

# 社会化推荐

张叶银

[Yeyin.zhang@renren-inc.com](mailto:Yeyin.zhang@renren-inc.com)

# 推荐系统

- 给用户推荐他可能感兴趣的东西
  - 好友推荐，商品推荐，文档推荐，广告推荐
- Amazon, Netflix, Google, Facebook, Youtube.....

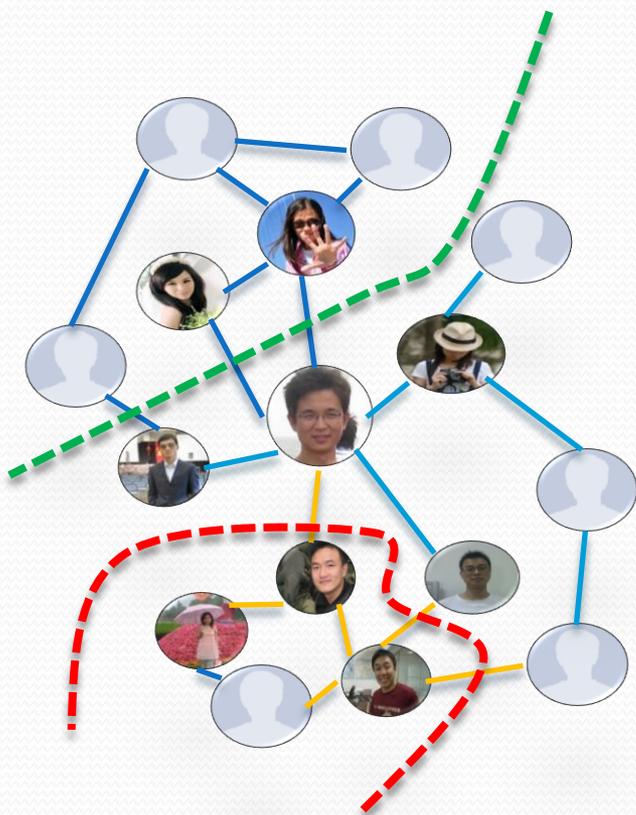
# 主流推荐算法

- 协同过滤
  - User-based, Item-based
- 内容过滤
  - 属性的相关性
- 社会网络
  - 图的方法

# 推荐系统评测

- 准确度
- 覆盖率
- 多样性
- 新颖性

# Social Graph



我们为你找到了49位可能认识的人

 陈箭 加为好友	 林豪华 加为好友	 尚美美 加为好友	 任思睿 加为好友	 高金来 加为好友
 瑞华华 加为好友	 周彦伟 加为好友	 詹捷 加为好友	 刘春文 加为好友	 国力群 加为好友

为你推荐

			
拍客/结婚偶遇开 主 分享:66456	当B-BOX牛人遇上 诈骗集团打来的电 话 分享:57979	《你猜你猜你猜猜 猜》夜店女王素颜 太惊人... 分享:68530	《索命DV》延边 大学探险社团鬼怪 集体失踪死... 分享:28231

为你推荐

	
---	---

为你推荐

【某A的折纸课堂】如何折出世界上最强的三大纸飞机之二：复仇者  
白:王念州 Vico

# 好友推荐

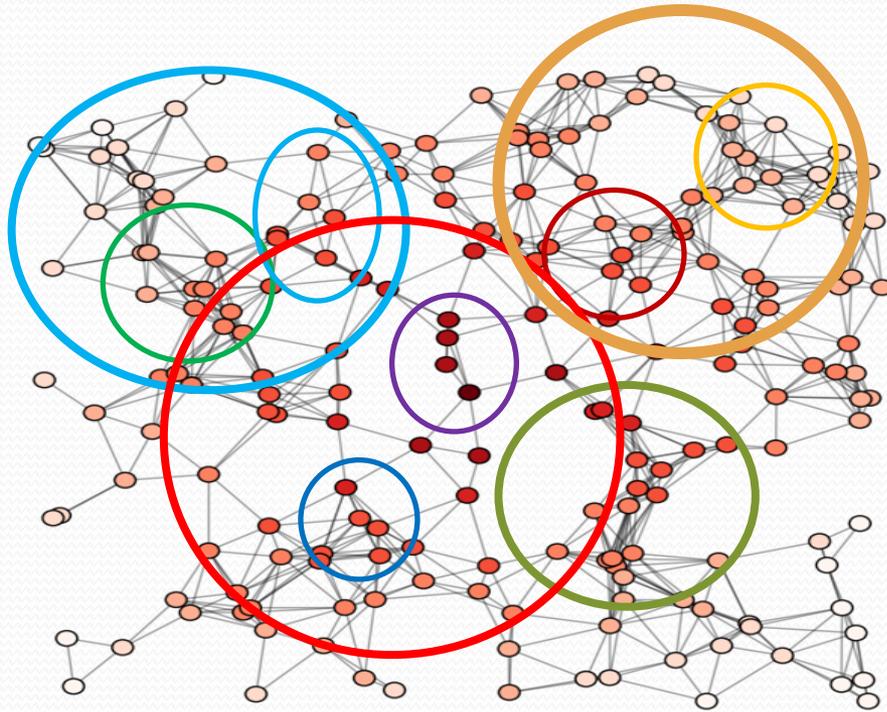
- 关键问题
  - 用户间的亲密度
    - 共同好友数目
    - 个人资料相似程度
    - 用户互动频度
    - 用户的兴趣

# People you may know

- 二度好友 Friends of friends

$$\text{similarity}(user1, user2) = |friendset1 \cap friendset2|$$

# Circles



- 部门 » 公司
- 班级 » 学校
- 实验室 » 班级
- ...

好友簇

Strong ties



Weak ties

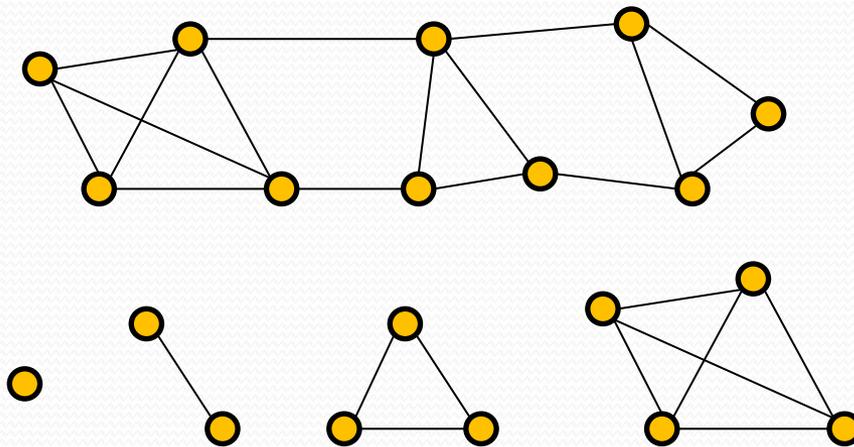
# Strong ties

- Community Detection

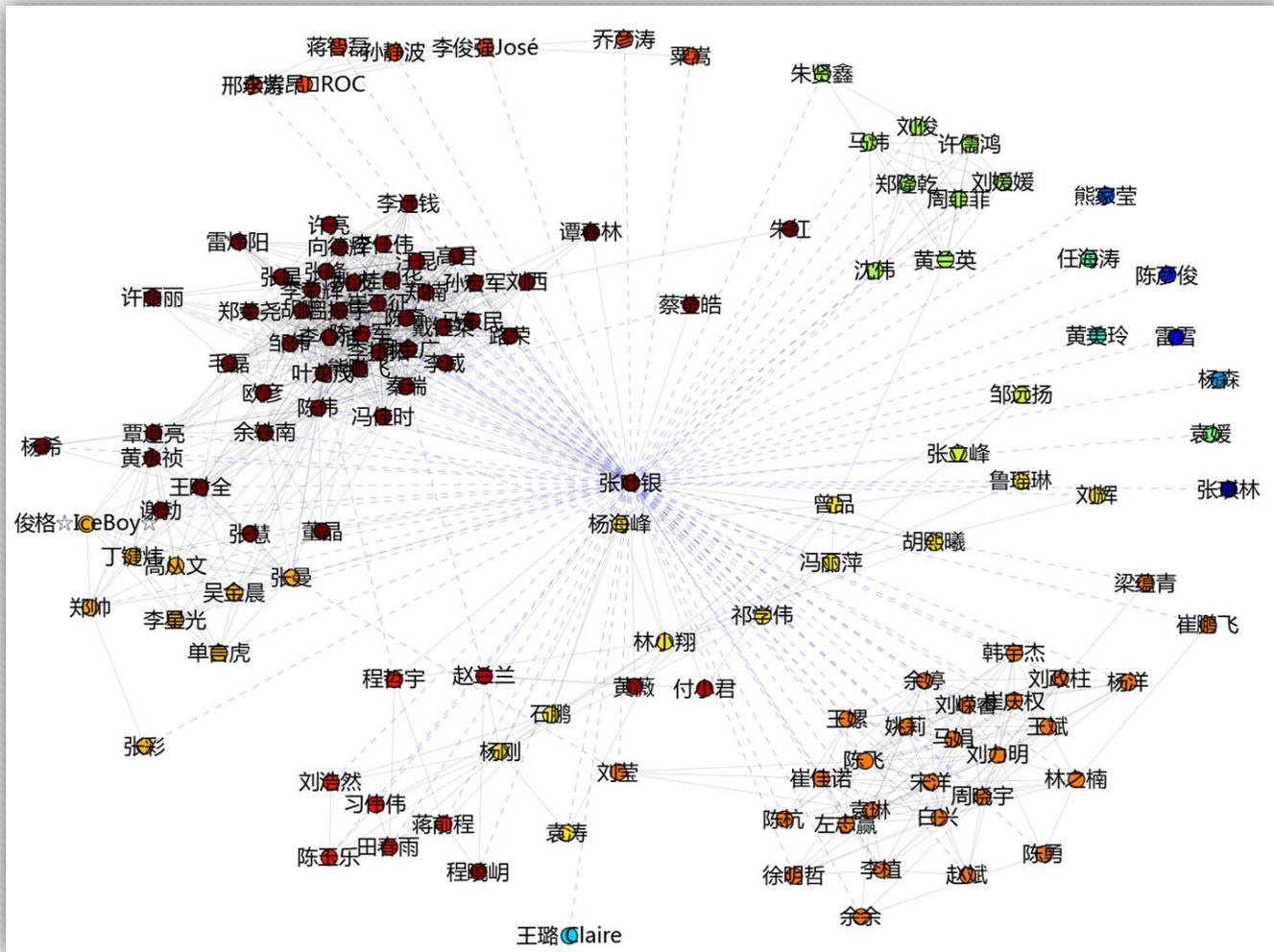
- 图的连通性

- Clique:

- 例如:  $V' \subseteq V$ , 对于任意的  $u, v \in V', u \neq v, uv \in E$



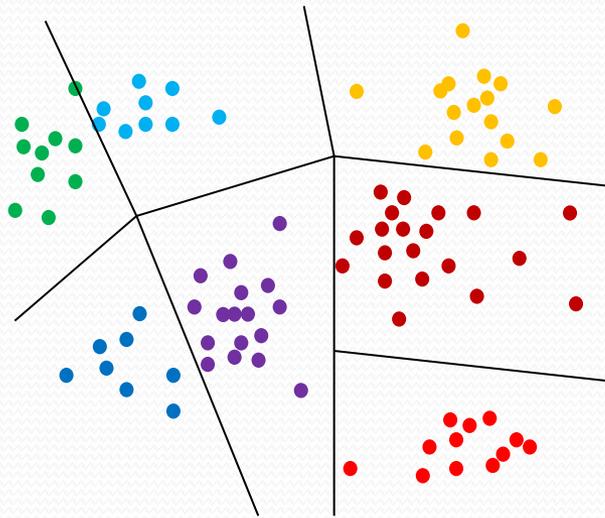
# 圈子



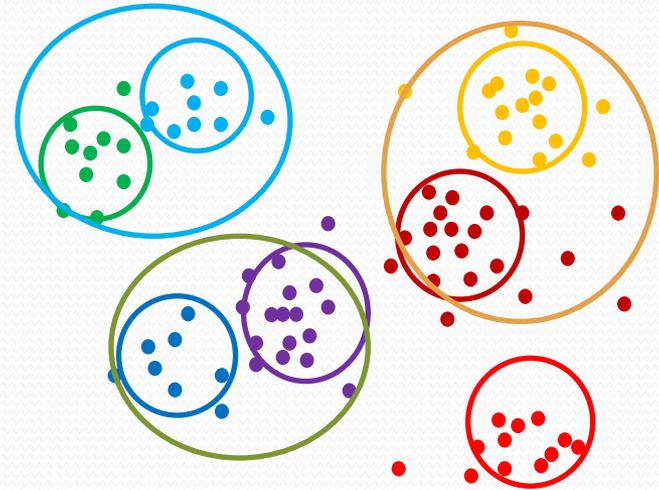
# 应用

- 推荐
- 隐私控制
- 新鲜事定制

# To weak ties : unsupervised learning

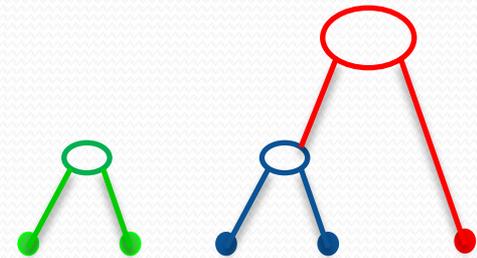
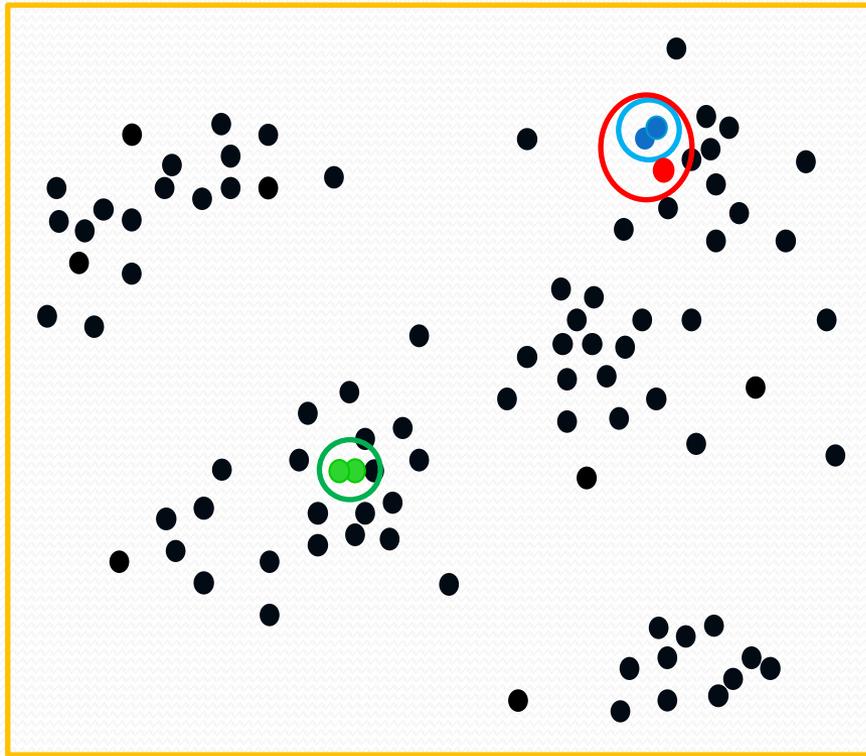


flat clustering



hierarchical clustering

# Hierarchical agglomerative Clustering



# Hierarchical agglomerative clustering

**Method:** Merge the nearest clusters until a single cluster is left

*Procedure HAC* (N points, stop criterion)

{

(1) Initialize n points as n cluster centers;

(2) Iterate over centers until stop criterion is satisfied:

a. Compute pair-wise similarity between

any two centers  $sim(c_i, c_j)$

b. Find the nearest pair of centers

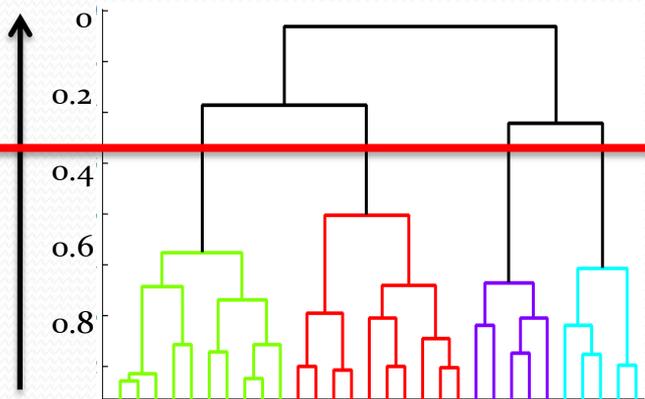
c. Merge the two centers

$$\langle i, j \rangle \leftarrow \arg \max_{i, j} sim(c_i, c_j)$$

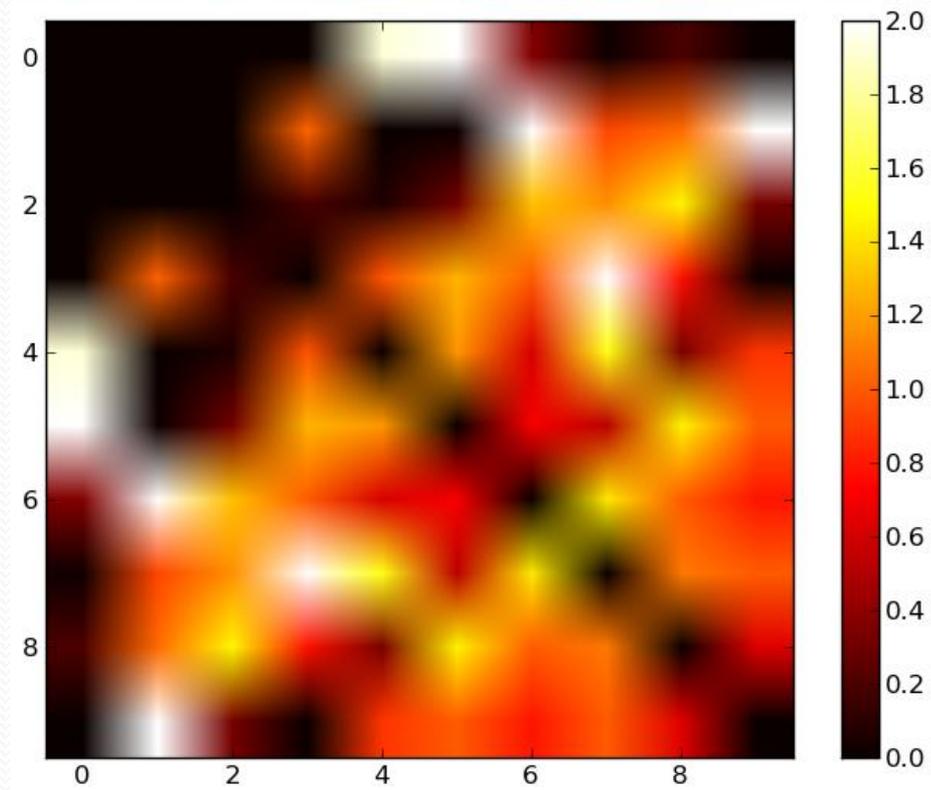
(3) Output the hierarchical clusters.

}

Monotonic



# Pair-wise distances



**Symmetric similarity matrix**

# To weak ties

- 相似性度量

$$\text{similarity}_{ij} = \left\| v_j \cap v_i \right\| \cdot \sum_k \sum_l R_{kl}$$

# To weak ties

- 选择

- 特征:

- 用户交互:  $r_{ij} = \sum_{k \in A} a^k \cdot n_{ij}^k$

- 用户关系:

$$c_i^o =$$

$$\arg \max_{c_i^o} \{ \sum r_{jk} \mid c_i^o \in C_i^o, j, k \in v_i \}$$

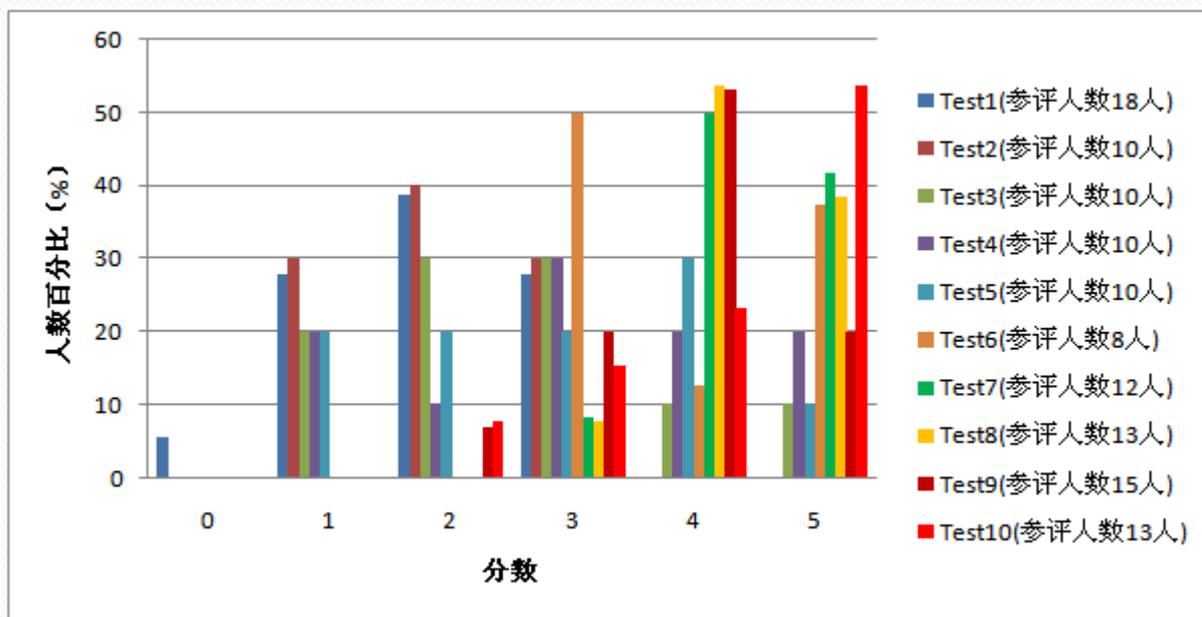
# 合并

$$o_i^r(U_i^r, C_i^r) = o_j^{r-1}(U_j^{r-1}, C_j^{r-1}) \cup o_k^{r-1}(U_k^{r-1}, C_k^{r-1}),$$

其中,  $U_i^r = U_j^{r-1} \cup U_k^{r-1}, C_i^r = C_j^{r-1} \cap C_k^{r-1}$

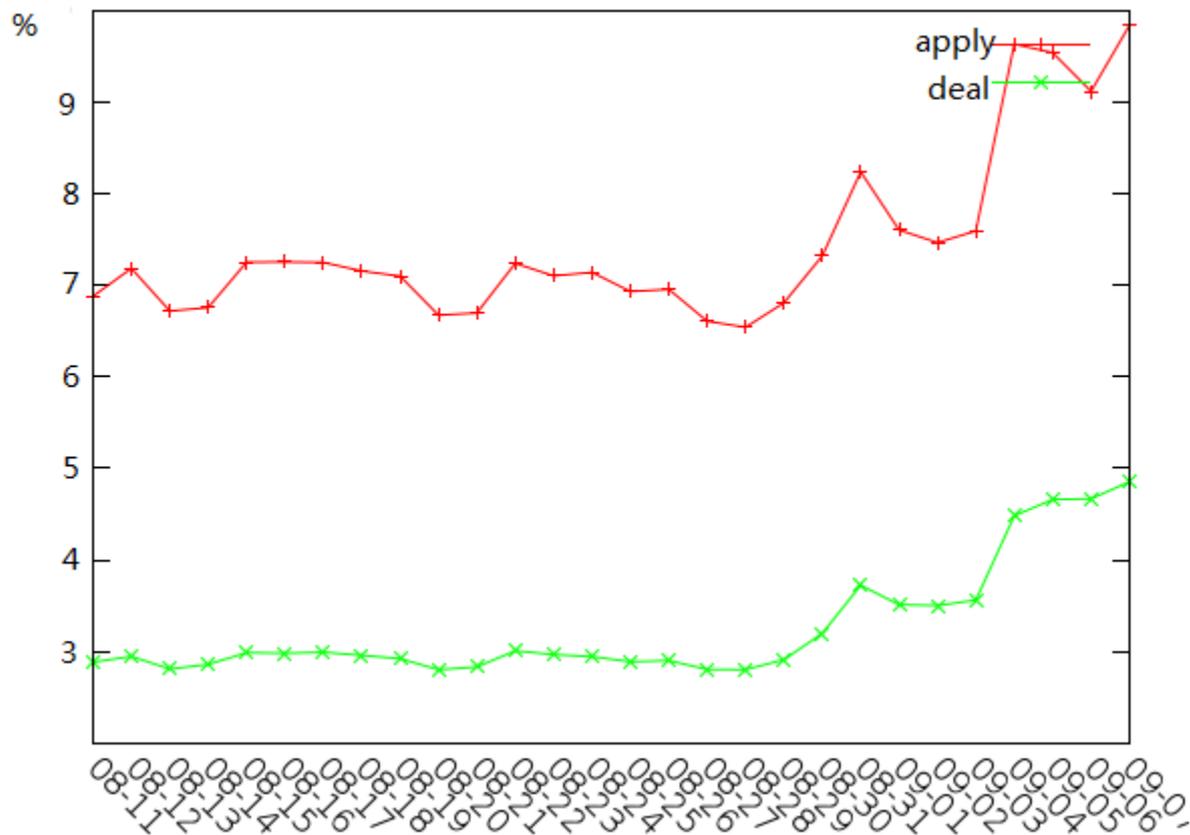
# 评测

- 缺少标定的数据，主观评价为主



# 线上数据

Home Recommend Stats Graph

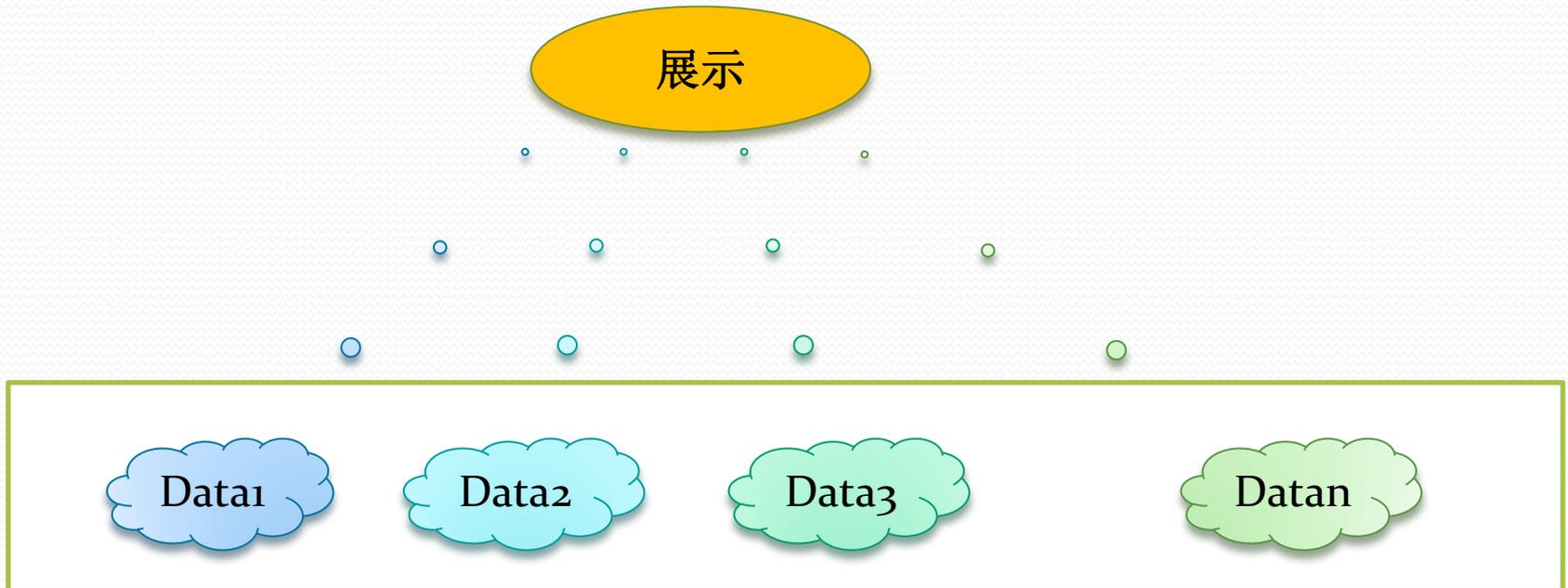


# 个性化推荐

- 用户偏好
- 用户兴趣
- 用户成长

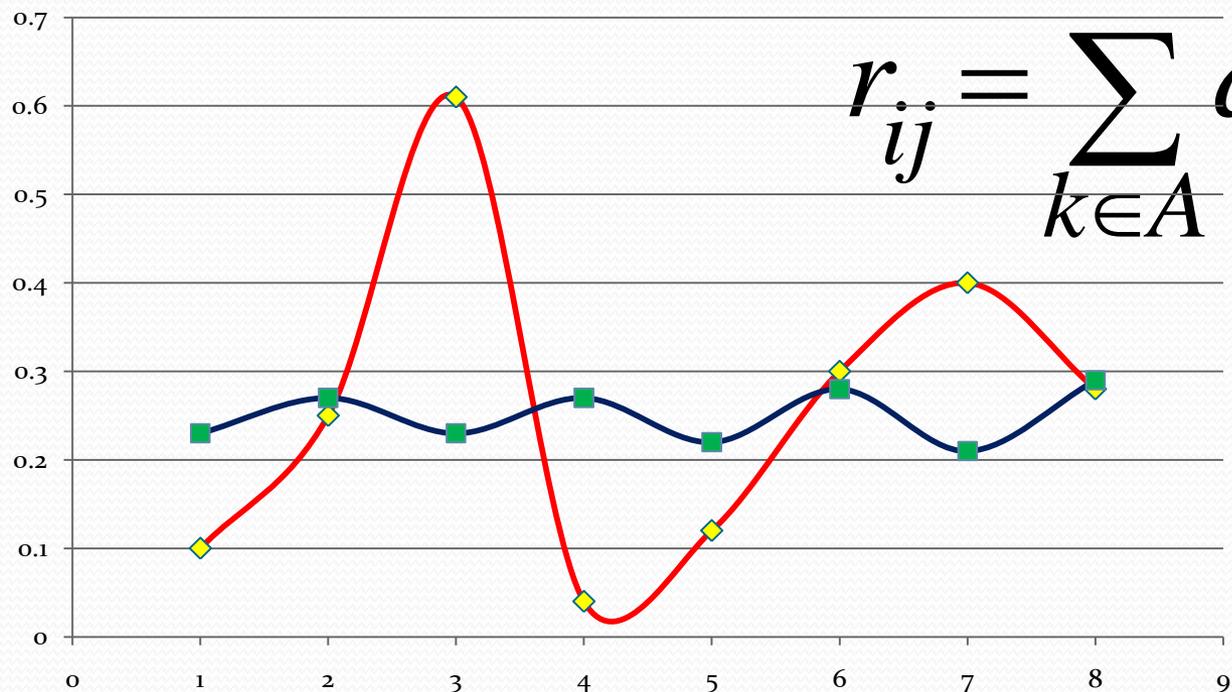
# 偏好

- 圈子
  - 年龄，学校，地域，性别.....



# 偏好的度量

- 用户行为
  - 访问个人主页、相册，分享，留言，评论等



# 量化

- 信息熵
  - 年龄，学校，地域，性别
    - $x = (0.1, 0.25, 0.61, 0.04) \approx \text{nonuniform}$
    - $y = (0.23, 0.27, 0.23, 0.27) \approx \text{uniform}$

$$H[x] = -\sum_x p(x) \log p(x)$$

- $H(x) = -0.1 \cdot \log 0.1 - 0.25 \cdot \log 0.25 - 0.61 \cdot \log 0.61 - 0.04 \cdot \log 0.04 = 1.007$
- $H(y) = -0.23 \cdot \log 0.23 - 0.27 \cdot \log 0.27 - 0.23 \cdot \log 0.23 - 0.27 \cdot \log 0.27 = 1.383$

# 好友推荐

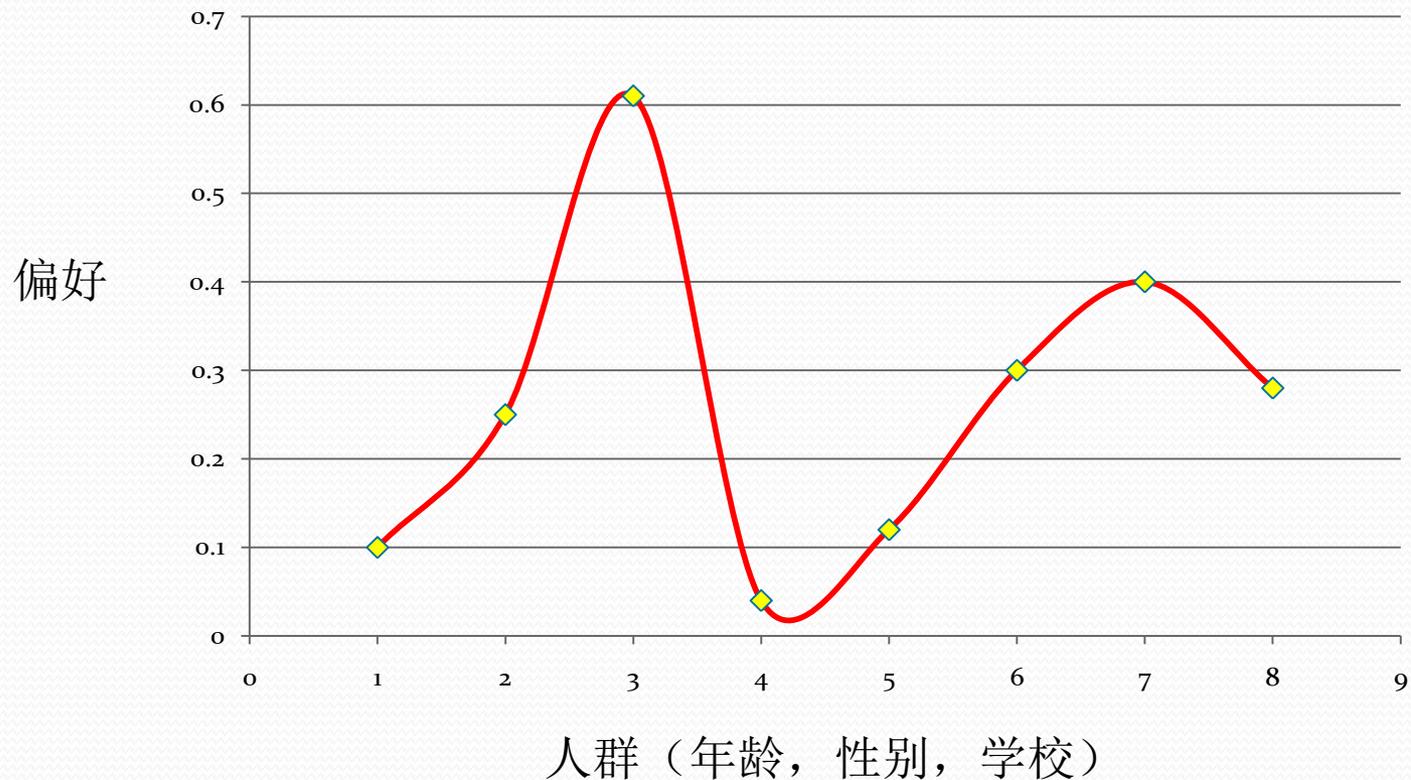
- 首页28个推荐位，从推荐源数据中随机选取

## Ranking VS Sampling

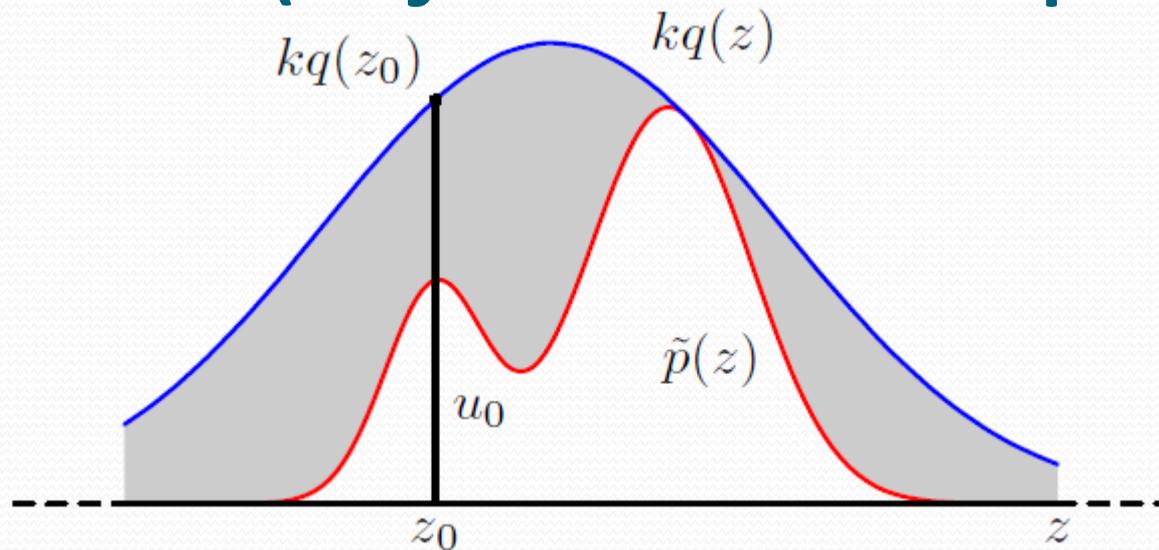
- Ranking
  - 亲密度
  - 影响力
- Sampling
  - 偏好
  - 推荐人群的多样化
  - 新颖性

# 采样

- 偏好分布密度函数 $P(x)$ 未知



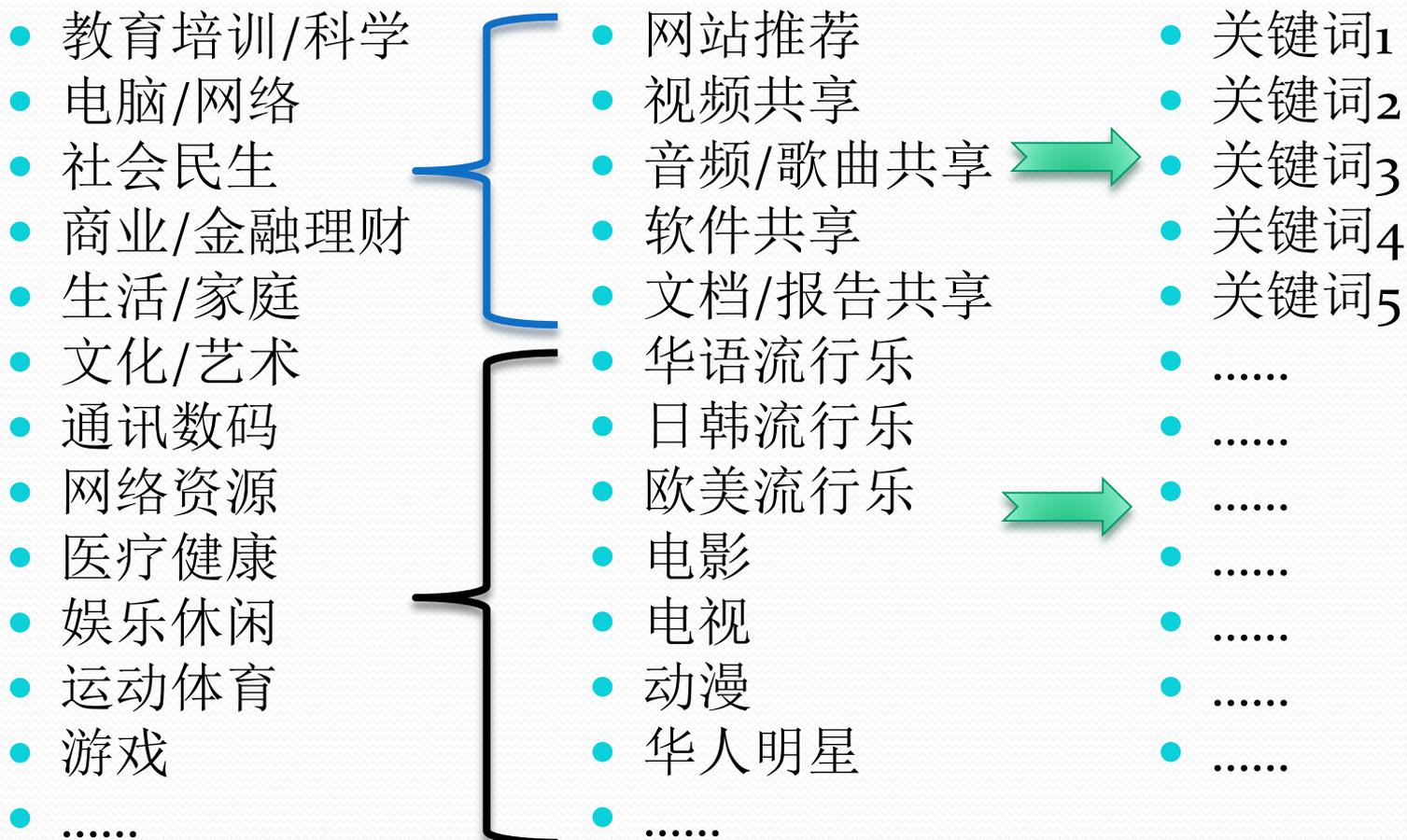
# 取舍抽样(rejection sampling)



$$p(\text{accept}) = \int \{ \tilde{p}(z) / kq(z) \} \cdot q(z) dz$$

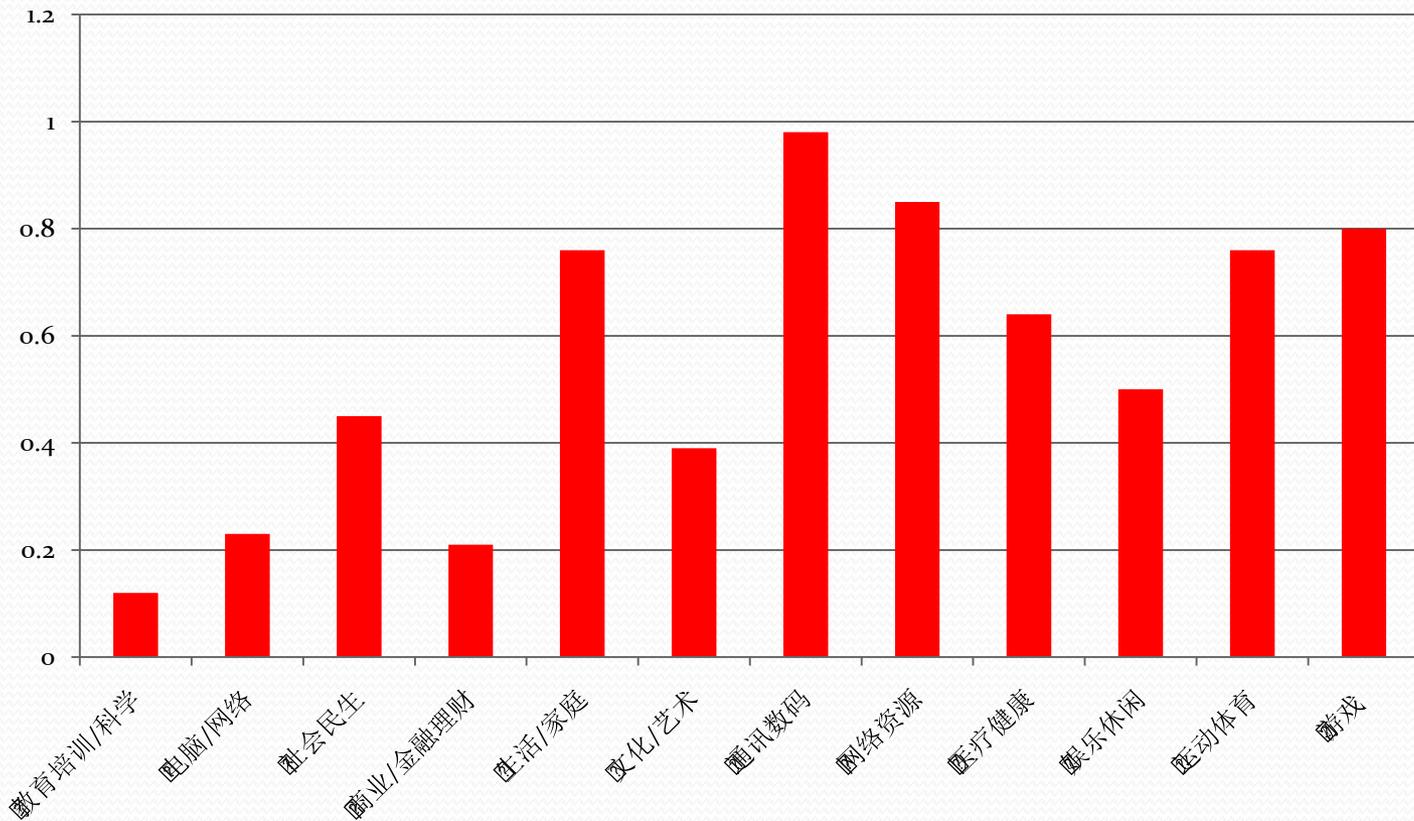
$$= \frac{1}{k} \int \tilde{p}(z) dz$$

# 用户兴趣

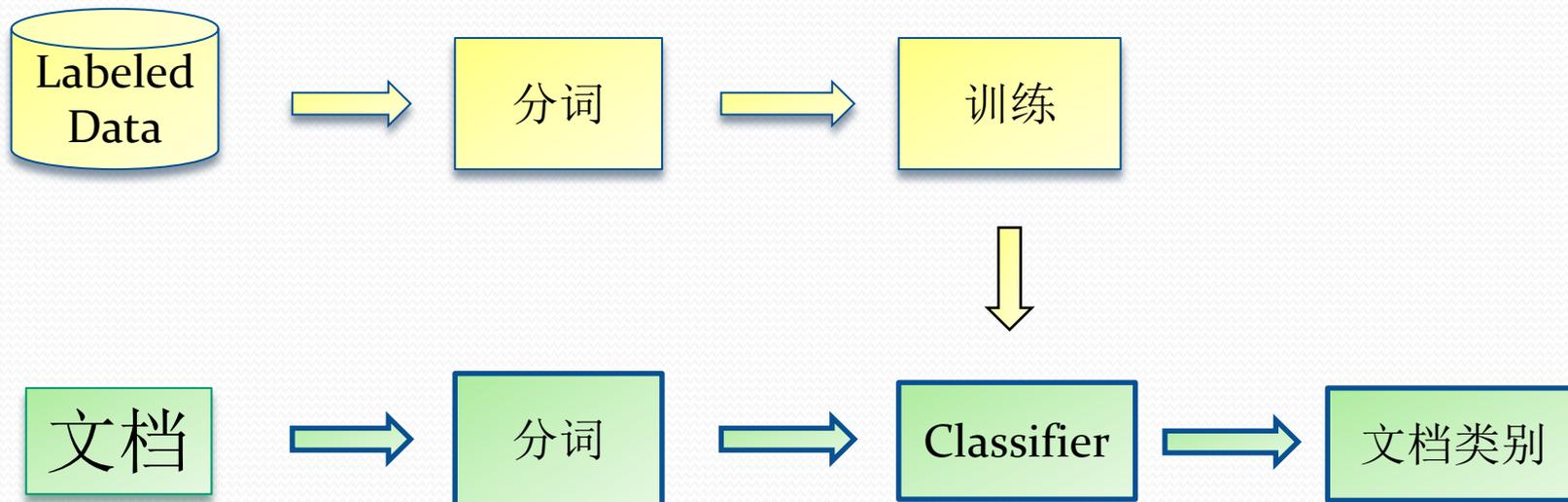


# User access

- Access log



# 文档分类 Supervised learning



# Naïve Bayes

- Bayes

$$p(\mathbf{C} | \vec{x}) = \frac{p(\vec{x} | C_k) p(C_k)}{p(\vec{x})}$$

Diagram illustrating the Naïve Bayes formula and its components:

- likelihood**:  $p(\vec{x} | C_k)$
- prior**:  $p(C_k)$
- posterior**:  $p(\mathbf{C} | \vec{x})$

Bayesian Network structure showing a parent node  $C$  (white circle) and three child nodes  $X_1$ ,  $X_2$ , and  $X_3$  (yellow circles).

```
graph TD; C((C)) --> X1((X1)); C --> X2((X2)); C --> X3((X3));
```

- Conditional independent

$$p(C_k | x_1, x_2) \propto p(x_1, x_2 | C_k) p(C_k)$$
$$\propto p(x_1 | C_k) p(x_2 | C_k) p(C_k)$$

# 内容过滤

- 公共主页
  - 相关性

$$\cos(U, P) = (\sum_i U_i)^T \cdot (\sum_j P_j) / \|\sum_i U_i\| \|\sum_j P_j\|$$

# 用户成长

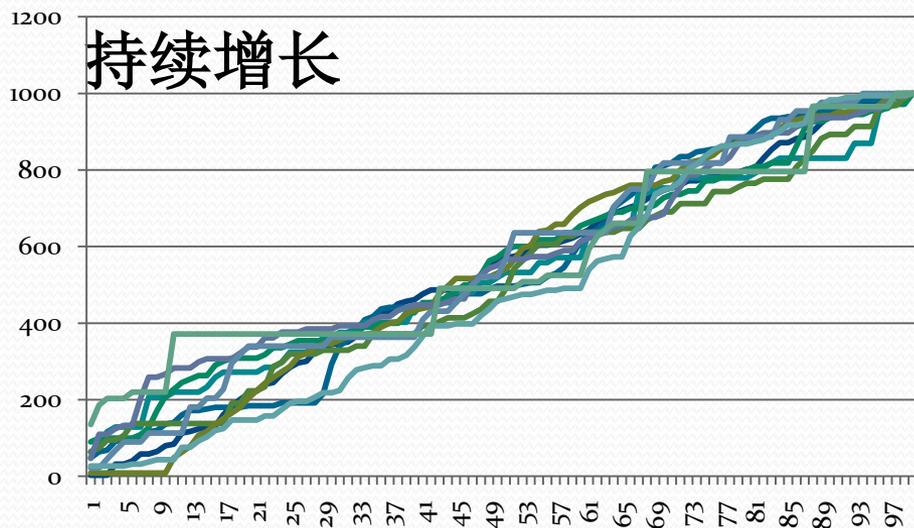
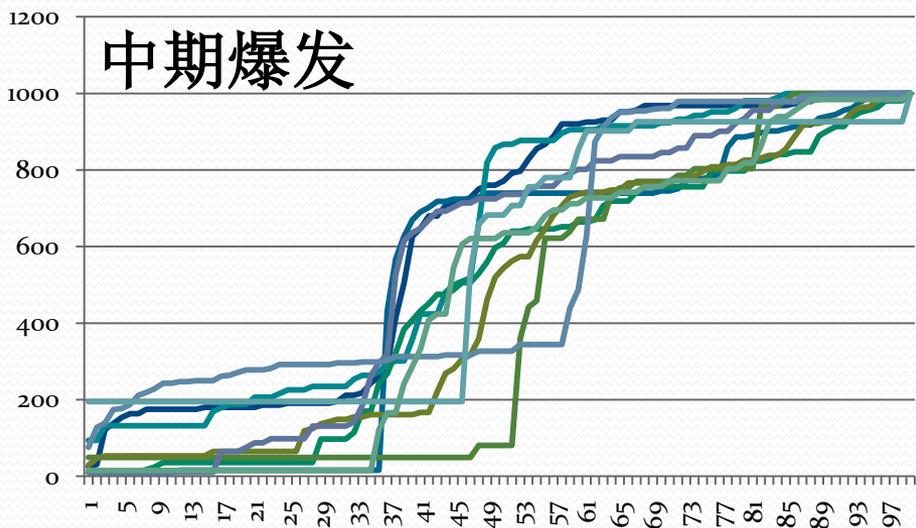
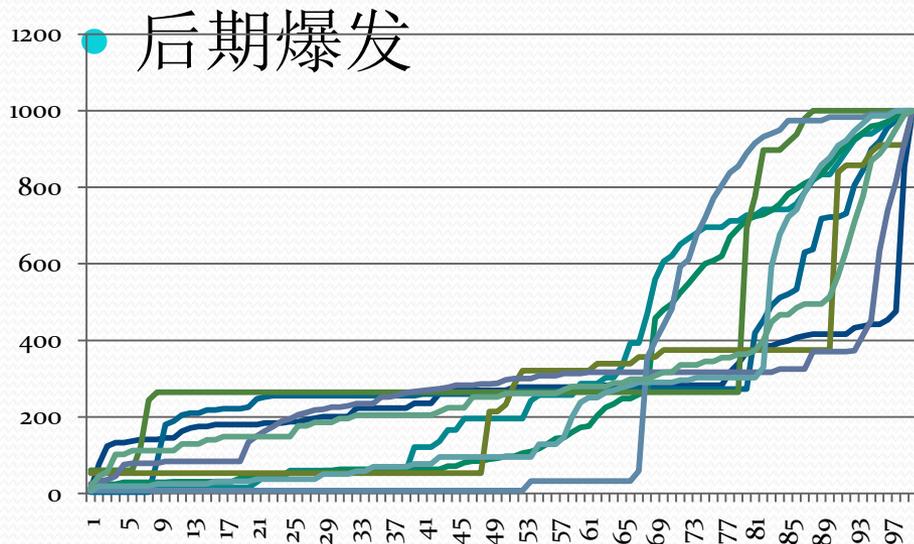
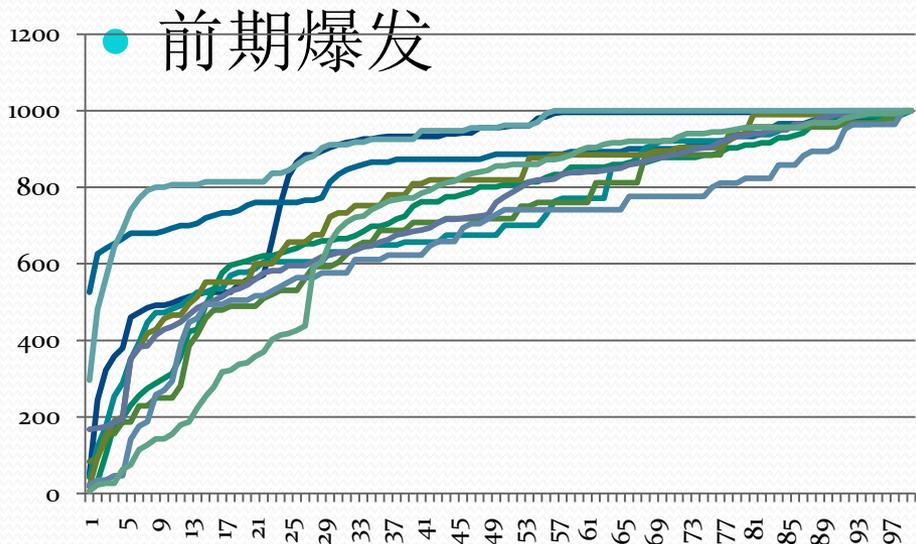
- 问题1

- 用户在什么阶段加好友最多？
- 爆发，还是持续增长？

- 问题2

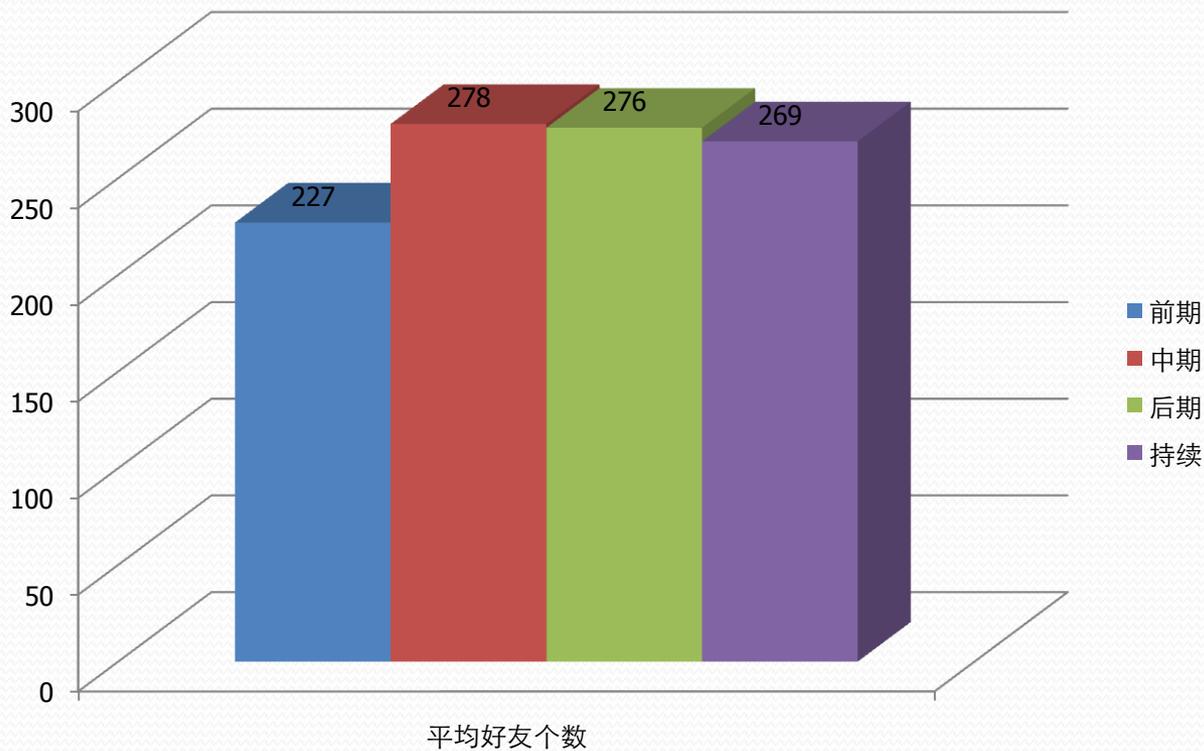
- 用户在什么阶段内，加什么样的好友？
- 资料？二度好友？陌生人？

# 好友增长类型



# 不同类型增长好友数

前期增长最终达到的好友数目较低



# 推荐

- 好友，公共主页，小站，音乐，日志，视频，小组.....
- 年龄 $a$ ，好友数目 $n$ ，兴趣 $i$ ，活跃度 $d$ ，推荐 $r$

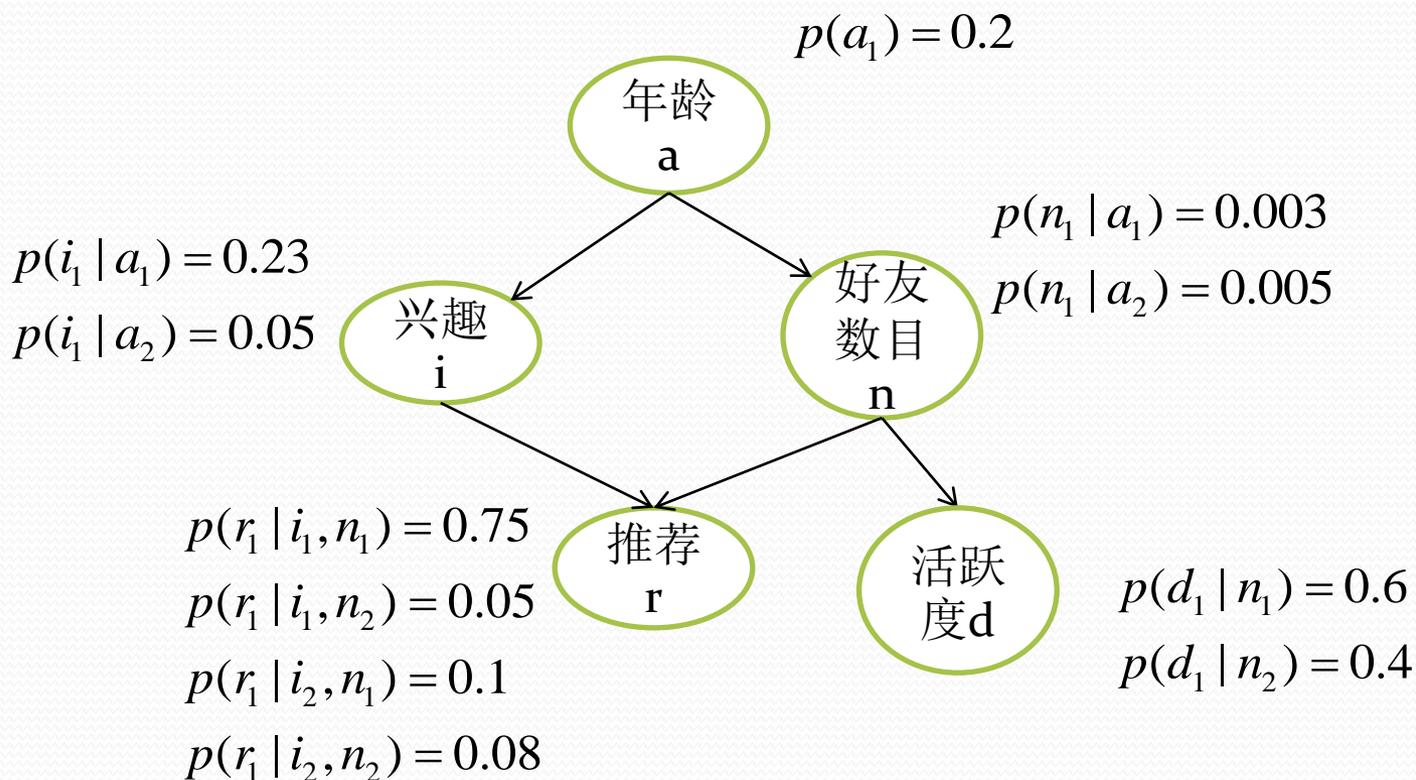
$$p(r | a, n, d)$$

$$= \frac{p(r, a, n, d)}{p(a, n, d)} = \frac{\sum_i p(r, a, i, n, d)}{\sum_{i, r} p(r, a, i, n, d)}$$

2<sup>5</sup>

# 贝叶斯推理

- 先验知识：用户成长模型，用户兴趣，用户偏好



# Factorization

$$p(a, i, n, r, d)$$

$$= p(a) p(i | a) p(n | a) p(r | i) p(r | n) p(d | n)$$



- 谢谢！