

淘宝网语义分析产品、技术介绍

王天舟（空海）

淘宝网-交易线-语义分析



个人介绍

- 花名空海



淘宝网中的文本

□ 淘宝主站：

- 30亿店铺、宝贝浏览
- 10亿计的在线宝贝数
- 千万量级交易笔数

□ 文本数据：

- 用户评论
- 商品标题、详情页
- 用户query数据
- SNS、论坛等其他数据

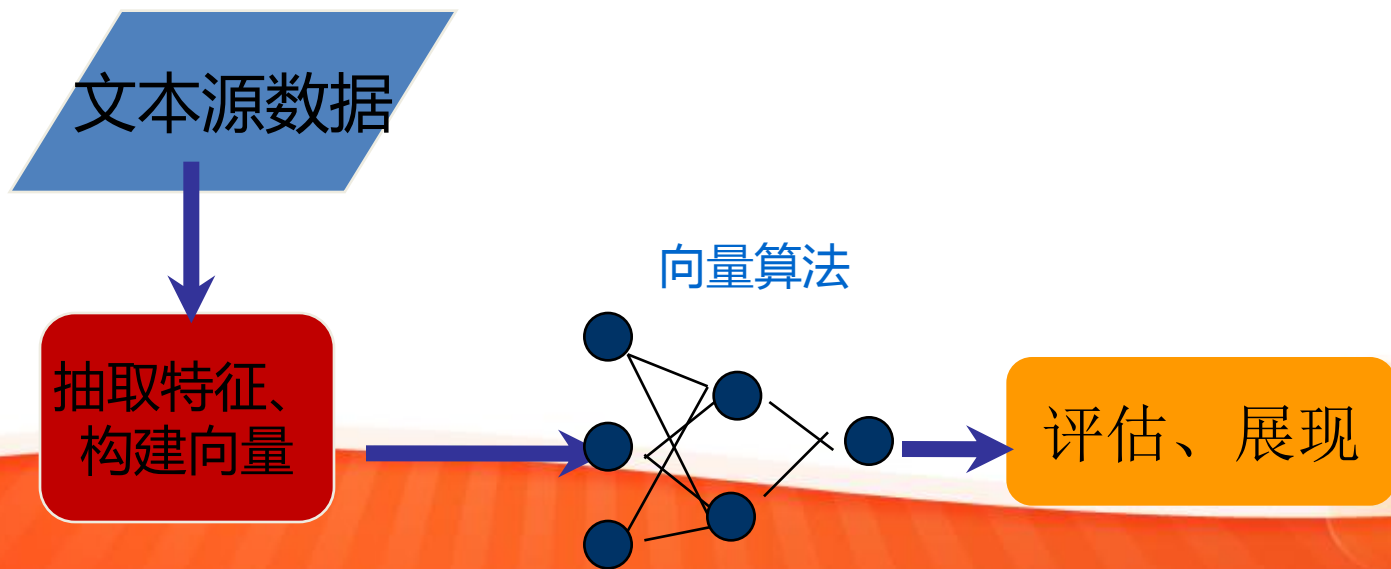


海量文本数据带来的价值和挑战

有限框架数据的补充、用户UGC信息

如何结构化、高效计算

呈现文本数据挖掘文本价值



语义分析平台架构总览

文本
数据

评价

商品标题

详情页



SNS/分享数据

特征计
算层

分词/新词

词之间相关

序列标注

句法分析

实体语义

聚类
算法

分类
算法

层次聚类

谱系聚类

KNN

Kmeans

SVM/ANN/决策树/贝叶斯

产品

大家印象

北极圈

U站推荐

UE反馈



今天的话题

- 相关业务场景: 标签、内容、分类打标
- 相关算法和问题
- 文本技术拓展



文本标签服务

内容相关、相似
框架提取

文本分类



单品标签 “大家印象”

宝贝详情	评价详情(2521)	成交记录 (2251件)
------	------------	--------------

宝贝与描述相符 4.7 分 共20378次打分 店铺评价

大家印象:

质量好 (677)版型很好 (440)穿上好看 (230)快递不错 (213)服务不错 (190)码子很准 (189)
衣衣很舒服 (112)物流慢 (97)颜色不喜欢 (66)

☒ 全部 ☐ 好评 (2449) ☐ 中评 (51) ☐ 差评 (21) ☐ 追加 (71) ☒ 有内容的评价 推荐排序



焦点宝宝5...



面料蕾丝都不错，挺漂亮，我要的是x码的，长袖，下面可能因臀围较宽99cm，所以比较撑衣服，下面的裙子和上面连接部分有凸起的轧边儿，有些不美观，准备让妈妈帮忙处理下看看，还有一种办法就是自己瘦身。

[2012.05.30] 颜色分类:蓝色长袖(热销款) 尺码:XL码

有用(0)



tqh68...

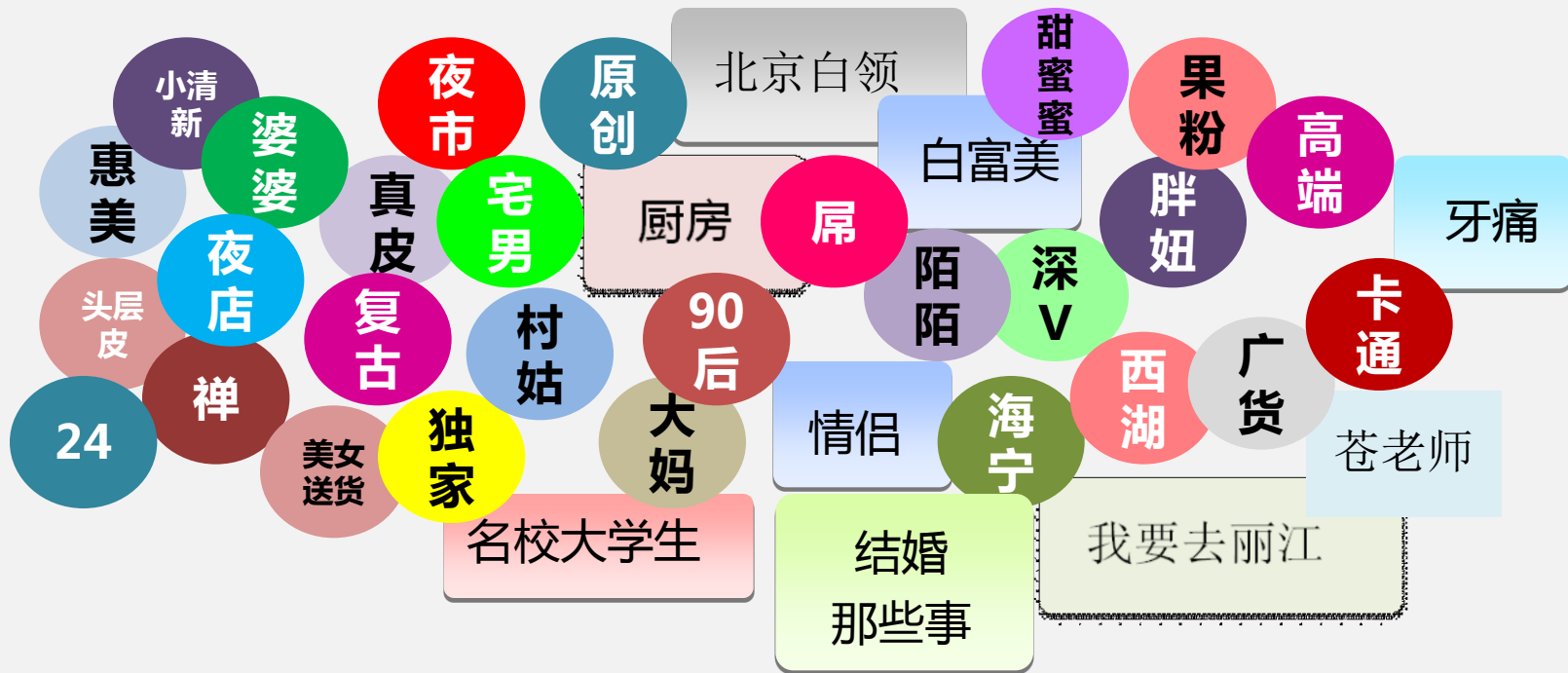


很漂亮的裙子，非常喜欢！朋友也很喜欢！超显瘦哦。卖家也很贴心，耐心！谢谢！全5分

[2012.05.21] 颜色分类:蓝色短袖(热销款) 尺码:XL码

有用(0)





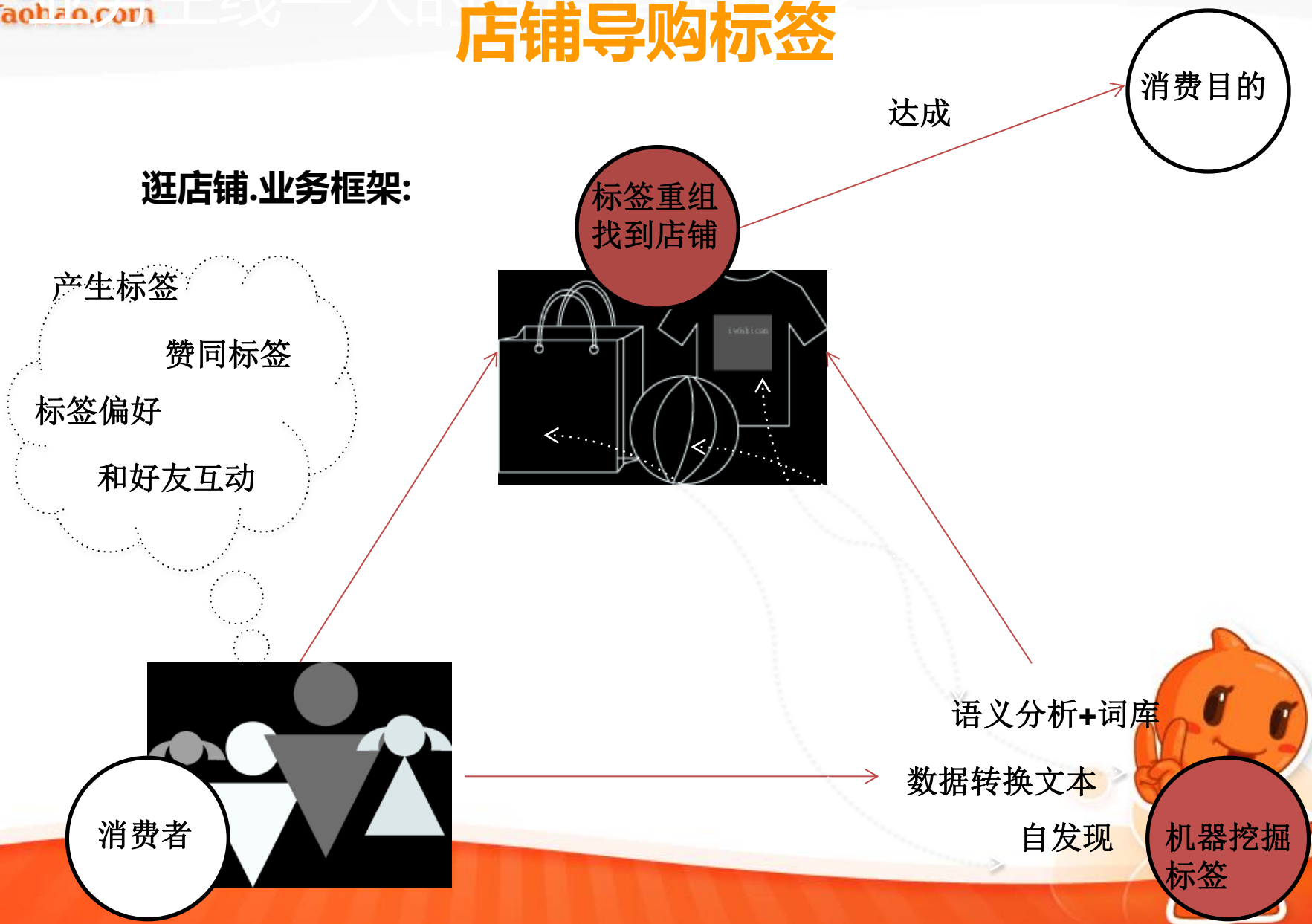
怎么样的标签的规模是合适的，长尾标签怎么处理？
需要BI提供数据模型，对标签的规模、数量提供界定范围；
根据标签覆盖的商品数来确定需要增加的标签。



业务主线——人的维度组织推荐

店铺导购标签

逛店铺.业务框架:



标签当中的问题

标签的来源

标签的属性词义

标签的关系



标签来源、分词问题

□一元分词

夏季 新款 女装 雪纺 连衣裙

□二元分词

夏季 新款 女装 雪纺 连衣裙

□CRF分词

夏季 新款 女装 雪纺 连衣裙
B E B E B E B E B M E



标签中的短语、新词来源

- 互信息、聚合度、左右熵发现二元
- 前缀树发现长字符串模式

1. Accessor Variety 语义实体的上下文独立性()

$$AV(ab)=\min(|XL|,|XR|)$$

其中 $XL=\{x|xab\text{为文档中的连续汉字串}\}$

其中 $XR=\{x|abx\text{为文档中的连续汉字串}\}$

$|XL|,|XR|$ 分别为集合 XL, XR 包含的元素个数

AV值越高说明ab上下
文独立性越强，越
可能成为实体

2.熵

衡量后缀的混乱度

$$H(X)=\sum_{x\in X}p(x)\log\frac{1}{p(x)}$$

$$H(X)=x\times\frac{1}{x}\times\log(x)\equiv\log(x)$$



序列标注问题

均可以看做概率图模型的不同表现形式

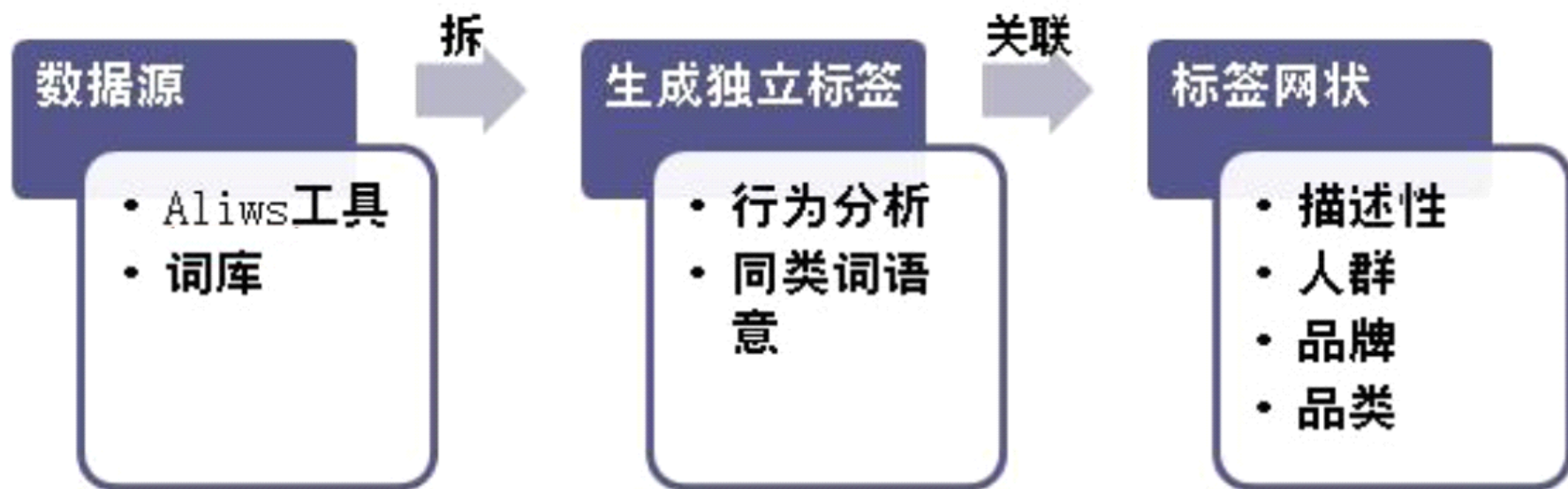
隐马尔可夫模型 (Hidden Markov Model , HMM)

最大熵模型 (Maximum Entropy Model , MEM)

条件随机场 (conditional random fields , CRF)



实体识别的标注问题



异类词（行为分析）

通过商品对应的tag与买家关联；通过tag对应买家重合度算出tag关联度；

同类词（语意分析）

先词库分析，计算词与词之间的相同出现的概率；

9: 佛珠手链是一种流行饰品



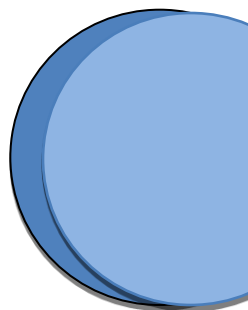
标签间的关系

手机套和手机

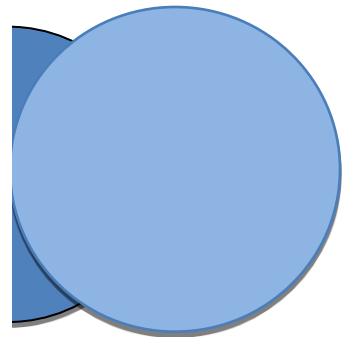
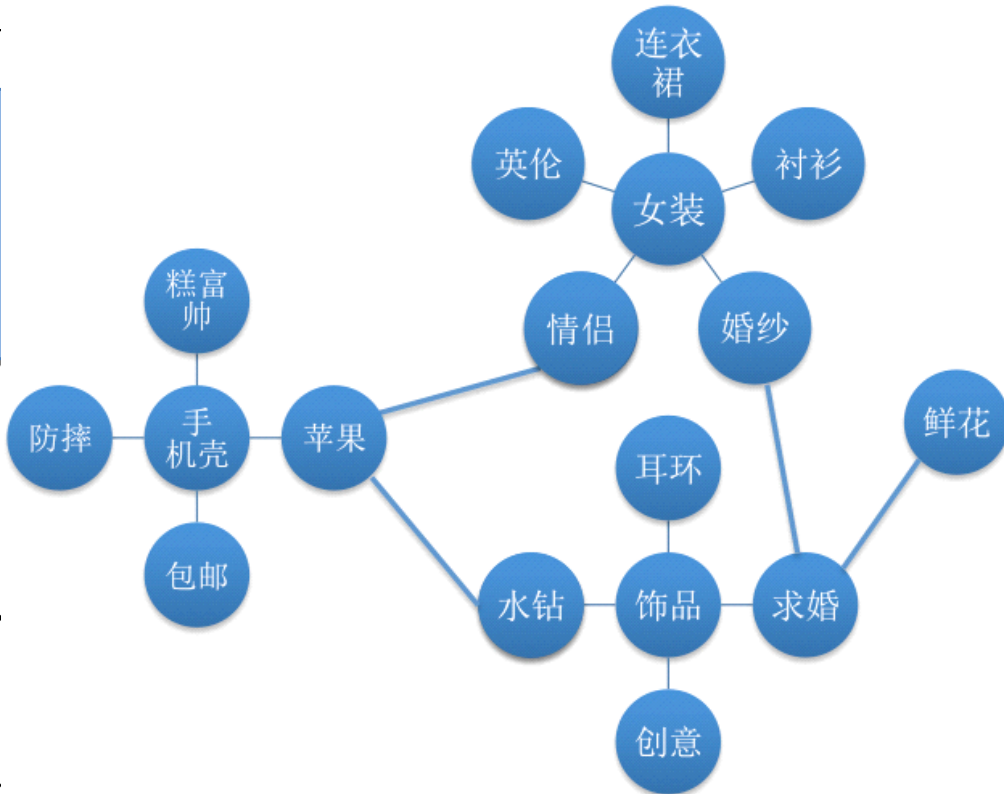
手机壳和手机

连衣裙和女

婚纱和女士



合并



关联

共同的购买人

标签的结构化.

类目、来源、变化、内容





产品框架

业务场景

首页

频道

搜索

List

场景、主题页

...

标签系统

管理层:

类目、来源、变化速度、内容...

功能层:

标签产生

商品打标

标签关联

模型层:

筛选
剔除
合并

打标商品范围确定

商品重合比例
人群选择
人群重合比例

数据源

搜索

交易

用户特征

标题

详情

评价

资讯

专辑

日记

百科



其他文本标签问题

- 标签质量判定
- 标签排序、相关性展现
- 标签合并去重



内容相关、提取问题

U站内容推荐

相似Query查询

特定内容提取



如何分析相关内容

- 人的行为
- 词之间相关性
- 句子、段落之间相关性



同义词、近义词、词之间关系

- 基于统计
- 基于词法分析
- 基于行为

用语
找出

Hub

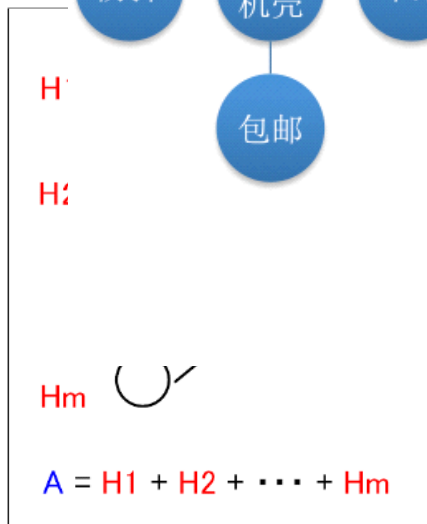
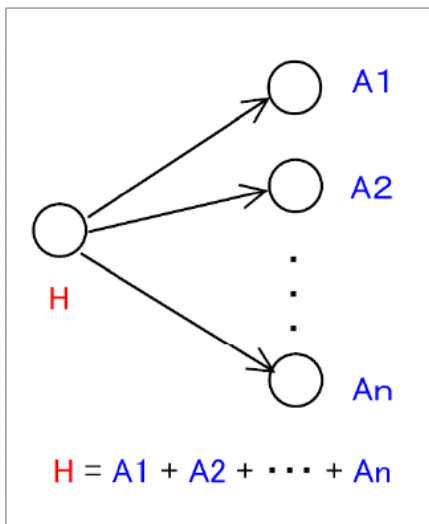
Authority

← 弓

← 被

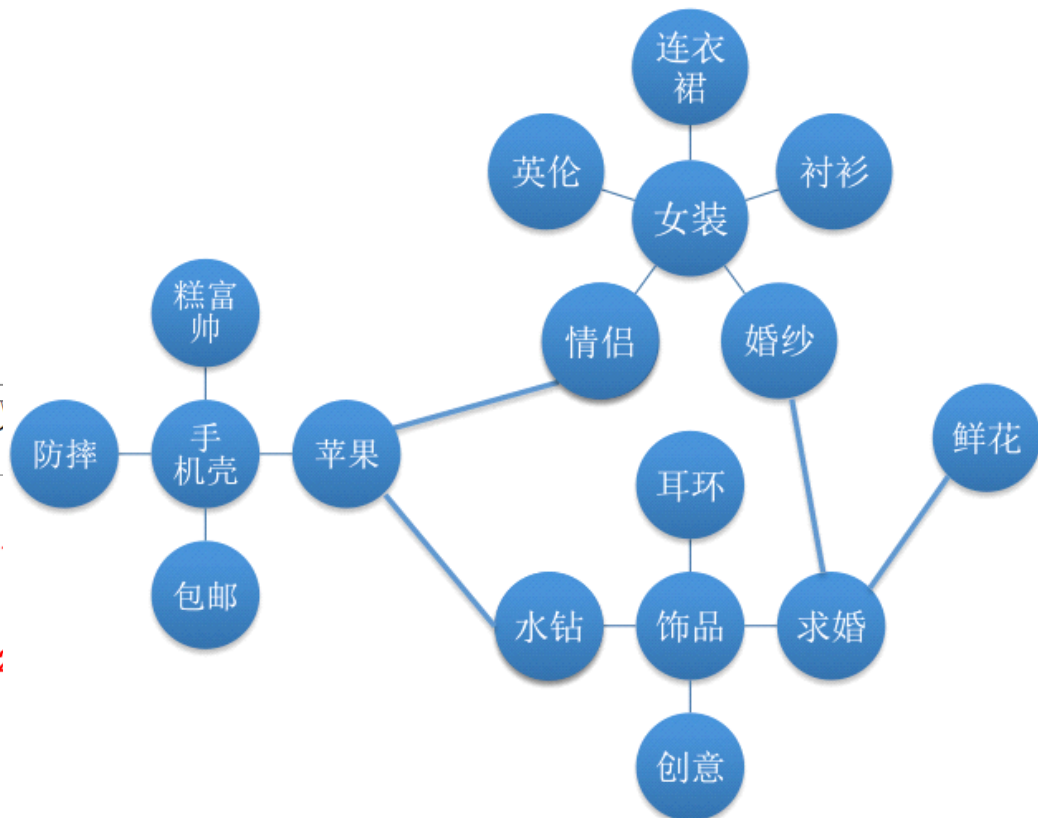
weigh

$p_{att,i}$



$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$

$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$



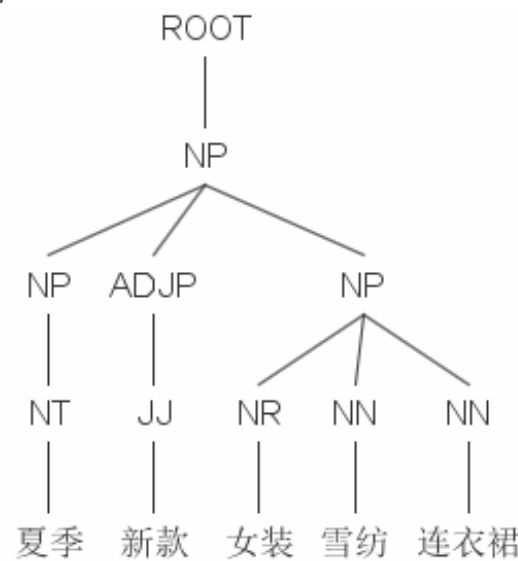
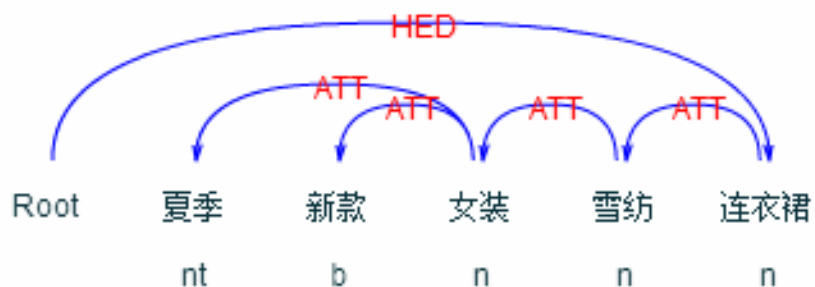
句法分析

□ 短语句法分析

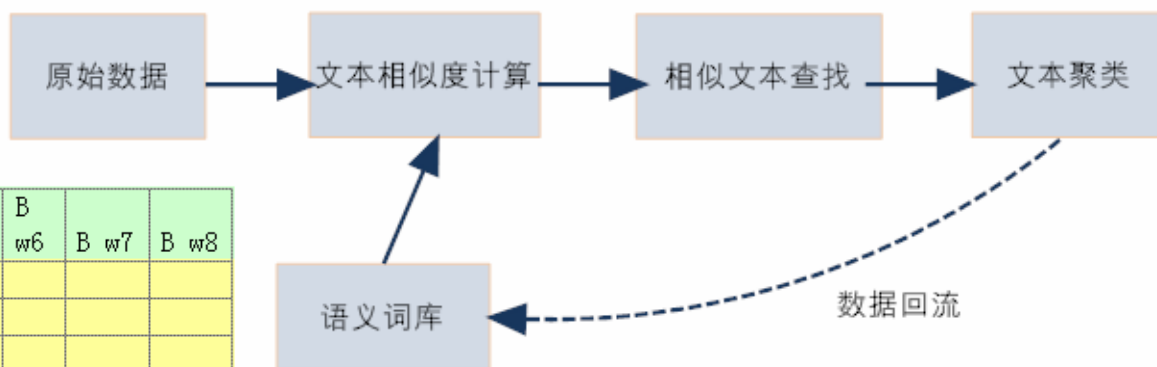
head-rule, lexicalize, grammar-based un-lexicalize

□ 依存句法分析

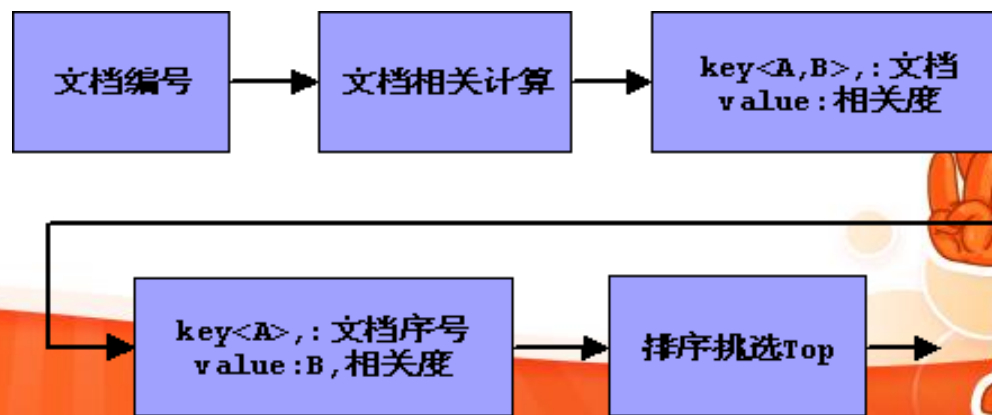
memory based , classifier-based,
feature-verification



语义推荐



	B w1	B w2	B w3	B w4	B w5	B w6	B w7	B w8
A w1								
A w2								
A w3								
A w4								
A w5								
A w6								
A w7								
A w8								
A w9								



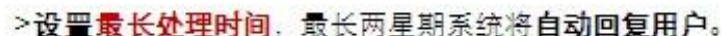
内容提取问题

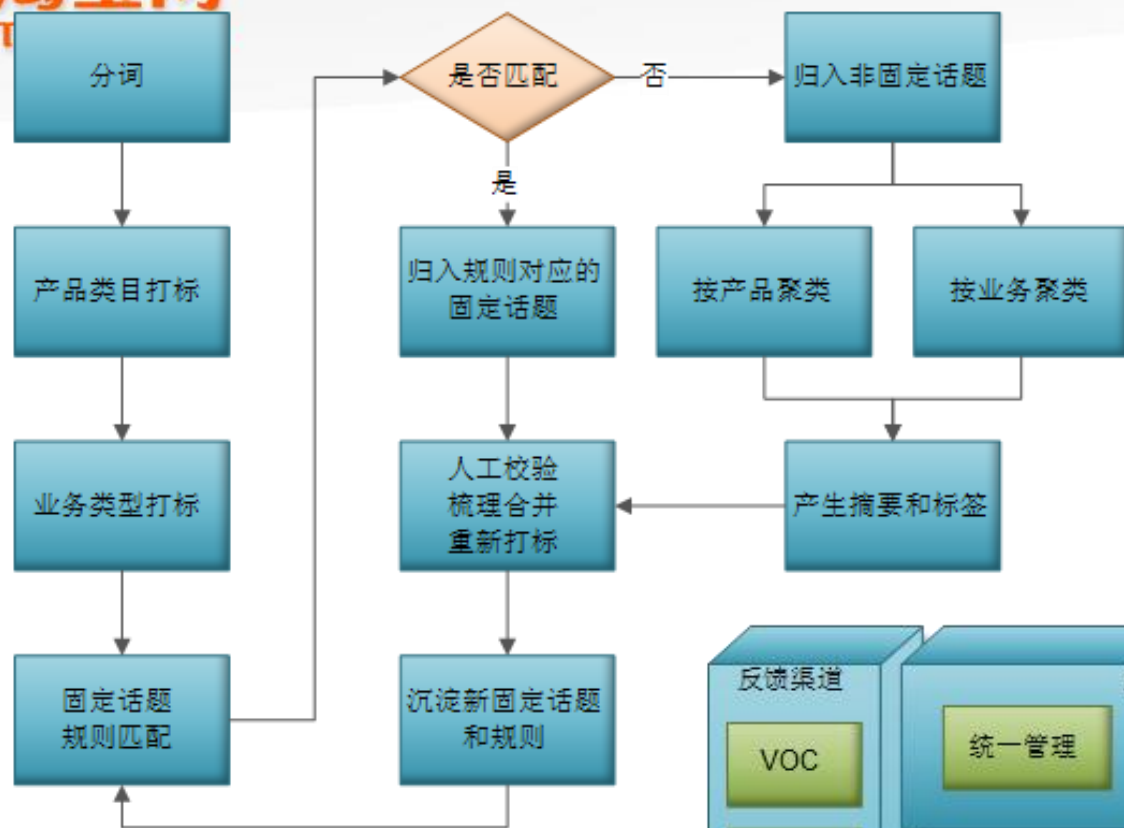
- 关键词、中心词
- 特定场景地址、礼物



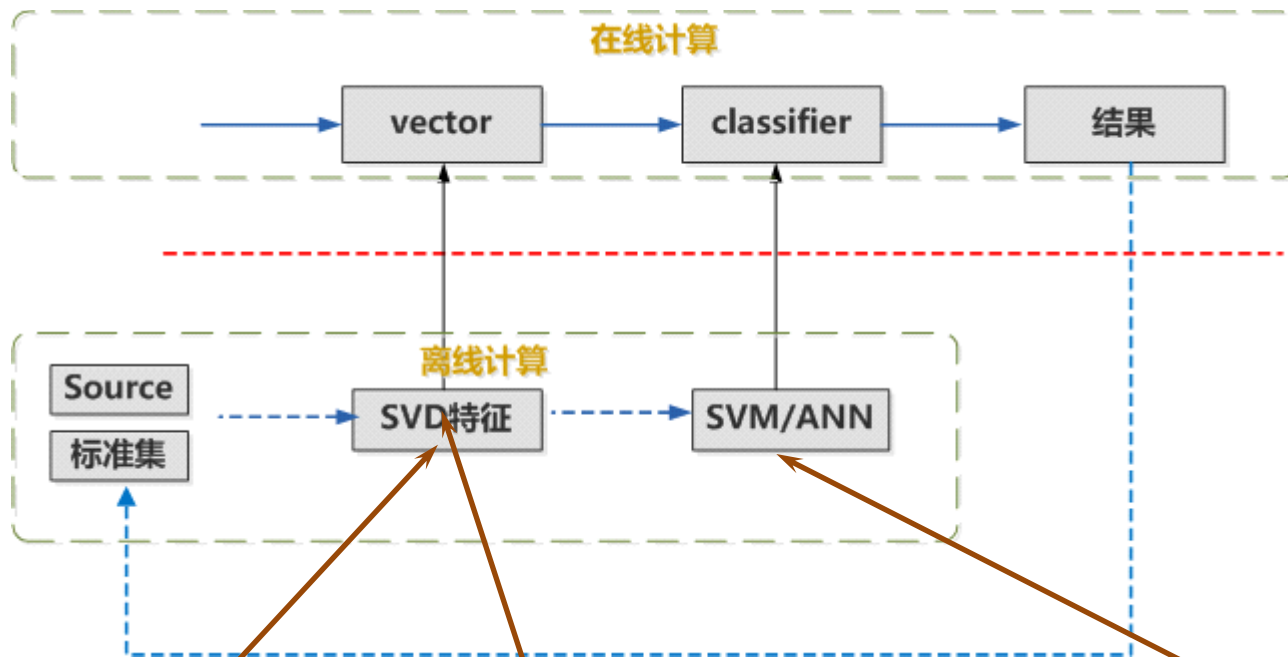
分类打标问题







分类打标问题

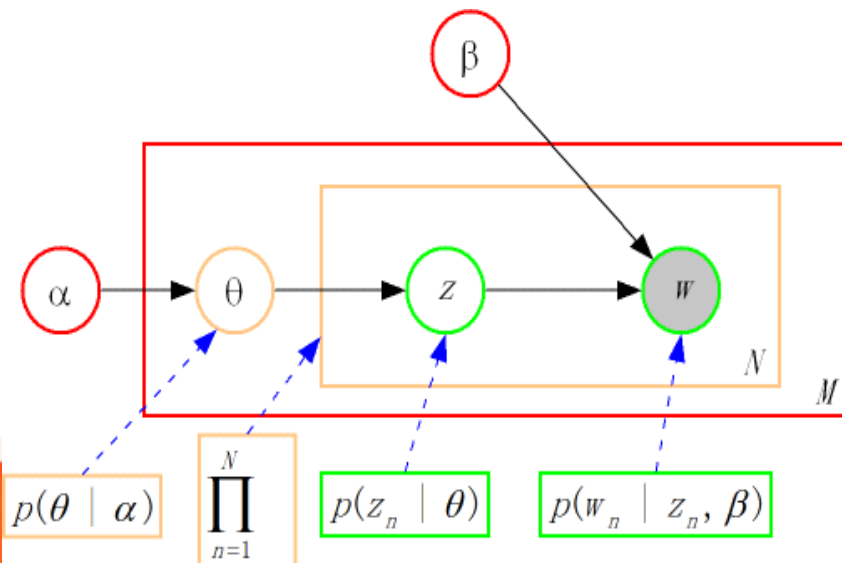
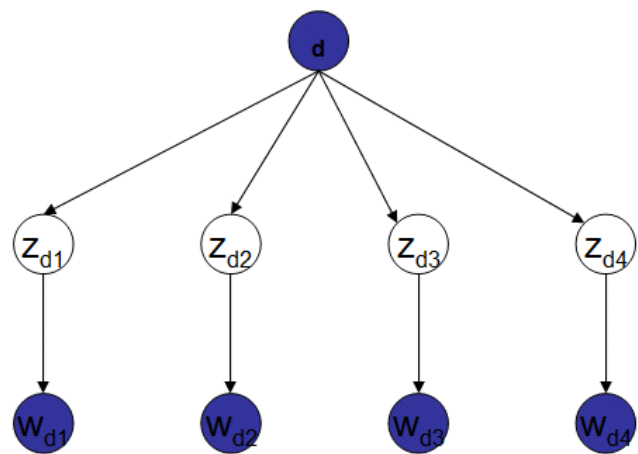
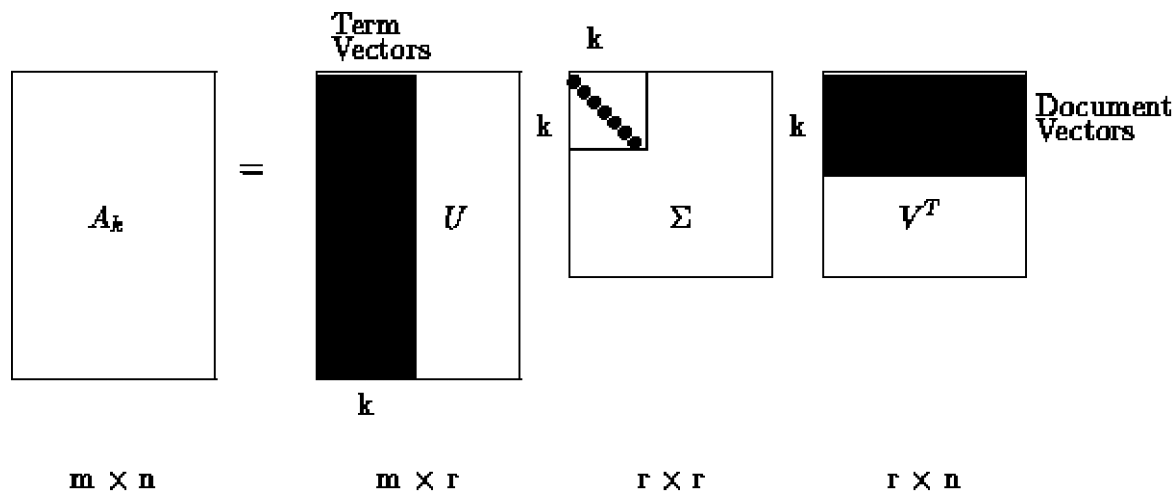


SVD, LDA

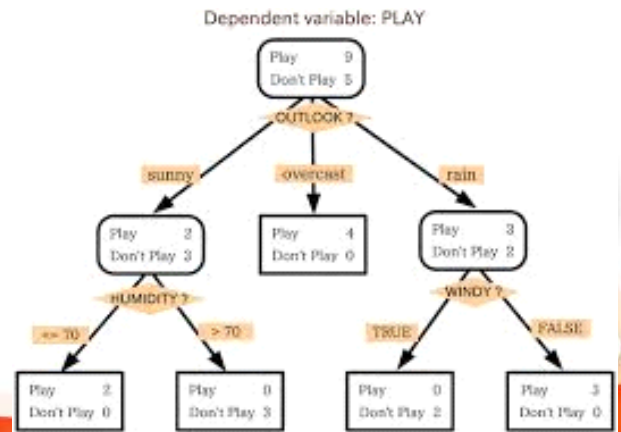
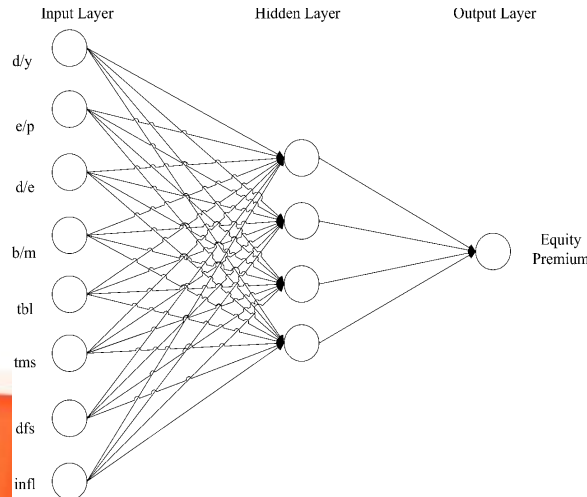
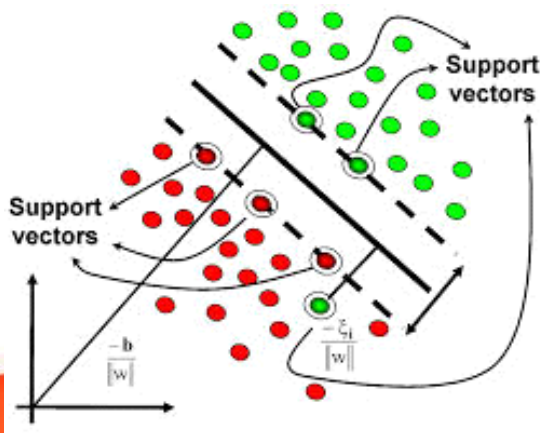
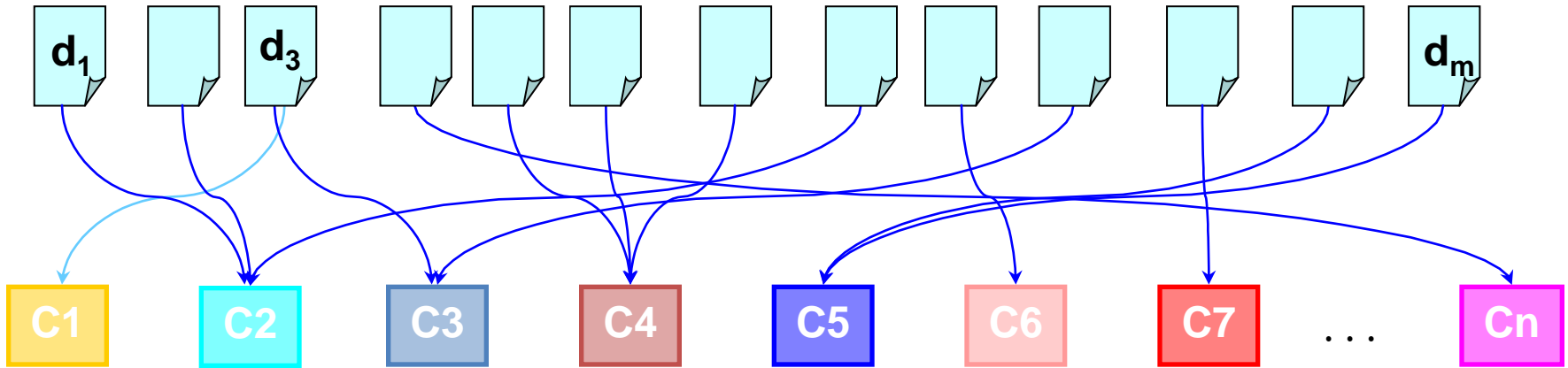
有效特征选取
、组合。特征
计算

分类器设计
。SVM、
ANN、贝叶
斯、决策树





机器分类问题



效果--
[同义] 作用 影响
[下位] 反馈 关系 害处 好处 后效 效应 效应 余波
[角色] 反作用
[上位] 反程度 结果
[原因] 影响

穿着--

[同义] 衣着 衣著 服饰 衣饰 衣装 着装 穿戴 穿章 行头
[下位] 领巾 帽子 包头 面纱 背带 吊带 兜肚 耳帽 发箍 发网 盖头 护耳 护颈 护腿 护膝
[角色] 围嘴 戏装 鞋 胸围 眼镜 眼罩 腰带 腰封 衣服 髭髭
[上位] 衣物
[部分] 筒子
[属性] 挺脱
[功能] 着装

衣服--

[同义] 服装 衣衫 衣裳 衣 装 衣袂
[下位] 上衣 下装 袍子
· 内衣 外衣
· 春装 夏装 秋装 冬装
· 中装 西服 查曼多 德瑞莎 韩服 和服
· 男装 女装 童装
· 百家衣 百衲衣 便服 便服 布衣 蝉衣 成服 成衣 单衣 登山服 法衣 防化服 工
· 航天服 号衣 华服 嫁衣 接衫 礼服 列宁服 露脐装 罗衫 迷彩服 偏衫 奇装异服
· 时装 寿衣 睡衣 俗装 素服 套装 透视装 晚装 舞服 戏衣 学生装 血衣 泳衣 浴
[上位] 衣着
[部分] 衣服的部分
[材料] 衣料
[是...的地点] 揣 笔挺 褴褛 单薄 郎当 麻花
[属性] 合身 肥瘦 服色
[拥有] 上身 绌绌 撞衫
[发出动作] 翻改 更衣 和衣 宽衣解带 置装
[接受动作]

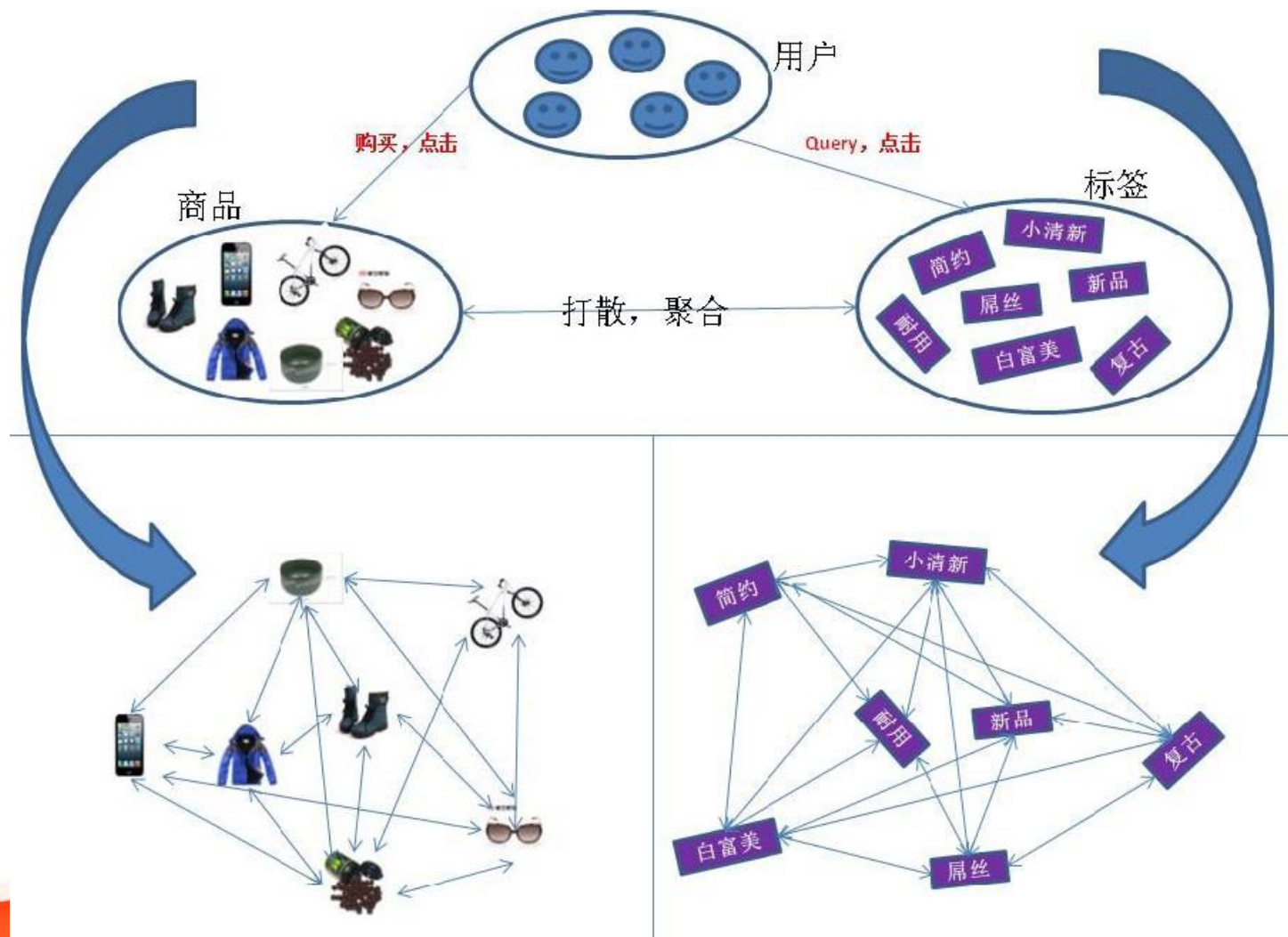
卖家提供

新词、短
语
发现

和已有词
去除重

后产出找
的的词

小结



文本技术相关介绍



数据聚类、天然类目、类目团簇

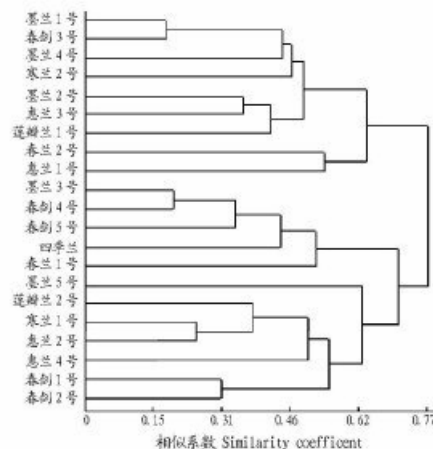
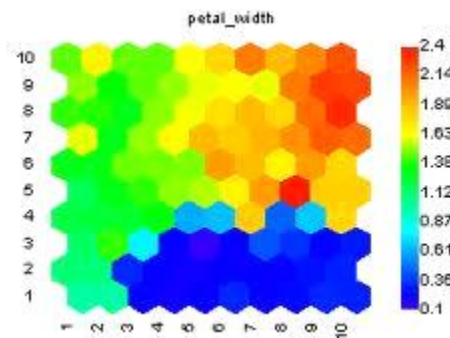
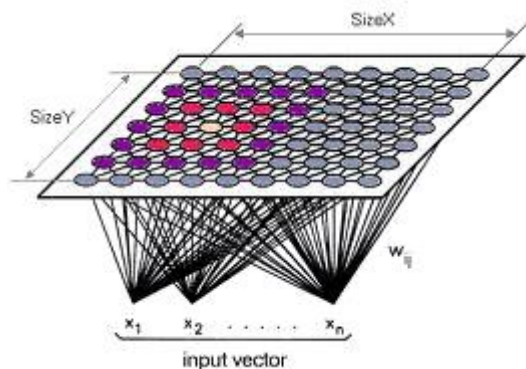
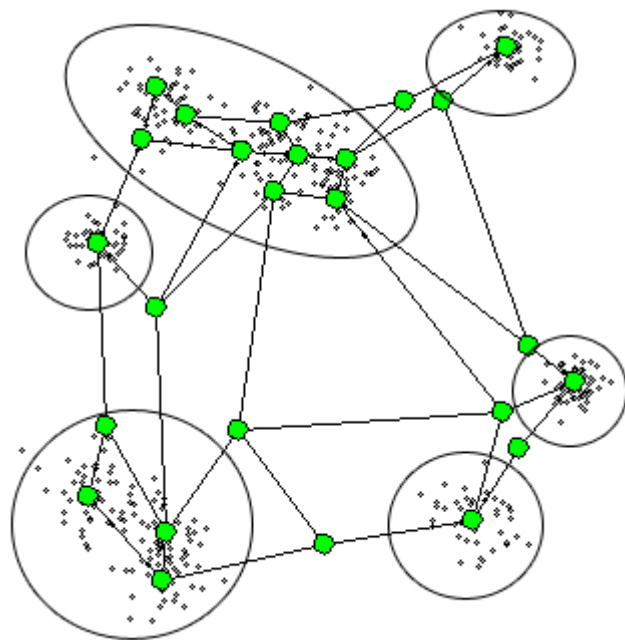
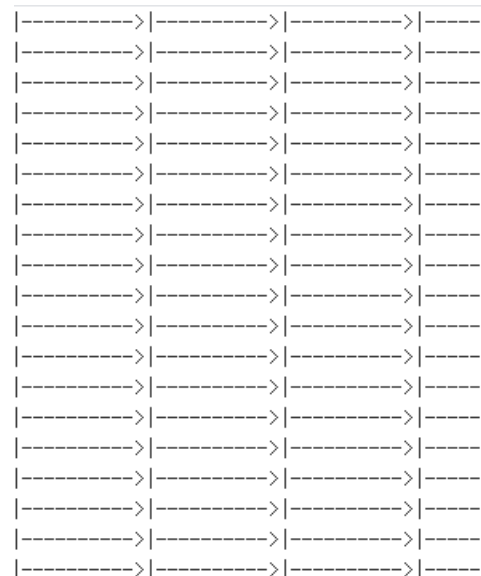


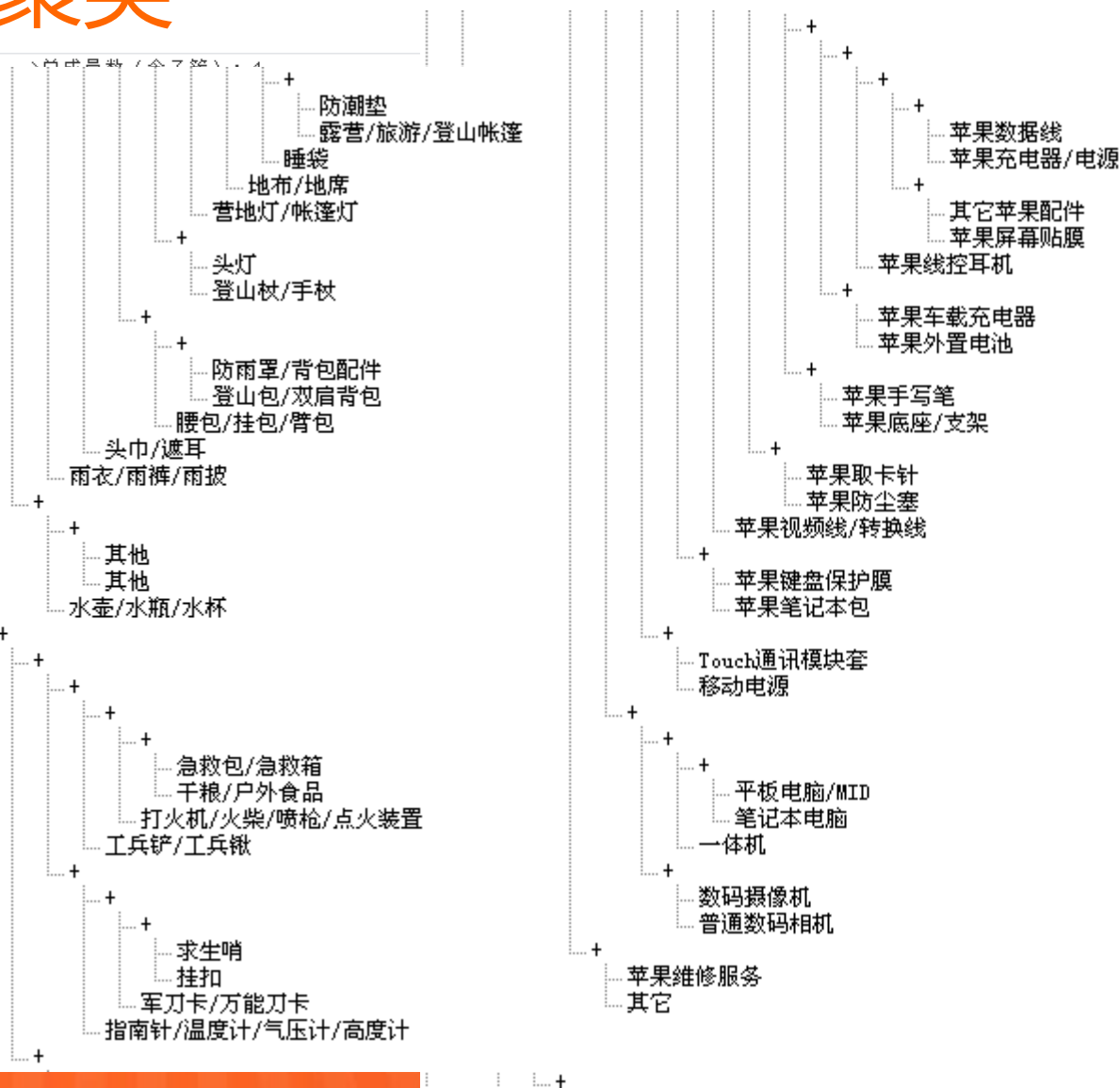
图 2 ERAPD 聚类分析结果



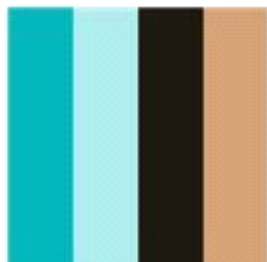
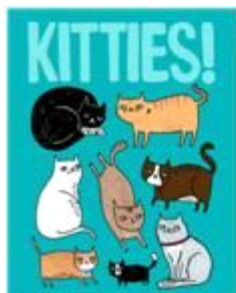
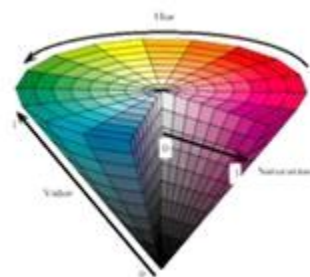
谱系、层次聚类



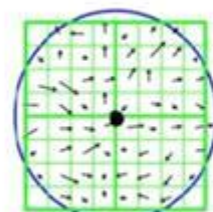
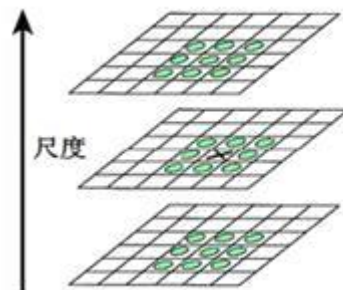
变 ??37天甩14斤 ??谁瘦了都很美 ?? ?? ?? 注
多姐妹的支持和帮顶, 在这里小妹谢了 瘦身专家
变 ??37天甩14斤 ??谁瘦了都很美 ?? ?? ?? 注
多姐妹的支持和帮顶, 在这里小妹谢了 瘦身专家
变 ??37天甩14斤 ??谁瘦了都很美 ?? ?? ?? 注
多姐妹的支持和帮顶, 在这里小妹谢了 瘦身专家
变 ??37天甩14斤 ??谁瘦了都很美 ?? ?? ?? 注
多姐妹的支持和帮顶, 在这里小妹谢了 瘦身专家



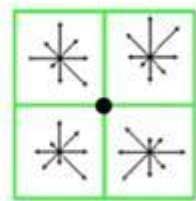
图片语义应用



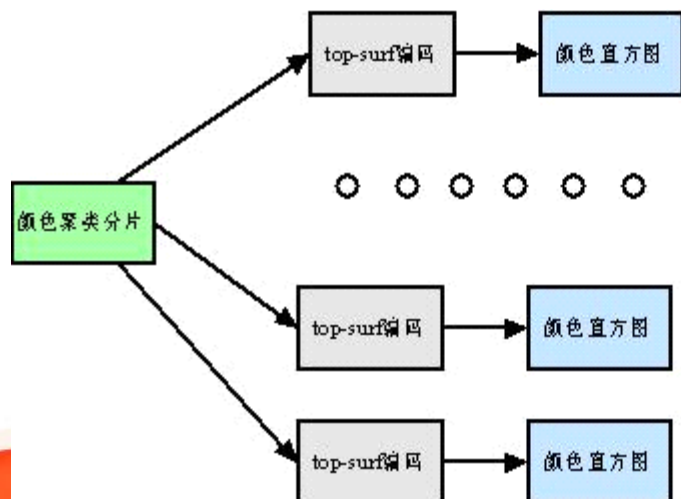
0.60== 135070, 3.39, 179.25, 184.17 =====
 0.18== 39755, 171.14, 230.79, 232.35 =====
 0.13== 29245, 29.35, 24.46, 16.35 =====
 0.10== 22548, 210.60, 158.77, 115.21 =====



领域梯度方向



关键点特征向量



➡ 选出结果



谢谢

