# Spotting Fake Reviewer Groups in Consumer Reviews

**Arjun Mukherjee**
Department of Computer Science
University of Illinois at Chicago

**Bing Liu**
Department of Computer Science
University of Illinois at Chicago

**Natalie Glance**
Google Inc.

# Fake/Spam/Biased Reviews

- Online **reviews** play an important role in **Decision making** process

- Review Spamming, Motivations:
  - Fame
  - Financial gain

- To **promote** or **demote** other products

- Opinion spamming is now a business
  - people get paid to write fake reviews
  - So they **write many reviews about many products**, such collective behavior can give them away

- This study: Focusing on **Spammer Groups** instead of **individual** reviews/reviewers

2

# Dataset

- Created a Labeled dataset for group opinion spam
- Refers to prior studies and argue that in absence of labeled data, best option is to create one based on human expert
- Dataset Stats:
  - Amazon Dataset from 2006 (updated on 2010)
  - Only manufactured products (53K reviewer with 110K reviews on 39K products)
  - Attributes: Title , content , star rating , posting date and helpful feedbacks

- 1) **Candidate groups**: Using Frequent Itemset Mining (FIM)
  - On list of reviewer ids per product
  - All groups with min_sup =3 and 2 items
    - Groups with at least 2 reviewers who worked at least on 3 products

# Dataset

- 2) Opinion Spam signals:
  - Provided the list spam signals from prior research and websites:
    - (i) having zero caveats, (ii) full of empty adjectives, (iii) purely glowing praises with no downsides, (iv) being left within a short period of time of each other
  - Access to review Database

- Judges: employees from Rediff shopping (4) and eBay.in (4)
  - Spent 8 weeks to label 2431 groups.
  - **Spamicity Rate (SR)**
    - 1: spam , 0.5 borderline , 0: no-spam
      - 8 votes, average of all would be the SR for review.
      - average of reviews SR => group spamicity

# SPAMMING BEHAVIOR INDICATORS

For modeling or learning, a set of effective spam indicators or features is needed

# Spamming behavior indicators

1. Group spam behavior indicators
   - Group time window (GTW)

$$GTW(g) = \max_{p \in P_g}(GTW_P(g, p)),$$

$$GTW_P(g, p) = \begin{cases} 0 & \text{if } L(g, p) - F(g, p) > \tau \\ 1 - \dfrac{L(g, p) - F(g, p)}{\tau} & \text{otherwise} \end{cases},$$

   - Group Deviation (GD)

$$GD(g) = \max_{p \in P_g}(D(g, p)),$$

$$D(g, p) = \frac{|r_{p,g} - \bar{r}_{p,g}|}{4},$$

# Spamming behavior indicators

- Group spam behavior indicators
  - Group Content Similarity (GCS)

$$GCS(g) = \max_{p \in P_g}(CS_G(g, p)),$$

$$CS_G(g, p) = \underset{m_i, m_j \in g, i < j}{avg}\left(cosine(c(m_i, p), c(m_j, p))\right),$$

  - Group Member Content Similarity (GMCS)

$$GMCS(g) = \frac{\sum_{m \in g} CS_M(g, m)}{|g|},$$

$$CS_M(g, m) = \underset{p_i, p_j \in P_g, i < j}{avg}\left(cosine(c(m, p_i), c(m, p_j))\right)$$

# Spamming behavior indicators

- Group spam behavior indicators
  - Group Early Time Frame (GETF)

$$GETF(g) = \max_{p \in P_g}(GTF(g, p)),$$

$$GTF(g, p) = \begin{cases} 0 & \text{if } L(g,p) - A(p) > \beta \\ 1 - \dfrac{L(g,p) - A(p)}{\beta} & \text{otherwise} \end{cases},$$

  - Group Size Ratio (GSR)

$$GSR(g) = \underset{p \in P_g}{avg}(GSR_P(g, p)),$$

$$GSR_P(g, p) = \frac{|g|}{|M_p|},$$

# Spamming behavior indicators

- Group spam behavior indicators
  - Group Size (GS)

$$GS(g) = \frac{|g|}{\max(|g_i|)}$$

  - Group Support Count (GSUP)

$$GSUP(g) = \frac{|P_g|}{\max(|P_{g_i}|)}$$

# Spamming behavior indicators

2. **Individual Spam Behavior Indicators**
    ○ Individual Rating Deviation (IRD):

$$IRD(m, p) = \frac{|\, r_{p,m} - \bar{r}_{p,m}\,|}{4},$$

    ○ Individual Content Similarity (ICS)

$$ICS\,(m, p) = avg\,(cosine\,(c(m, p))$$

# Spamming behavior indicators

2.  **Individual Spam Behavior Indicators**
    - Individual Early Time Frame (IETF)

$$IETF(m,p) = \begin{cases} 0 & \text{if } L(m,p) - A(p) > \beta \\ 1 - \dfrac{L(m,p) - A(p)}{\beta} & \text{otherwise} \end{cases},$$

    - Individual Member Coupling in a group (IMC)

$$IMC(g,m) = \underset{p \in P_g}{avg}\left( \frac{|(T(m,p) - F(g,p)) - avg(g,m)|}{L(g,p) - F(g,p)} \right),$$

$$avg(g,m) = \frac{\sum\limits_{m_i \in G-\{m\}}(T(m_i,p) - F(g,p))}{|g|-1},$$

This behavior measures how closely a member works with the other members of the group. If a member m almost posts at the same date as other group members, then m is said to be tightly coupled with the group

# Empirical Analysis

# Statistical validation

- Spamicity threshold : 0.5 => 62% non-spam and 38% spam groups
- Feature effectiveness:
  - $$Eff(f) \equiv P(f > 0 \mid Spam) - P(f > 0 \mid Non - spam),$$

  $$P(f > 0 \mid Spam) = \frac{|\{g \mid f(g) > 0 \wedge g \in Spam\}|}{|\{g \mid g \in Spam\}|}$$

  $$P(f > 0 \mid Non - spam) = \frac{|\{g \mid f(g) > 0 \wedge g \in Non - spam\}|}{|\{g \mid g \in Non - spam\}|}$$

  - Using Fisher's exact test, it is reported that spam groups are more likely to exhibit feature.
    - null hypothesis rejected with p< 0.0001

# Behavioral Distribution

- ## Position
  - for a given cumulative percentage cp, the corresponding feature value xn for non-spam groups is less than xs for spam groups
- ## Steep initial jumps
  - very few groups obtain significant feature values
- ## Gaps
  - The separation margin refers to the relative discriminative potency



Figure 4: Behavioral Distribution. Cumulative % of spam (solid) and non-spam (dashed) groups vs. feature value

14

# MODELING RELATIONS

# MODELING RELATIONS

Better not to follow the classic approach : Classification

1. training and testing instances are not independently and identically drawn from some distribution (groups share members)
2. Group features only summarize the group behaviors (avg/sum)
   a. lead to loss of information
3. It is difficult to include the effect of Products!

So, they propose a more effective model to address the above concerns and also cover three binary relations:

Group Spam–Products
Member Spam–Products,
and Group Spam–Member Spam.

# Group Spam—Products Model

- The relation among groups and products they target.
  - (i) spam contribution to p by each group reviewing p and
  - (ii) "spamicity" of each such group

$$w_1(p_i, g_j) = \frac{1}{5}[GTW_P(g_j, p_i) + D(g_j, p_i) + GTF(g_j, p_i) + CS_G(g_j, p_i) + GSR_P(g_j, p_i)],$$

$$W_{PG} = [w_1(p_i, g_j)] \; |P| \times |G| \tag{16}$$

$W_{PG}$ denotes the corresponding contribution matrix.

$$s(p_i) = \sum_{j=1}^{|G|} w_1(p_i, g_j) s(g_j); \quad V_P = W_{PG} V_G,$$

$$s(g_j) = \sum_{i=1}^{|P|} w_1(p_i, g_j) s(p_i); \quad V_G = W_{PG}^T V_P$$

# Member Spam-Product Model

- IRD (individual rating deviation of m towards p)
- ICS (individual content similarity of reviews on p by m)
- IETF (individual early time frame of spam infliction by m towards p)

$$w_2(m_k, p_i) = \frac{1}{3}[IRD(m_k, p_i) + ICS(m_k, p_i) + IETF(m_k, p_i)],$$

$$W_{MP} = [w_2(m_k, p_i)]_{|M| \times |P|} \tag{19}$$

We sum the individual contribution of each member w2, weighted by its spamicity:

$$s(m_k) = \sum_{i=1}^{|P|} w_2(m_k, p_i) s(p_i); \quad V_M = W_{MP} V_P$$

$$s(p_i) = \sum_{k=1}^{|M|} w_2(m_k, p_i) s(m_k); \quad V_P = W_{MP}^T V_M$$

# Group Spam–Member Spam Model

- IMC (degree of m's coupling in g),
- GS (size of g with which m worked), and
- GSUP (number of products towards which m worked with g)

$$w_3(g_j, m_k) = \frac{1}{3}[IMC(g_j, m_k) + \boxed{(1 - GS(g_j))} + GSUP(g_j)],$$

$$W_{GM} = [w_3(g_j, m_k)]_{|G| \times |M|}$$

for large groups, the individual contribution of a member diminishes. Hence we use 1-GS(gj) to compute w3.

$$s(g_j) = \sum_{k=1}^{|M|} w_3(g_j, m_k) s(m_k); \quad V_G = W_{GM} V_M,$$

$$s(m_k) = \sum_{j=1}^{|G|} w_3(g_j, m_k) s(g_j); \quad V_M = W_{GM}^T V_G.$$

# GSRank: Ranking Group Spam

**Algorithm**: GSRank

Input: Weight matrices $W_{PG}$, $W_{MP}$, and $W_{GM}$
Output: Ranked list of candidate spam groups

1. Initialize $V_G^0 \leftarrow [0.5]_{|G|}$; $t \leftarrow 1$;
2. Iterate:
   i. $V_P \leftarrow W_{PG} V_G^{(t-1)}$; $V_M \leftarrow W_{MP} V_P$;
   ii. $V_G \leftarrow W_{GM} V_M$; $V_M \leftarrow W_{GM}^T V_G$;
   iii. $V_P \leftarrow W_{MP}^T V_M$; $V_G^{(t)} \leftarrow W_{PG}^T V_P$;
   iv. $V_G^{(t)} \leftarrow V_G^{(t)} / \| V_G^{(t)} \|_1$;
   until $\| V_G^{(t)} - V_G^{(t-1)} \|_\infty < \delta$
3. Output the ranked list of groups in descending order of $V_G^*$

**Complexity**: linear in the number of candidate groups discovered by FIM

$$O(t(|G|(|M|+|P|) + |M||P|))$$

# EXPERIMENTAL EVALUATION

- We first split 2431 groups:
  - The development set, D with 431 groups (randomly sampled) for parameter estimation
    - for GTW and GETF ,
      - using a greedy hill climbing search to maximize the log likelihood of the set D
      - $\tau = 2.87$
        the time interval beyond which members in a group are not likely to be working in collusion
      - $\beta = 8.86$
        denotes the time interval beyond which reviews posted are not considered to be "early" anymore

  - The validation set, V with 2000 groups for evaluation.

- All evaluation metrics averaged over 10-fold cross validation (CV)

# Ranking Experiments :: baselines

1. Using **regression**
   - The problem of ranking spammer groups can be seen as:
     - optimizing the spamicity of each group as a regression target
   - the support vector regression (SVR) system in SVMlight is used

2. and **Learning to Rank**
   - we treat each feature f as a ranking function
   - The rank produced by each feature is based on a certain spamicity dimension
   - None of the ranks may be optimal. A learning to rank method basically learns an optimal ranking function using the combination of f1...f8
   - Each group is vectorized with (represented with a vector of) the 8 group spam features

# Ranking Experiments, cont.

- Normalized Discounted Cumulative Gain (NDCG) as our evaluation metric
- GSRank performs the best at all top rank positions except at the bottom,
  - which are unimportant because they are most likely to be non-spam (since in each fold of cross validation, the test set has only 200 groups and out of the 200 there are at most 38% spam groups)

# Ranking Experiments, cont.

- we also experimented with the following baselines:
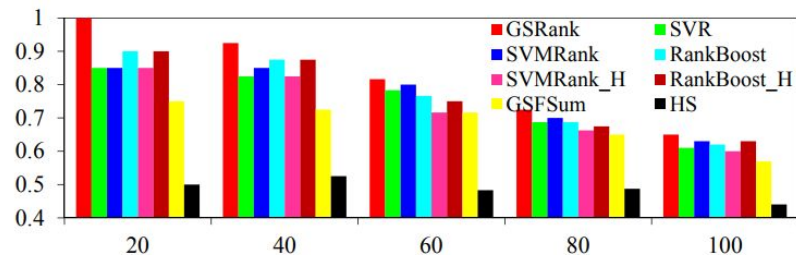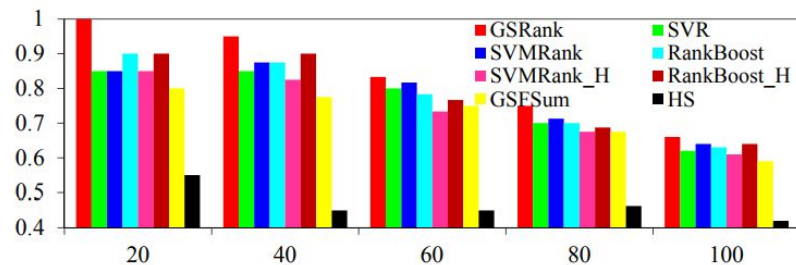  - Group Spam Feature Sum (GSFSum)
    - to rank the groups in descending order of the sum of all feature values
  - Helpfulness Score (HS)
    - HS uses the mean helpfulness score (percentage of people who found a review helpful) of reviews of each group to rank groups in ascending order of the scores
  - Heuristic training rankings (H)
    - three heuristic rankings using feature mixtures

$$h_1(g) : G \rightarrow \mathbf{R}^+, h_1(g) = GCS(g) + GMCS(g)$$
$$h_2(g) : G \rightarrow \mathbf{R}^+, h_2(g) = GS(g) + GSUP(g) + GTW(g)$$
$$h_3(g) : G \rightarrow \mathbf{R}^+, h_3(g) = GSR(g) + GETF(g) + GD(g)$$



(a) The spamicity threshold of $\xi = 0.5$

(b) The spamicity threshold of $\xi = 0.7$

**Figure 6: Precision @ $n$ = 20, 40, 60, 80, 100 rank positions.**
All the improvements of GSRank over other methods are statistically significant at the confidence level of 95% based on paired $t$-test.

|  | $\xi = 0.5$ | $\xi = 0.7$ |
|---|---|---|
| Spam | 38% | 29% |
| Non-spam | 62% | 71% |

# Classification

- If a spamicity threshold is applied to decide spam and non-spam groups, supervised classification can also be applied
- features that we consider in learning:
  - Group Spam Features (GSF)
  - Individual Spammer Features (ISF)
  - Linguistic Features of reviews (LF)
        (word and POS (part-of-speech) n-gram features )

- AUC (Area Under the ROC Curve) is employed for classification evaluation

# Classification

| Feature Settings | SVM | LR | SVR | SVM Rank | Rank Boost | SVM Rank_H | Rank Boost_H | GS Rank |
|---|---|---|---|---|---|---|---|---|
| GSF | 0.81 | 0.77 | 0.83 | 0.83 | 0.85 | 0.81 | 0.83 | **0.93** |
| ISF | 0.67 | 0.67 | 0.71 | 0.70 | 0.74 | 0.68 | 0.72 | |
| LF | 0.65 | 0.62 | 0.63 | 0.67 | 0.72 | 0.64 | 0.71 | |
| GSF + ISF + LF | 0.84 | 0.81 | 0.85 | 0.84 | **0.86** | 0.83 | 0.85 | |

(a)    The spamicity threshold of $\xi = 0.5$

| Feature Settings | SVM | LR | SVR | SVM Rank | Rank Boost | SVM Rank_H | Rank Boost_H | GS Rank |
|---|---|---|---|---|---|---|---|---|
| GSF | 0.83 | 0.79 | 0.84 | 0.85 | 0.87 | 0.83 | 0.85 | **0.95** |
| ISF | 0.68 | 0.68 | 0.73 | 0.71 | 0.75 | 0.70 | 0.74 | |
| LF | 0.66 | 0.62 | 0.67 | 0.69 | 0.74 | 0.68 | 0.73 | |
| GSF + ISF + LF | 0.86 | 0.83 | 0.86 | 0.86 | **0.88** | 0.84 | 0.86 | |

(b)    The spamicity threshold of $\xi = 0.7$

**Table 2: AUC results of different algorithms and feature sets.**

# The End