# Generating ordered list of Recommended Items: A Hybrid Recommender System of Microblog

Yingzhen Li and Ye Zhang

School of Mathematics & Computational Science, Sun Yat-sen University, Guangzhou, China

## Introduction

- A solution of KDD Cup 2012, track 1 task, which requires predicting users a user might follow in Tencent Microblog.
- Tencent Microblog has some special properties which we'll introduce.
- The system consists of:
  - ▷ keyword analysis
  - ▷ user taxonomy
  - ▷ item ranking
  - ▷ (potential)interests extraction
  - ▷ item recommendation(grading process)

## Speciality of Tencent Microblog

- It attracted a lot of registered users and became one of the dominant microblog platforms in China based on the large user group of its instant messaging service QQ.
- It is embedded in Tencent's other leading platforms.
- It considers those who frequently write(including retweet and comment) or read microblog messages - no matter on the website or other associated platforms - as active users.
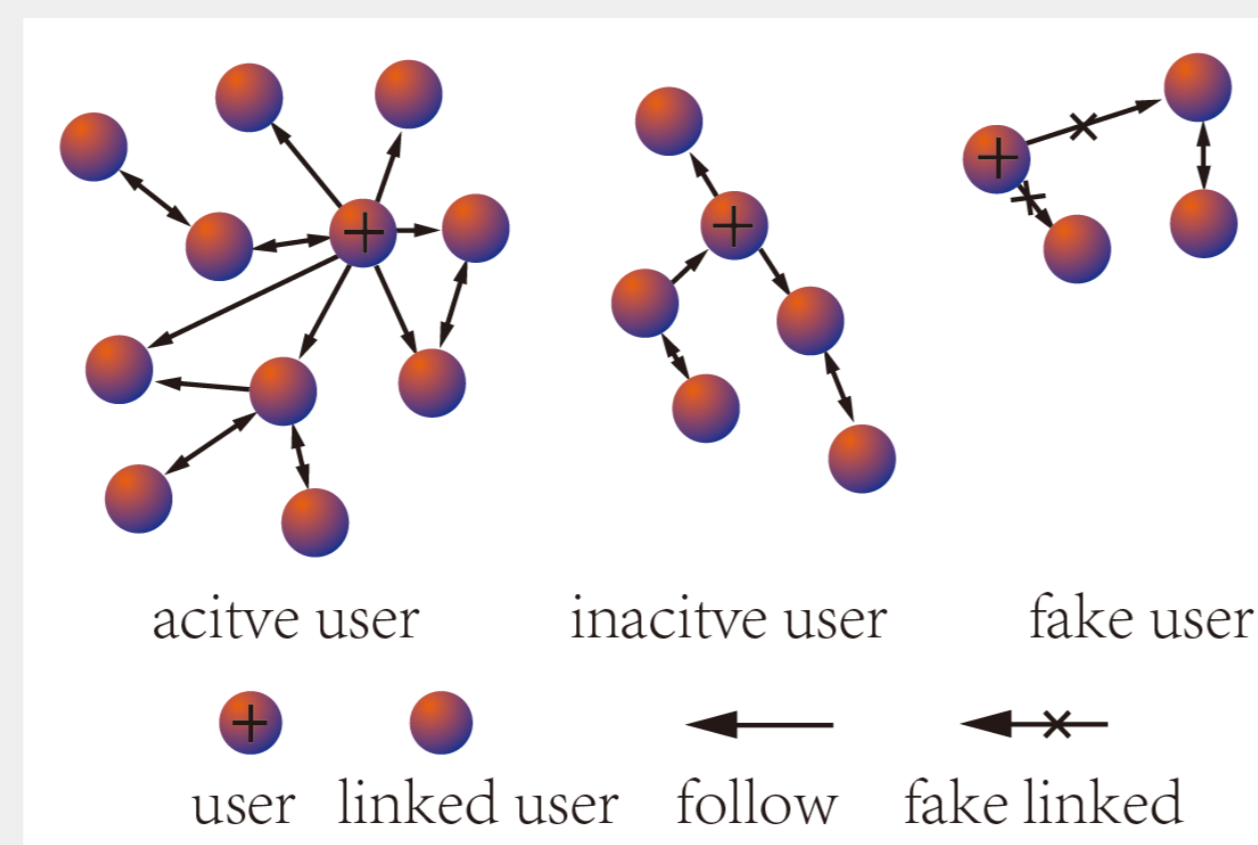
## Keyword Analysis

- Applying association rule algorithm to find them directly in the huge keyword set is unrealistic.
- We parallel this process by adopting revised FDM and insert the ambiguous keywords into different classes simultaneously.
- The necessity and sufficient condition for a frequent itemset is the frequency of all its subsets.
- User set: $U = \{u_1, u_2, ..., u_m\}$
- $u_j$'s keyword set: $K_j = \{k_{j1}, k_{j2}, ..., k_{jn_j}\}$ with weights $\mathcal{W}_j = \{w_{j1}, w_{j2}, ..., w_{jn_j}\}$
- keyword_class $= \{class_1, class_2, ..., class_N\}$, $class_i = \{k_{i1}, k_{i2}, ..., k_{im}\}$.

## User Taxonomy

- We divide the users into 3 groups by **user_class($u_j$)**:
  - ▷ Active - lots of tweets/interactions;
  - ▷ Inactive - few messages/interactions;
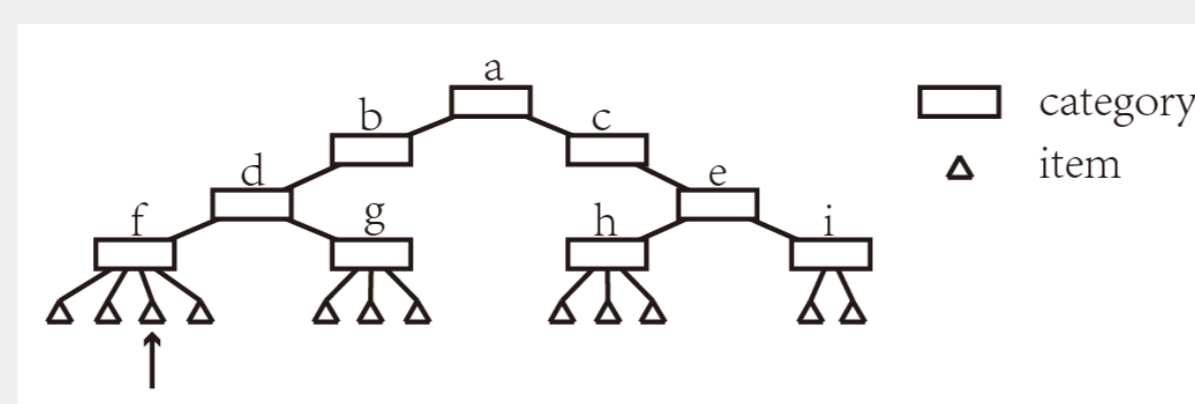  - ▷ **Fake** - they don't login the platform directly, and their messages are synchronized from related platforms.
-

$$user\_class(u_j) = \begin{cases} active, & act(u_j) \geq min\_activeness \\ inactive, & 0 \leq act(u_j) < min\_activeness \\ fake, & act(u_j) = 0 \end{cases}$$

acitve user   inacitve user   fake user

user   linked user   follow   fake linked

- $act(u_j) = tweet \times is\_fake(u_j)$
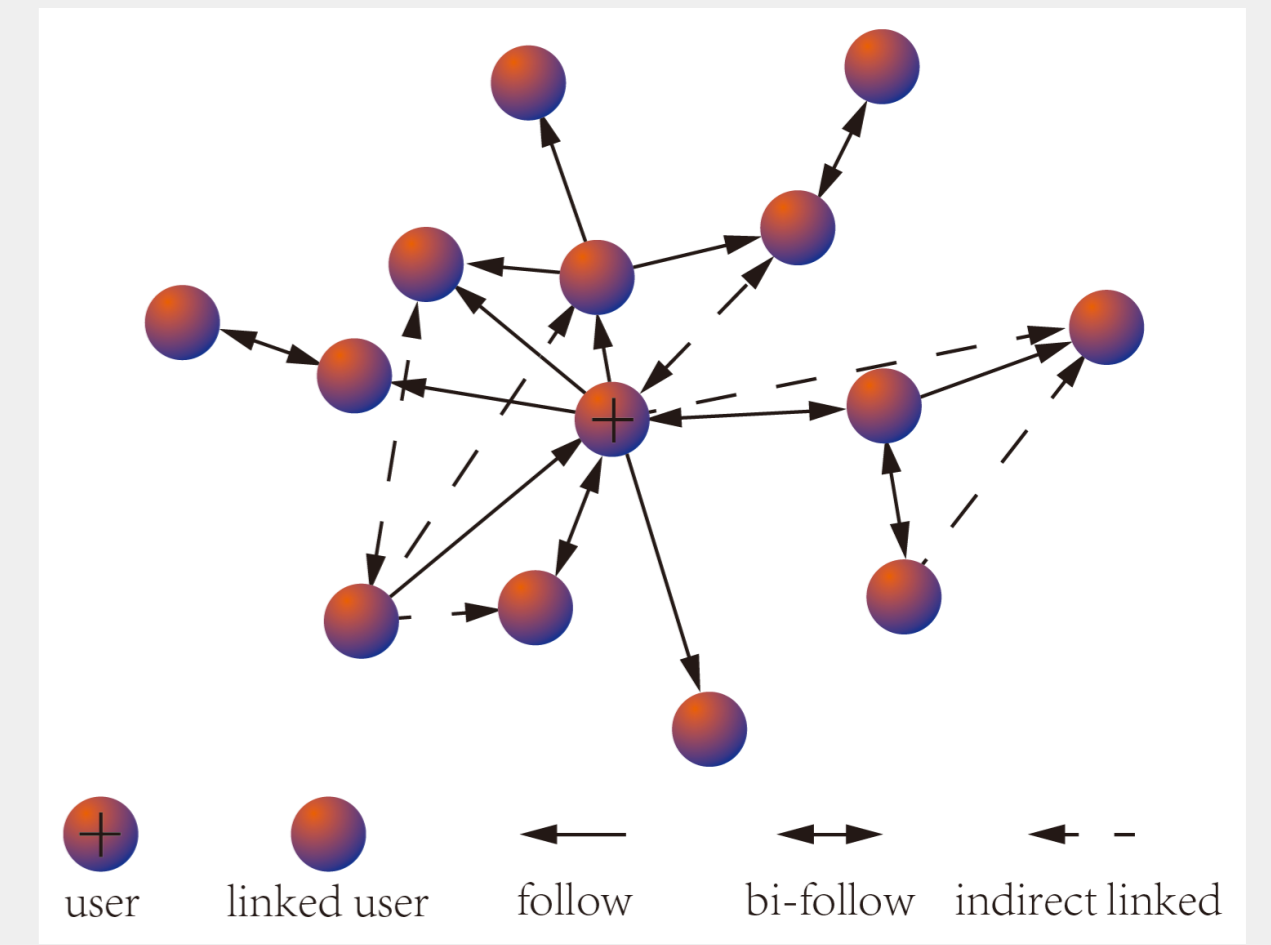- $is\_fake(u_j) = \dfrac{1 + sgn(at + retweet + comment - min\_action)}{2}$

## Item Ranking

- Items are organized in different categories of professional domains by Tencent to form a hierarchy.
  - ▷ The pointed item belongs to the category a.b.d.f.

category
item

- Item set: $I = \{i_1, i_2, ..., i_n\}$ ($i_k \in h_k$)
- Category set: $H = \{h_1, h_2, ..., h_n\}$
- Counts the number of $i_k$'s followers and return its ranking in $h_k(I)$:
  - ▷ The rank of $i_k$ in $h_k$: $hot_k = get\_hot\_rank(i_k, h_k)$
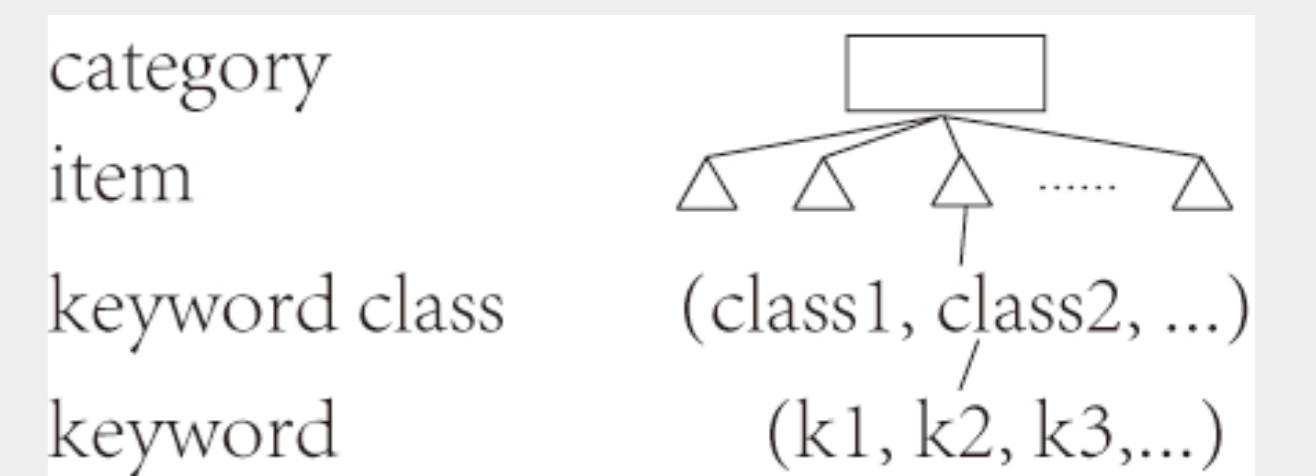  - ▷ The rank of $i_k$ in $I$: $HOT_k = GET\_HOT\_RANK(i_k)$

## (Potential)Interest Extraction

- $key\_class(u_j) = \{class_{ji}\}$
  - ▷ $\overline{W}_{ji} = \sum\limits_{k_l \in K_j \cap class_{ji}} w_l$
- $potential\_key(u_j) = \{class_{ji}\} = \bigcup\limits_{u_k \in related\_users(u_j)} key\_class(u_k)$
  - ▷ $\widetilde{W}_{ji} = \sum\limits_{\substack{u_k \in related\_users(u_j), \\ class_{kl_k} class_{ji}}} \overline{W}_{kl_k} fami(u_j, u_k)$
  - ▷ $fami(u_j, u_k) = \omega_1 f(at) + \omega_2 f(retweet) + \omega_3 f(comment)$
- $interests(u_j) = key\_class(u_j) \cap potential\_key(u_j) = \{class_{jl}\}$
  - ▷

$$W_{jl} = \begin{cases} \overline{W}_{jl}, & class_{jl} \in key\_class(u_j) \\ \widetilde{W}_{jl}, & class_{jl} \in potential\_key(u_j) \\ \frac{1}{2}(\overline{W}_{jl} + \widetilde{W}_{jl}), & class_{jl} \text{ in both sets} \end{cases}$$

user   linked user   follow   bi-follow   indirect linked

## Item Recommendation(Grading Process)

- $KH(h_k) = \{class_{kj}\}$
  - ▷ $\hat{W}_{kp} = average(\overline{W}_{jl_j})$

    $class_{jl_j} = class_{kp} \in key\_class(i_j)$
- $fond(u_j, h_k) = g(class\_weight(u_j) \cdot class\_weight(h_k), 100)$
- $grade(u_j, i_k) = 2fond(u_j, h_k)(\alpha_1 hot_k + \alpha_2 sim(u_j, i_k)) - 1$
  - ▷ $sim(u_j, i_k) = n(|class\_weight(u_j) - class\_weight(i_k)|)$
  - ▷ $\alpha_1 + \alpha_2 = 1, \alpha_i \geq 0$ (obtained in the training process)
- $n(x)$ and $g(x, y)$ are the normalization functions.
- The top-3 items are picked out for recommendation.

| category | |
| item | |
| keyword class | (class1, class2, ...) |
| keyword | (k1, k2, k3,...) |

user → keyword analysis → ? → KH mapping → category (hierarchy) → similarity & ranking → item

interest (keyword class)

## Results: Training Results of the Parameters

- $\omega_i = \frac{1}{3}$
- $\alpha_1$ reflects the inclination of accepting popular items.

| class | user | followee | interaction | keyword | $\alpha_1$ |
|---|---|---|---|---|---|
| active | 3919 | 46 | 87 | 10 | 0.33 |
| inactive | 1194 | 27 | 42 | 8 | 0.18 |
| fake | 825 | 18 | 2 | 5 | / |

## Result: Evaluation of the Performance

- Evaluation metric: average precision (which KDD Cup's organizers adopted):

$$AP@3(u_j) = \sum_{i=1}^{3} p(i)\Delta r(i)$$

  - ▷ $p(i)$ is the precision of the $i^{th}$ recommended item,
  - ▷ $\Delta r(i)$ is the change in the recall from $i - 1$ to $i$.

| | active | inactive | fake | total |
|---|---|---|---|---|
| **MAP@3** | 0.41066 | 0.46879 | 0.33606 | 0.41198 |

| $u_j$ | user class | item | accepted | **AP@3** |
|---|---|---|---|---|
| 2071402 | active | 1606902 | 1606902 | 0.83 |
| | | 1760350 | 1774452 | |
| | | 1774452 | | |
| 942226 | inactive | 1606902 | 1606902 | 1.00 |
| | | 1606609 | | |
| | | 1774452 | | |
| 193889 | fake | 1760642 | 1774862 | 0.33 |
| | | 1774684 | | |
| | | 1774862 | | |

## About

- Yingzhen Li
  Senior Year Undergraduate
  Department of Mathematics
  liyzhen2@mail2.sysu.edu.cn
  http://www.yingzhenli.net/

- Ye Zhang
  Senior Year Undergraduate
  Department of Mathematics
  zhangye5@mail.sysu.edu.cn
  http://dantepy.yslsg.org/