



北京邮电大学

BEIJING UNIVERSITY OF POSTS AND TELECOMMUNICATIONS



计算机应用编程 实验四

找你邻居家的妹

熊永平@网络技术研究院

ypxiong@bupt.edu.cn

周一10:10-12:00@3-134

2015.12.20

课程表

学 期	秋 季 学																				寒 假					
年 份	二 〇 一 五 年															二 〇 一 六 年										
月 份	九 月				十 月				十 一 月				十 二 月				一 月				二 月					
周 次	〇	一	二	三	四	五	六	七	八	九	十	十一	十二	十三	十四	十五	十六	十七	十八	十九	二十	廿一	廿二	廿三	廿四	廿五
星期一	31	7	研究生 新生上课	21	22	23	12	19	26	2	9	16	23	30	7	14	21	28	4	11	18	25	1	春节	15	22
星期二	1	本科生 新生报到	5	22	29	国庆节 假期	13	20	27	3	10	17	24	1	8	15	22	29	5	12	19	26	2	9	16	23
星期三	2		6	23	30		14	21	28	4	11	18	25	2	9	16	23	30	6	13	20	27	3	10	17	24
星期四	3	10*	7	24		8		22	29	5	12	19	26	3	10	17	24	31	7	14	21	28	4	11	18	25
星期五	4	研究生 开学典礼	8	25		9		23	30	6	13	20	27	4	11	18	25	元旦	8	15	22	29	5	12	19	26
星期六	5	12	9	26		10		24	31	7	14	21	28	5	12	19	26		9	16	23	30	6	13	20	27
星期日	6	1	10	中 国 民 主 义 青年 节	11		25	1	8	15	22	29	5	12	19	26	3	10	17	24	31	7	14	21	28	

注:

以党政办

科

制日

课程介绍

实验一讨论课

实验二讨论课

实验三讨论课

实验四讨论
论结束

实验三：发现相关词

实验目标

➤ 目标

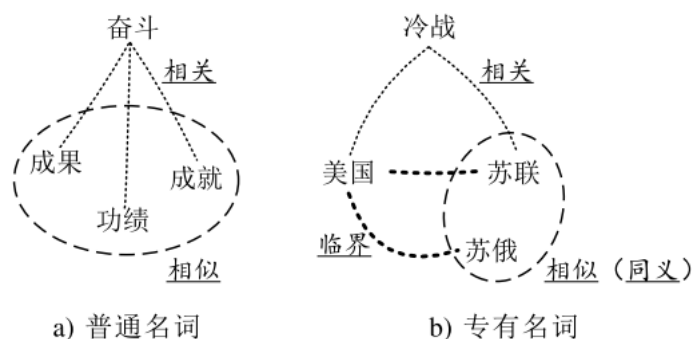
- ❑ 设计一个分析程序，实现
- ❑ 对一个大规模文本语料库分词
- ❑ 计算词之间的距离
- ❑ 寻找最相关的前20个词对

➤ 编程技能

- ❑ Java语言练习
- ❑ NLP
- ❑ 统计距离
- ❑ 概率模型
- ❑ 数据挖掘

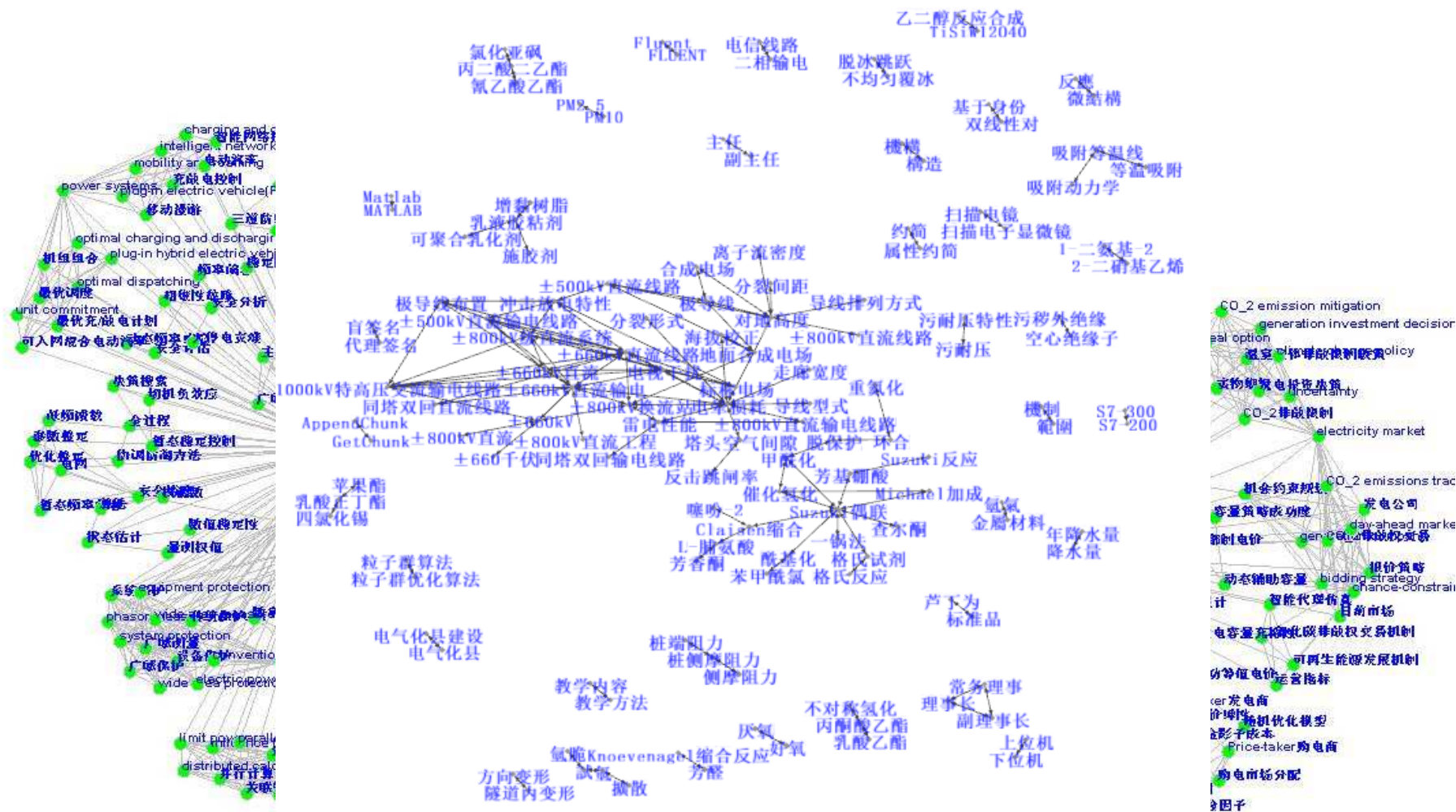
相关词是什么？

- 共同出现，co-occurrence
- 同时发生、存在或出现
- 在一个文本语料库中，两个术语超出正常频率的同时发生、存在或出现。
- 在语言学里，共现可以解释为语义上相似的指标或者是固定搭配（成语）表达。

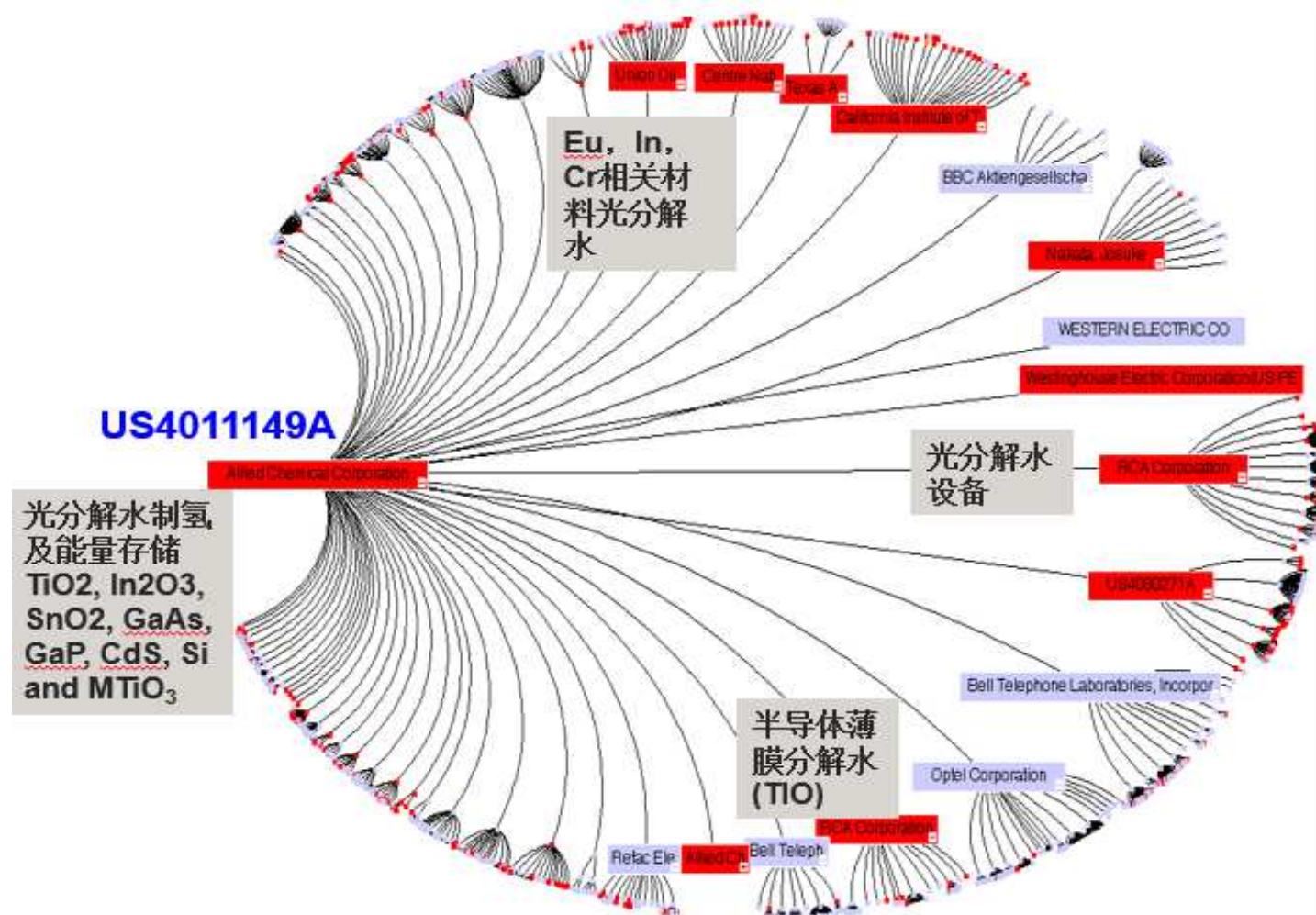


© 2015 Pearson Education, Inc. or its affiliate(s). All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or by any information storage or retrieval system, without prior written permission from Pearson Education, Inc. or its affiliate(s).

Downloaded from <http://ajphaphysocpharm.sagepub.com/> at 11:05 11 November 2014



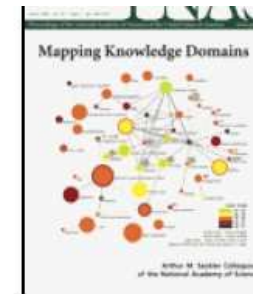
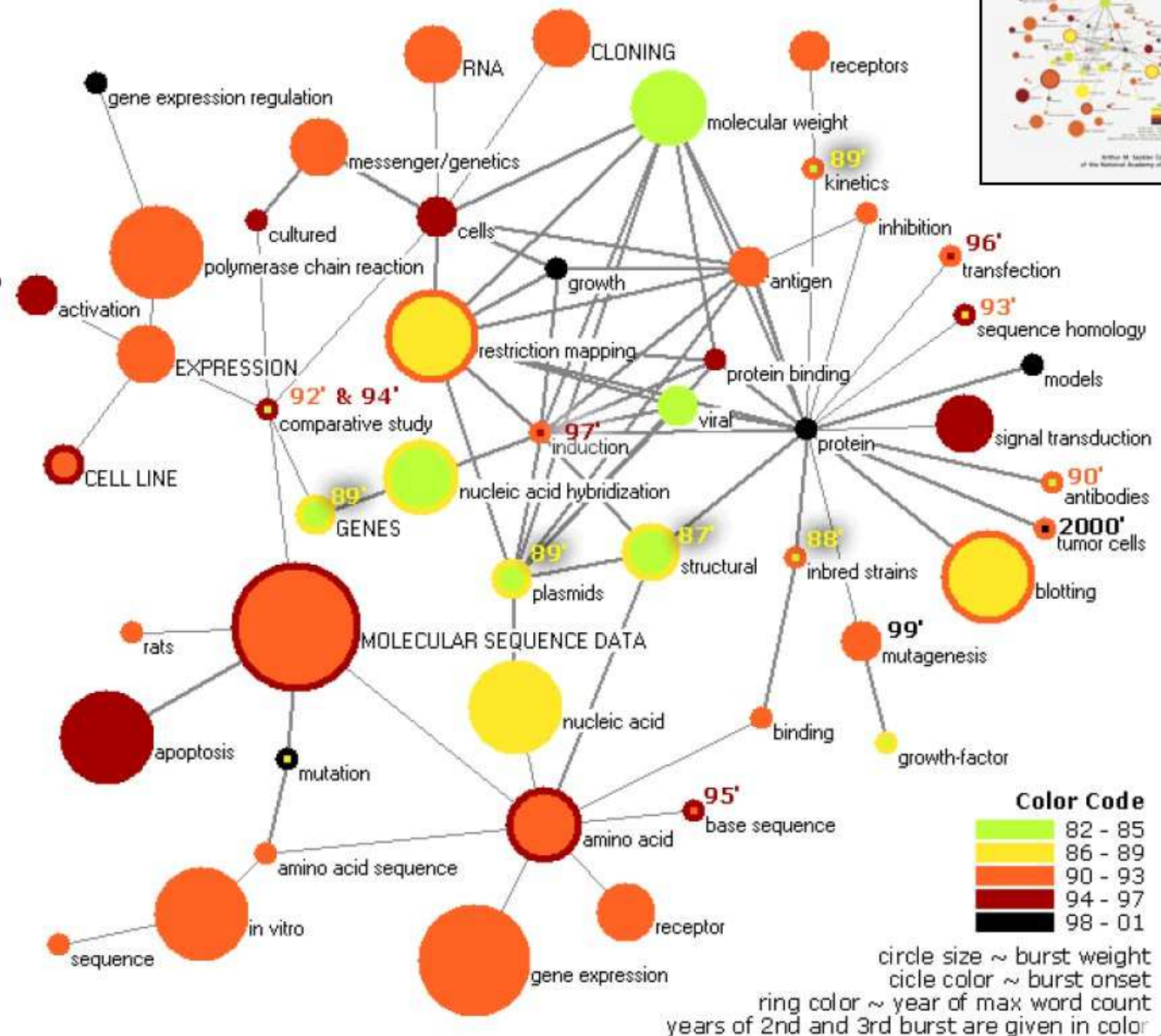
相关词意义



Mapping Topic Bursts

Co-word space of the top 50 highly frequent and bursty words used in the top 10% most highly cited PNAS publications in 1982-2001.

*Mane & Börner. (2004)
PNAS, 101(Suppl. 1):
5287-5290.*



中文分词

➤ 作用

- ❑ 找出单个词汇

➤ 中文文本处理的第一步

- ❑ 我国科学家近日研制出一套水下反恐监控系统
- ❑ 我国/科学家/近日/研制/出/一套/水下/反恐/监控/系统

➤ 分词方法

- ❑ **基于词典的方法**：给出一部词典，根据这部词典进行匹配
- ❑ 无词典的方法：不需要词典，根据某种人工构词规则或者统计规则从字生成词

中文分词实际挑战

➤ 中文分词歧义

❑ 交集型

- “部分居民生活水平”：分居、居民、民生、生活、
- “我们小组合成氢气”：我们/小组/合成/氢气或我们/小/组合/成/氢气

❑ 组合型

- “老人家”：老人、老人家
- 他/从/马/上/下/来；我/马上/就/来/了

➤ 未登录词

- ❑ 专有名词（人名、地名、机构名、译名、术语等）、新词



本实验不考虑

中文分词方法1

➤ 正向最大匹配(基于词典的方法)

最大匹配法的分词实现很简单，并且可以满足一些对分词准确率要求不高的检索系统，该方法在早期的分词系统中被广泛使用

0 1 2 3 4 5 6

他	说	的	确	实	在	理
---	---	---	---	---	---	---

指针位置	剩余词串	首字	最大匹配词条
0	他说的确实在理	他	他
1	说的确实在理	说	说
2	的确实在理	的	的确
4	实在理	实	实在
6	理	理	理

中文分词方法2

➤ 逆向最大匹配(基于词典的方法)

0 1 2 3 4 5 6

他	说	的	确	实	在	理
---	---	---	---	---	---	---

指针位置	剩余词串	尾字	最大匹配词条
6	他说的确实在理	理	在理
4	他说的确实	实	确实
2	他说的	的	的
1	他说	说	说
0	他	他	他

几个著名的中文分词器

➤ Paoding Analysis（庖丁解牛）

- ❑ Lucene御用中文分词
- ❑ 48个java文件，6895 行，使用不用的 Knife 切不同类型的流

➤ Imdict

- ❑ imdict智能词典，使用ICTCLAS HHMM隐马尔科夫模型中文分词
- ❑ 20个java文件，2399行

➤ mmseg4j

- ❑ 用 Chih-Hao Tsai 的 MMSeg 算法 实现的中文分词器
- ❑ 23个java文件，2089行

➤ ik

- ❑ 采用了特有的“正向迭代最细粒度切分算法”
- ❑ 多子处理器分析模式
- ❑ 22个java文件，4217行

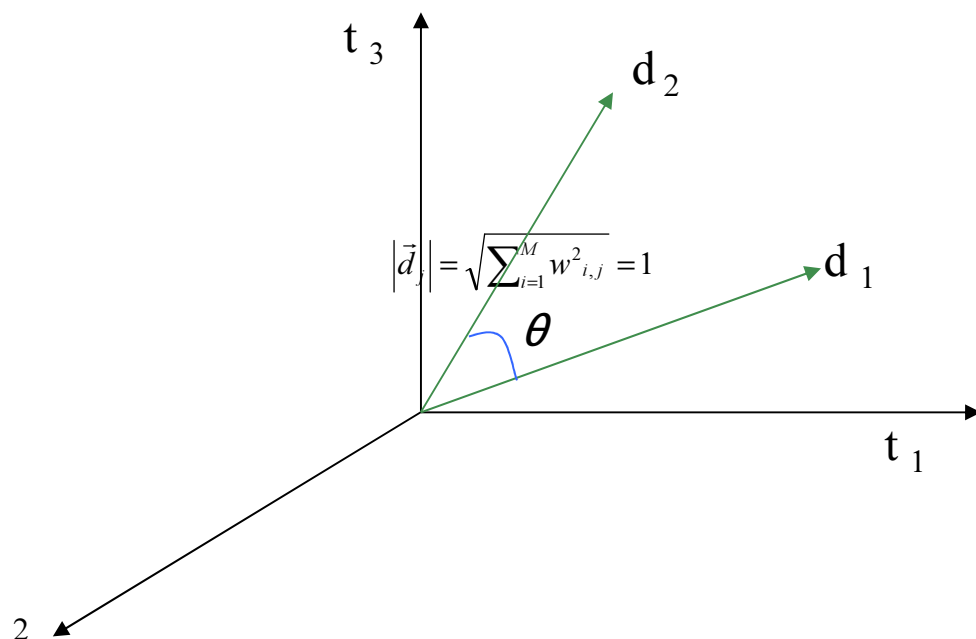
建立统计窗口

同时，除了解发文量较高的研究机构分布之外，还需进一步分析相对活跃的机构之间的合作关系。研究人员可以据此开展访问交流，或合作研究，**从而进一步推动研究的进展。**如图 7 所示，**东京大学不仅在动画相关领域中发文量较多，**在国际合作中也表现得相对活跃。另外，斯坦福大学及多伦

以标点符号分隔的一句话作为统计窗口

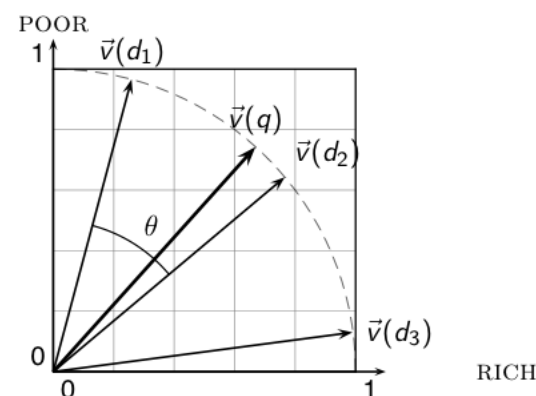
普通方法（1）： 向量空间模型

Cosine similarity 余弦相似度



$$\text{sim}(d_j, d_k) = \frac{\vec{d}_j \cdot \vec{d}_k}{|\vec{d}_j| |\vec{d}_k|} = \frac{\sum_{i=1}^M w_{i,j} w_{i,k}}{\sqrt{\sum_{i=1}^M w_{i,j}^2} \sqrt{\sum_{i=1}^M w_{i,k}^2}}$$

- 向量 d_1 和 d_2 的 “closeness” 可以用它们之间的夹角大小来度量
- 具体的，可用cosine of the angle θ 来计算向量相似度.
- 向量按长度归一化
Normalization



构建向量

假定生成了N个统计窗口
每个词对应一个N个元素的向量
计算所有两两词对之间的Cos距离

	窗口1	窗口2	窗口3	窗口4
词I	0	1	0	1	
词j	0	1	1	0	

普通方法 (2) : 互信息排序

信源

- 信源：信息来源。
 - 离散信源
 - 连续信源
- 通常用随机变量 x 来表示一个离散信源
- 思考：当信源发出某个信号 x 时，它提供了多少信息？
- 进一步思考：什么样的事件信息更多？

自信息

- 事件的不确定性决定了信息的多少
- 我们用自信息 $I(x)$ 来表示事件 x 所具有的信息，那么 $I(x)$ 应该是 $p(x)$ 的一个函数。
- 思考： $I(x)$ 应该有哪些性质？
 - ❑ 非负性： $I(x) \geq 0$
 - ❑ $P(x)=0, I(x)=+\infty; P(x)=1, I(x)=0$
 - ❑ 单调性：若 $p(x) > p(y)$ ，那么 $I(x) < I(y)$
 - ❑ 若 $p(x,y) = p(x)p(y)$ ，那么 $I(xy) = I(x) + I(y)$

$$I(x) = C * \log \frac{1}{p(x)}$$

最优编码

- 实例：假定有一个房间中有时没有人，有时甲在房间中，有时乙在房间中，有时甲乙都在房间中，房间状态服从下面的概率分布

房间状态	房间没有人	甲在房间	乙在房间	甲乙均在房间
概率	0.5	0.125	0.125	0.25

- 某人受命监视房间，每五分钟记录一次房间状态，并经一个通讯设备将房间状态发送出去。
- 问题：状态如何编码？

最优编码

- 一种可行的定长编码方案：用00 表示没有人在房间中，01 表示甲在房间中，10 表示乙在房间中，11 表示甲乙两人均在在房间中。
- 按照这样的编码，发送一个消息所需要的码的长度为2，平均发送一个消息需要2个二进制位。
- 思考：有没有一种更短的编码方式？

消息	编码
房间没有人	0
甲在房间	110
乙在房间	111
甲乙均在房间	10

最优编码

- 如果消息 x 的概率为 $p(x)$ ，则给其分配一个长度为 $\lceil -\log_2 p(x) \rceil$ 个二进制位的编码
- 平均发送一个消息所需要的编码的长度（以二进制位衡量）：

$$-\sum p(x) \log_2 p(x)$$

熵

定义1 熵 设 X 是取有限个值的随机变量，它的分布密度为

$$p(x) = P\{X=x\}, \quad \text{且 } x \in X$$

则， X 的熵定义为

$$H(X) = - \sum_{x \in X} p(x) \log_a p(x)$$

规定 $0 \log_a 0 = 0$

通常 $a=2$ ，此时熵的单位为比特。

熵的基本性质：

1. $H(X) \geq 0$ ，等号表明确定场(无随机性)的熵最小。
2. $H(X) \leq \log |X|$ ，等号表明等概场的熵最大。



熵描述了随机变量的不确定性。

熵

- 假定一种语言P有6个字母p、t、k、a、i、u，每个字母的概率为：

P	<i>p</i>	<i>t</i>	<i>k</i>	<i>a</i>	<i>i</i>	<i>u</i>
概率	1/8	1/4	1/8	1/4	1/8	1/8

则随机变量 P 的熵为：

$$\begin{aligned} H(P) &= - \sum_{i \in \{p, t, k, a, i, u\}} p(i) \log p(i) \\ &= - \left[4 \times \frac{1}{8} \log \frac{1}{8} + 2 \times \frac{1}{4} \log \frac{1}{4} \right] \\ &= 2 \frac{1}{2} \text{ bit} \end{aligned}$$

语言的字母熵

联合熵与条件熵

- ◆ **定义2 联合熵** 设 X 、 Y 是两个离散型随机变量，它们的联合分布密度为 $p(x,y)$ ，则 X,Y 的联合熵定义为：

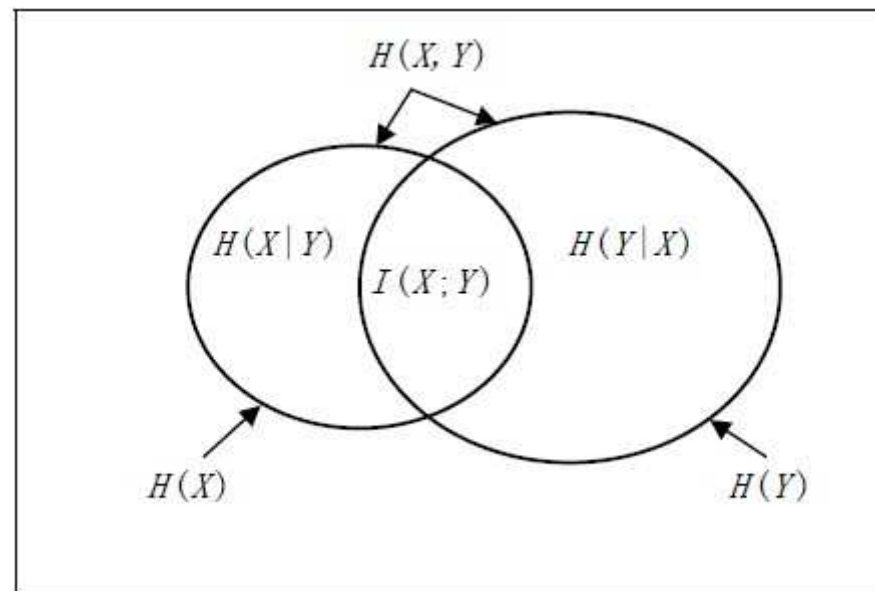
$$H(X,Y) = - \sum_{x \in X} \sum_{y \in Y} p(x,y) \log p(x,y)$$

- ◆ **定义3 条件熵** 设 X 、 Y 是两个离散型随机变量，它们的联合分布密度为 $p(x,y)$ ，则给定 X 时 Y 的条件熵定义为：

$$\begin{aligned} H(Y|X) &= - \sum_{x \in X} p(x) H(Y|X=x) \\ &= \sum_{x \in X} p(x) \left[- \sum_{y \in Y} p(y|x) \log p(y|x) \right] \\ &= - \sum_{x \in X} \sum_{y \in Y} p(x,y) \log p(y|x) \end{aligned}$$

- ◆ **链式规则** $H(X,Y) = H(X) + H(Y|X)$

- $H(X)$ 和 $H(X|Y)$ 的差称为互信息，一般记作 $I(X;Y)$ 。 $I(X;Y)$ 描述了包含在 X 中的有关 Y 的信息量，或包含在 Y 中的有关 X 的信息量。



互信息

➤ 如何量化互信息？

$$I(X;Y) = H(X) - H(X|Y)$$

$$= H(X) + H(Y) - H(X,Y)$$

$$= \sum_x p(x) \log \frac{1}{p(x)} + \sum_y p(y) \log \frac{1}{p(y)} + \sum_{x,y} p(x,y) \log p(x,y)$$

$$= \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)}$$

点间互信息

- 更为常用的是两个具体事件之间的
 - 我们常说的“互信息”，一般称之为点间互信息。
- 点间互信息：事件 x, y 之间的互信息定义为：

$$I(x, y) = \log \frac{p(x, y)}{p(x)p(y)}$$

- 点间互信息度量两个具体事件之间的相关程度
 - 当时 $I(x, y) \gg 0$ ， x 和 y 高度相关。
 - 当时 $I(x, y) = 0$ ， x 和 y 高度相互独立。
 - 当时 $I(x, y) \ll 0$ ， x 和 y 呈互补分布

计算方式

同时，除了解发文量较高的研究机构分布之外，还需进一步分析相对活跃的机构之间的合作关系。研究人员可以据此开展访问交流，或合作研究，从而进一步推动研究的进展。如图 7 所示，东京大学不仅在动画相关领域中发文量较多，在国际合作中也表现得相对活跃。另外，斯坦福大学及多伦

$P(x,y)$ 计算所有两两词对同时出现在一个滑动窗口内的概率

$P(x)$ 计算单个词出现的概率

高阶方法（1）： 频繁模式挖掘

关联规则简介

- 关联规则反映一个事物与其他事物之间的相互依存性和关联性。如果两个或者多个事物之间存在一定的关联关系，那么，其中一个事物就能够通过其他事物预测到。
 - 典型的关联规则发现问题是对超市中的货篮数据（Market Basket）进行分析。通过发现顾客放入货篮中的不同商品之间的关系来分析顾客的购买习惯。
-

什么是关联规则挖掘

- 关联规则挖掘
 - 首先被Agrawal, Imielinski and Swami在1993年的SIGMOD会议上提出
 - 在事务、关系数据库中的项集和对象中发现频繁模式、关联规则、相关性或者因果结构
 - **频繁模式**: 数据库中频繁出现的项集
 - 目的: 发现数据中的规律
 - 超市数据中的什么产品会一起购买? — 啤酒和尿布
 - 在买了一台PC之后下一步会购买?
 - 哪种DNA对这种药物敏感?
 - 我们如何自动对Web文档进行分类?
-

频繁模式挖掘的重要性

- 许多重要数据挖掘任务的基础
 - 关联、相关性、因果性
 - 序列模式、空间模式、时间模式、多维
 - 关联分类、聚类分析
 - 更加广泛的用处
 - 购物篮分析、交叉销售、直销
 - 点击流分析、DNA序列分析等等
-

关联规则基本模型

- 关联规则基本模型
 - Apriori算法
 - Fp-Tree算法
-

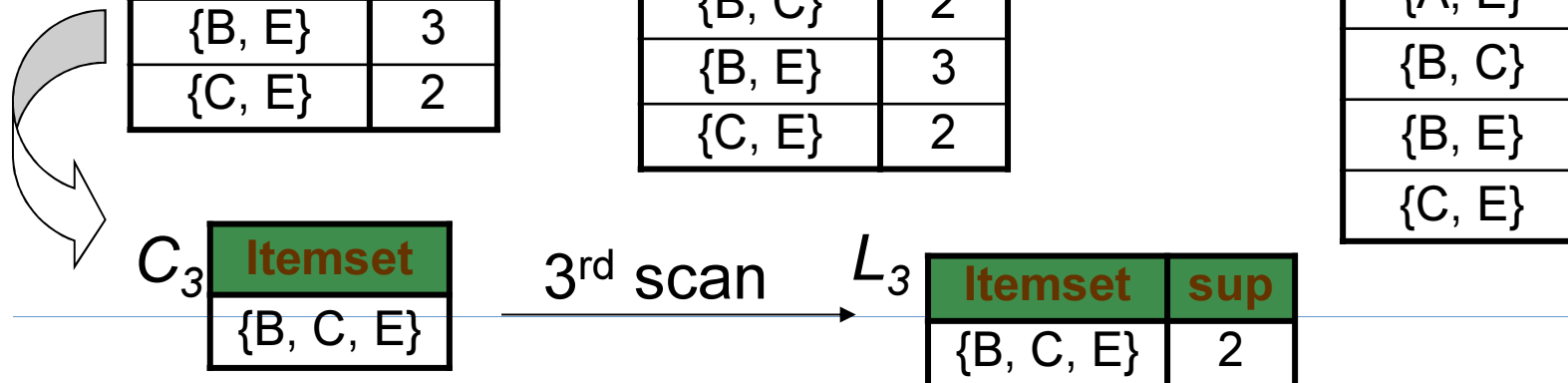
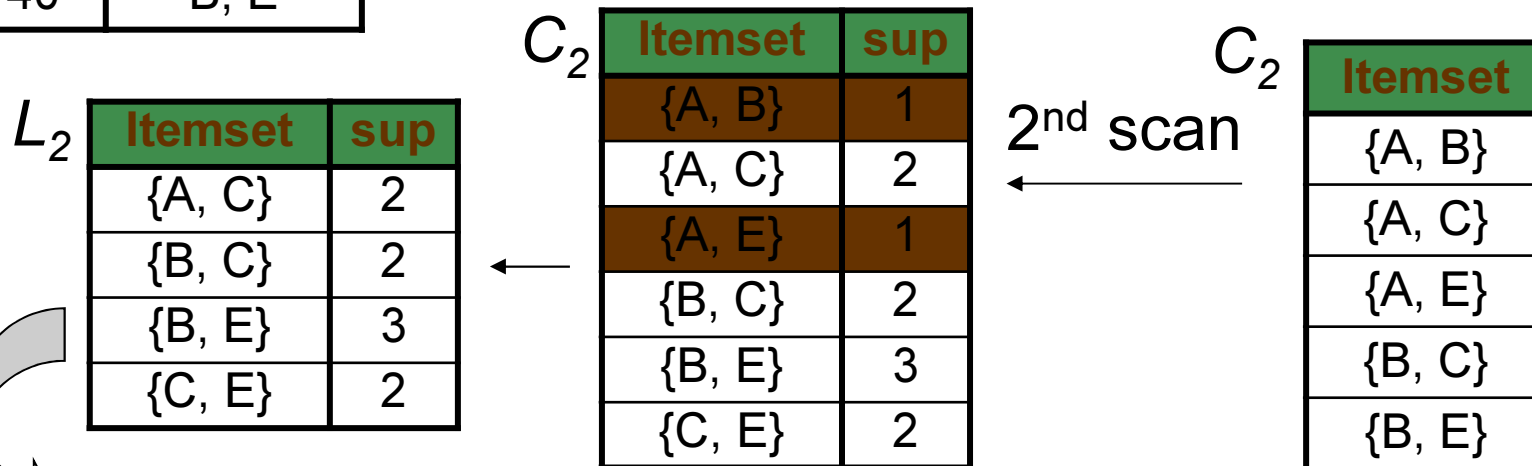
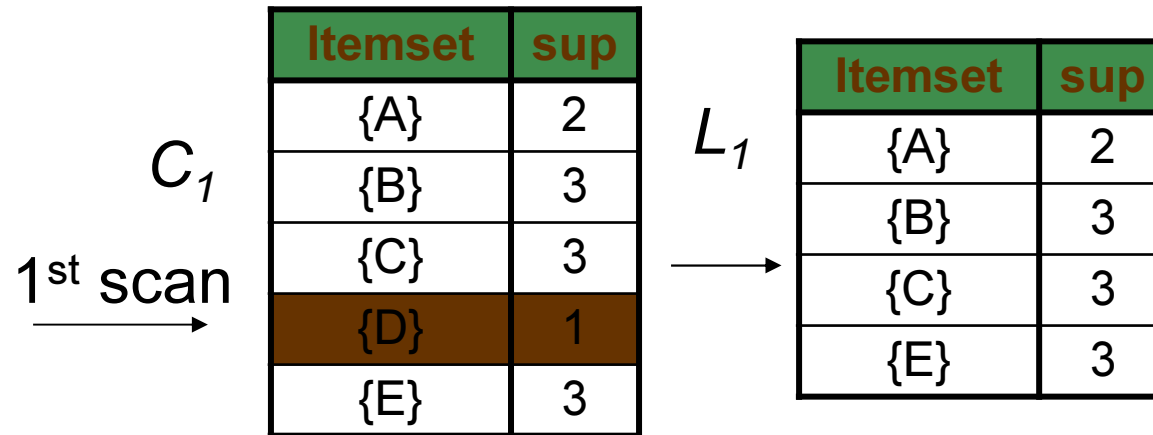
Apriori算法

- (1) $L_1 = \{\text{频繁1项集}\}$;
 - (2) for($k=2; L_{k-1} \neq \emptyset; k++$) do begin
 - (3) $C_k = \text{apriori_gen}(L_{k-1})$; //新的潜在频繁项集
 - (4) for all *transactions* $t \in D$ do begin
 - (5) $C_t = \text{subset}(C_k, t)$; //t中包含的潜在频繁项集
 - (6) for all *candidates* $c \in C_t$ do
 - (7) $c.\text{count}++$;
 - (8) end;
 - (9) $L_k = \{c \in C_k \mid c.\text{count} \geq \text{minsup}\}$
 - (10) end;
 - (11) $\text{Answer} = \bigcup_k L_k$
-

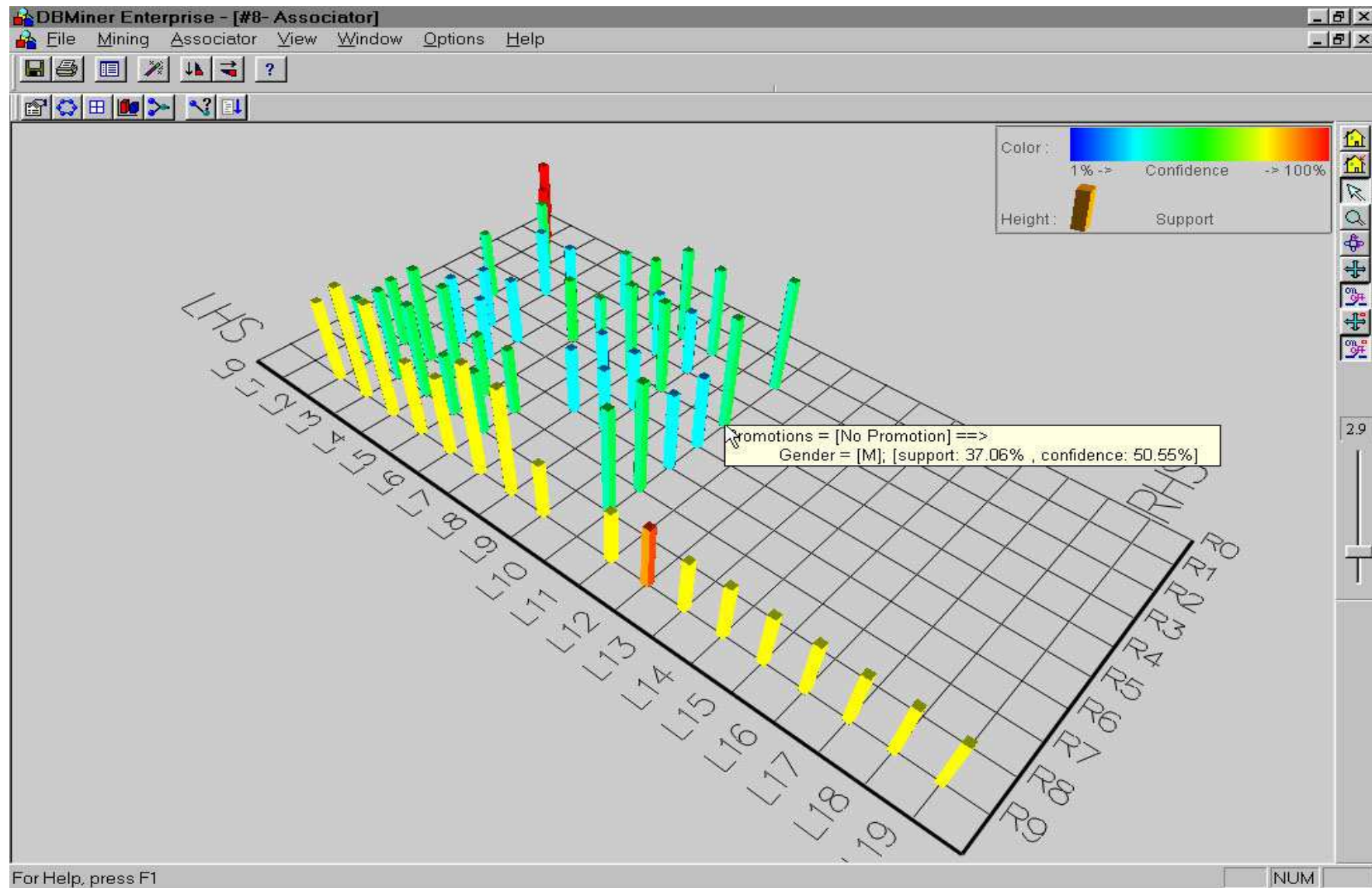
实例

Database TDB

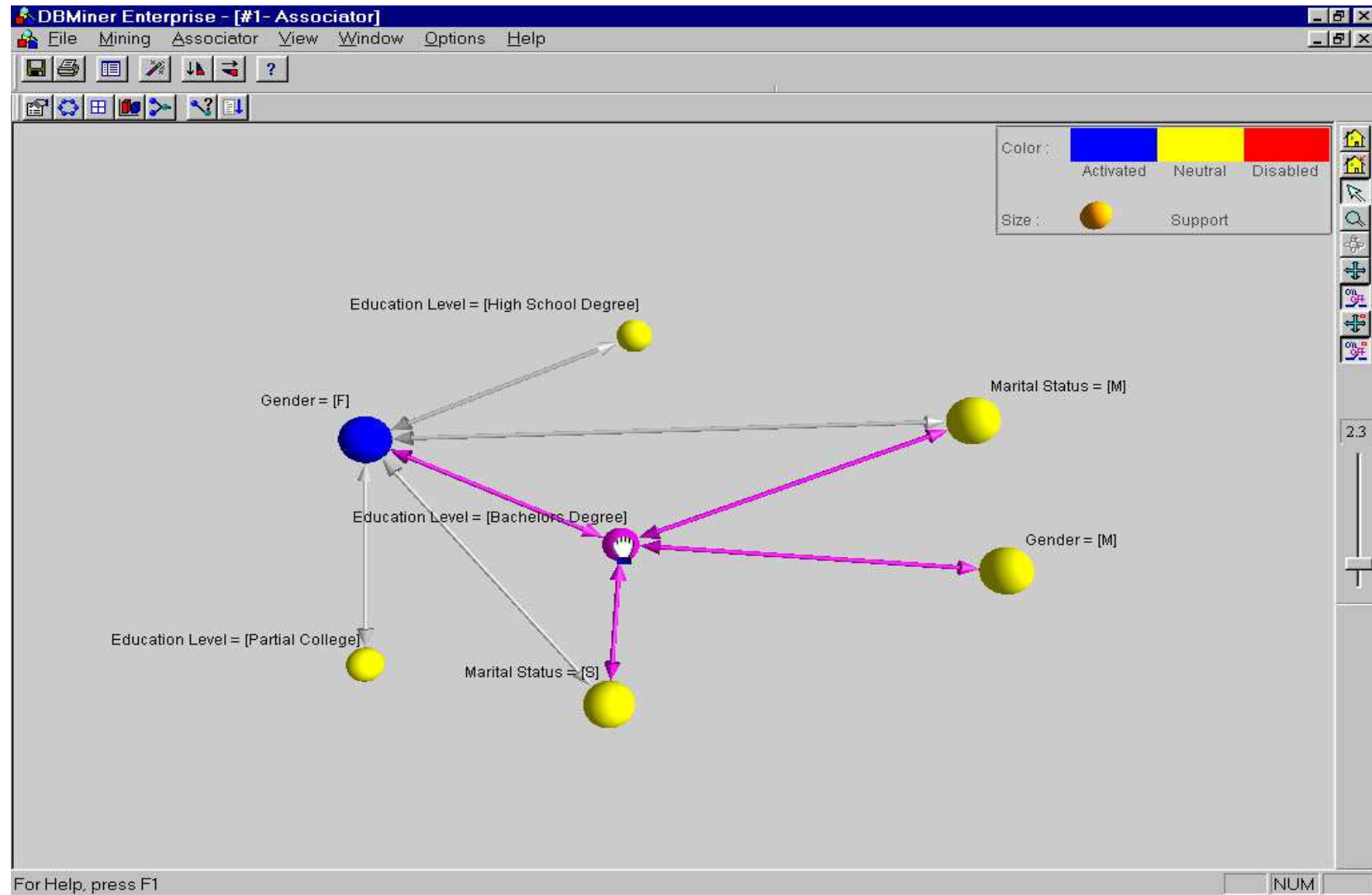
Tid	Items
10	A, C, D
20	B, C, E
30	A, B, C, E
40	B, E



Visualization of Association Rules: Pane Graph



Visualization of Association Rules: Rule Graph



挖掘频繁集 不用生成候选集

- 用 Frequent-Pattern tree (FP-tree) 结构压缩数据库,
 - 高度浓缩, 同时对频繁集的挖掘又完备的
 - 避免代价较高的数据库扫描
 - 开发一种高效的基于FP-tree的频繁集挖掘算法
 - 采用分而治之的方法学: 分解数据挖掘任务为小任务
 - 避免生成关联规则: 只使用部分数据库!
-

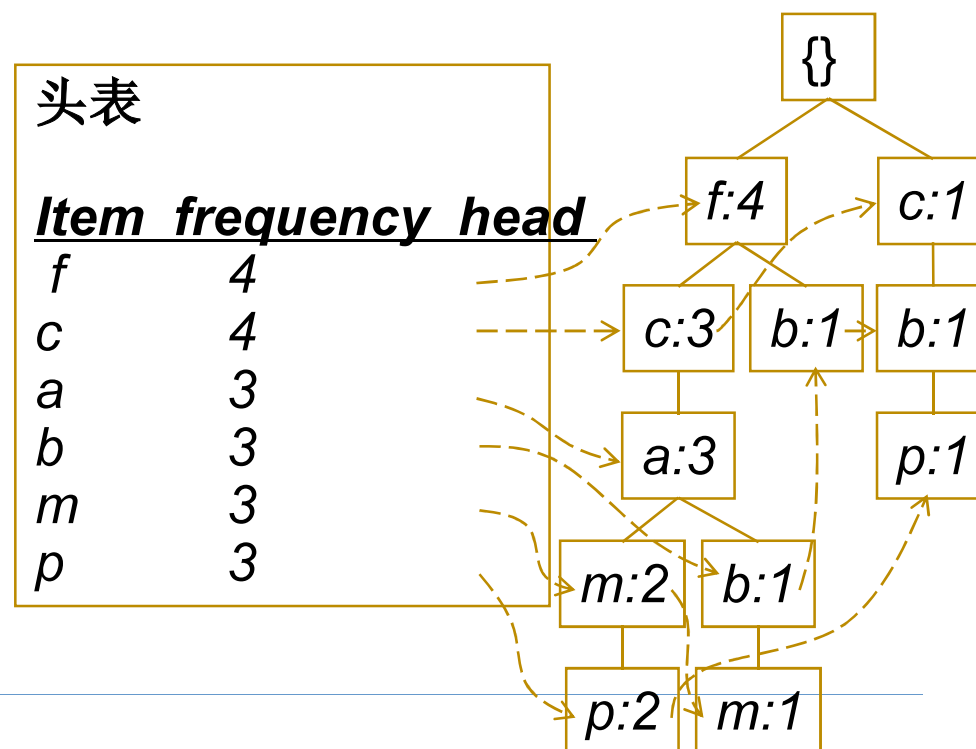
建立 FP-tree树

<i>TID</i>	<i>Items bought</i>	<i>(ordered) frequent items</i>
100	{f, a, c, d, g, i, m, p}	{f, c, a, m, p}
200	{a, b, c, f, l, m, o}	{f, c, a, b, m}
300	{b, f, h, j, o}	{f, b}
400	{b, c, k, s, p}	{c, b, p}
500	{a, f, c, e, l, p, m, n}	{f, c, a, m, p}

最小支持度 = 0.5

步骤:

1. 扫描数据库一次, 得到频繁 1-项集
2. 把项按支持度递减排序
3. 再一次扫描数据库, 建立FP-tree



高阶方法 (2) :

深度学习

Word2Vector

词向量

➤ 自然语言中的词语在机器学习中表示符号

□ One-hot Representation

例如：

- “话筒” 表示为 [0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 ...]
- “麦克” 表示为 [0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 ...]
- 实现时就可以用0,1,2,3,...来表示词语进行计算，这样“话筒”就为3，“麦克”为8.

□ 存在两个问题

- 维度比较大,尤其是用于 Deep Learning 的一些算法时
 - 词汇鸿沟：任意两个词之间都是孤立的，不能体现词和词之间的关系
-

词向量

➤ Distributional Representation

❑ 词表示为：

- $[0.792, -0.177, -0.107, 0.109, 0.542, \dots]$ ，常见维度50或者100

❑ 解决“词汇鸿沟”问题

- 可以通过计算向量之间的距离（欧式距离、余弦距离等）来体现词与词的相似性

➤ 如何训练这样的词向量

❑ 没有直接的模型可训练得到

❑ 可通过训练语言模型的同时，得到词向量

语言模型

➤ 判断一句话是不是正常人说出来的，用数学符号描述为

- 给定一个字符串" w_1, w_2, \dots, w_t ", 计算它是自然语言的概率 $p(w_1, w_2, \dots, w_t)$ ，一个很简单的推论是

$$p(w_1, w_2, \dots, w_t) = p(w_1) \cdot p(w_2 | w_1) \cdot p(w_3 | w_1, w_2) \cdot \dots \cdot p(w_t | w_1, w_2, \dots, w_{t-1})$$

- 例如，有个句子"大家,喜欢,吃,苹果"

- $P(\text{大家, 喜欢, 吃, 苹果}) = p(\text{大家})p(\text{喜欢}|\text{大家})p(\text{吃}|\text{大家, 喜欢})p(\text{苹果}|\text{大家, 喜欢, 吃})$

- 简单表示为 $p(s) = p(w_1, w_2, \dots, w_T) = \prod_{i=1}^T p(w_i | \text{Context}_i)$

➤ 计算 $p(w_i | \text{Context}_i)$ 问题

word2vec原理

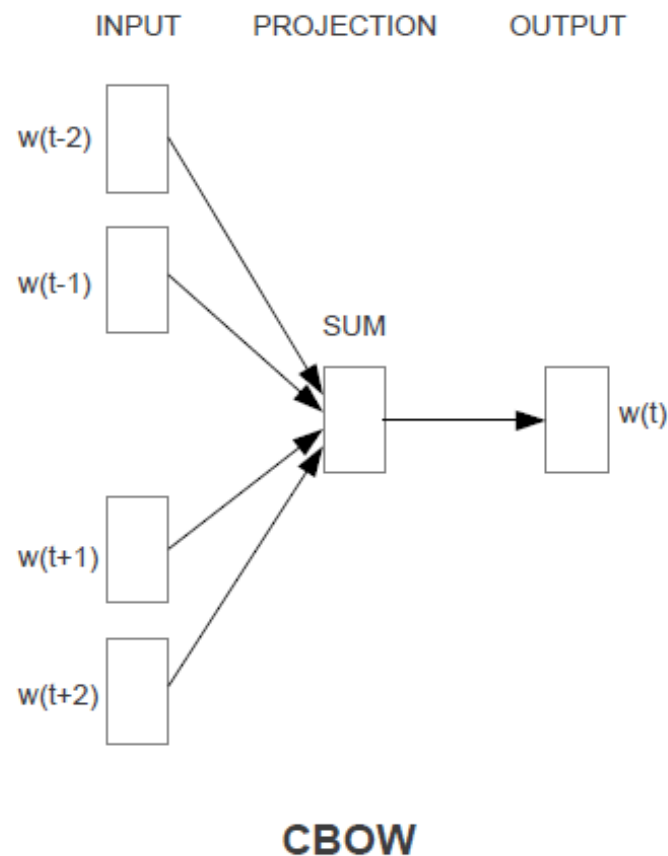
➤ 两种模型，两种方法

模型	CBOW		Skip-Gram	
方法	Hierarchical Softmax	Negative Sampling	Hierarchical Softmax	Negative Sampling

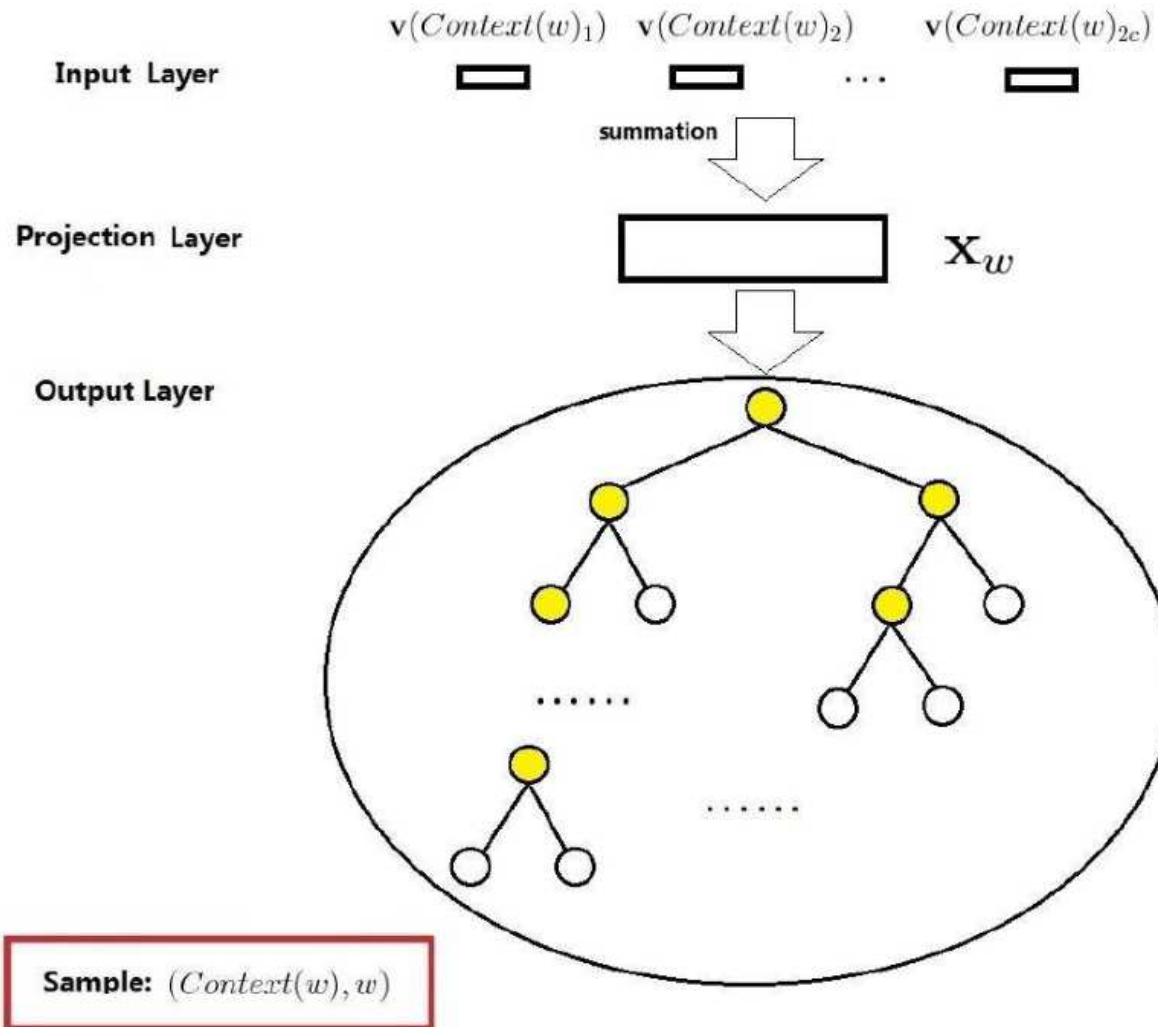
CBOW模型+Hierarchical Softmax方法

➤ CBOW模型

- ❑ INPUT:输入层
- ❑ PROJECTION:投影层
- ❑ OUTPUT:输出层
- ❑ $w(t)$:当前词语 (向量)
- ❑ $w(t-2), w(t-1), w(t+1), w(t+2)$:当前词语的上下文
- ❑ SUM:上下文的累加和



CBOW模型+Hierarchical Softmax方法（续）



CBOW模型+Hierarchical Softmax方法（续）

1. 输入层: 包含 $Context(w)$ 中 $2c$ 个词的词向量 $\mathbf{v}(Context(w)_1), \mathbf{v}(Context(w)_2), \dots, \mathbf{v}(Context(w)_{2c}) \in \mathbb{R}^m$. 这里, m 的含义同上表示词向量的长度.
2. 投影层: 将输入层的 $2c$ 个向量做求和累加, 即 $\mathbf{x}_w = \sum_{i=1}^{2c} \mathbf{v}(Context(w)_i) \in \mathbb{R}^m$.
3. 输出层: 输出层对应一棵二叉树, 它是以语料中出现过的词当叶子结点, 以各词在语料中出现的次数当权值构造出来的 Huffman 树. 在这棵 Huffman 树中, 叶子结点共 $N (= |\mathcal{D}|)$ 个, 分别对应词典 \mathcal{D} 中的词, 非叶子结点 $N - 1$ 个 (图中标成黄色的那些结点).

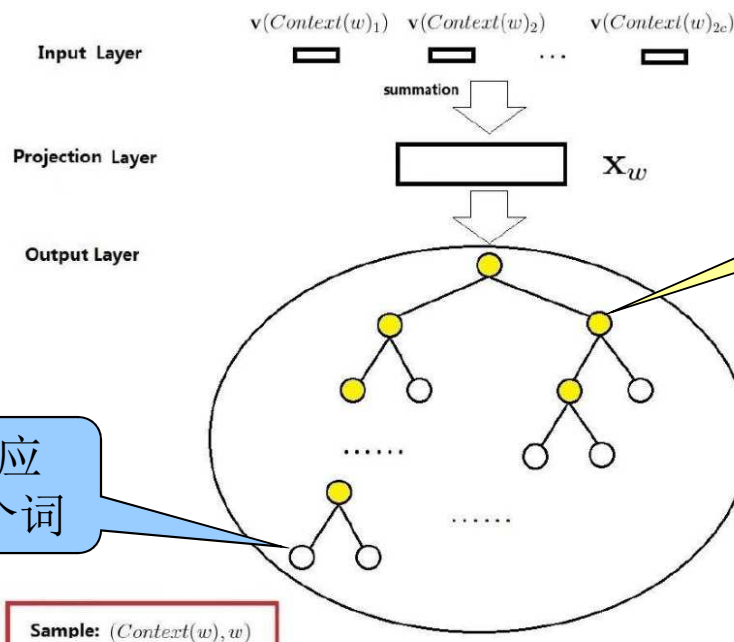
目标: $p(w_i | Context_i)$

为什么建哈夫曼树？



叶子结点对应词典中的一个词

非叶子结点为LR分类器



CBOW模型+Hierarchical Softmax方法（续）

➤ 句子：我,喜欢,观看,巴西,足球,世界杯

➤ w =足球

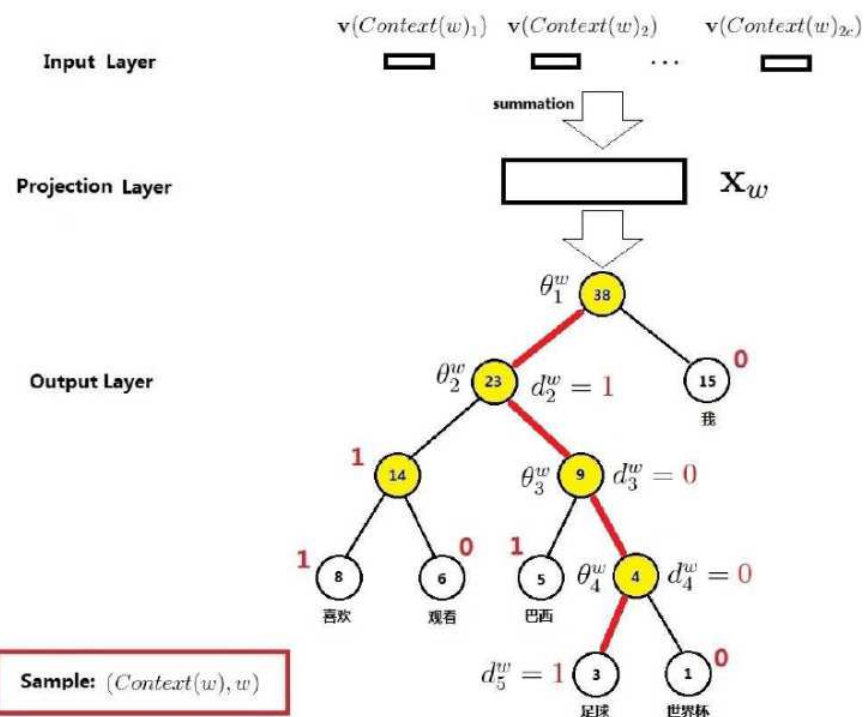
d_j^w : 编码（1或0）

二分类

d_j^w : 正负类（1：负类,0：正类）

θ_j^w : 非叶子节点向量

θ_j^w : 类别向量



CBOW模型+Hierarchical Softmax方法（续）

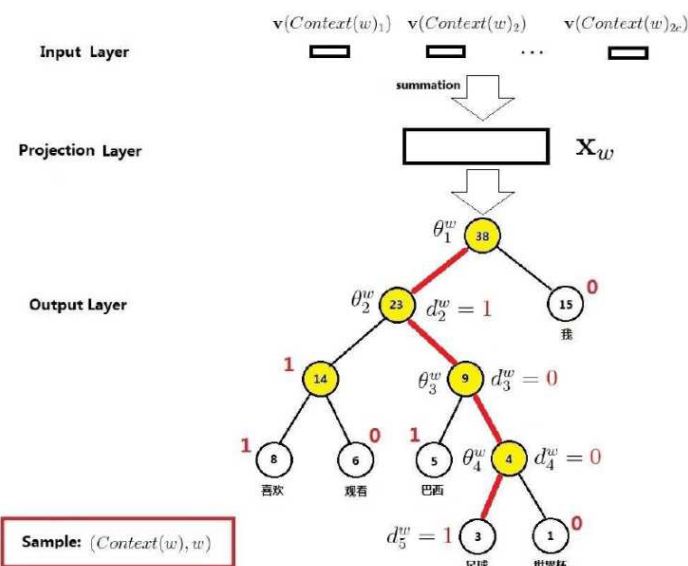
- 正类概率: $\sigma(\mathbf{x}_w^\top \theta) = \frac{1}{1 + e^{-\mathbf{x}_w^\top \theta}}$,
- 负类概率: $1 - \sigma(\mathbf{x}_w^\top \theta)$,
- "足球" 叶子节点经过4次二分类，每次分类结果对应的概率为

第 1 次: $p(d_2^w | \mathbf{x}_w, \theta_1^w) = 1 - \sigma(\mathbf{x}_w^\top \theta_1^w)$;

第 2 次: $p(d_3^w | \mathbf{x}_w, \theta_2^w) = \sigma(\mathbf{x}_w^\top \theta_2^w)$;

第 3 次: $p(d_4^w | \mathbf{x}_w, \theta_3^w) = \sigma(\mathbf{x}_w^\top \theta_3^w)$;

第 4 次: $p(d_5^w | \mathbf{x}_w, \theta_4^w) = 1 - \sigma(\mathbf{x}_w^\top \theta_4^w)$,



- 由Context("足球")预测"足球"出现的概率

$$p(\text{足球} | \text{Context}(\text{足球})) = \prod_{j=2}^5 p(d_j^w | \mathbf{x}_w, \theta_{j-1}^w).$$

CBOW模型+Hierarchical Softmax方法（续）

- 对于词典中的每个词 w 有 l^w 结点个数

$$p(w | Context(w)) = \prod_{j=2}^{l^w} p(d_j^w | X_w, \theta_{j-1}^w)$$

- 其中 , $p(d_j^w | X_w, \theta_{j-1}^w) = \begin{cases} \sigma(X_w^T \cdot \theta_{j-1}^w), d_j^w = 0; \\ 1 - \sigma(X_w^T \cdot \theta_{j-1}^w), d_j^w = 1. \end{cases}$

- 或者表示为 $p(d_j^w | X_w, \theta_{j-1}^w) = [\sigma(X_w^T \cdot \theta_{j-1}^w)]^{1-d_j^w} \cdot [1 - \sigma(X_w^T \cdot \theta_{j-1}^w)]^{d_j^w}$

- 对于由 s 个句子组成的语料库 C 有

$$L(X, \theta) = \prod_{s \in C} \prod_{w \in s} p(w | Context(w)) = \prod_{s \in C} \prod_{w \in s} \prod_{j=2}^{l^w} p(d_j^w | X_w, \theta_{j-1}^w)$$

- 取对数似然函数

参数1

$$\log L(X, \theta) = \sum_{s \in C} \sum_{w \in s} \sum_{j=2}^{l^w} \log p(d_j^w | X_w, \theta_{j-1}^w)$$

参数2

$$= \sum_{s \in C} \sum_{w \in s} \sum_{j=2}^{l^w} [(1 - d_j^w) \cdot \log \sigma(X_w^T \cdot \theta_{j-1}^w) + d_j^w \cdot \log(1 - \sigma(X_w^T \cdot \theta_{j-1}^w))]$$

CBOW模型+Hierarchical Softmax方法（续）

➤ 梯度下降法进行求解

- 令 $f(w, j) = -(1 - d_j^w) \cdot \log \sigma(X_w^T \cdot \theta_{j-1}^w) - d_j^w \cdot \log(1 - \sigma(X_w^T \cdot \theta_{j-1}^w))$
- $f(w, j)$ 关于 θ_{j-1}^w 和 X_w 的梯度分别为

$$\frac{\partial f(w, j)}{\partial \theta_{j-1}^w} = -[1 - d_j^w - \sigma(X_w^T \cdot \theta_{j-1}^w)] \cdot X_w$$

$$\frac{\partial f(w, j)}{\partial X_w} = -[1 - d_j^w - \sigma(X_w^T \cdot \theta_{j-1}^w)] \cdot \theta_{j-1}^w$$

- 更新公式

$$\theta_{j-1}^w := \theta_{j-1}^w - \eta \cdot \frac{\partial f(w, j)}{\partial \theta_{j-1}^w}$$

$$V(\tilde{w}) := V(\tilde{w}) - \eta \cdot \sum_{j=2}^{l^w} \frac{\partial f(w, j)}{\partial X_w}, \tilde{w} \in \text{Context}(w)$$

Word2vec效果

- 训练数据集：经过分词后的新闻数据，大小184MB
 - ❑ 查看"中国"，"钓鱼岛"，"旅游"，"苹果"几个词语的相似词语如下所示

请输入词语<exit退出>:中国

大陆	0.66763467
中共	0.57856727
共产党	0.56305367
解放军	0.55761635
台湾	0.5368497
反攻	0.5271177
日本	0.5103535
王文莹	0.49295437
内地	0.48557448
对岸	0.48428434

请输入词语<exit退出>:钓鱼岛

钓鱼台	0.6219264
钓岛	0.6123347
南海	0.6018163
领土	0.51753837
领海	0.4928774
岛屿	0.4853142
舰队	0.47854927
渔权	0.47229362
主权	0.46729872
东海	0.4613399

请输入词语<exit退出>:旅游

观光	0.65619475
景点	0.60212
陆客	0.59477097
旅行	0.5677106
游憩	0.557839
赏樱	0.5571045
游玩	0.52199984
观光客	0.51974636
行程	0.51943743
参观	0.5077874

请输入词语<exit退出>:苹果

三星	0.7224437
微软	0.7101249
Apple	0.66682446
iPhone5	0.62071097
Google	0.597368
iPadmini	0.5609188
新机	0.559093
库克	0.5589176
宏达电	0.555409
产品	0.55437565

Word2vec效果

➤ 向量加减法

- ❑ "中国+北京-日本", "中国+北京-法国", "家庭+孩子-学校"

请输入词语<exit退出>:中国,北京,日本

中国 + 北京 - 日本 =

东京	0.6188478
外海	0.4843038
广东	0.45942163
安倍	0.44709134
参拜	0.44158116
派出	0.43384194
釜山	0.4310615
海域	0.42759195
湖北	0.4272585
官邸	0.42487407

请输入词语<exit退出>:中国,北京,法国

中国 + 北京 - 法国 =

巴黎	0.7990697
伦敦	0.7614885
义大利	0.7026581
纽约	0.6534368
德国	0.6426668
首都	0.6354907
英国	0.6339937
西班牙	0.6326488
柏林	0.628184
加拿大	0.6255576

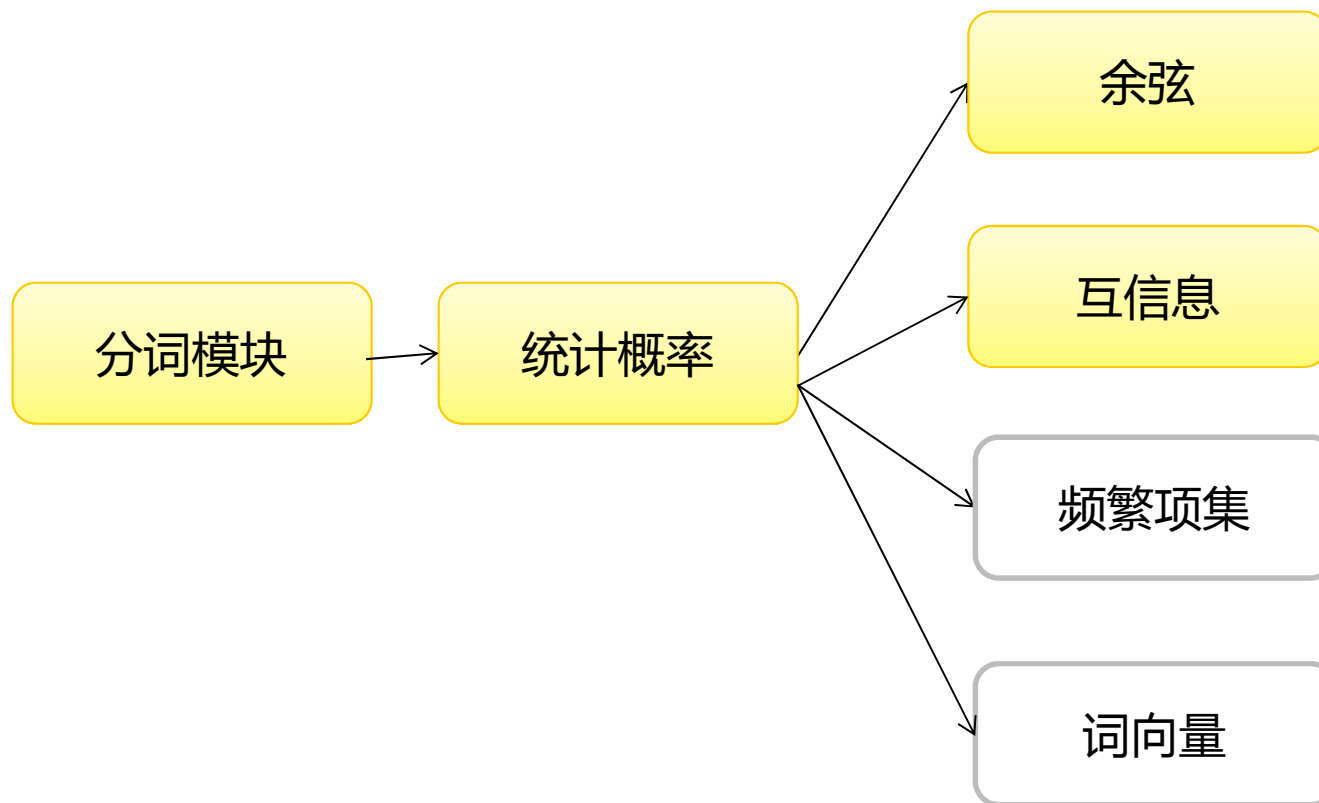
请输入词语<exit退出>:家庭,孩子,学校

家庭 + 孩子 - 学校 =

老师	0.948024
国中	0.76681674
学生	0.73207045
高中	0.7140167
小孩	0.70768154
小朋友	0.692751
上课	0.6889297
校方	0.687891
家长	0.66356415
高职	0.65407264

具体工程实现

总体设计



实验说明

➤ 语料库

- ❑ 大约1G电子书

➤ 运行要求

- ❑ `./simword article.txt result.txt`
- ❑ `Article.txt`
 - 原始文本文件
- ❑ `Result.txt`
 - 第1行输出: **Cosine**
 - 第2-21行输出: 前20个相关性最高的词对
 - 第22行输出: **MutualInfo**
 - 第23-42行输出: 前20个最相关的词对

The End

