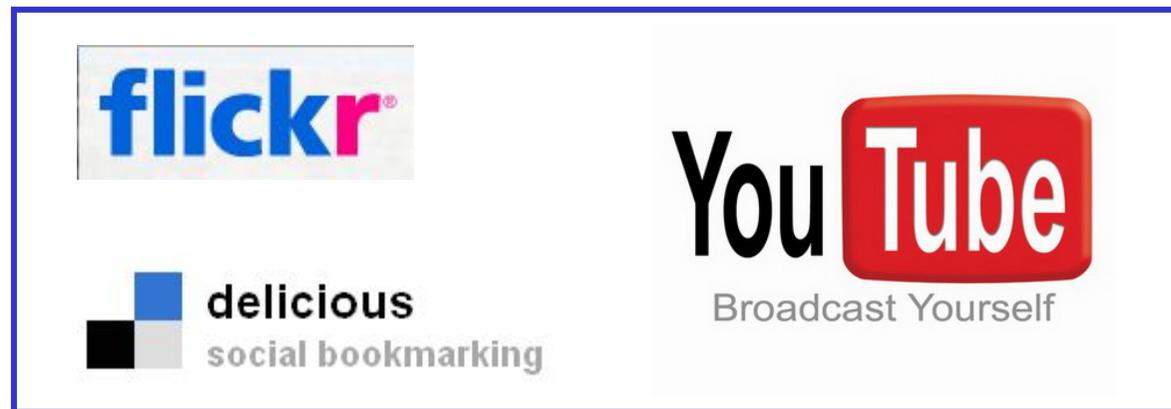
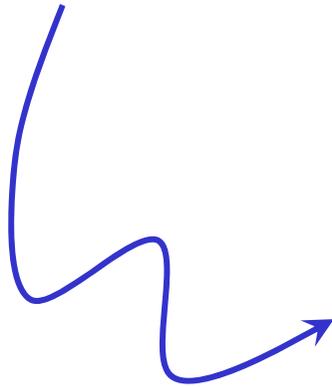


SpamClean: Towards Spam-free Tagging Systems

Ennan Zhai, Huiping Sun, Sihan Qing, Zhong Chen
Peking University

The Problem: Tag Spam

First of all, I have a question that what is the problem of *tag spam* in the current popular tagging systems?



The Problem: Tag Spam

Figure depicts the results for searching tag “iphone” in MyWeb’s site.

When we click this link, we will find the following page

Tag: **iphone**

Bookmarks 1 - 20 of about 3,320

Viewing Tag: **iphone** [Clear \[x\]](#)

> Dig Deeper: [apple](#), [ipod](#), [podcasting](#), [apple blog](#), [apple information](#), [get a iphone](#), [iphone auction](#), [iphone games](#), [iphone linux](#), [iphone rumors](#), [iphone superstore](#), [ipod auction](#), [ipod business](#), [ipod downloads](#), [ipod information](#), [ipod nano](#), [ipod player review](#), [ipod search](#), [apple business](#), [apple mac ipod](#) [\[view all\]](#)

[Apple May Shoot Own iPod Soaps - Smart House](#) 1 save Save
Tags: [iphone](#), [ipod](#), [podcasting](#), [shuffle](#)
Shared by: [Mike j](#) on 9/3/2007 at 6:30 PM - [Details](#)

[iPhone Nano? Not Likely](#) 1 save Save
Tags: [apple blog](#), [apple information](#), [apple mac](#)
Shared by: [Leland](#) on 9/3/2007 at 6:30 PM - [Details](#)

[GPS Robot Boats to Race Across Atlantic](#) 1 save Save
Tags: [apple](#), [apple blog](#), [apple business](#), [apple](#)
Shared by: [Leland](#) on 9/3/2007 at 3:40 PM - [Details](#)

[Apple May Introduce New iPod on Wednesday - Glas...](#) 1 save Save
Tags: [iphone](#), [ipod](#), [podcasting](#), [shuffle](#)
Shared by: [Mike j](#) on 9/3/2007 at 3:15 PM - [Details](#)

[Apple iPhone's Disassembling: Smoke, Sparks ...](#) 1 save Save
... begging all was pretty cool - all necessary software is copied and installed, and we have started iPhone's disassembling process. F... [more...](#)
Tags: [apple](#), [hacks](#), [iphone](#), [ipod](#), [mac](#),
Shared by: [askripko](#) on 9/3/2007 at 2:42 PM - [Details](#)

[Counterfeit Chocolate from China Has Worms](#) 1 save Save
Tags: [apple](#), [apple blog](#),
Shared by: [Leland](#) on 9/3/2007 at 2:23 PM - [Details](#)

The Problem: Tag Spam

Clicking the link leads to the site depicted in this Figure, which is not related to iphones.

We can also observe that this site has been assigned many other popular but irrelevant tags.



Small **ROBOTIC BOATS** from all by the world are set to race each other next year across the Atlantic, some 4,000 miles from Brittany, France, to the Caribbean. The race, called Microtransat 2008, was conceived by a computer scientist at the University of Wales, Aberystwyth. Entry boats must be "fully autonomous" (they can use GPS), self-sufficient in terms of energy (via solar panels) and no longer than 13 feet.

Original post by *Mike*

[iPod Copying Software.](#)

Copy all your iPod content back into iTunes. PC or Mac. Try it free

Ads by Google

everything ipod ipod photos ipod link iphone hack ipod work iphone news ipod video 60 gig iphone price drop iphone business iphone rumors buy ipod iphone linux itunes podcasting iphone search iphone link iphone user ipod bose mac ipod auction apple the apple store ipod ibook iphone ipod nano ipod stories songs apple reveal apple information apple work ipod update apple company iphone deals itunes download apple superstore apple search ipod player review iphone software apple mac ipod iphone information macworld ipod downloads itunes search ipod generation ipod stereo iphone games iphone wifi iphone blog ipod information

The Problem: Tag Spam

That is the problem of tag spam!

Definition on Tag Spam: The erroneous or misleading tags that are generated by some malicious users to confuse the normal users in the systems [Heymann, IEEE Internet Computing'07].



Existing Solutions

Detection-based:

[Krause, AIRWeb'08],
[Kyriakopoulou, RSDC'08],
[Gkanogiannis, RSDC'08];

Demotion-based:

[Koutrika, AIRWeb'07],
[Heymann, IEEE Internet
Computing'07];

Interface-based:



CAPTCHAs



Overview of SpamClean

SpamClean encompasses two key mechanisms:

- **Experience Mechanism;**
- **Socially-enhanced Mechanism.**



Experience Mechanism

After a tag (e.g., Sea) search, client Alice computes the experiences of Alice with respect to other users in the system based on cosine technique. Then, Alice ranks each search result based on the average of experiences of Alice with respect to all the **annotators** of this result.



**Annotator is the user who
annotates the result with tag Sea**

How to compute experience ?

$$E_{A,B} = \frac{\sum_{r_j \in R} \left(\sum_{t_i \in C_{r_j}} |N(t_i, r_j)| \right)^2}{\sqrt{\sum_{r_j \in R} \left(\sum_{t_i \in T_A(r_j)} |N(t_i, r_j)| \right)^2} \sqrt{\sum_{r_j \in R} \left(\sum_{t_i \in T_B(r_j)} |N(t_i, r_j)| \right)^2}} \quad (1)$$

where

- R : The set of resources shared by A and B in common.
- r_j : The j th resource of the common resource set R .
- C_{r_j} : The set of the tags annotated by A and B in common to the resource r_j .
- t_i : The i th tag of the tag set.
- $T_x(r_j)$: The set of tags annotated by the user x to the resource r_j .
- $N(t_i, r_j)$: The set of annotation that annotated r_j with t_i .
- $|N(t_i, r_j)|$: The size of $N(t_i, r_j)$.

Experience Mechanism

Incentive: Users who do not annotate correctly and actively will find the quality of the estimate they compute noticeably degraded. Indeed, a client can benefit from SpamClean by annotating honestly.

Socially-enhanced Mechanism

We utilizes friend-relationships, the social nature of tagging systems, to enhance SpamClean, since the client's social friends can share their previous experiences and help improve both the performance and convergence of SpamClean.

Socially-enhanced Mechanism

After Alice ranks the search results

- Alice finds out all the resources which have the ranking scores lower than 0.5 in the search results;
- Alice demands his friends to compute ranking scores for those resources;
- If more than half of the friends return the ranking scores of the above resource higher than 0.5, SpamClean re-locates the position of the resource based on the new ranking score by computing the average of these friends' returned scores higher than 0.5. Otherwise, the client maintains the original rank unchanged.

Socially-enhanced Mechanism

- **The example about socially-enhanced mechanism of SpamClean please see our paper.**
- **The practical issue on unreliable friends please see our paper.**

Evaluation

Search Models:

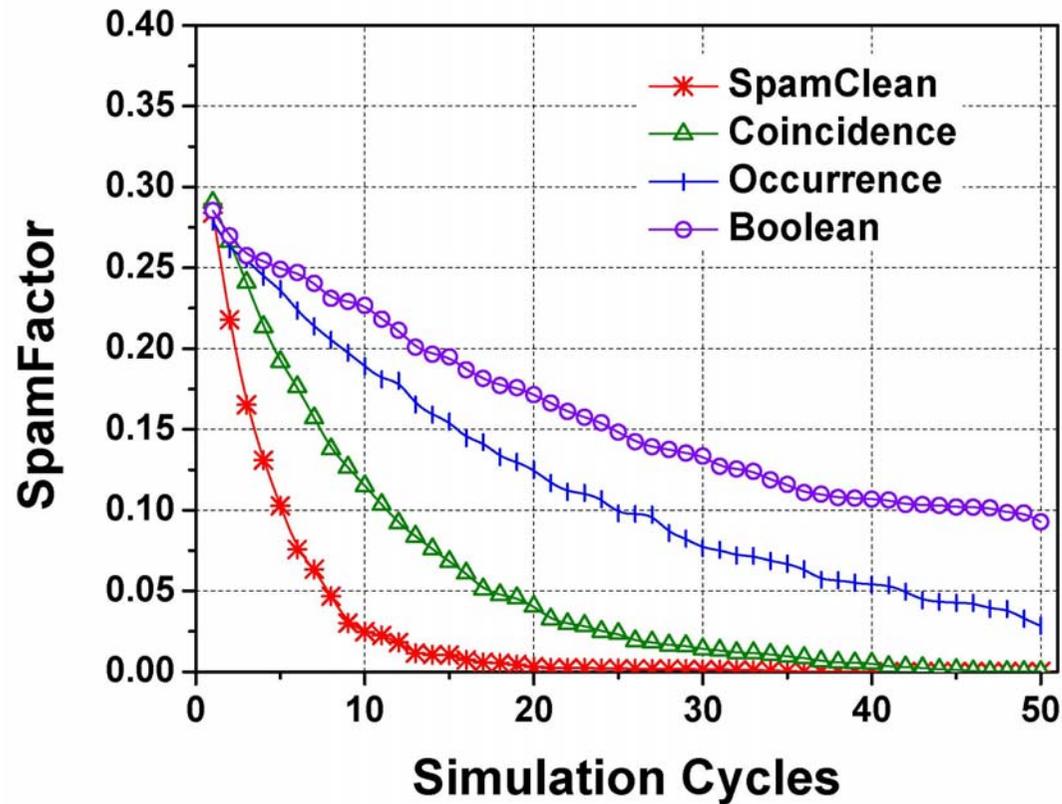
- **Boolean Model:** randomly ranks the results;
- **Occurrence Model:** returns the top number results;
- **Coincidence Model:** assigns each user a global score which is the sum of the same annotations between this user and the other users in the system; then the system ranks the search results based on the average of all the annotators' scores of each result.

Evaluation

Threat Models:

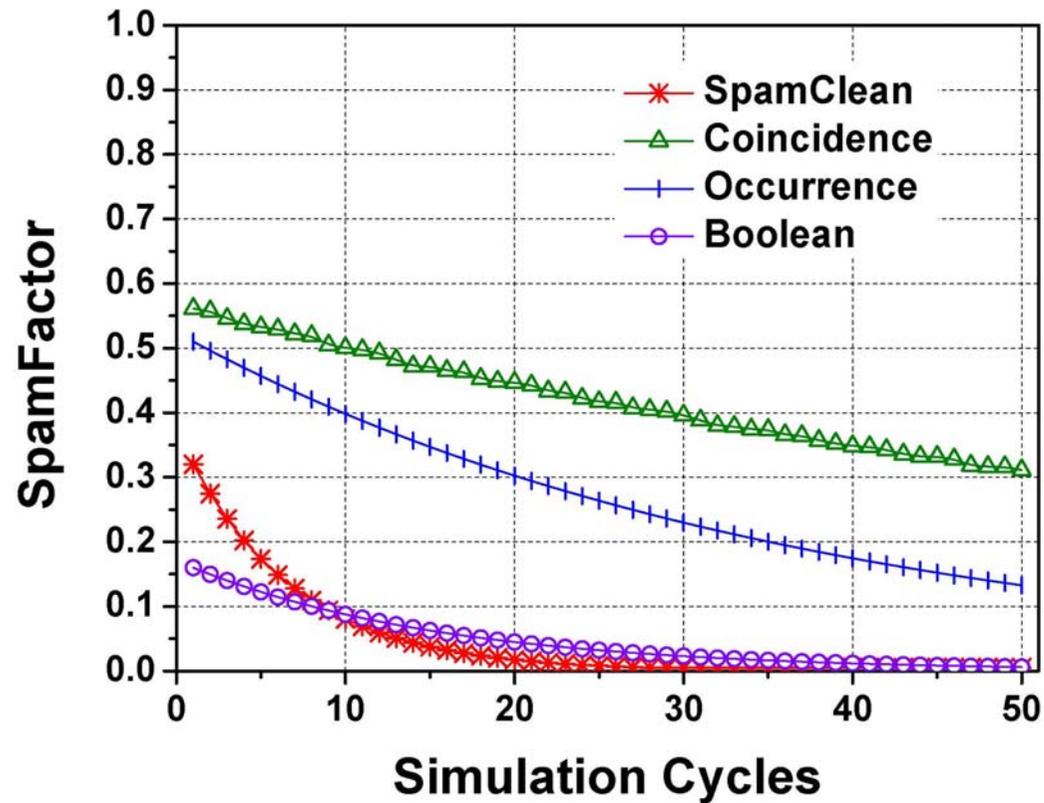
- **Random Attacks:** randomly annotate misleading tags to the resources in the system;
- **Collusive Attacks:** collusively annotate resources with the same misleading tags;
- **Tricky Attacks:** annotate resources with both correct and misleading tags. This attack could make some anti-spam scheme unusable.

Evaluation



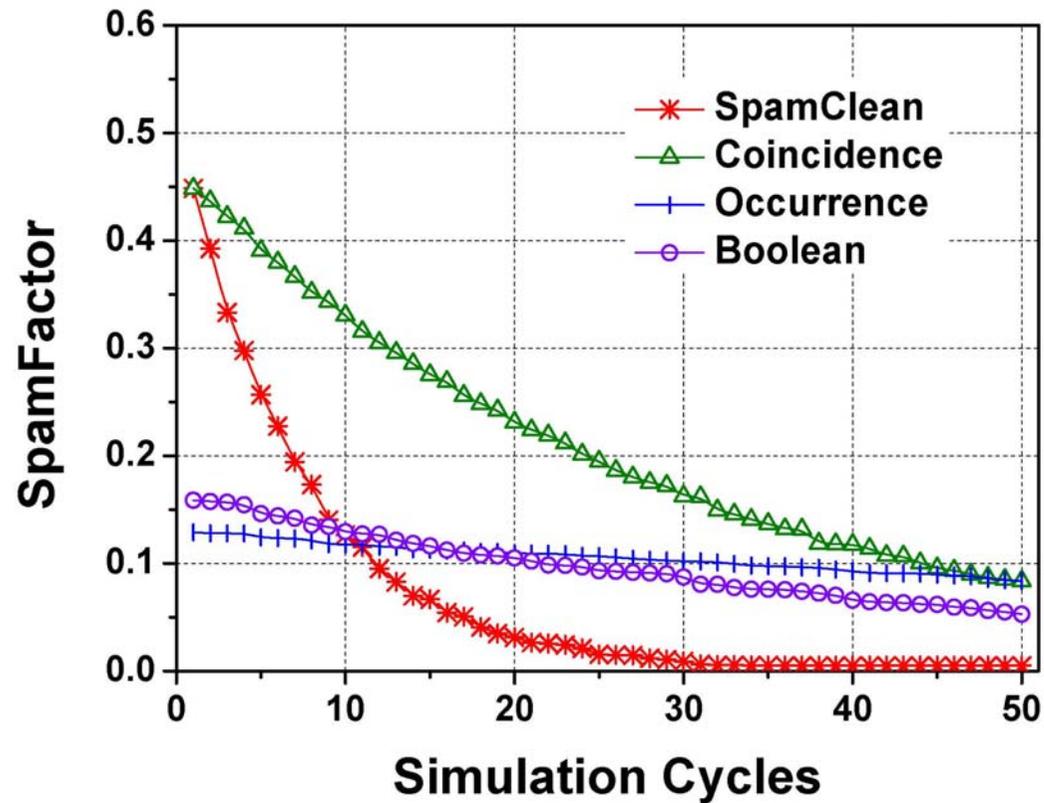
The impact of random attack under 20% random attackers (more detail see paper)

Evaluation



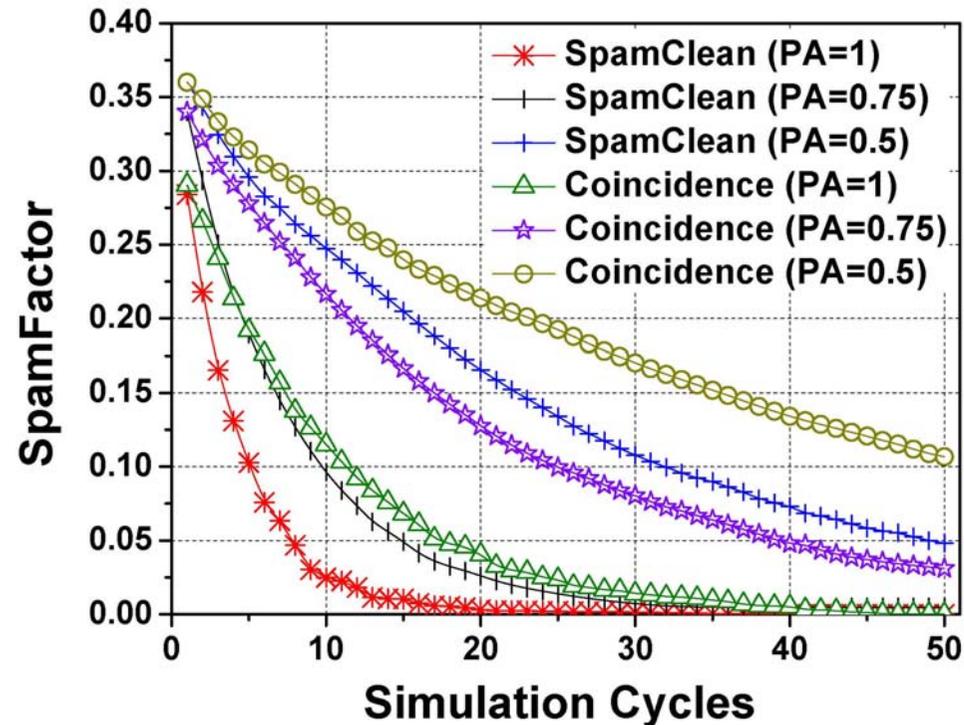
The impact of collusive attack under 20% collusive attackers (more detail see paper)

Evaluation



The impact of tricky attack under 20% tricky attackers (more detail see paper)

Evaluation



The impact of cooperation of users, PA means the probability that users share their tags (more detail see paper)

Conclusions and Discussions

SpamClean is a novel social experience-based scheme towards spam-free and personalized tag search results in the tagging systems.

Discussions:

How does SpamClean defense Sybil attacks?

How does SpamClean defense DoS attacks?

(The detail please see our paper)

Q & A

