

AIRWeb'07

Third International Workshop on
Adversarial Information Retrieval on the Web

Web Spam Challenge Track I

<http://webspam.lip6.fr/>

Supported by the EU PASCAL Network of Excellence Challenge Program



Web Spam Challenge Overview

- Shared dataset
 - Page links and contents (77M in 11,000 hosts)
 - 500 GB full version, 11 GB summary version (comp.)
- Training labels (~5,500 hosts)
 - And sets of pre-computed features (optional)
- Testing labels (~2,300 hosts)
 - Labeled by participants



Dataset: WEBSpAM-UK2006

- 1) Crawling of base data**
- 2) Design of guidelines and interface**
- 3) Labeling by volunteers**



Crawling of base data

1) Crawling of base data

- Laboratory of Algorithms on the Web, University of Milan: P. Boldi, M. Santini and S. Vigna, using UbiCrawler
- May 2006, .UK domain

2) Design of guidelines and interface

3) Labeling by volunteers



Guidelines and interface

1) Crawling of base data

2) Design of guidelines and interface

- Available at <http://www.yr-bcn.es/webspam/>
- Examples and guidelines

3) Labeling by volunteers



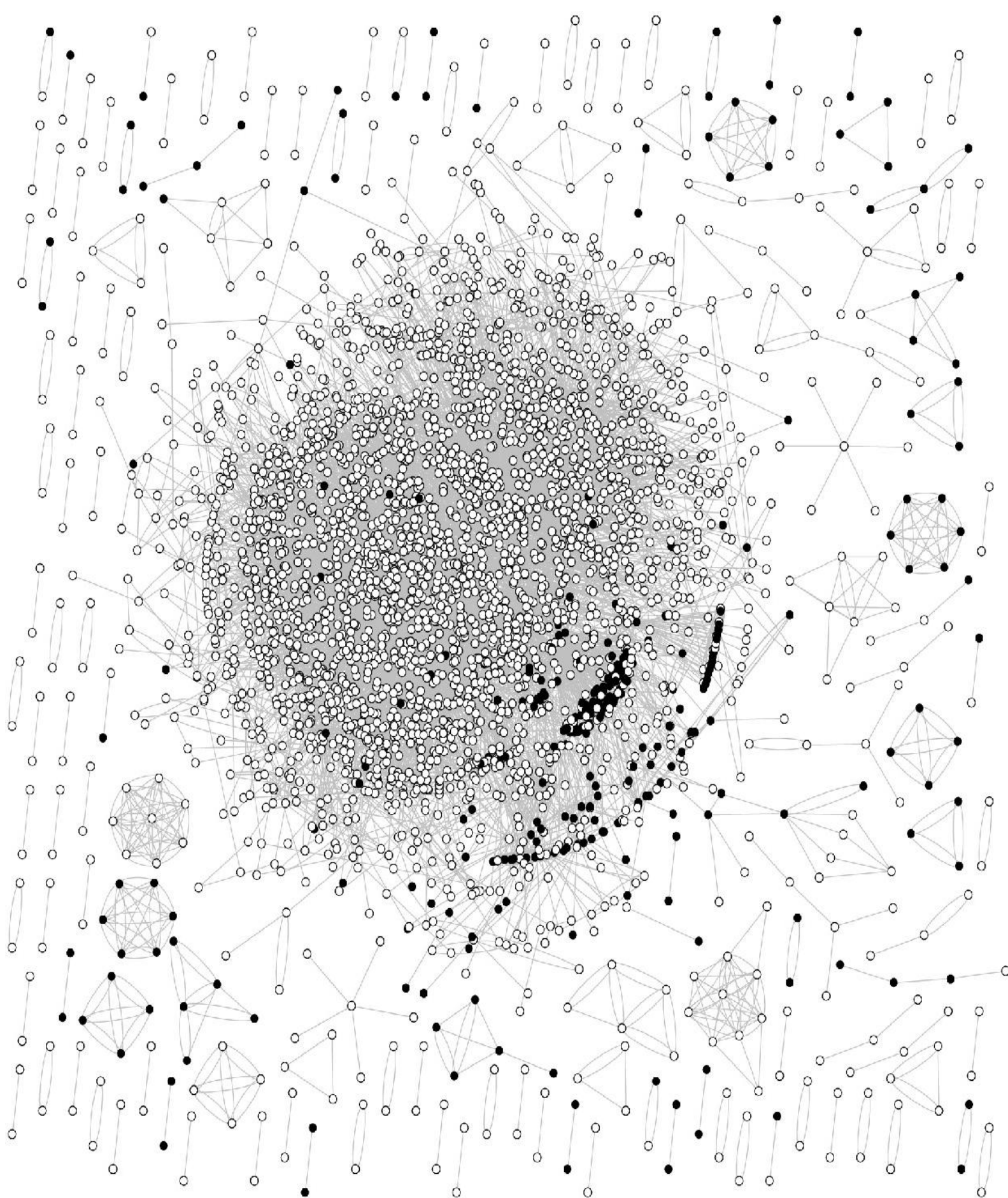
Labeling by volunteers

- 1) Crawling of base data
- 2) Design of guidelines and interface
- 3) Labeling by volunteers
 - Host-level labeling as **Normal**/**Borderline**/**Spam**
 - 20 volunteers, paired at random for each host
 - In general their labels agree



>200 man-hours later ...

- 6,500 host labels
 - 62% normal
 - 11% borderline
 - 22% spam
 - 5% can't classify
- More details:
 - C. Castillo, D. Donato, L. Becchetti, P. Boldi, S. Leonardi, M. Santini, S. Vigna: “A Reference Collection for Web Spam”, SIGIR Forum Dec/2006.





Base features (optional)

- Link- and content-based features proposed in the literature, mostly from:
 - Z. Gyöngyi, H. García-Molina and J. Pedersen: “Computing Web Spam with TrustRank”, VLDB 2004
 - A. Ntoulas, M. Najork, M. Manasse and D. Fetterly: “Detecting Web Spam Pages through Content Analysis”, WWW 2006
 - C. Castillo, A. Gionis, D. Donato, V. Murdock and F. Silvestri: “Know your Neighbors: Web Spam Detection Using the Web Topology”, SIGIR 2007



Rules (webspam.lip6.fr)

- Submissions (max. 2 per team)
 - Spam/Nospam labels and/or “spamnicity” predictions for each host
- Metrics (computed using KDD-CUP’s “perf”)
 - F1 measure
 - Combines precision and true positives for a specific threshold
 - Area under ROC curve (for “spamnicity” only)
 - ROC curve = false positives vs. true positives for different thresholds



Rules (webspam.lip6.fr)

- Test set labeled by participants
 - Considers all hosts with at least two labels from {normal, borderline, spam}
 - Ground-truth spamicity is average label: normal=0.0, borderline=0.5, spam=1.0
- Remove hosts with ground-truth spamicity=0.5 and evaluate using F1 and AUC (there can be more than one winner)



Rules (webspam.lip6.fr)

- Base evaluation is removing hosts with spamicity=0.5 in the test set
- Two scenarios for these hosts
 - Include them as normal, then test
 - Include them as spam, then test
- Significant difference (for ranking submissions): higher than the difference obtained using these scenarios; otherwise, consider it a tie



Participating Teams

- Tony Abou-Assaleh and Tapajyoti Das
- András A. Benczúr, István Bíró, Károly Csalogány, Miklós Kurucz and Tamás Sarlós
- Gordon Cormack
- Pascal Filoche, Tanguy Urvoy, Chauveau Emmanuel and Lavergne Thomas
- Dennis Fetterly, Steve Chien, Marc Najork, Mark Manasse and Alexandros Ntoulas
- Guanggang Geng, Chunheng Wang, Xiaobo Jin, Qiudan Li and Lei Xu