# Social media content: Applications

Saptarshi Ghosh

Department of CSE, IIT Kharagpur

Social Computing course, CS60017

# Content posted on OSM

- User generated content (UGC)
- Crowdsourced from the user population
- Huge volume, posted with high velocity
- Variety of content: text, images, videos, …

- Large variation in quality
  - News articles, celebrity / expert posts, conversational chatter, spam, abusive and hate speech, fake news, …

# Few applications

- Classifying different types of information
- Sentiment analysis
- Filtering harmful content
- Clustering similar information
- Event detection and tracking
- Summarization
- Expert / important user identification
- Social search and recommendation
- Handling content in different languages
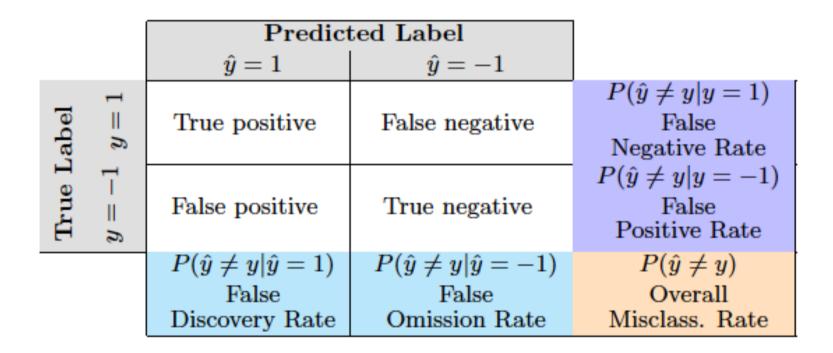
# Classification

- Many aspects along which OSM content can be classified

- Type of content:
  - During a disaster: situational information / sentiment and opinion
  - Political leaning: democratic-leaning / republican-leaning

- Credibility: rumor vs. true information, genuine vs. fake news

# Classification

- Supervised classification
  - Set of example items in each category known – training set
  - Extract features from the items
  - Learn a predictive function or model from the features
  - Apply model on a testing set to test performance – items for which categories are known, but not used for training

- Selecting training set and testing set
  - Cross validation
  - Held-out testing set

# Desirable property of classifiers

- Accuracy: measured using confusion matrix

| | Predicted Label | | |
|---|---|---|---|
| | $\hat{y} = 1$ | $\hat{y} = -1$ | |
| True Label $y = 1$ | True positive | False negative | $P(\hat{y} \neq y \mid y = 1)$ False Negative Rate |
| $y = -1$ | False positive | True negative | $P(\hat{y} \neq y \mid y = -1)$ False Positive Rate |
| | $P(\hat{y} \neq y \mid \hat{y} = 1)$ False Discovery Rate | $P(\hat{y} \neq y \mid \hat{y} = -1)$ False Omission Rate | $P(\hat{y} \neq y)$ Overall Misclass. Rate |

# Desirable properties of classifiers

- FAT:  Fairness, Accountability, Transparency

- Challenges
  - Some features may be sensitive (should not be used to discriminate), e.g., race, gender
  - Non-sensitive features may be correlated with sensitive features
  - Training set may be biased, and the bias may be inherited by the classifier
  - Misclassification rate may be different for different types of instances

# Classification

- Primary challenge: feature extraction and selection
- More features might not always guarantee better classification performance: feature selection

- Recent emphasis on neural network / deep learning techniques
  - Simplifies feature extraction
  - Reduced explainability, transparency

# Sentiment analysis

- Special type of classification
- Usually 3 classes: positive, neutral, negative

- Many applications:
  - Understanding general opinion about a product / movie
  - Predicting election outcomes

- What features can be used?

# Filtering harmful content

- Harmful content: spam, abusive and hate posts, rumors, fake news, …

- What features can be used?
  - Text features
  - User features
  - Network features
  - Temporal features

# Examples of rumor

40 # Crocodile out of the reserves. #Chennai people please be safe. #ChennaiFloods #ChennaiRainsHelp

More than 40 crocodile escaped from park at Chennai due to overflow of water. On the roads of ecr side. Vellachery #ChennaiFloods

# Examples of rumor and denials

40 # Crocodile out of the reserves. #Chennai people please be safe. #ChennaiFloods #ChennaiRainsHelp

No. the crocodiles have NOT escaped from the Madras Crocodile Bank. It's a hoax, so please don't panic #ChennaiFloods

More than 40 crocodile escaped from park at Chennai due to overflow of water. On the roads of ecr side. Vellachery #ChennaiFloods

Stop spreading rumors like crocodiles on the loose etc … #ChennaiFloods #ChennaiRainsHelp

# Clustering

- Unsupervised version of classification
  - Group similar items together …
  - … so that elements within a cluster are more similar to each other, than elements in different clusters


- Applications
  - Cluster similar OSM posts into stories, so that it is sufficient for human to check stories

# Clustering

- Two broad types

- Hard clustering: each item belongs to only one cluster

- Soft clustering: an item can simultaneously belong to multiple clusters with varying degrees

- Analogous to finding partitions / overlapping communities in networks

# Topic modeling: soft clustering

- Identifies "topics" for a given set of documents

- Very simply
  - Topic: a cluster of words which frequently occur together
  - A document assigned multiple topics with varying degrees

- Actually
  - Each topic is a distribution over all distinct terms
  - Each document assigned a distribution over all topics

# Topic modeling: soft clustering

- Examples of topics identified from social media posts during an earthquake
  - {tsunami, disaster, relief, earthquake}
  - {dead, bodies, missing, victims}
  - {aid, help, money, relief}

# Summarization

- Summarizing a single document vs. summarizing a set of documents vs. summarizing a stream of documents

- Types of summarization
  - How is the summary generated: Extractive vs. Abstractive
  - Incremental summarization or update summarization: a set of documents already read, and set of new documents

# Summarization

- Application of both clustering and classification

- Clustering: group similar documents, choose representative from each cluster

- Classification: separate out different types of documents, summarize each type separately

# Event detection and tracking

- New event detection
  - Given an incoming stream of documents, check each to see whether it is a new story
  - Check whether a document is 'sufficiently' different from previous ones, according to some similarity metric

- Event tracking
  - Follow the evolution of an event / topic
  - Detect sub-events

# New event detection: Possible methods

- Cluster documents, check if new document sufficiently close to cluster representative / center

- Look for keyword bursts:
  - Frequency of a keyword sharply increases, compared to historical running average
  - Need to distinguish between events in physical world and Twitter memes like #musicmonday or #followfriday

# Identify influential users / experts

- Several metrics of user influence
  - #followers, PageRank, #times retweeted in Twitter, …
  - Topic-specific expertise

- Experts in specific scenarios
  - Community leaders during emergencies [Tyshchuk, ASONAM 2013]
  - Geographically 'local' sources [Yardi, ICWSM 2007]

# Search and Recommendation

- Help users discover interesting content, friends, groups

- Basis: friends likely to have similar tastes

- Recommend friends, groups to join [Chen, WWW09], resources [Konstas, SIGIR09], tags [Sen, WWW09][Song, SIGIR08]

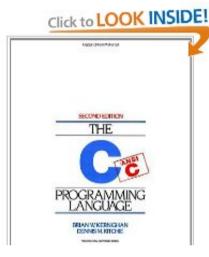- Personalized answers to queries [Xu, SIGIR08] [Bao, WWW07] [Mislove, HotNets06]

# Recommendation algorithms

- Two broad types
    - Collaborative filtering
    - Content-based filtering

- Hybrid schemes also used

# Collaborative Filtering

- Input:
  - Data on users' past behavior, or preferences for items
  - Typically, a user-item matrix where entries are ratings
- Idea:
  - For user u, identify users with similar interests, recommend to u the items that they liked
  - For a user who has liked an item, recommend other similar items
- No "understanding" of items / users required
- Challenges: scalability, scarcity
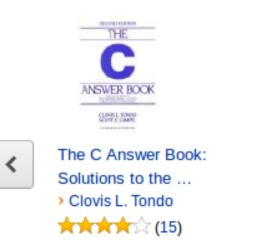
# Recommendation of books in Amazon

Click to **LOOK** INSIDE!

SECOND EDITION
THE
**C**
PROGRAMMING
LANGUAGE
BRIAN W KERNIGHAN
DENNIS M. RITCHIE

## C Programming Language (2nd Edi

Brian W. Kernighan ⌄ (Author), Dennis M. Ritchie (Author)
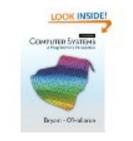
★★★★☆ ⌄ (367 customer reviews)

**Buy New**
$52.49 & **FREE Shipping**. Details

**In Stock.**
Ships from and sold by **Amazon.com**. Gift-wrap available.

**Want it Tuesday, June 4?** Order within **55 hrs 50 mins** and choose **One-Day Shipping** at checkout. Details
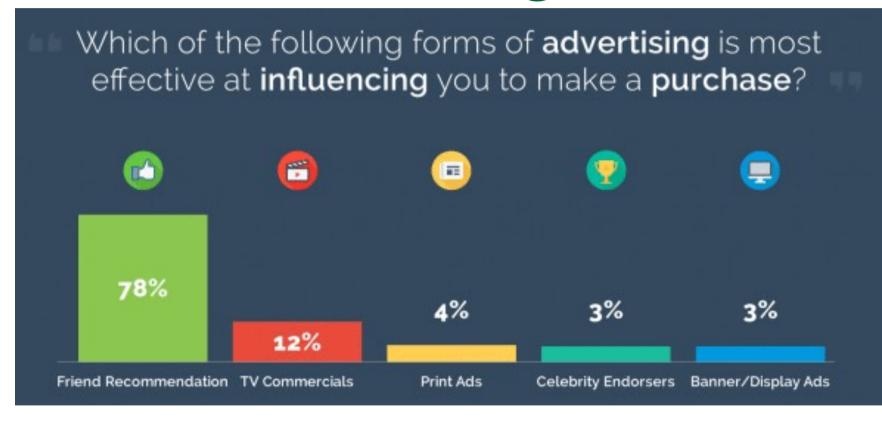
## Customers Who Bought This Item Also Bought

SECOND EDITION
THE
**C**
ANSWER BOOK
CLOVIS L. TONDO
SCOTT E. GIMPEL

LOOK INSIDE!
COMPUTER SYSTEMS
A Programmer's Perspective
Bryant · O'Hallaron

LOOK INSIDE!
Programming in C

‹

The C Answer Book:
Solutions to the …
› Clovis L. Tondo
★★★★☆ (15)

Computer Systems: A
Programmer's …
› Randal E. Bryant
★★★★☆ (23)

Programming in C (3rd
Edition)
› Stephen G. Kochan
★★★★☆ (70)

# Social recommendations: special case of collaborative filtering



Which of the following forms of **advertising** is most effective at **influencing** you to make a **purchase**?

| Friend Recommendation | TV Commercials | Print Ads | Celebrity Endorsers | Banner/Display Ads |
|---|---|---|---|---|
| 78% | 12% | 4% | 3% | 3% |

# Content-based filtering

- Input:
  - Data on users' past behavior, or preferences for items
  - Some information about the items (keywords, attributes)
- Idea:
  - Learn a profile / representation of a user, and recommend matching items
  - Recommend items that are similar to those that a user liked in the past
- Requires an "understanding" of users and items

# Evaluation of RS

- Accuracy / Relevance
- Diversity, novelty, serendipity (trade-off with relevance)
- Privacy
- Trust and explainability
- Fairness (unbiased)

- Related terms: filter bubbles, echo chambers, segregation or polarization

# Multi-lingual content

- Increased use of non-English languages
- Code mixing
- Transliteration

भूकंप पीड़ितों को खाना-टेंट चाहिये

भरतीय रेलवे ने 1 लाख पानी की बोतले भेजी है। धन्यवाद @sureshpprabhu जी #NepalEarthquake

100 feet statue for Modi Temple in UP. 30Cr to be spent. Who is malik? Why Modi's temple? Whose ACCHE DIN??