# T-61.6060 Presentation

# Tag and document relations

Aki Saarinen, Laszlo Kozma, Stevan Keraudy
26.2.2008

# Objectives

- Grouping of similar documents together
- Grouping of similar tags together

# Outline of the methodology

1. Fetching of data and saving it to XML-files

2. Preprocessing + feature extraction with a Python script

3. (Dimension reduction)

4. Visualization with Self-Organizing Maps

5. Tag prediction with various methods

# Data source: Newsvine.com

- Newsvine.com is a community news website
- Created in 2005
- Based in Seattle
- Owned by msnbc.com
- Articles can be written by anyone
- Any kind of topic can be covered
- Around 1000 articles per day

# 1. Data Fetching

- Retrieve all articles from December 2007
- Script to retrieve data written in Java
- Store data into XML files
- 1 XML file per article

# Fetching algorithm

**For each day**

Go to page that lists all articles
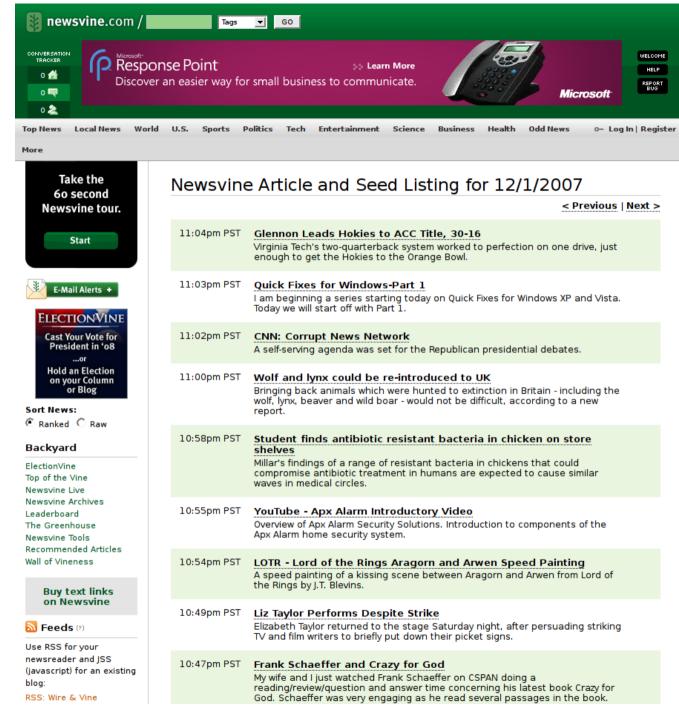
Get the list of articles

**For each article**

Get meta information

Go to the article page

Get tags, category and text

**If article not empty (to reduce spam)**

Write information to XML file

Clinton: Obama Mailings Are Deceptive

# Information retrieved

- For each document, we retrieved:
  - Date and time
  - Title
  - Summary
  - Tags
  - Category
  - Text

- We basically only used the tags for analysis, other information could have been used too but we decided to leave it out because of a relatively tight schedule

# Figures about the collected data

- **Newsvine.com**
  - 31 days: 2007 December 1$^{st}$ to 31$^{st}$
  - 27 779 documents -> average 896 docs/day
  - 45 039 unique tags
  - 31 360 tags used only once

# Tag histogram for Newsvine

# Figures about the collected data

- As a curiosity, we also fetched data from Slashdot, even though most of the analysis was done for Newsvine

- **Slashdot.org**

  - 6 months (2007-07 to 2007-12)

  - 4017 documents -> average 22 docs / day

  - 5123 unique tags

  - 4052 tags used only once


  - *Free tip: don't fetch ~1000 pages/day from a single machine or you will get banned :)*

# 2. Preprocessing and feature extraction

- Read data from XML-files

  - Drop tags which are used less than N times

  - Ignore option for the most general tags (categories)


- Generate 2 matrices to be processed with Matlab

  - Tags used in documents (docs=rows, tags=columns)

  - Tags occurring together (tags=rows, tags=columns)

# Generated matrices

- Document-tag (=tags used in documents)

|  | Gw-bush | Elections | Cooking | Internet | Web2.0 | ... |
|---|---|---|---|---|---|---|
| 20071201-1 | 1 | 1 | 0 | 0 | 0 | |
| 20071201-2 | 0 | 0 | 0 | 1 | 1 | |
| ... | | | | | | |

- Tag-tag (=tags occurring together)

|  | Gw-bush | Elections | Cooking | Internet | Web2.0 | ... |
|---|---|---|---|---|---|---|
| Gw-bush | 600 | 100 | 1 | 2 | 0 | |
| Elections | 100 | 520 | 0 | 10 | 0 | |
| Cooking | 1 | 0 | 80 | 0 | 0 | |
| Internet | 2 | 10 | 0 | 400 | 200 | |
| Web2.0 | 0 | 0 | 0 | 200 | 350 | |
| ... | | | | | | |

# Normalization for tag-tag-matrix

- Original matrix

|  | Gw-bush | Elections | Cooking | Internet | Web2.0 | ... |
|---|---|---|---|---|---|---|
| Gw-bush | 600 | 100 | 1 | 2 | 0 | |
| Elections | 100 | 520 | 0 | 10 | 0 | |
| Cooking | 1 | 0 | 80 | 0 | 0 | |
| Internet | 2 | 10 | 0 | 400 | 200 | |
| Web2.0 | 0 | 0 | 0 | 200 | 350 | |
| ... | | | | | | |

- Each cell divided by total count of matching tag

|  | Gw-bush | Elections | Cooking | Internet | Web2.0 | ... |
|---|---|---|---|---|---|---|
| Gw-bush | 600/600 | 100/600 | 1/600 | 2/600 | 0 | |
| Elections | 100/520 | 520/520 | 0 | 10/520 | 0 | |
| Cooking | 1/80 | 0 | 80/80 | 0 | 0 | |
| Internet | 2/400 | 10/400 | 0 | 400/400 | 200/400 | |
| Web2.0 | 0 | 0 | 0 | 200/350 | 350/350 | |
| ... | | | | | | |

# 3. Dimension reduction

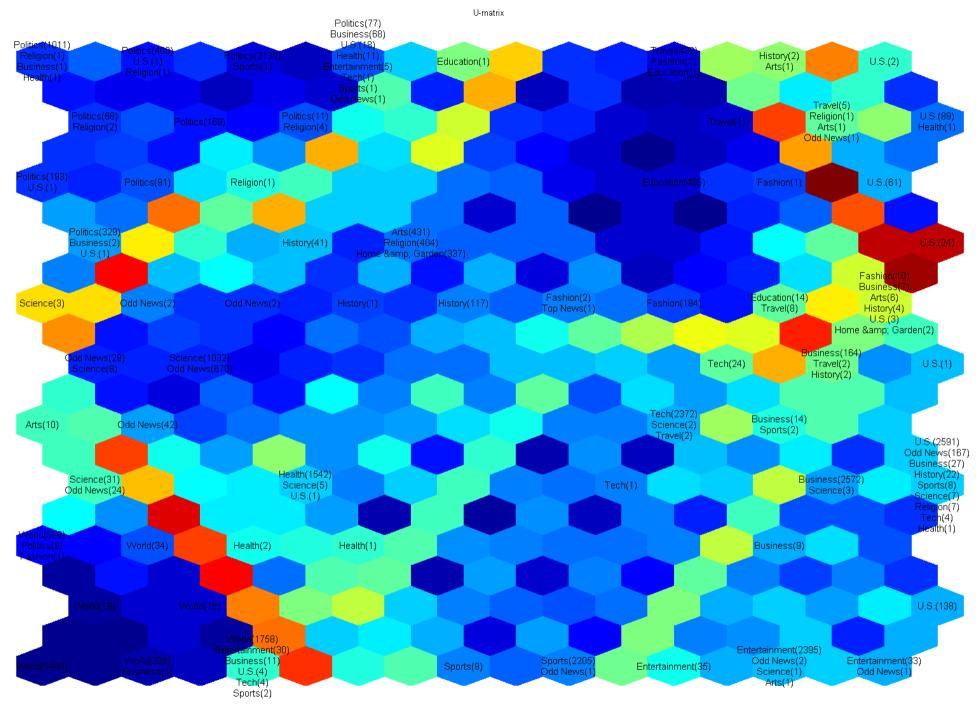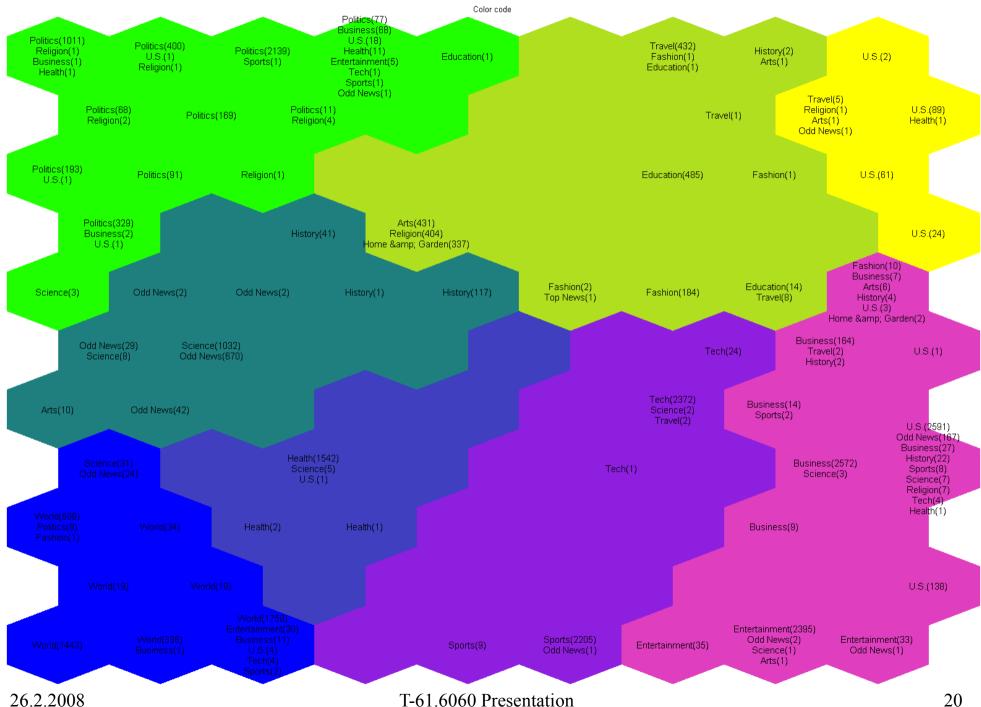- Generated document-tag-matrices are quite big
  - Document-tag-matrix size for whole Dec 2007 is 27779 x 45039
  - Need for dimension reduction
- Achieved using Principle Component Analysis (PCA)
  - We used for example top-100 or top-200 components
  - SOM could do it too, but takes a very long time to process, PCA is convenient
  - Even PCA is not doable for the whole month, half a month used

# 4. Visualization with SOMs

- Self-Organizing Maps

  - Artificial neural networks that can map multidimensional data to (usually) 2 dimensions

  - Map seeks to preserve topological properties of the input space

  - Results can be inspected visually

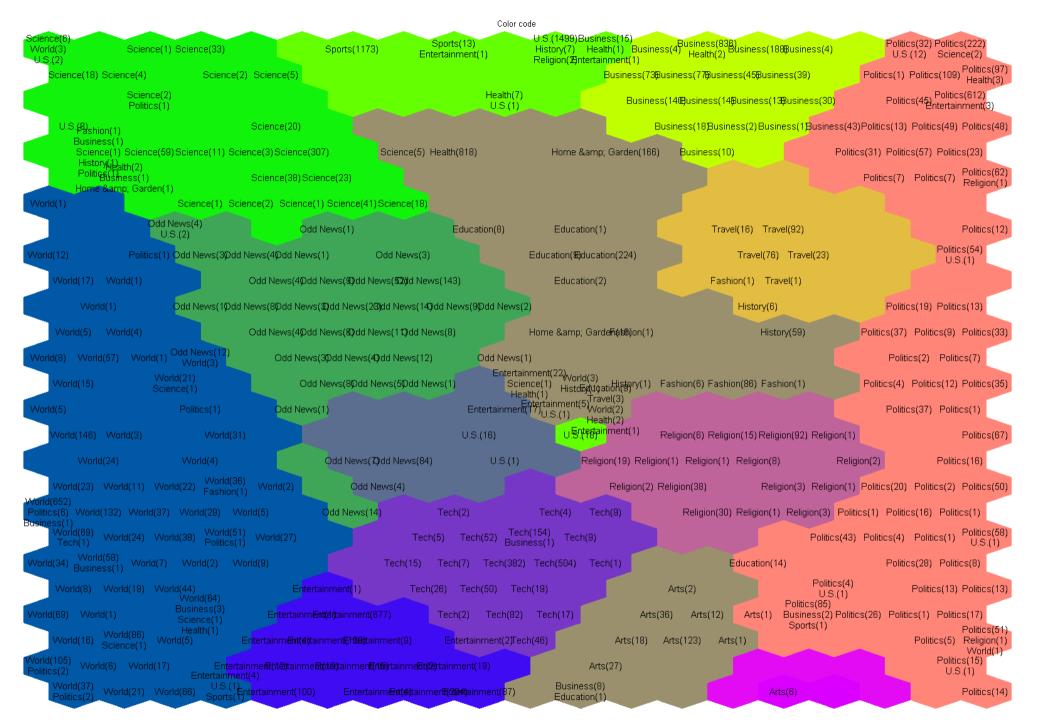  - U-matrices and k-means clustering shown

# SOM #1

- Data from Newsvine.com
  - Whole December 2007
- Document-tag-matrix used as input to SOM
  - Only tags which occur >= 40 times are used
  - No ignore list (category tags are there, too)
  - 27779 documents, 423 tags
  - No need to do PCA (matrix size: 27779 x 423)
  - Labeled with categories from the site

U-matrix

# SOM #2

- Data from Newsvine.com
    - December 1$^{st}$ to December 15$^{th}$
- Document-tag-matrix used as input to SOM
    - Only tags which occur more than 3 times are used
    - No ignore list (category tags are there, too)
    - 14679 documents, 4951 tags
    - Need to reduce with PCA, top-100 components used
    - Labeled with categories from the site

# SOM #3

- Data from Newsvine.com
  - Whole December 2007
- Normalized tag-tag-matrix used as input to SOM
  - Only tags which occur >= 40 times are used
  - No ignore list (category tags are there, too)
  - 423 tags
  - No need to do PCA (matrix size: 423 x 423)

# SOM #4

- Data from Slashdot.org

  - From July 2007 to December 2007

- Normalized tag-tag-matrix used as input to SOM

  - Only tags which occur >= 5 times are used

  - No ignores for categories

  - 268 tags

  - No need to do PCA

T-61.6060 Presentation

# 5. Tag prediction

- Tag-tag co-occurrence matrix revisited

|  | Gw-bush | Elections | Cooking | Internet | Web2.0 | ... |
|---|---|---|---|---|---|---|
| Gw-bush | 600 | 100 | 1 | 2 | 0 | |
| Elections | 100 | 520 | 0 | 10 | 0 | |
| Cooking | 1 | 0 | 80 | 0 | 0 | |
| Internet | 2 | 10 | 0 | 400 | 200 | |
| Web2.0 | 0 | 0 | 0 | 200 | 350 | |
| ... | | | | | | |

|  | Gw-bush | Elections | Cooking | Internet | Web2.0 | ... |
|---|---|---|---|---|---|---|
| Gw-bush | 600/600 | 100/600 | 1/600 | 2/600 | 0 | |
| Elections | 100/520 | 520/520 | 0 | 10/520 | 0 | |
| Cooking | 1/80 | 0 | 80/80 | 0 | 0 | |
| Internet | 2/400 | 10/400 | 0 | 400/400 | 200/400 | |
| Web2.0 | 0 | 0 | 0 | 200/350 | 350/350 | |
| ... | | | | | | |

# Tag-tag-matrix

- Very sparse, feasible to store for whole month as adjacency lists

- Original matrix symmetrical, normalized matrix not

- M ("obama", "politics") = 0.99

- M ("politics", "obama") = 0.02

- Possible to find partial orders between tags (see Mannila, et. al ...)

- General term > specific term

# Tags frequently occurring together

- For each tag i find tag j for which $M(i, j)$ maximal (or higher than a threshold)

- If $M(j, i)$ not too high, then j probably more general

- If both are high, terms probably closely related

- Before processing remove diagonal elements

- Note: results can be just as good as the quality of tags

# Example 1

- Co-occurrance high both for f(B|A) and f(A|B)

- Best match for each term

| term A | term B | f(B\|A) | f(A\|B) |
|---|---|---|---|
| free-premium-search-engine | search-engine | 1 | 0,69 |
| dogs | pets | 0,52 | 0,51 |
| united-kingdom | uk-news | 0,59 | 0,52 |
| aids | hiv | 0,47 | 0,68 |
| pkk | turkey | 1 | 0,35 |
| ron | paul | 1 | 0,59 |
| chavez | venezuela | 0,81 | 0,39 |
| energy | nuclear | 0,34 | 0,34 |
| au-news | australia | 1 | 0,45 |
| playstation | xbox | 0,64 | 0,47 |
| benazir-bhutto | pakistan | 0,84 | 0,31 |
| on | trial | 0,32 | 0,37 |
| wii | nintendo | 0,64 | 0,68 |
| based | home | 1 | 0,38 |
| bosnia-news | bosnia | 1 | 0,84 |
| pets | animals | 0,59 | 0,47 |
| slain | taylor | 0,56 | 0,86 |
| gay | lesbian | 0,43 | 0,74 |
| maryland | dcmetro | 0,66 | 0,82 |
| dcmetro | maryland | 0,82 | 0,66 |
| plane | crash | 0,38 | 0,32 |
| veterinary | animals | 0,57 | 0,33 |
| animals | pets | 0,47 | 0,59 |
| homosexual | gay | 0,77 | 0,38 |
| climate-change | global-warming | 0,73 | 0,62 |

# Example 2

- Co-occurrance high both for f(B|A) and f(A|B)

- All matches for each term which are higher than specified threshold

| term A | term B | f(B\|A) | f(A\|B) |
|---|---|---|---|
| mobile-phone | cell-phone | 0,58 | 0,43 |
| mobile-phone | at-t | 0,58 | 0,56 |
| mobile-phone | nokia | 0,69 | 0,44 |
| mobile-phone | verizon | 0,62 | 0,59 |
| free-premium-search-engine | search-engine | 1 | 0,69 |
| search-engine | engine | 0,69 | 1 |
| uranium | fuel-cycle | 0,62 | 0,68 |
| uranium | mining | 0,65 | 0,63 |
| mitt-romney | mike-huckabee | 0,35 | 0,3 |
| nuclear | energy | 0,34 | 0,34 |
| mp3 | download | 0,33 | 0,56 |
| dogs | pets | 0,52 | 0,51 |
| dogs | animals | 0,45 | 0,35 |
| dogs | cats | 0,39 | 0,61 |
| climate | conference | 0,31 | 0,41 |
| nfc-east | philadelphia-eagles | 0,7 | 0,67 |
| nfc-east | philadelphia | 0,72 | 0,63 |
| college-basketball | top-25 | 0,88 | 0,84 |
| board-of-trade | chicago-board | 1 | 1 |
| united-kingdom | uk-news | 0,59 | 0,52 |
| cell-phone | mobile-phone | 0,43 | 0,58 |
| cell-phone | at-t | 0,43 | 0,56 |
| cell-phone | nokia | 0,46 | 0,39 |

# Example 3

- (specific, general)-pairs
- f(B|A) large
- f(A|B) small

| term A | term B | f(B\|A) | f(A\|B) |
|---|---|---|---|
| music | entertainment | 0,7 | 0,11 |
| tech | technology | 0,81 | 0,02 |
| mobile-phone | technology | 0,77 | 0,01 |
| john-edwards | politics | 0,97 | 0,02 |
| dennis-kucinich | politics | 0,91 | 0,01 |
| election-08 | politics | 0,95 | 0,01 |
| cheney | politics | 0,77 | 0,01 |
| investment | business | 0,78 | 0,01 |
| company | marketing | 0,72 | 0,22 |
| company | business | 0,86 | 0,01 |
| uranium | nuclear | 0,74 | 0,09 |
| mitt-romney | politics | 0,93 | 0,03 |
| biology | science | 0,82 | 0,03 |
| mp3 | music | 0,76 | 0,17 |
| mp3 | entertainment | 0,76 | 0,03 |
| bill-clinton | politics | 0,92 | 0,01 |
| britney-spears | entertainment | 0,87 | 0,01 |
| college-basketball | sports | 1 | 0,09 |
| board-of-trade | business | 1 | 0,01 |
| cell-phone | technology | 0,71 | 0,01 |
| movies | entertainment | 0,82 | 0,06 |
| venezuela | world-news | 0,8 | 0,02 |
| bhutto | pakistan | 0,75 | 0,22 |

# Applications of generalization

- Build trees of terms, proceeding from specific to general
- Automatically finding categories
- Predicting new tags

# Tag prediction, method #1

- For a set of tags, what are the most plausible next tags
  - Simple method:    $t\_i$ in tags
  - Predict $t\_o$, such that   $\sum_{i=1}^{n} M(t_0, t_i)$  is maximized
- Examples
  - ("john-mccain", "mitt-romney", "obama")
    --> "bill-richardson", "mccain", "rudy-giuliani", "fred-thompson"
  - ("nokia", "mobile-phone")
    --> "technology", "cell-phone", "verizon", "at-t"
  - ("gordon-brown", "george-w-bush")
    --> "uk-news", "politics", "world-news"

# Tag prediction, method #2

- Until now we relied on tag co-occurence
  - But "george-w-bush" and "g-w-bush" probably don't occur together for the same document, one writer only uses one of these tags in one article.
- They probably still have similar co-occurence patterns with other terms ("politics", "iraq-war", "united-states", ...)
- **Idea**: cross-correlation between rows of the tag-tag matrix
  - Matching terms will not necessarily be related

# Examples for method #2

- **"finland"**
  - reuters, britain, philippines, somalia, africa, germany, russia, mexico, algeria, lebanon, ireland, brazil, cuba, oil, violence

- **"president-bush"**
  - republicans, constitution, democrats, white-house, budget, congress, george-w-bush, bush, taxes, senate, foreign-policy, election-2008, george-bush

- **"beer"**
  - alcohol, coffee, party, food, children, charity, new-year, relationships, baby, odd, holidays, theft, ireland, medical, santa-claus, christmas

# Questions?