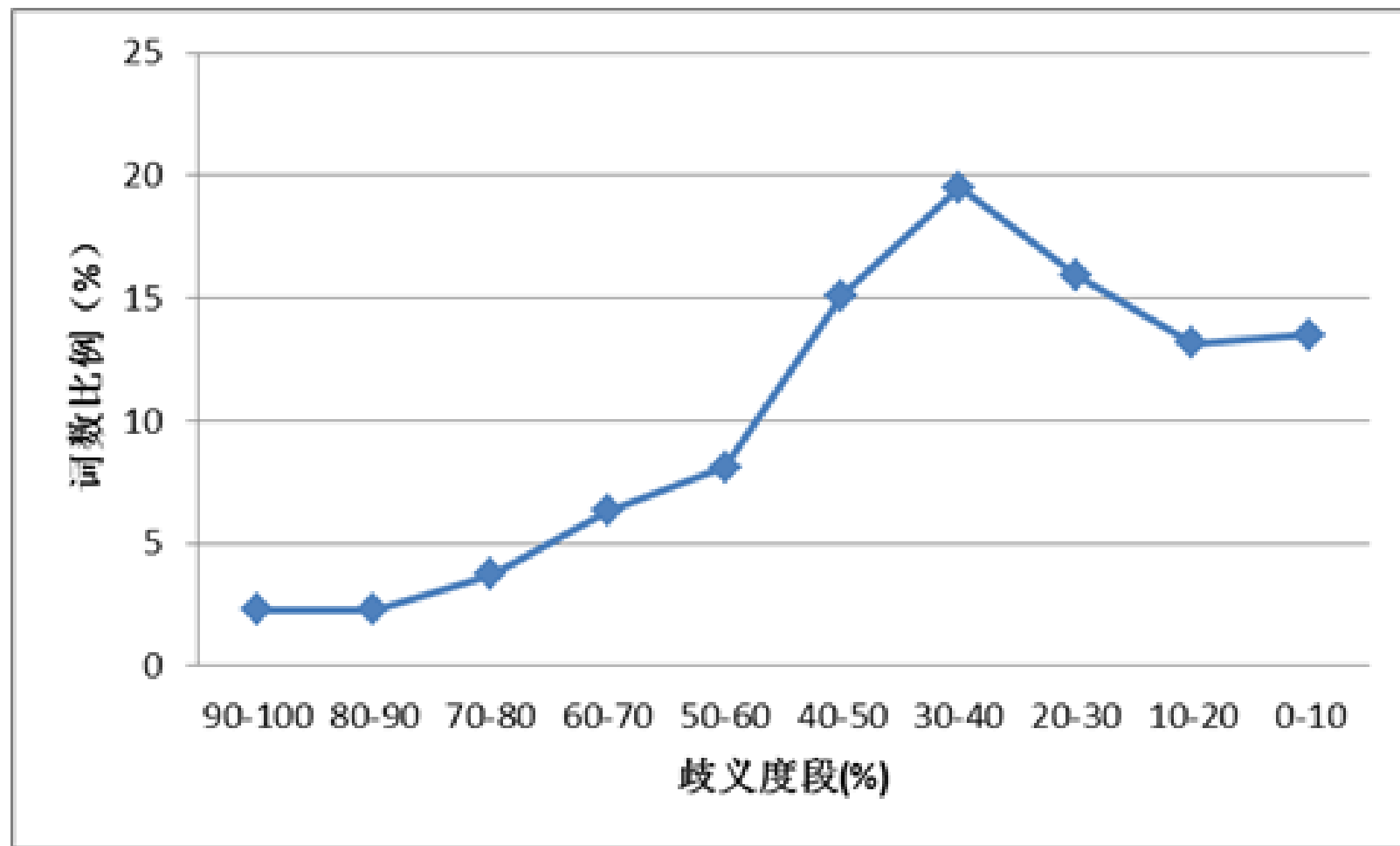
同义近义词的歧义度分布



同类词的歧义度分布



跨义类词的歧义度分布



同义类词与跨义类词歧义度成因不同

○ 同义类词

- 造成歧义度的原因为义项间共同义素的多少。
- 消歧模式应以描写义项间的相同和不同义素为主。

○ 跨义类词

- 造成歧义度的原因是义项间意义引申模式或隐喻、暗喻的概念映射模式。

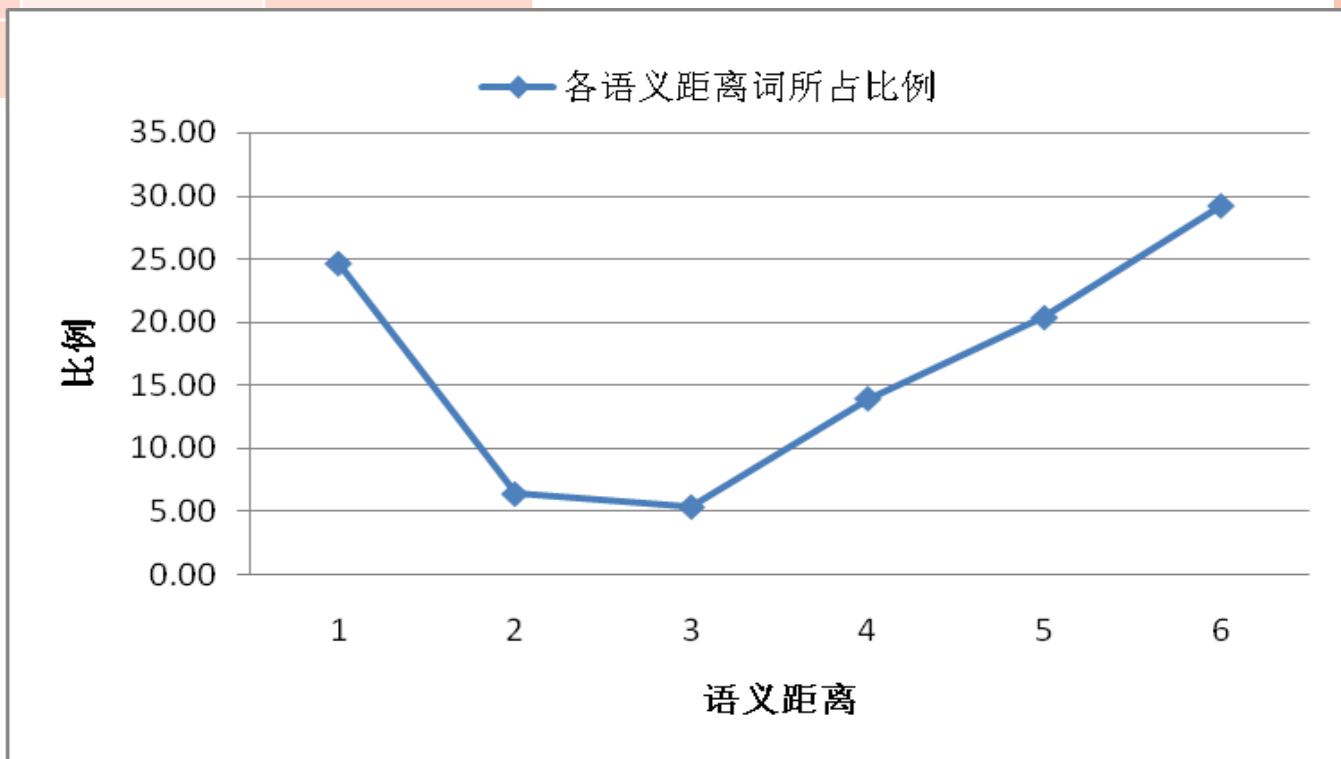
词义消歧的词典问题

- 词义划分有很强的主观性。
- 义项区分粒度是词义消歧中的一个基础、核心问题。
- 目前既没有确定义项粒度的方法，也没有形成取舍的标准。
- 分析论和功能论无法单独解释义项粒度现象。



“现汉”有部分词义项粒度过小

语义距离	词数	比例(%)	累加(%)
1	265	24.65	24.65
2	69	6.42	31.07
3	58	5.40	36.47
4	150	13.95	50.42
5	219	20.37	70.79
6	314	29.21	100
合计	1075		



举例

- **【大雨】** ①指24小时内雨量达25—50毫米的雨。②指下得很大的雨。
- **【封面】** ①线装书指书皮里面印着书名和刻书者的名称等的一页。②新式装订的书刊指最外面的一层，用厚纸、布、皮等做成。③特指新式装订的书刊印着书刊名称等的第一面。也叫封一。
- **【砖头】** ①不完整的砖②砖。
- **【工钱】** ①做零活儿的报酬②工资。
- **【职工】** ①职员和工人。②旧时指工人。
- **【匪徒】** ①强盗。②危害人民的反动派或坏分子。



调整义项粒度对歧义度的影响

- (一) 合并语义距离为1的义项对歧义度的影响
- 1、多义词减少24.1%
- 研究范围内语义距离为1的词有259个，占双义项名词数的24.1%，也就是合并这部分义项后，多义词会减少24.1%。
- 2、歧义度为100%的词减少79.22%
- 歧义度为100%的词为词义消歧中最难处理的部分，双义项名词中共有77个，其中语义距离为1的部分有61个，占到总数的79.22%，如果合并这部分义项，将同时排除这部难以辨析的词。
- 3、总体平均歧义度降低6个百分点
- 双义项名词平均歧义度为46.20%，语义距离为1的词平均歧义度为64.76%，如果排除这一部分，双义项名词歧义度将降低6个百分点，变为40.54%。
- (二) 合并同义类词义项对歧义度的影响
- 1、多义词减少42.79%
- 本文研究范围内同义类词和跨义类词总数为1075个，其中同义类词数量为460个，如果合并同义类词多义词数量将减少42.79%。
- 2、歧义度为100%的词减少94.81%
- 歧义度为100%的同义类词有73个，占全部歧义度为100%的同义类词和跨义类词的94.81%。
- 3、总体平均歧义度降低10.52%
- 同义类词的平均歧义度为60.26%总体上是比较高的，如果去掉同义类词部分，平均歧义度将由46.2%降为35.68%，降低10.52个百分点。



谢谢!

