













- Search refers to retrieving relevant information from web
  - Goal oriented
- Mining refers to discovery of useful, hitherto unknown information from web usage data and web contents
  - Opportunistic























- Click-stream analysis
- Recommend items
   Weighted average of rating



Weighted average of neighbour's ratings



















## Impact of social network analysis

- Marketing strategies
- Influence of one community over another analyzed
- Analysis of Enron Email corpus
   Who was passing information to whom?







Text Mining is not Information Extraction

- Information Extraction yields informative chunks of data extracted from text
- Uses Natural Language Processing Techniques
- IR can answer user posed queries
- Does not address the issue of establishing "interestingness" of information

### Text Mining applications commerce Find "important" business houses from web houses Frequently occurring names may not be "interesting" since well known Boolly interesting – identify names that

- Really interesting identify names that are worth keeping a watch on
- How???????
  - – analyze views, reviews, news, .....



### Social Impact [<u>Narin et</u> <u>al.1997]</u>

- Technology industry relies more heavily than ever on governmentsponsored research results.
- Examined the science references on the front pages of American patents
- 73.3 percent had been written at public institutions - universities, government labs and other public agencies, both in the United States and abroad.













- Refers to identifying user opinions about various objects on web
- Available in the form of
  - o Explicit User responses
  - Implicit responses
  - Text documents



Three categories of opinions
 positive, neutral and negative





- o Small but useful overall positive
- o Reliable but expensive negative



### Extracting Product Features

- Pre-processing
  - Stemming
  - Stop word deletion
  - Fuzzy matching for word variant matches (auto-focus / autofocus/auto focus)
- Nouns and Noun Phrases from each sentence

# Frequent Feature Generation

- Association rule mining for frequent item sets
- Features of a product are
  - k-item sets (k words occurring together)
  - Occur in at least c% of reviews











- Intra-sentence conjunction
  - Battery life is very long positive or negative??
     The camera takes great pictures and has a long battery life long is positive for battery
- And positive conjugation in a sentence
- But negative conjugation in a sentence
- The camera takes great pictures but has a short battery life

# Pseudo Intra-sentence conjunction

- Missing and
- Use opinion from other sentences like The camera has a long battery life which is great
- Contradictory reviews majority counts









- To identify semantic orientation of word w
- $P(I(w)=L|A_k)_{(m)} = P(I(w)=L|U_T(A_{k,T}))_{(m)}$ 
  - A<sub>K,T</sub> represents the labels assigned by A<sub>k</sub> to neighbors of *w* connected to it by relationship (synonym/antonym) T
  - Ak =  $\{w_i|L_i\}$  is the set of labeled word

### **Opinion Words**

Table 1. Examples of opinion-bearing/non-opinion-

bearing words			
Adjectives	Final score	Verbs	Final score
Careless	0.63749	Harm	0.61715
wasteful	0.49999	Hate	0.53847
Unpleasant	0.15263	Yearn	0.50000
Southern	-0.2746	Enter	-0.4870
Vertical	-0.4999	Crack	-0.4999
Scored	-0.5874	combine	-0.5852









- MPQA and NTCIR-6
- SVM based system

or not

- Opinionated sentence recognition
   Each sentence classified as opinionated
- Opinion holder extraction
  - Opinion holder one word or multiple words in opinionated sentence

#### Features used

- Unigram of tokens with two attributes lemma and POS
- Sentence represented in terms of unigram - tf \* idf
  - o tf frequency of term in a sentence
  - *idf* inverse of term frequency in all sentences in training data
- SVM algorithm for learning a two-class problem



- POS tag
- Morphological features of word
- Entity types















