# Learning Influence Propagation on Personal Blogs

Il-Chul Moon, Dongwoo Kim, Yohan Jo, Alice H. Oh

yohan.jo@kaist.ac.kr

KAIST
Daejeon, Korea

SMS LAB
SYSTEMS MODELING SIMULATION

u&i lab
users and information

KAIST

# In this presentation, I will talk about

- how to measure a blog's influence to readers

- how to predict the potential influence of a new blog post well

# Blog Influence

# Blog Influence



**Hits Measured by WordPress**
- 3,783,718 hits since April 1, 2006

2,500 visits / day

# Blog Influence

# Blog Influence

Power of stimulating the readers to express their thoughts in response

The influence is reflected,

- quantitatively, by the **network position**
  e.g. number of people influenced

- qualitatively, in the **content**
  e.g. similarity of the topics

# Content Analysis

Apply the author-topic model to the blogs [Rosen-Zvi, 2004]

| Topic | Highly Related Words |
|:---:|:---:|
| 1 | gold, wave, market, term, short, cycle |
| 2 | money, tax, government, fund, pay, financial |
| ... | |
| 50 | school, student, university, education, class |

| Blog | Topic | | | | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | 1 | 2 | 3 | 4 | ... | 50 |
| 1 | 0.011 | 0.324 | 0.003 | 0.008 | | 0.003 |
| 2 | 0.250 | 0.007 | 0.012 | 0.009 | | 0.011 |
| ... | | | | | | |
| 4,165 | 0.009 | 0.015 | 0.003 | 0.010 | | 0.363 |

# Network Analysis

# Network Analysis



In-degree Centrality = 3

# Network Analysis



Out-degree Centrality = 1

# Network Analysis



Total-degree Centrality = 4

# Network Analysis



Betweenness Centrality = 16

# Network Analysis



Clustering Coefficient = 1/3

# "Influence Size"

Takes into account

- how many readers write posts in response

- how similar their topics are

$$S_i = \{B_j : \ B_j \text{ can reach } B_i \text{ by following links}\}$$

$$\text{InfluenceSize}(B_i) = \sum_{B_j \in S_i} \text{TopicSimilarity}(B_i, B_j)$$

$$= \sum_{B_j \in S_i} \frac{\mathbf{t}_i \cdot \mathbf{t}_j}{\|\mathbf{t}_i\|\|\mathbf{t}_j\|}$$

# Existing Measures of Blog Influence

- Number of Comments

  - Quantitative

- Digg Score (Digg.com)

  - Quantitative, partly qualitative

"How can we predict the potential influence of a new blog post?"

# Prediction of Blog Influence



The Post's Content

The Blog's Network

# Experiments

1. Content information and network information capture different aspects of blog influence

2. It is important to use both content information and network information for finding influential blogs

# Dataset

- Selected from TREC Blog08

  http://ir.dcs.gla.ac.uk/test_collections/blogs08info.html

| # of Blogs | # of Posts | # of Unique Words | Average # of Words/Post |
|:---:|:---:|:---:|:---:|
| 4,165 | 72,143 | 53,257 | 225.24 |

- ✓ Blogspot, LiveJournal
- ✓ Contains >50 words
- ✓ Has at least one link to another post
- ✓ Written in 2008
- ✓ Written in English

# Experiment I

" Do content information and network information play different roles in predicting blog influence? "

# Design

- **Task**

  Classify the blogs into three groups
  (non-influential, influential, and very-influential)

- **Ground Truth**

  Sort the blogs in the order of *influence size* and group them into three

# Classification

- **Linear SVM**

- **Features**

  - Content: Topic proportions (50)

  - Network: In-degree, out-degree, total-degree, betweenness, clustering coefficient (5)

- **Training**
  Train a classifier using 30% of the data blogs

- **Testing**
  Test with the rest of the data

# Feature Importance



| Group 1 | Group 2 | Group 3 |
|---|---|---|
| OutDegree | OutDegree | Topic 47 |
| TotalDegree | TotalDegree | Topic 14 |
| Betweenness | Topic 9 | Topic 15 |
| Clustering Coefficient | Topic 35 | Topic 43 |
| Topic 8 | Topic 13 | Topic 41 |
| Topic 28 | Topic 0 | Topic 38 |
| Topic 19 | Topic 24 | Topic 42 |
| InDegree | Topic 27 | Topic 24 |
| Topic 43 | Topic 16 | Topic 12 |
| Topic 23 | Topic 19 | Topic 22 |

## Most Important Features

# of
Instances

Lower Threshold
of Influence Size

Upper Threshold
of Influence Size

# Experiment II

" If we use both content information and network information, can we find potentially influential blogs better than when using only one of them? "

# Design

- **Task**

  Classify each blog whether it belongs to the top 10% influential blogs or not

- **Ground Truth**

  Top 10% influential blogs are obtained with regard to each influence measure (influence size, number of comments, and Digg score)
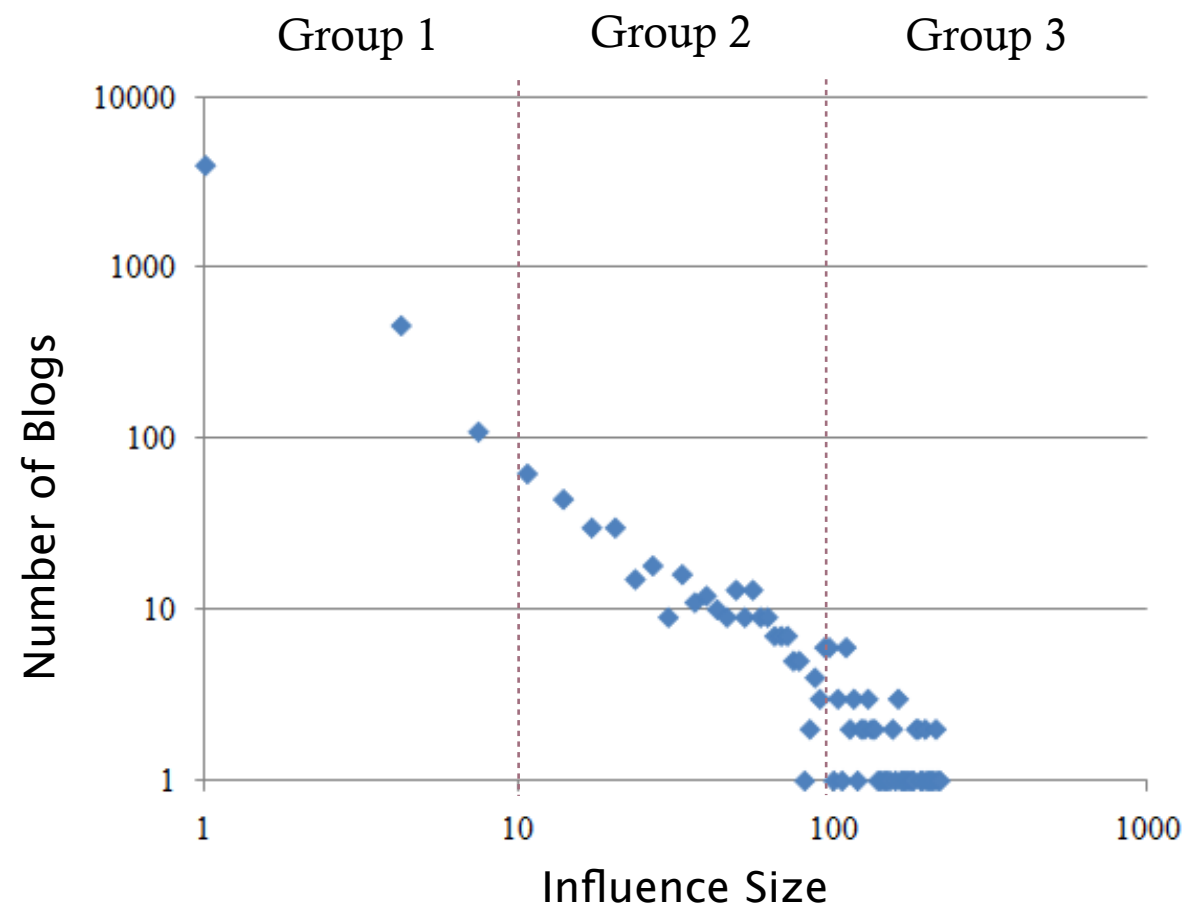
# Classification

- **Linear SVM**

- **Features**

  - Content: Topic proportions (50)

  - Network: In-degree, out-degree, total-degree, betweenness, clustering coefficient (5)

  - Content-Network: All (55)

- **Training**
  Train a classifier using 30% of the data blogs

- **Testing**
  Test with the rest of the data

# Prediction Result

| Features | Influence Size | Number of Comments | Digg Score |
|---|---|---|---|
| Content | 0.360 | 0.275 | 0.295 |
| Network | 0.726 | 0.308 | 0.239 |
| Content-Network | **0.727** | **0.322** | **0.308** |

## F1-Measure
(Harmonic mean of precision and recall)

# Prediction Result

| Features | Influence Size | | | | Number of Comments | | | | Digg Score | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | F1-Measure | Accuracy | Precision | Recall | F1-Measure | Accuracy | Precision | Recall | F1-Measure |
| Content | 0.841 | 0.303 | 0.445 | 0.360 | 0.631 | 0.172 | 0.691 | 0.275 | 0.688 | 0.197 | 0.583 | 0.295 |
| Network | 0.934 | 0.683 | 0.775 | 0.726 | 0.857 | 0.307 | 0.309 | 0.308 | 0.851 | 0.283 | 0.207 | 0.239 |
| Content-Network | 0.943 | 0.699 | 0.759 | 0.727 | 0.716 | 0.213 | 0.664 | 0.322 | 0.714 | 0.211 | 0.571 | 0.308 |

# Contributions

- We proposed a new measure of blog influence, *Influence Size*, which considers both content and network

- We showed that content information and network information play different roles in predicting blog influence

- We showed that it is important to use both content information and network information for finding influential blogs

# Future Work

- Evaluation of the influence measures

- Exploration of various methods for content analysis

# Thank You

*Presenter*

**Yohan Jo**

*yohan.jo@kaist.ac.kr*

*Users & Information Lab*
*Computer Science Dept.*
*Korea Advanced Institute of Science and Technology*