



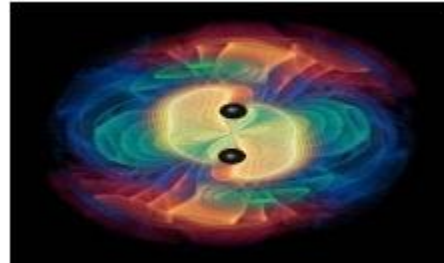
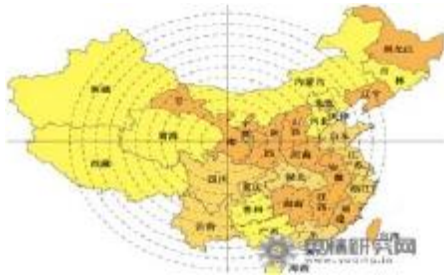
在线社会媒体计算的 机遇与挑战

王 娜

深圳大学信息工程学院

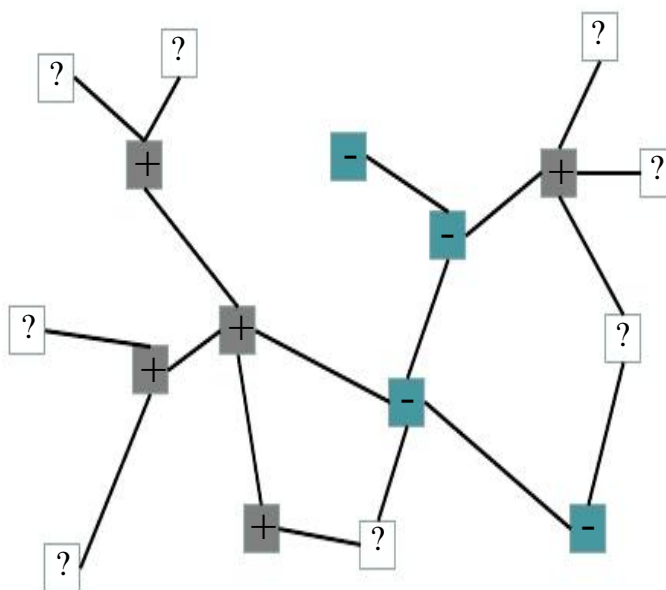
Social Networks

SN **bridges** our daily life and the **virtual** web space!

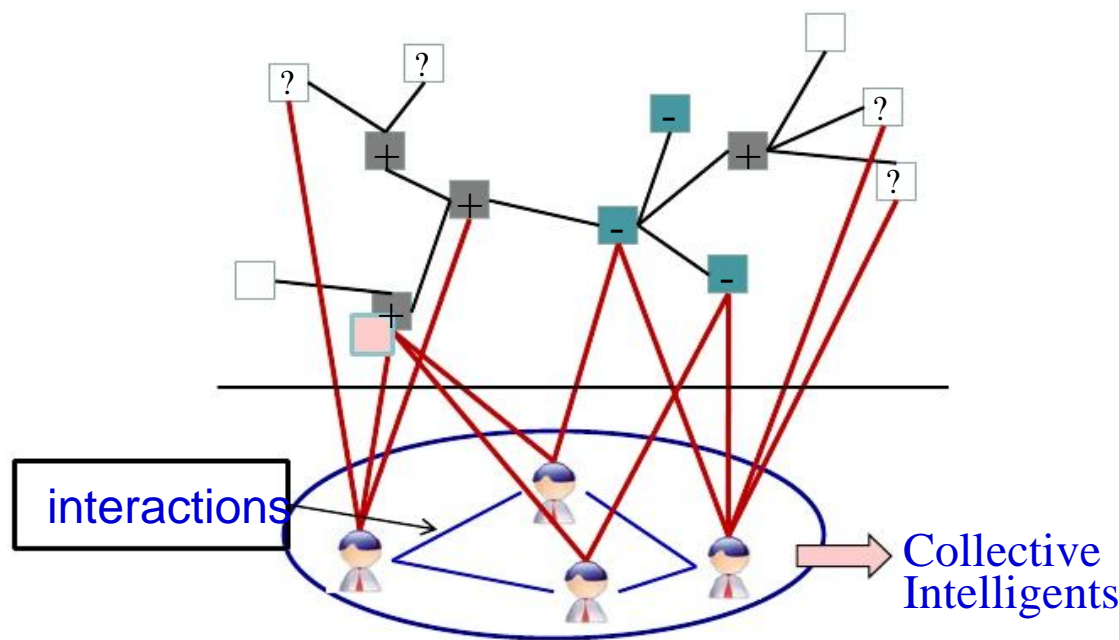


在线社会网络特征

- 多维网络 (Multi-Dimensional Networks)
- 多模网络 (Multi-Mode Networks)



Web 1.0



Social Web

在线社会网络大数据

■ 大量社交媒体应用，产生了巨量数据

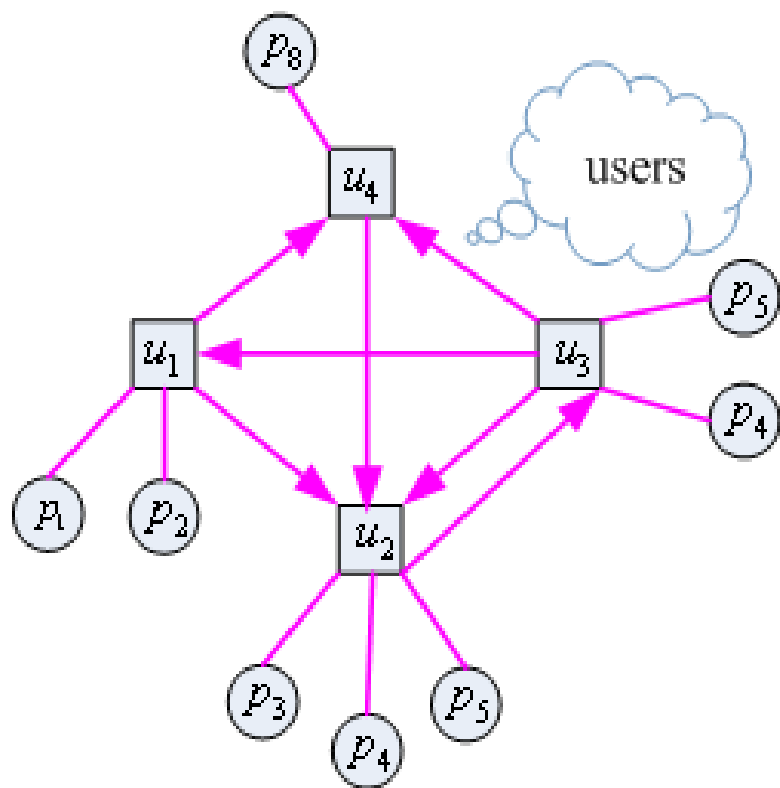
- 250 million tweets per day
- 3,000 photos in Flickr per minute
- 153 million blogs posted per year

■ 数据特点

- 高维（high dimension）
- 大规模（large-scale）
- 链接的（linked data）
- 噪声（noisy）



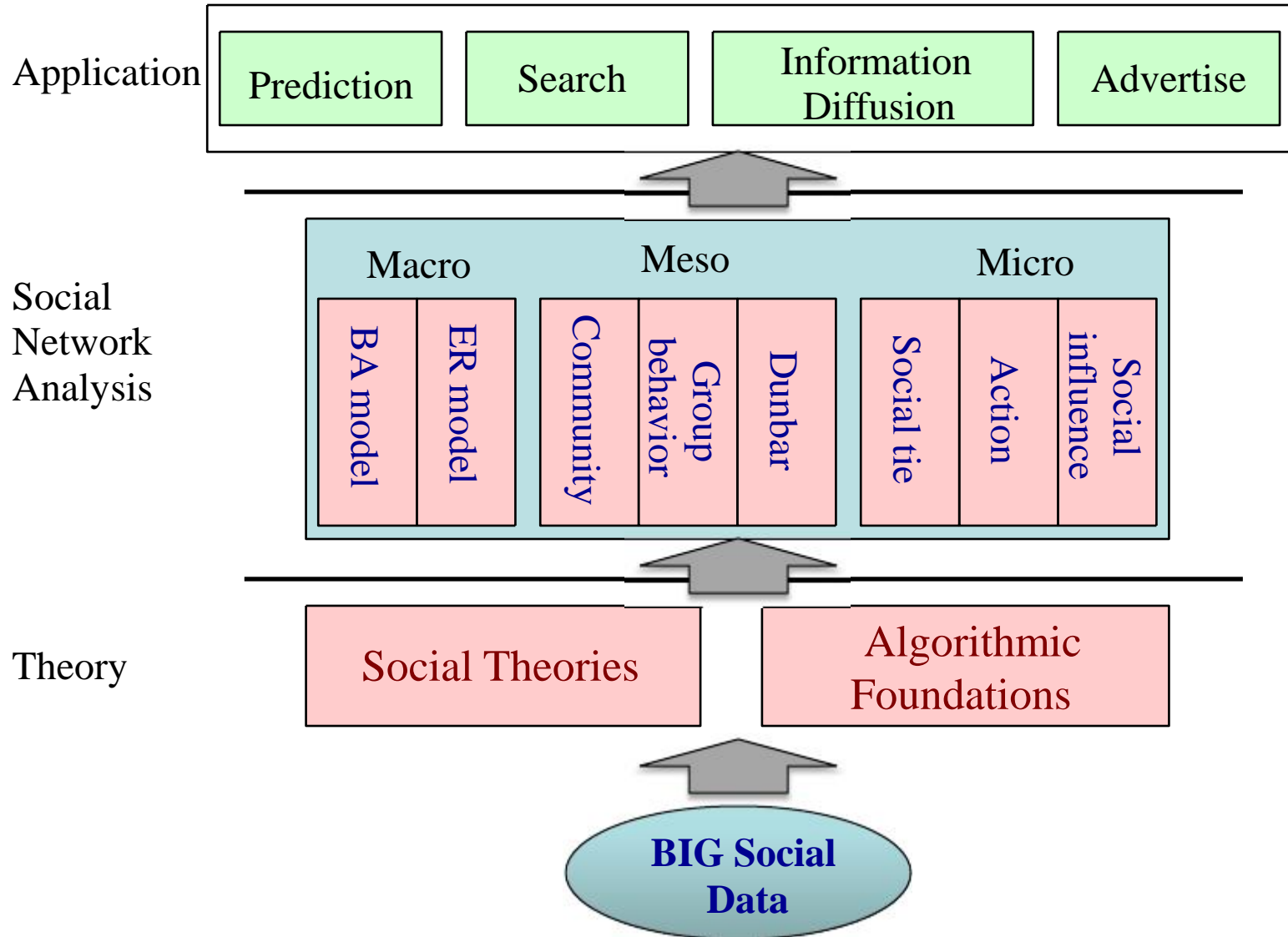
在线社交媒体—微博数据



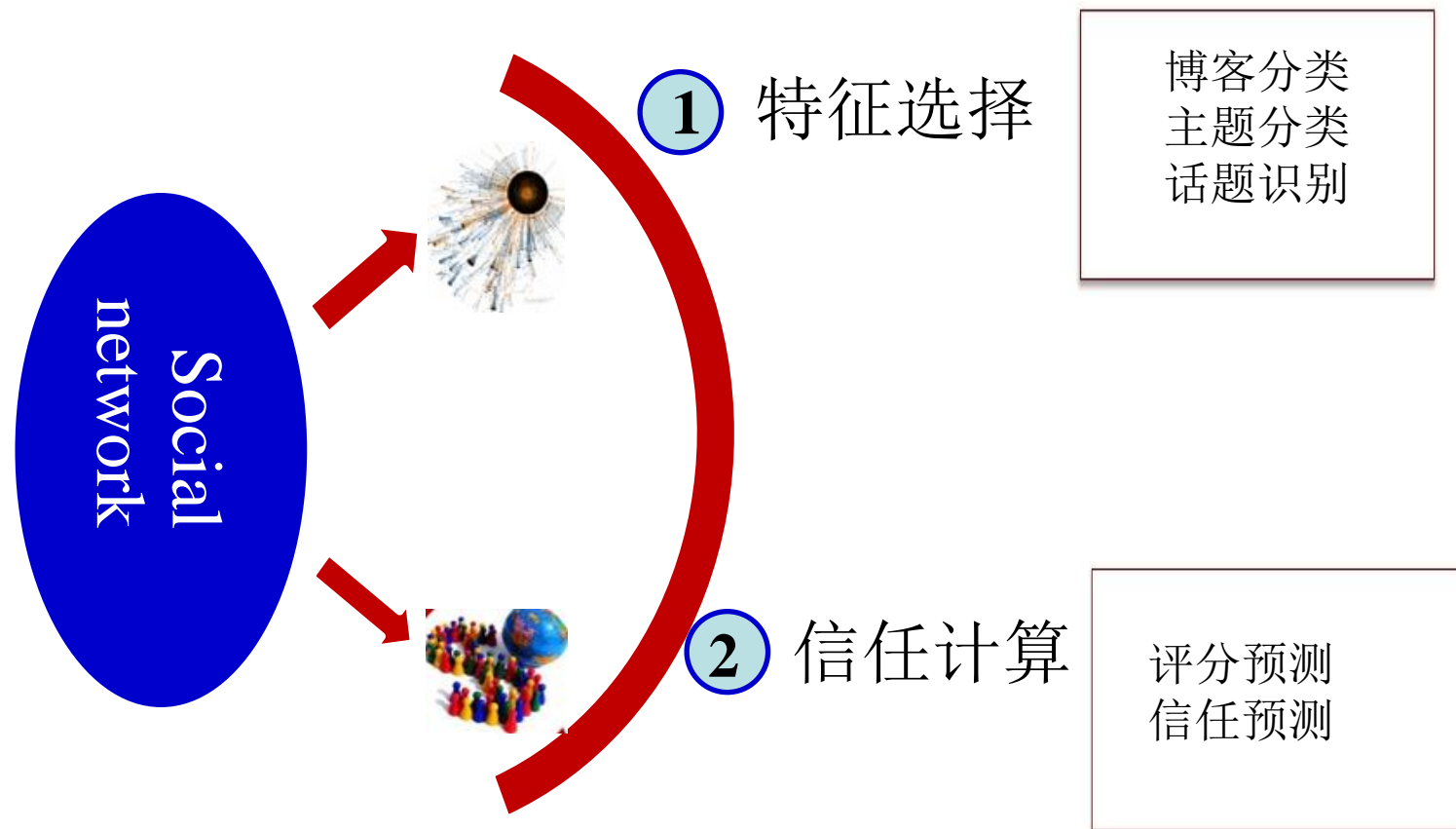
- 实体
 - 用户
 - 内容

- 关系
 - 用户—用户
 - 用户—内容
 - 内容—内容

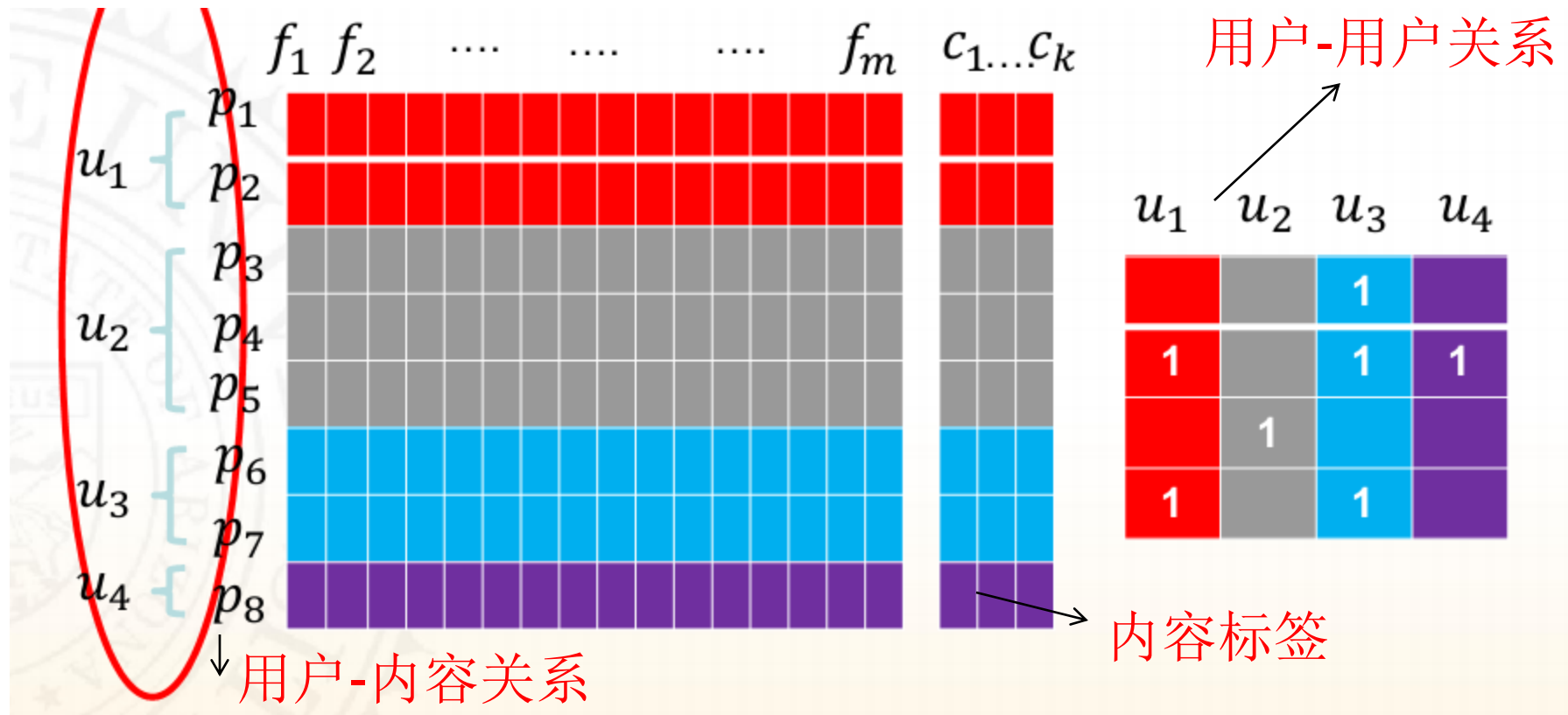
Core Research in Social Network



研究工作



特征选择—链接数据



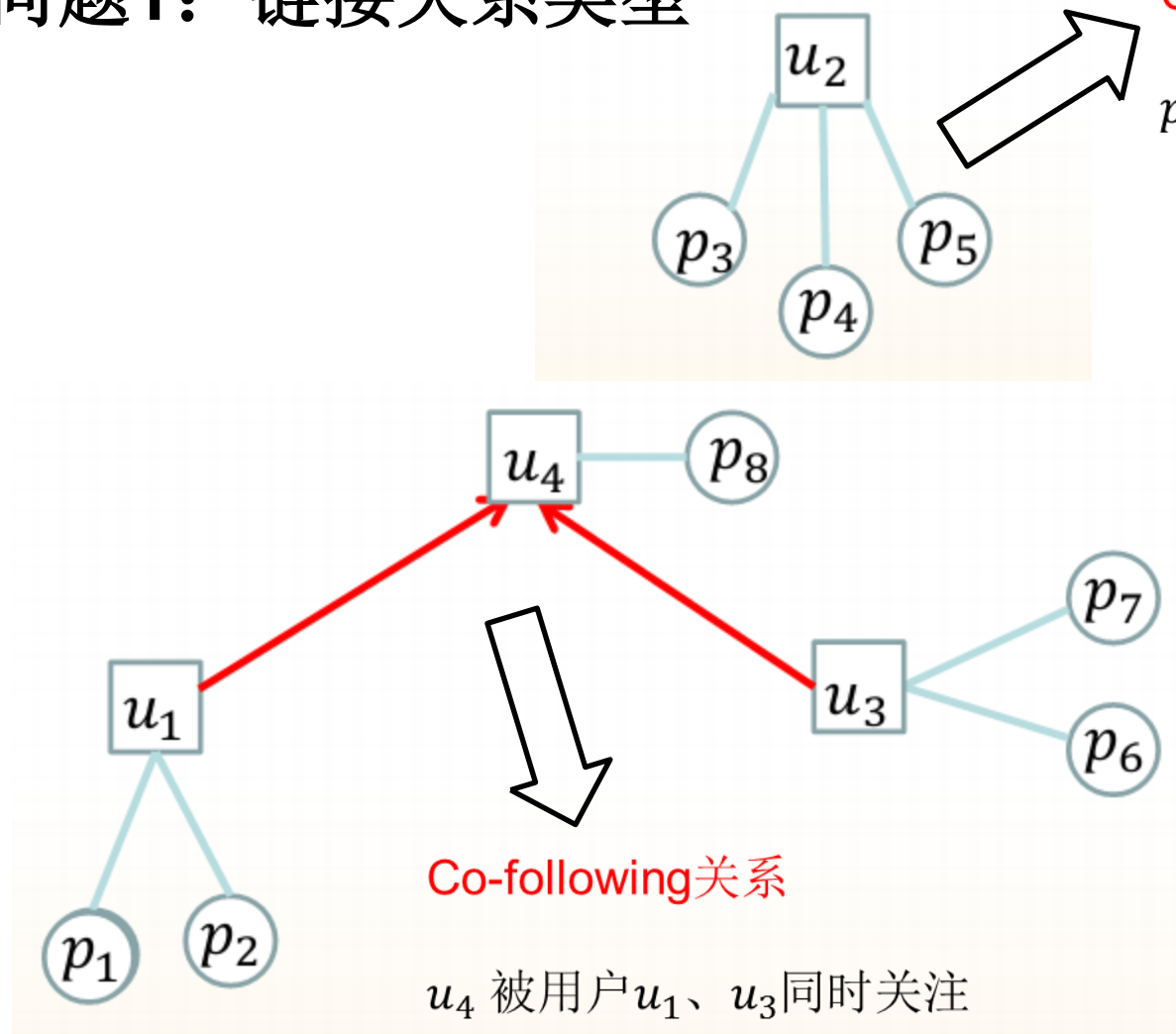
■ 社会学原理(Social Theories)

- **同质规则:**认为两个有链接关系的对象倾向于具有相同或相似的类别
- **合引规则:**与同一个对象有链接关系的两个对象也倾向于具有相似的类别标签

问题1：链接关系类型

Co-post 关系

p_3 、 p_4 、 p_5 由同一用户 u_2 发布



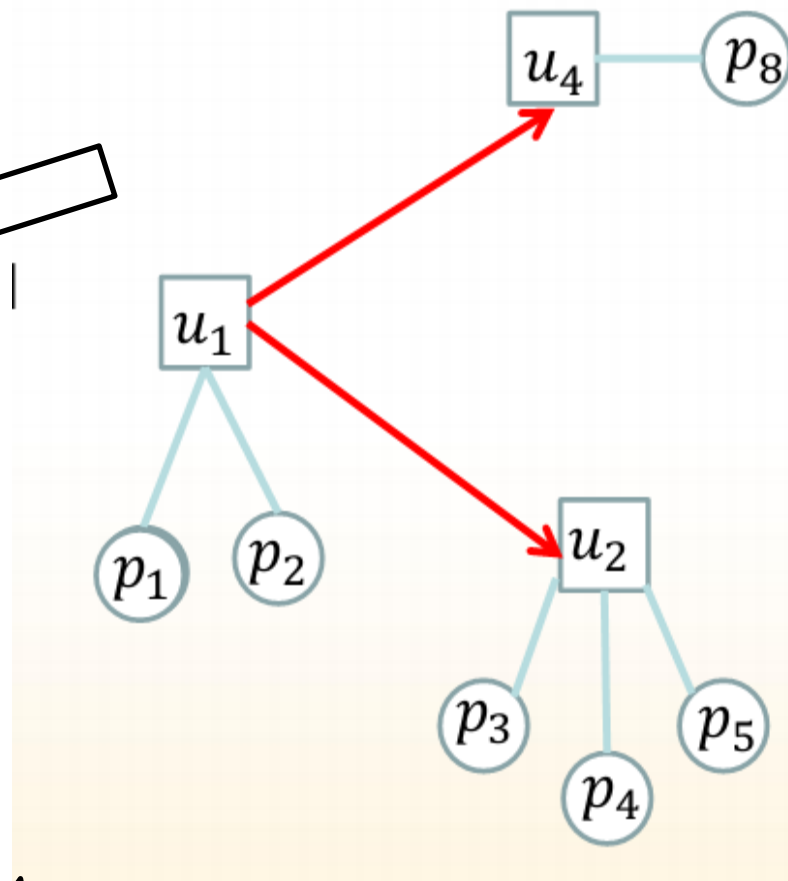
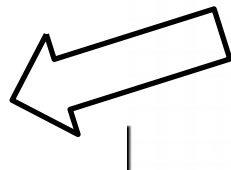
Co-following 关系

u_4 被用户 u_1 、 u_3 同时关注

问题1：链接关系类型

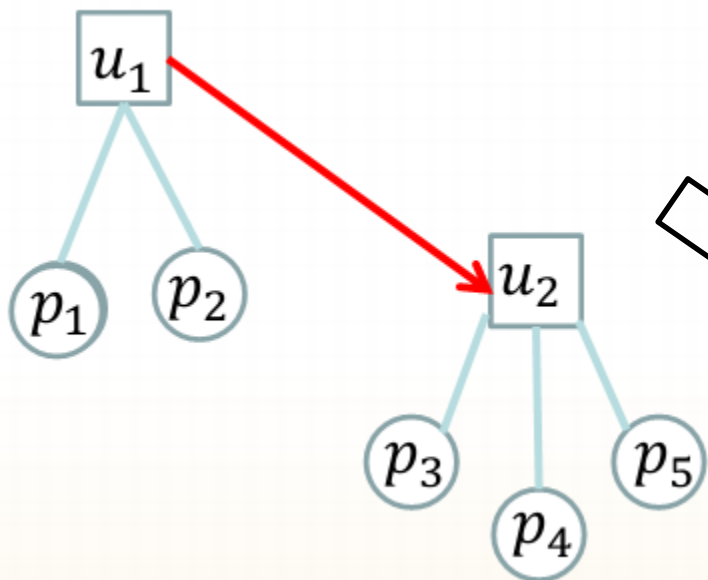
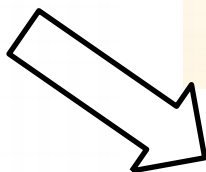
Co-followed 关系

u_1 同时关注 用户 u_2 、 u_4



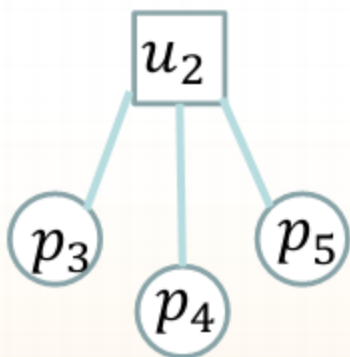
Following 关系

u_1 关注用户 u_2



问题2：关系利用—Co-post

为用户内容标签预测的函数

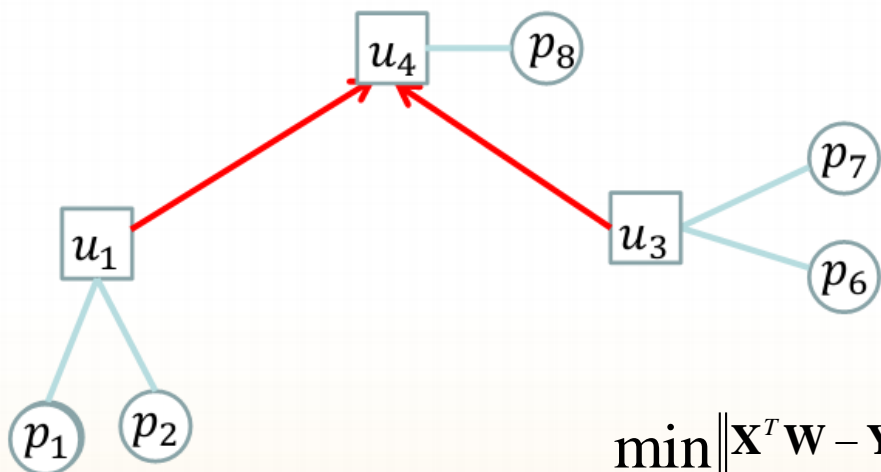


$$\min_{\mathbf{W}} \|\mathbf{X}^T \mathbf{W} - \mathbf{Y}\|_F^2 + \alpha \|\mathbf{W}\|_{2,1} + \beta \sum_{u \in \mathbf{U}} \sum_{\mathbf{f}_i, \mathbf{f}_j \in \mathbf{F}_u} \|T(\mathbf{f}_i) - T(\mathbf{f}_j)\|_2^2$$

$$= \min_{\mathbf{W}} \|\mathbf{X}^T \mathbf{W} - \mathbf{Y}\|_F^2 + \alpha \|\mathbf{W}\|_{2,1} + \beta \sum_i \sum_j \mathbf{A}(i, j) \|\mathbf{W}^T \mathbf{f}_i - \mathbf{W}^T \mathbf{f}_j\|_2^2$$

Co-post关系矩阵，如果内容 p_i 、 p_j 由同一个用户发布，那么 $\mathbf{A}(i, j) = 1$ ，否则为0.

问题2：关系利用—Co-following



$$\hat{T}(u_k) = \frac{\sum_{f_i \in F_k} T(f_i)}{|F_k|} = \frac{\sum_{f_i \in F_k} \mathbf{W}^T \mathbf{f}_i}{|F_k|}$$

定义用户主题兴趣：

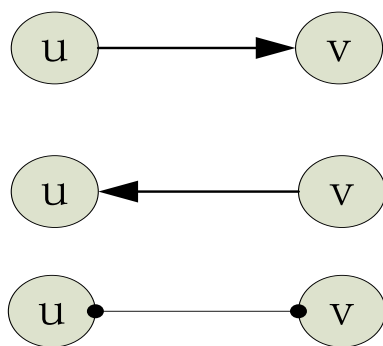
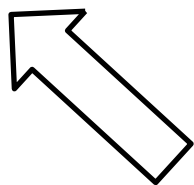
$$\begin{aligned} \min_{\mathbf{W}} & \left\| \mathbf{X}^T \mathbf{W} - \mathbf{Y} \right\|_F^2 + \alpha \left\| \mathbf{W} \right\|_{2,1} + \beta \sum_{u_k \in \mathbf{U}} \sum_{u_i, u_j \in F_k} \left\| \hat{T}(u_i) - \hat{T}(u_j) \right\|_2^2 \\ &= \min_{\mathbf{W}} \left\| \mathbf{X}^T \mathbf{W} - \mathbf{Y} \right\|_F^2 + \alpha \left\| \mathbf{W} \right\|_{2,1} + \beta \sum_{ij} \mathbf{FI}(i, j) \left\| \mathbf{W}^T \mathbf{F}\mathbf{H}(:, i) - \mathbf{W}^T \mathbf{F}\mathbf{H}(:, j) \right\|_2^2 \end{aligned}$$

Co-following关系矩阵，如果用户 u_i 与 u_j

同时关注至少一个用户. $\mathbf{FI}(i, j) = 1$

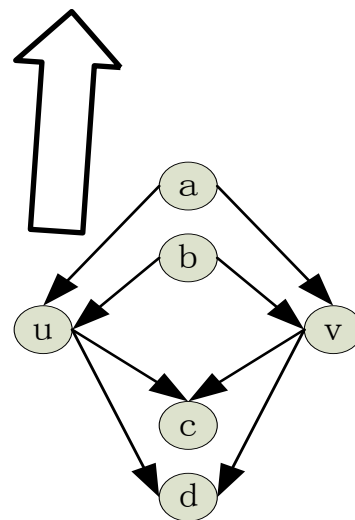
问题3：关系强度—局部相似性

用户间直接关注关系



(a) 传统的聚类假设

用户间局部链接关系



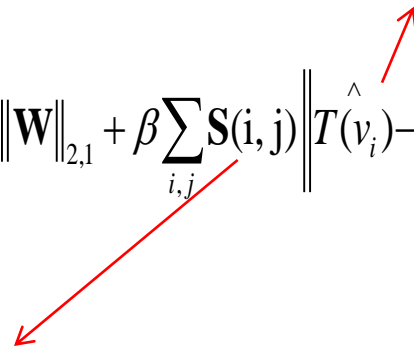
(b) 链接结构假设

目标函数



$$\min_{\mathbf{W}} \|\mathbf{X}^T \mathbf{W} - \mathbf{Y}\|_F^2 + \alpha \|\mathbf{W}\|_{2,1} + \beta \sum_{i,j} \mathbf{S}(i,j) \left\| \overset{\text{用户主题兴趣}}{\hat{T}(v_i)} - \overset{\text{用户主题兴趣}}{\hat{T}(v_j)} \right\|_2^2$$

用户关系强度矩阵



用户局部关系相似指标：(CN指标等)

$$score(u, v) = |\tau(u) \cap \tau(v)|$$

$$score(u, v) = |\tau(u) \cap \tau(v)| / |\tau(u) \cup \tau(v)|$$

$$score(x, y) = \frac{|\tau(x) \cap \tau(y)|}{\sqrt{k_x k_y}}$$

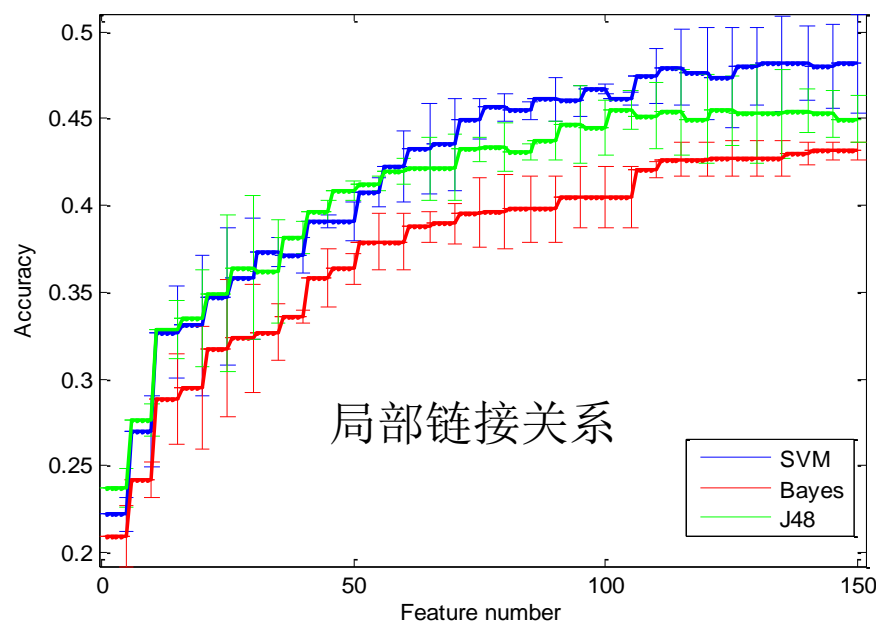
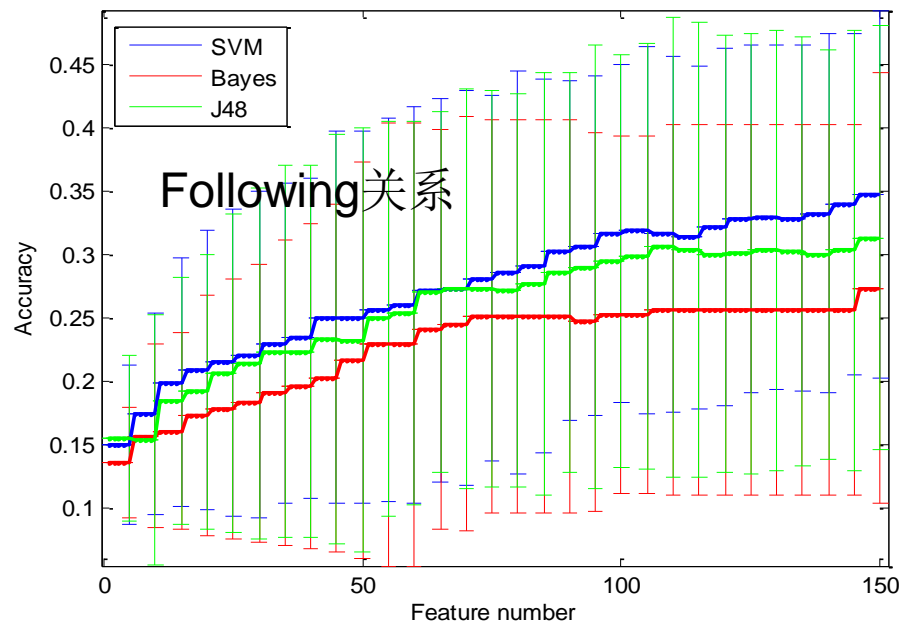
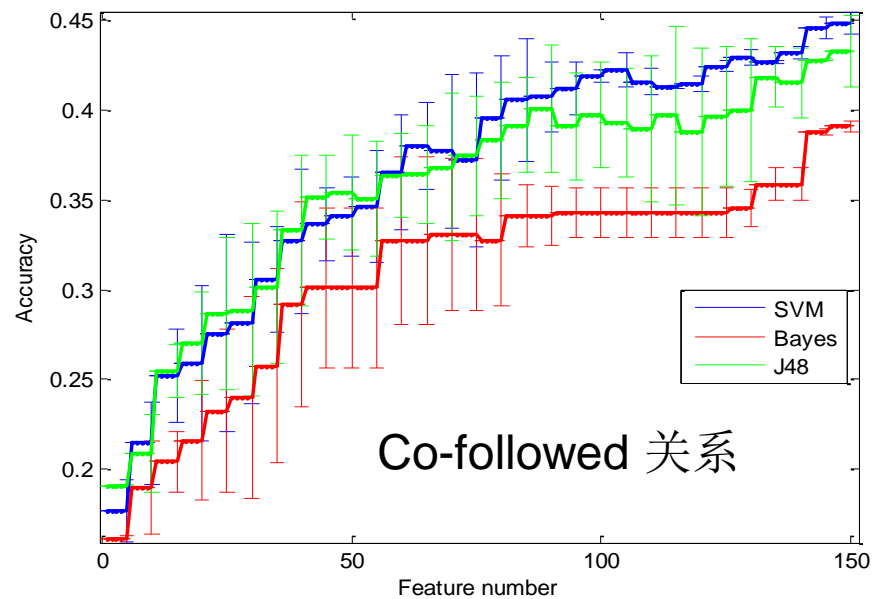
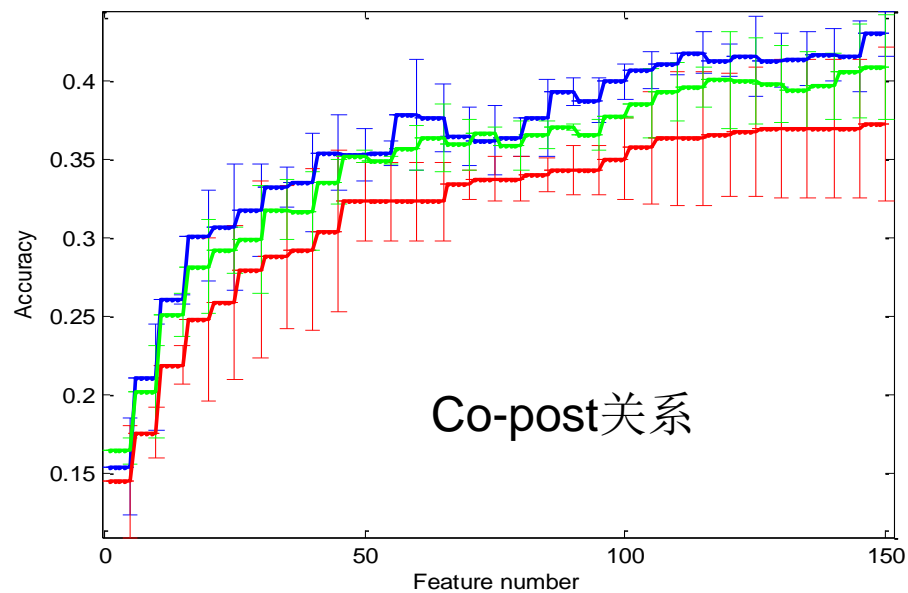
数据平台

BlogCatalog数据集是国外真实的一个博客平台，用户可以在预定义好的类别标签下记录并发表他们的博文。这些博文内容都经过了停用词的移除等相关预处理，去除了那些极不相关的特征词得到博文内容的属性特征。

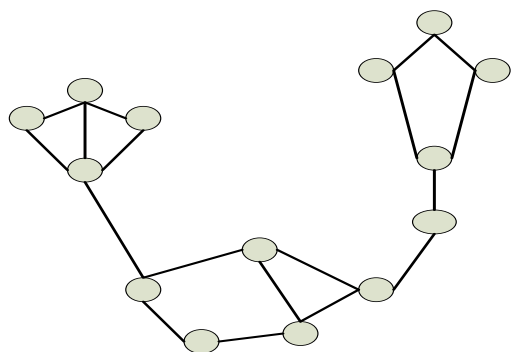
Table: Statistics of the BlogCatalog Dataset↵

# of <u>Users</u>	2242↵
# of <u>Posts</u>	3000↵
# Features after TFIDF	6000↵
# Following Relations	55356↵
# of Classed	15↵
Max#Followers	80↵
Min#Followers	1↵

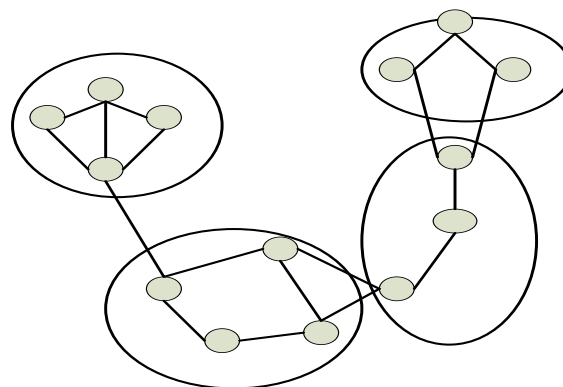
结论—BlogCatalog



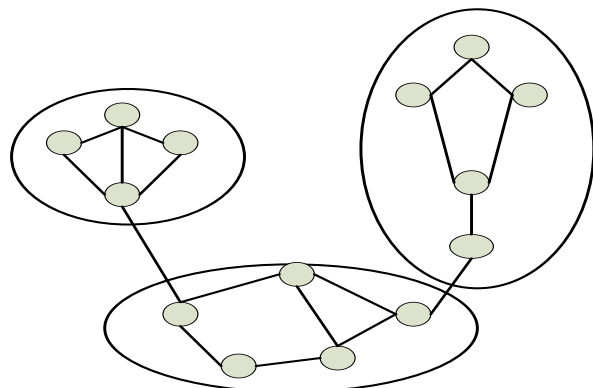
问题4：异质节点属性和链接关系利用



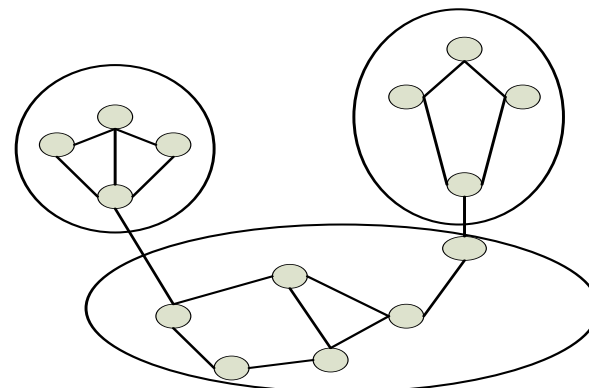
(a) 原始链接图



(b) 属性划分图



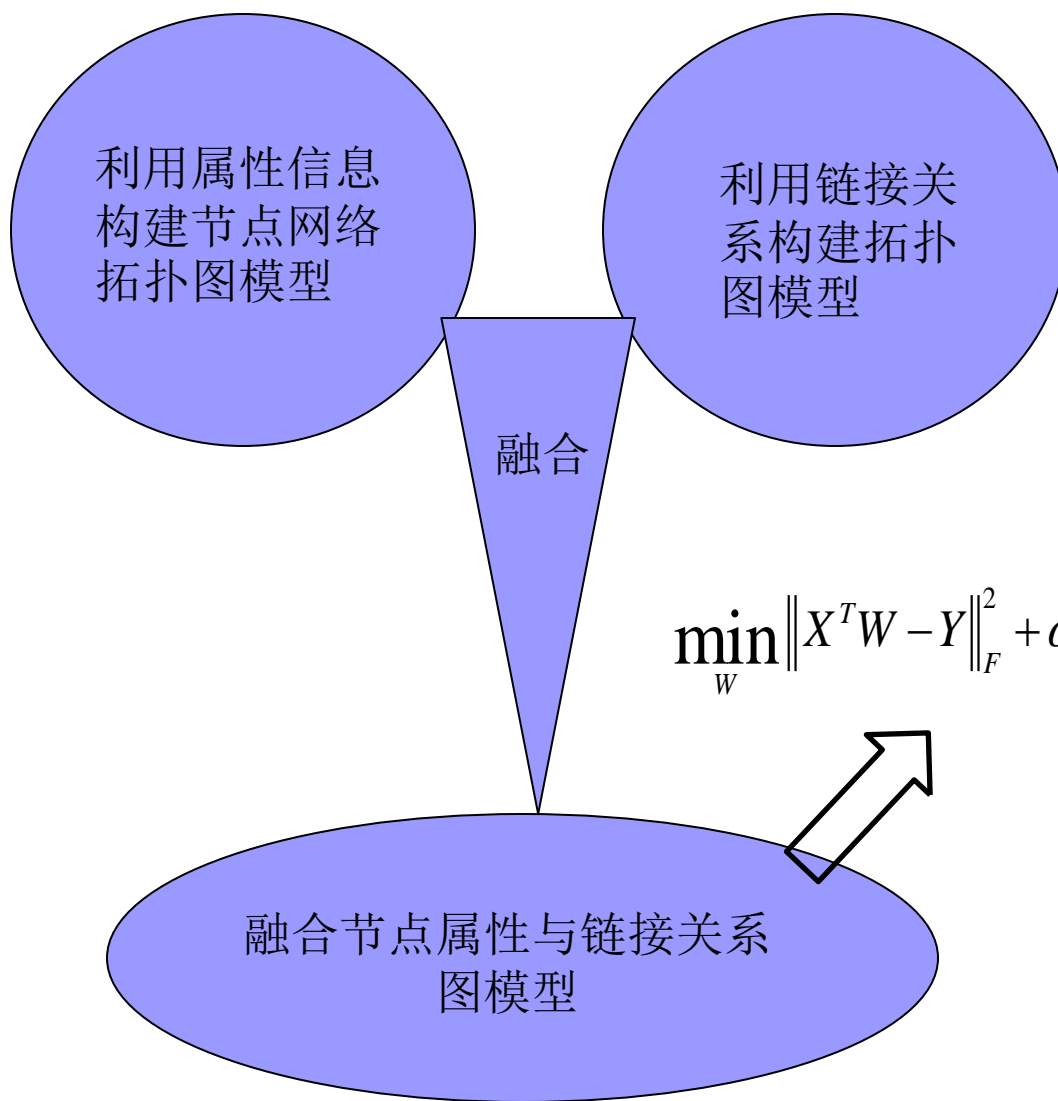
(c) 链接划分图



(d) 期待划分图

社区划分示意图

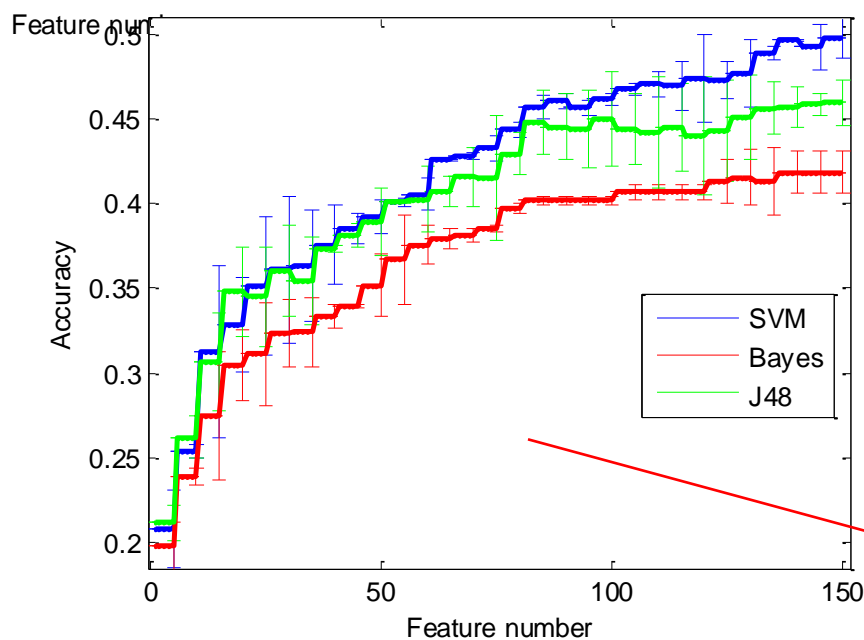
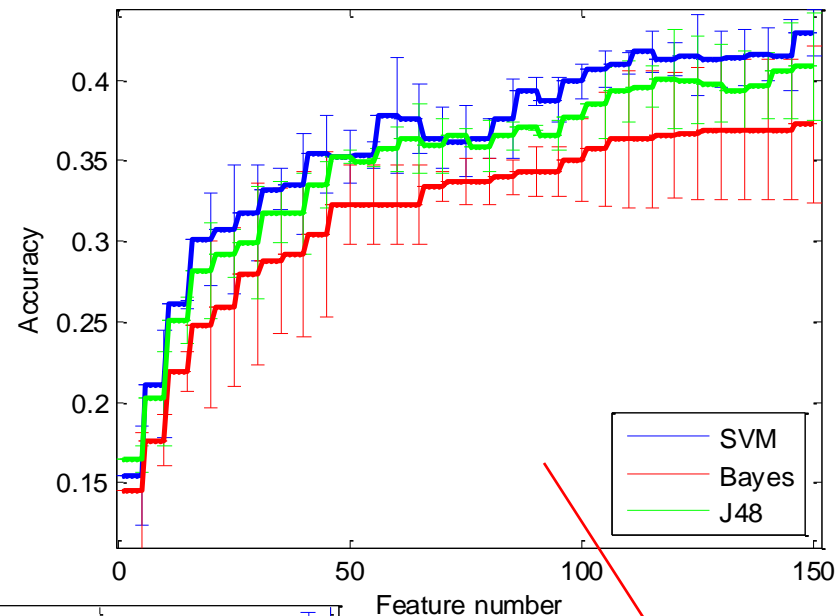
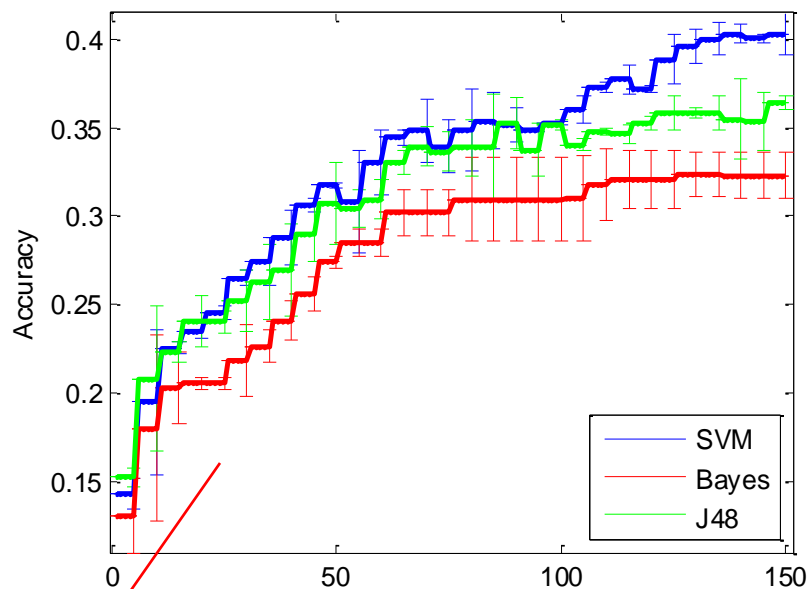
问题4：异质节点属性和链接关系利用



用户内容-内容相似关系强度

$$\min_W \|X^T W - Y\|_F^2 + \alpha \|W\|_{2,1} + \beta \sum_{i,j=1}^N \mathbf{M}(i,j) \|W^T \mathbf{x}_i - W^T \mathbf{x}_j\|_2^2$$

结论—BlogCatalog

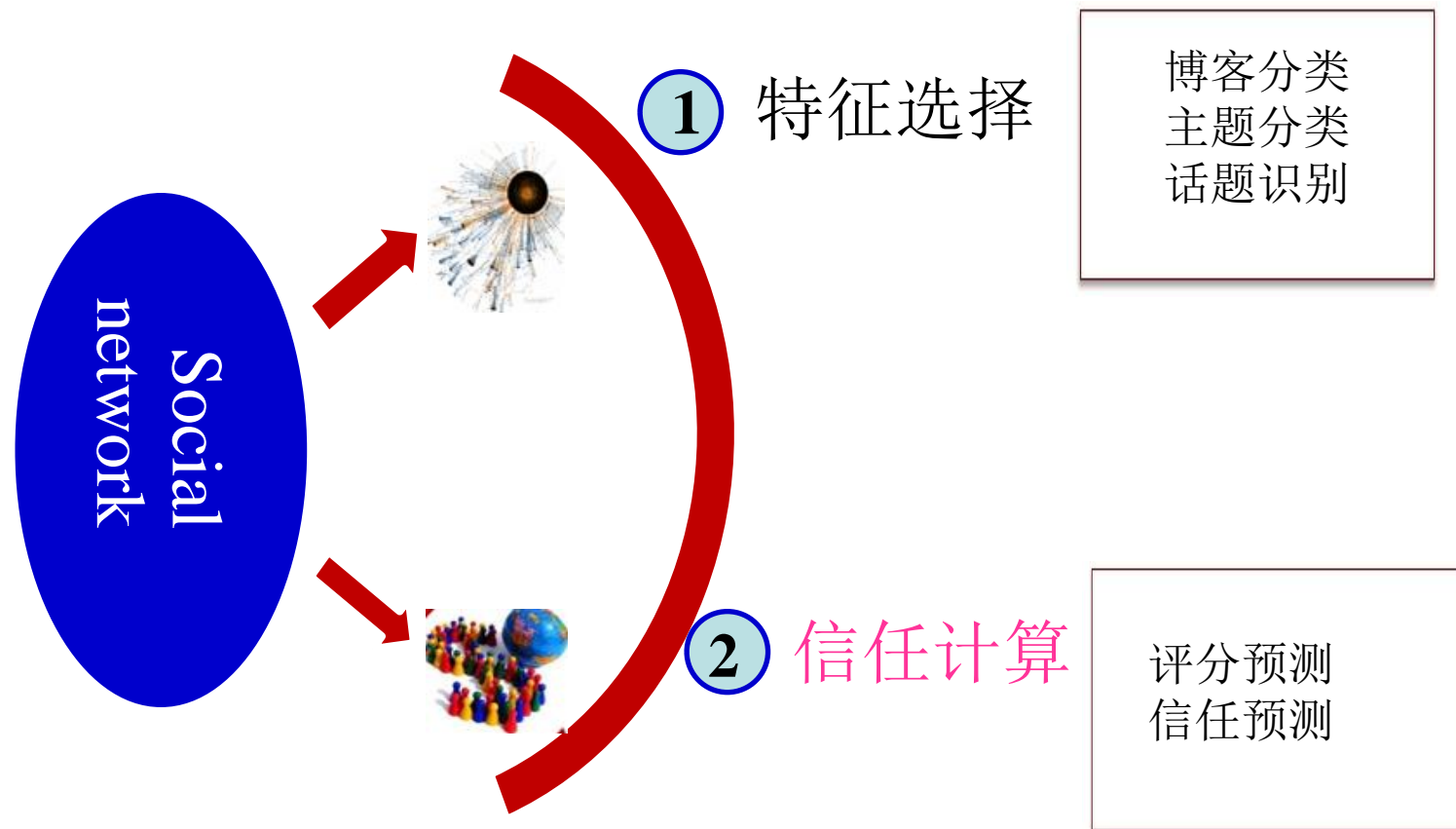


只利用属性信息

只利用链接关系

同时利用节点属性和链接关系

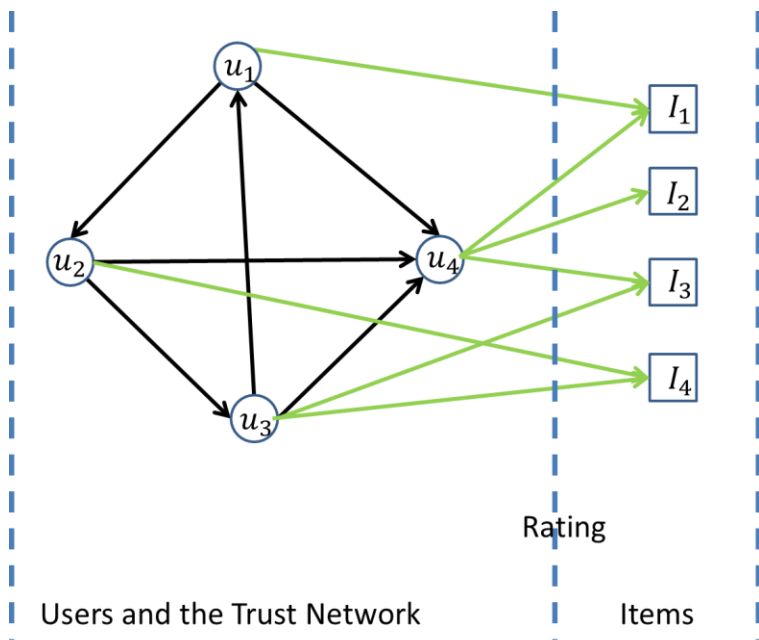
研究工作



信任计算

- 匿名性:允许用户匿名注册,一个用户可以拥有几个ID号。
 - 开放性:任何人在任何时间、任何地点都可以自由地加入或离开。
 - 用户生成内容:数据中包含了大量反映网民兴趣爱好、观点和价值观的信息。
-
- 由于互联网的开放性和匿名性,如何给在线社会网络中的用户提供合理的信任计算机制成为一个急需解决的问题。
 - 作为一个新的商业渠道,面向实际应用,如何利用大量的用户生成内容对在线社会网络进行挖掘分析从而取得更好的经济效益和社会效益也成为备受关注的一个焦点问题。
 - 这两方面直接关系到在线社会网络的安全性和实用性,因此对其研究既具有理论价值又具有实际意义

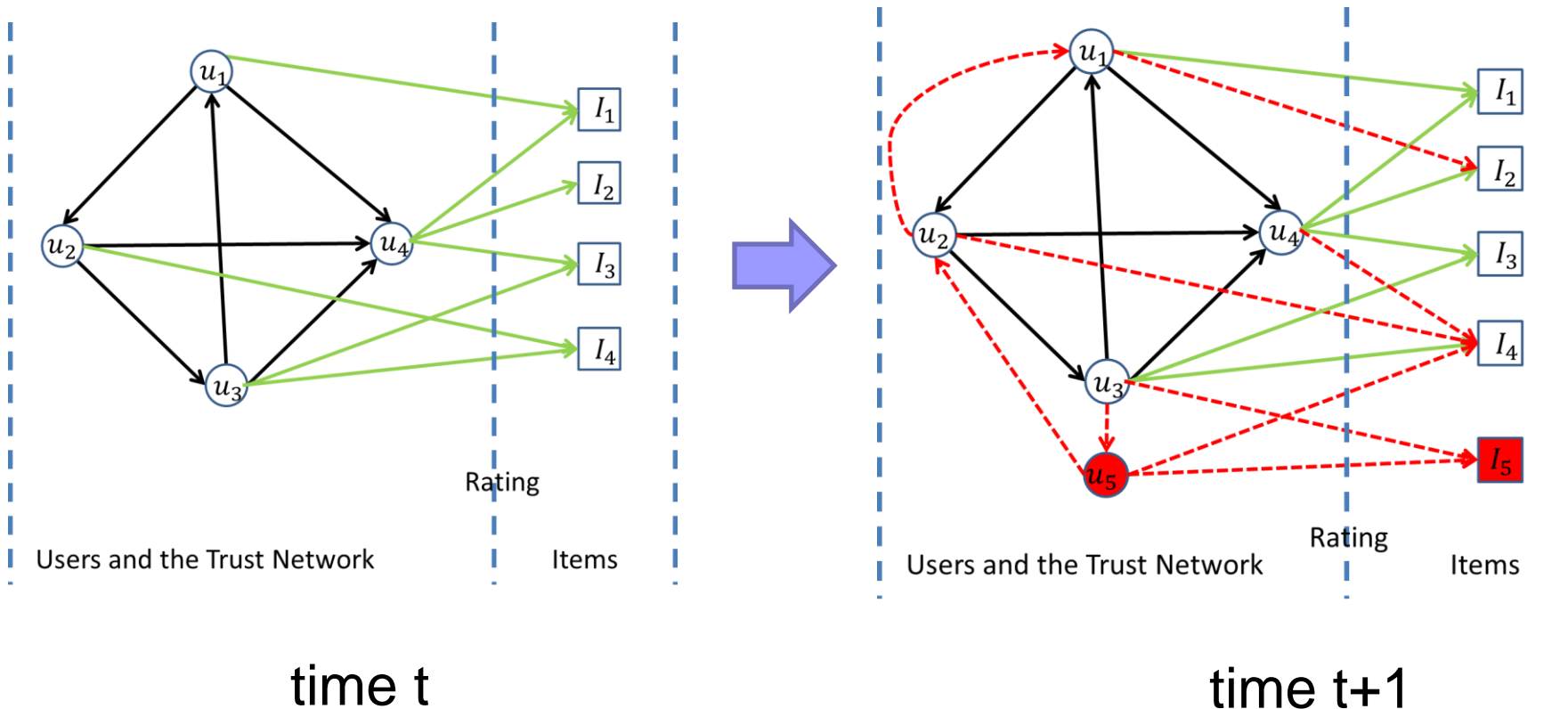
信任计算_动态演化



time t

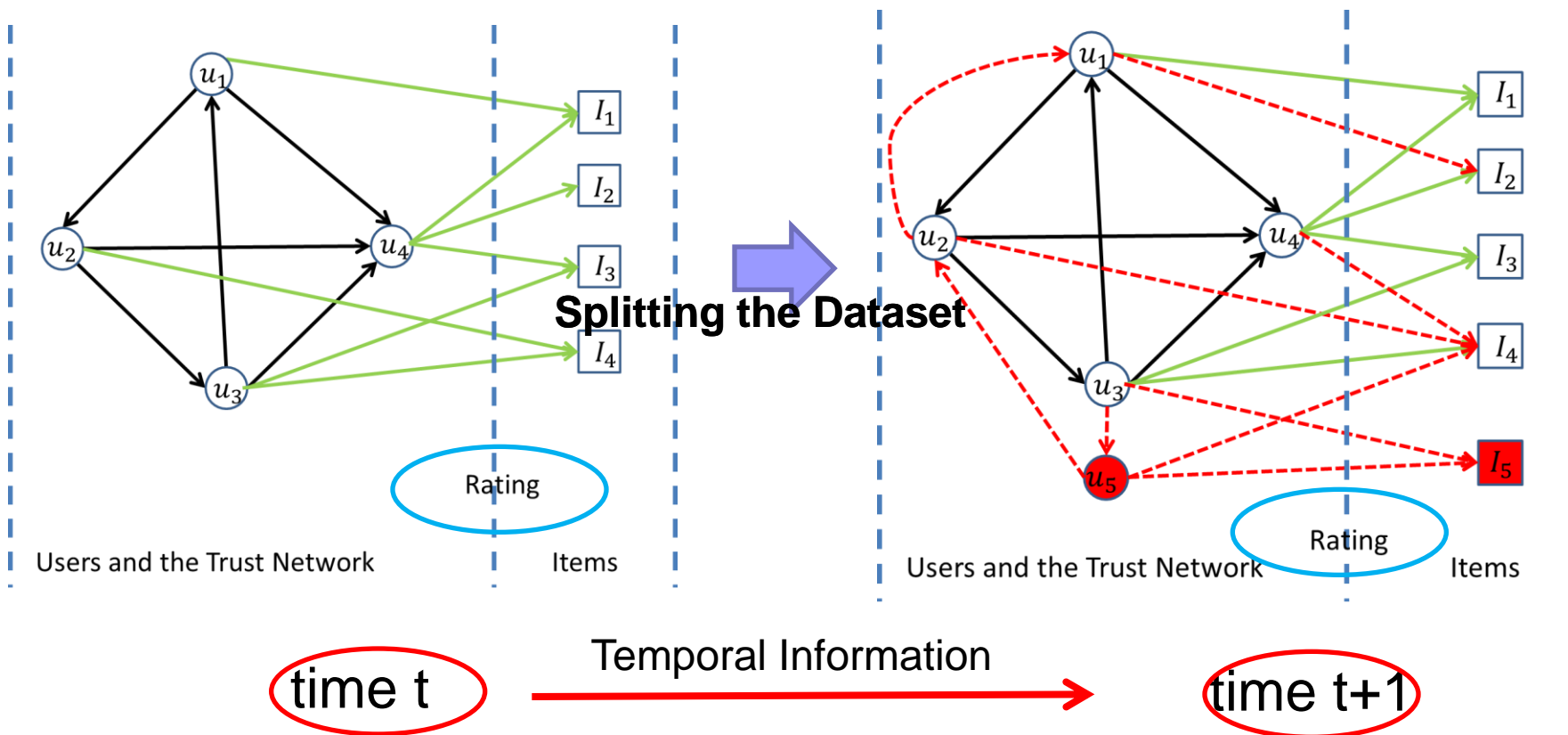
Epinions在线评分网络

信任计算_动态演化



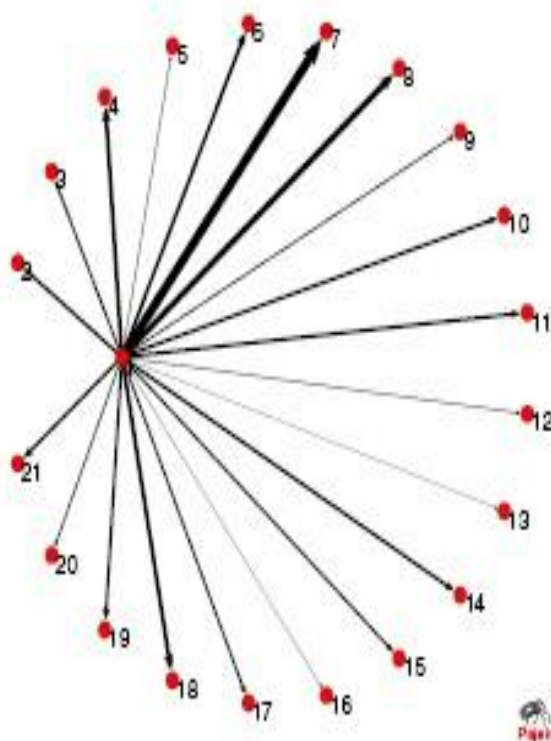
Epinions在线评分网络

信任计算_动态演化

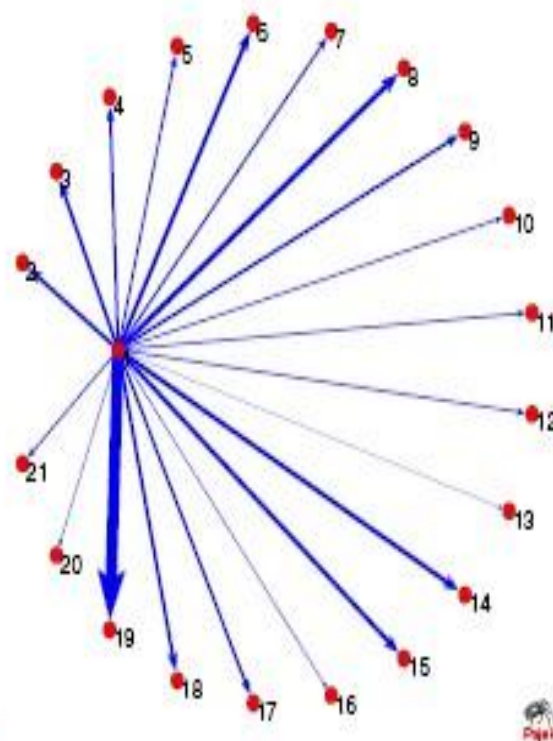


Epinions在线评分网络

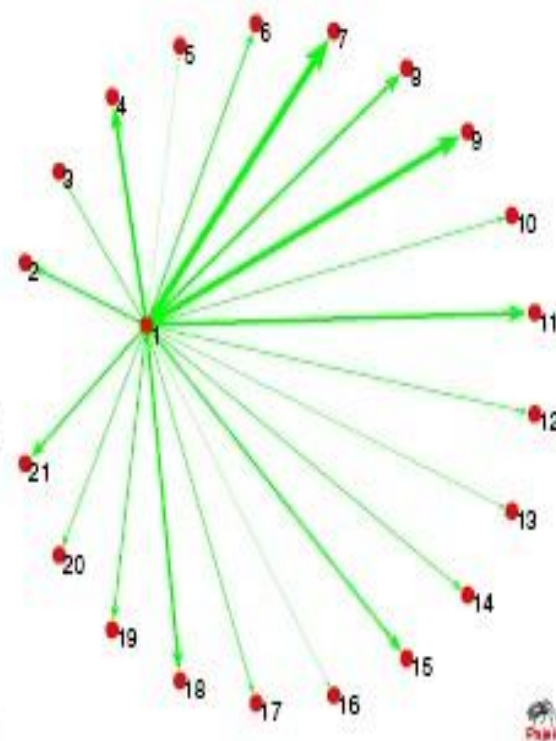
信任计算_多主题



(a) Single Trust



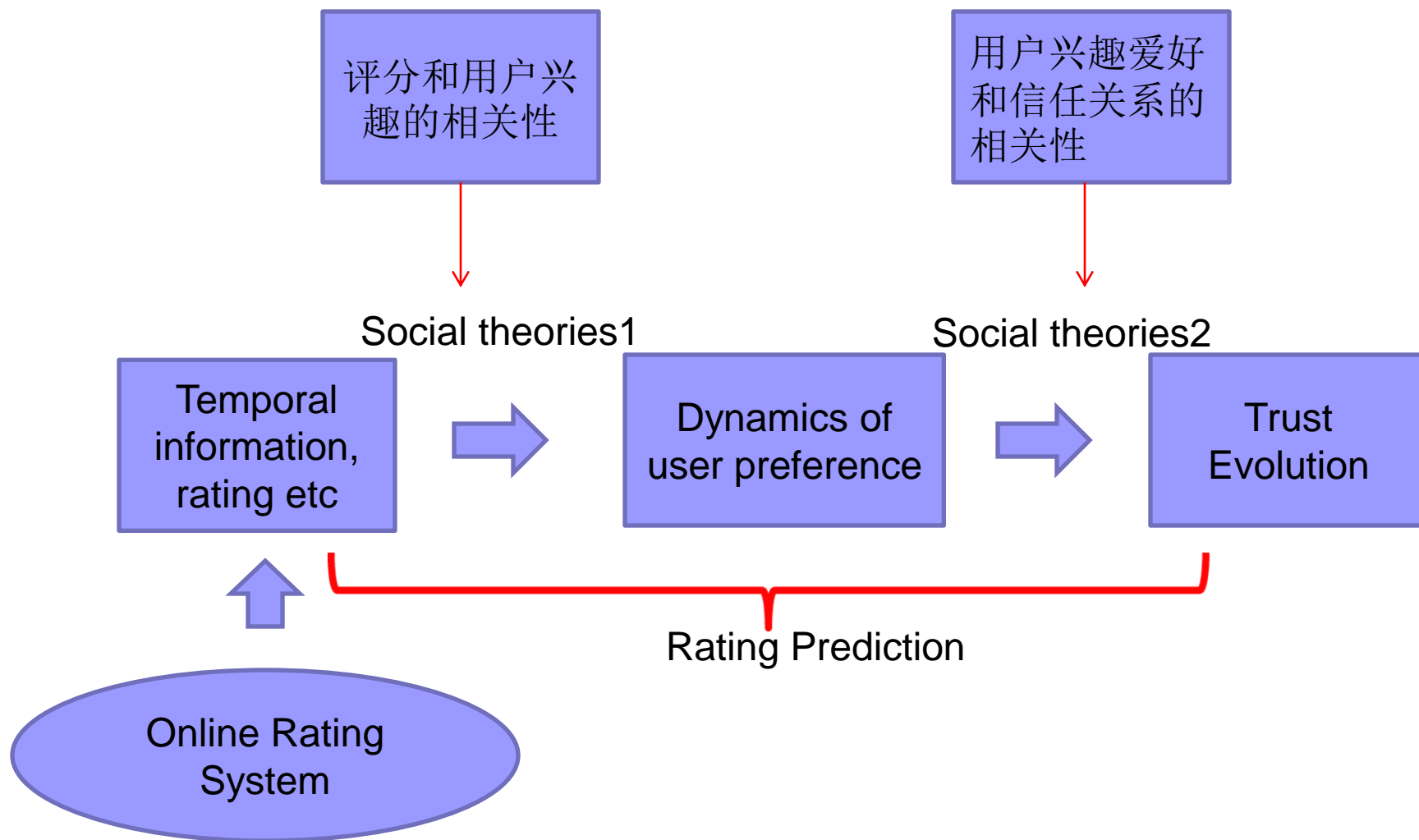
(b) Trust in Home & Garden



(c) Trust in Restaurants

Epinions在线评分网络

信任计算思想



信任计算模型

$$\begin{aligned}
 & \min_{\mathbf{p}_i^t \geq 0, \mathbf{q}_j \geq 0, \eta_i \geq 0} \sum_{t=1}^T \sum_{(i,j) \in \mathcal{O}_t} \left(r_{ij}^t - \alpha \sum_{k=1}^K \mathbf{q}_j(k) \mathbf{p}_i^t(k) \right) - \\
 & (1 - \alpha) \frac{\sum_{v \in N_i^t} \sum_{k=1}^K f(\mathbf{w}^\top \mathbf{s}_{ivk}^t + b_i) \mathbf{q}_j(k) r_{vj}}{\sum_{v \in N_i^t} \sum_{k=1}^K f(\mathbf{w}^\top \mathbf{s}_{ivk}^t + b_i) \mathbf{q}_j(k)} \Big)^2 \\
 & + \beta \left(\sum_{u_i \in \mathcal{U}_T} \sum_{t=t_{u_i}}^T \|\mathbf{p}_i^t\|_2^2 + \sum_{I_j \in \mathcal{I}_T} \|\mathbf{q}_j\|_2^2 \right. \\
 & \left. + \|\mathbf{w}\|_2^2 + \sum_{u_i \in \mathcal{U}_T} \|b_i\|_2^2 + \sum_{u_i \in \mathcal{U}_T} \|\eta_i\|_2^2 \right) \\
 & + \lambda \sum_{u_i \in \mathcal{U}_T} \sum_{t=t_{u_i}+1}^T \sum_{k=1}^K c(\mathbf{p}_i^t(k) - \mathbf{p}_i^{t-1}(k))
 \end{aligned}$$

Part 1: $\alpha \sum_{k=1}^K \mathbf{q}_j(k) \mathbf{p}_i^t(k)$
 Part 2: $\frac{\sum_{v \in N_i^t} \sum_{k=1}^K f(\mathbf{w}^\top \mathbf{s}_{ivk}^t + b_i) \mathbf{q}_j(k) r_{vj}}{\sum_{v \in N_i^t} \sum_{k=1}^K f(\mathbf{w}^\top \mathbf{s}_{ivk}^t + b_i) \mathbf{q}_j(k)}$
 Part 3: $f(\mathbf{w}^\top \mathbf{s}_{ivk}^t + b_i)$
 Part 4: $\lambda \sum_{u_i \in \mathcal{U}_T} \sum_{t=t_{u_i}+1}^T \sum_{k=1}^K c(\mathbf{p}_i^t(k) - \mathbf{p}_i^{t-1}(k))$

Part 1: Modeling Rating via User Preference

■ 用户兴趣爱好和商品特征相关性

$$- \hat{r}_{ij}^t = \mathbf{q}_j^\top \mathbf{p}_i^t = \sum_{k=1}^K \mathbf{q}_j(k) \mathbf{p}_i^t(k),$$

- \mathbf{p}_i^t 表示用户 i 在时间 t 时的兴趣爱好, \mathbf{q}_j 表示商品 j 的特征, k 为商品分类

Part 2: Modeling Rating via Trust Network

■ 利用信任网络计算商品评分

$$\hat{r}_{ij}^t = \frac{\sum_{v \in N_i^t} \sum_{k=1}^K w_{ivk}^t \mathbf{q}_j(k) r_{vj}}{\sum_{v \in N_i^t} \sum_{k=1}^K w_{ivk}^t \mathbf{q}_j(k)}$$

表示用户*i*和用户*v*
在主题分类*k*下的
信任值大小

早期评分的影响力会随着时间而衰减

$$r_{vj} = e^{-\eta_i(t-t_{vj})} r_{vj}^{t_{vj}}$$

Part 3: Modeling Trust and User preference

- 建模信任网络 and 用户兴趣爱好 的相关性

$$w_{ivk}^t = f(\mathbf{w}^\top \mathbf{s}_{ivk}^t + b_i),$$

\mathbf{s}_{ivk}^t 是用户间兴趣爱好相似矩阵, b_i 是用户偏差。

Part 4: Modeling Change of User Preference

■ 建模用户兴趣爱好的变化

$$\lambda \sum_{u_i \in \mathcal{U}_T} \sum_{t=t_{u_i}+1}^T \sum_{k=1}^K c\left(\mathbf{p}_i^t(k) - \mathbf{p}_i^{t-1}(k)\right)$$

c 函数描述用户兴趣爱好是如何变化, **λ**控制变化速度

实验数据

■ Epinions

- 商品评价网站



# of Users	22,166
# of Items	296,277
# of Categories	27
# of Ratings	922,267
# of Links	355,813
First Rating on	Jul 05 1999
Last Rating on	May 08 2011
Trust Network Density	0.0014
Clustering Coefficient	0.1518

■ Epinions分为11个时间戳



eTrust提高信任相关应用的性能

■ 评分预测 (Rating Prediction)

	\mathcal{K}	\mathcal{N}	$\mathcal{K} + \mathcal{N}$
Mean	1.1054	1.1562	1.1106
NN	1.1092	1.1566	1.1148
MF	1.0804	1.1472	1.0872
MF+NN	1.0675	1.1392	1.0747
mTrust	1.0566	1.1375	1.0646
eTrust	1.0299	1.0783	1.0347

\mathcal{K} 原有的无变化评分系统,
 \mathcal{N} 有新产品或者新用户加入
的评分演化系统(10.06%)

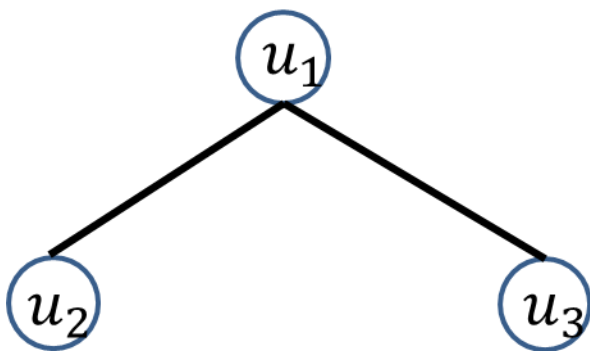
■ 信任预测 (Trust Prediction)

	$\mathcal{E}(\%)$	$\mathcal{N}(\%)$	$\mathcal{E} + \mathcal{N}(\%)$
Simi	48.41	28.94	43.93
TP	45.47	N.A.	35.01
Simi+TP	50.19	28.94	45.31
eTrust	55.07	33.83	50.18

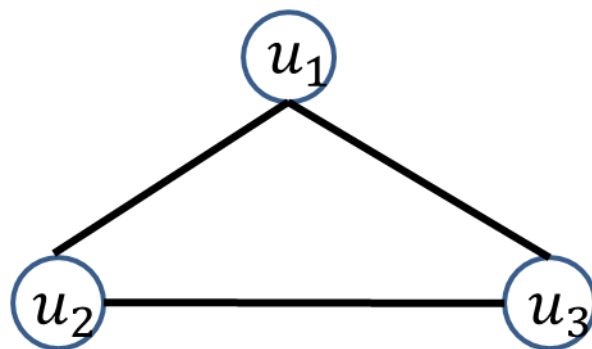
\mathcal{E} 现有用户的信任关系建立,
 \mathcal{N} 新加入新用户后的信任关
系预测(23.51%)

结论

- 开放三角（an open triad）的演化速度是闭合三角（a closed triad）的 **6.12 倍**

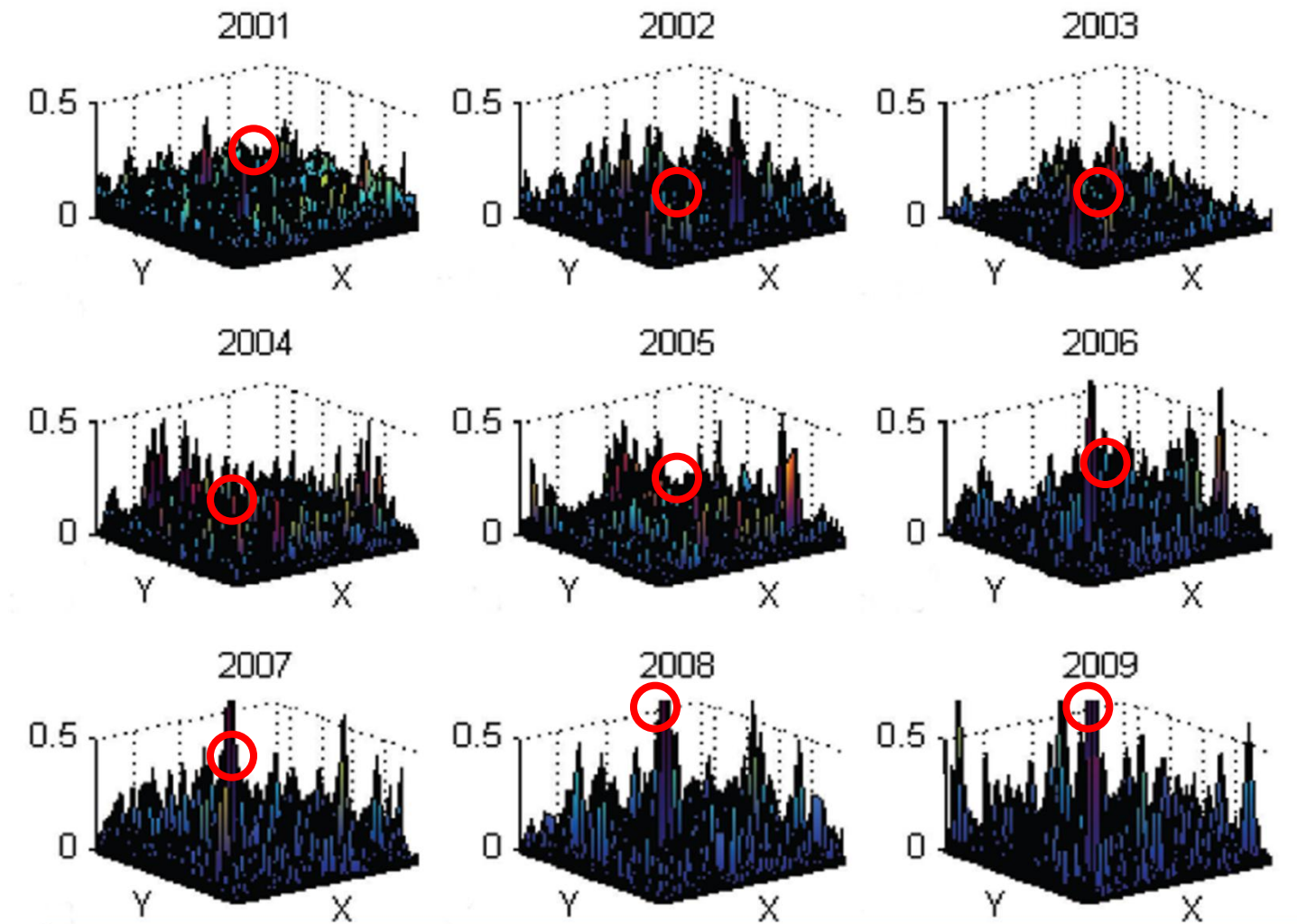


Open Triad

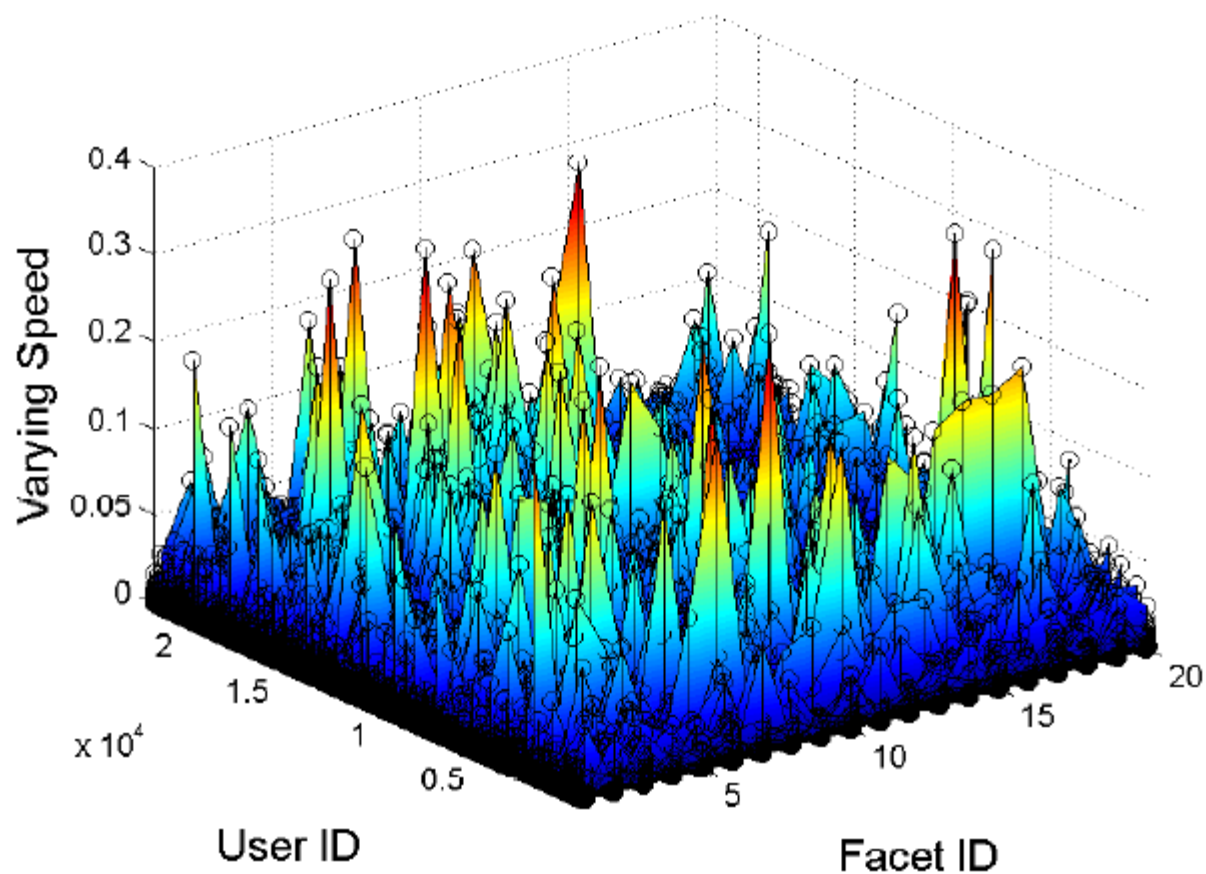


Closed Triad

用户兴趣爱好随时间变化图



用户和商品的变化速率图



未来工作

■ 寻找更多的实际应用

- Ranking evolution
- Recommendation systems
- Helpfulness prediction

SiteReco™

智能商品推荐



谢谢！

