# Introduction to Text Mining

Hongning Wang
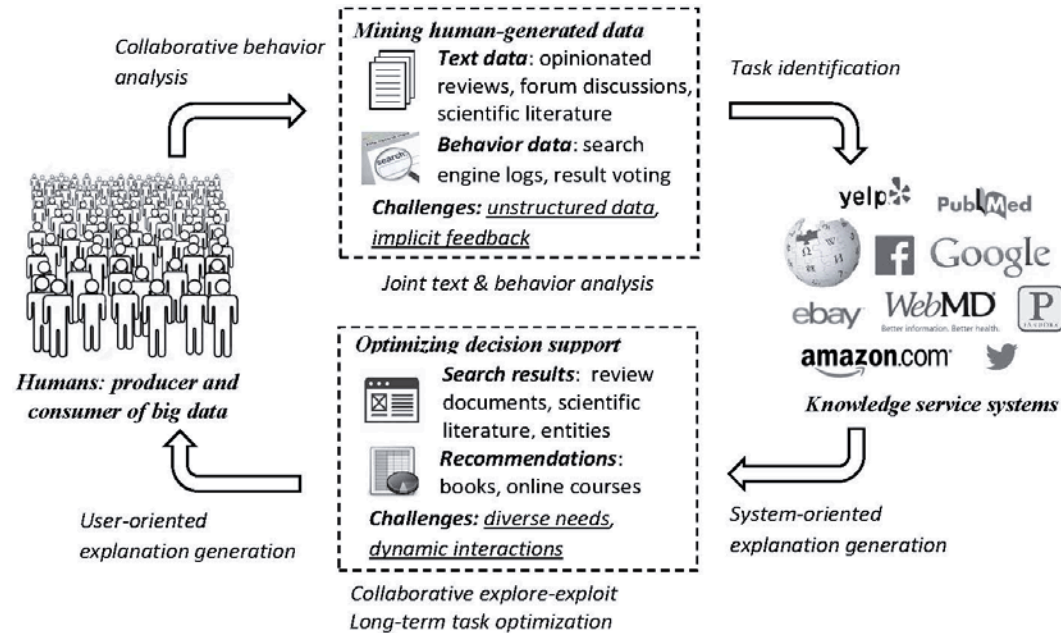
CS@UVa

# Who Am I?

- Hongning Wang
  - Assistant professor in CS@UVa since August 2014
  - Research areas
    - Information retrieval
    - Data mining
    - Machine learning

# Who Am I?

- Hongning Wang
  - Assistant professor in CS@UVa since August 2014

# What Am I Doing at UVa?

- Sentiment analysis with topic modeling

| | By: Kindle Customer | Date: June 25, 2014 |
|---|---|---|
| Topic | Samsung Galaxy Note 10.1 | Amazon Kindle Fire HDX |
| (+, battery) | Battery life is very good, it is easily an all day device with wifi on and high brightness while taking notes | Battery life is ok - probably need to recharge every other day with normal use |
| (-, battery) | My only issue is that it takes a long time to take a full charge and does not charge rapidly enough to use while charging, but the battery life is not bad | Everything works great, but the battery life is not nearly as long as advertised |
| (+, sound) | it has pretty good battery life, it also has an excellent quality sounding speakers, which i wasn't expecting on any tablet | Sound is really good (not home theater quality or anything) but better than any phone I've heard. |
| (-, sound) | The audio became occasionally inoperative and the headphone jack would crackle when using my ear buds | Users can get confused with volume buttons on the other side |
| (+, cpu) | quad core processor runs everything quickly and smoothly | The device features a fast 2.2GHz quad-core processor and 2GB of RAM for fast that run apps, games, and videos smoothly without an issues. |
| (-,cpu) | The OS was fast at first but as I added Apps it got slower and choppy | Compared to a galaxy note which is the same price, the Kindle HDX seems to have a slower processor |

those low standards $^{\{sound,-\}}$.

# What Am I Doing at UVa?

- Interactive online recommendation
  - Modeling recommendation as a two-party game
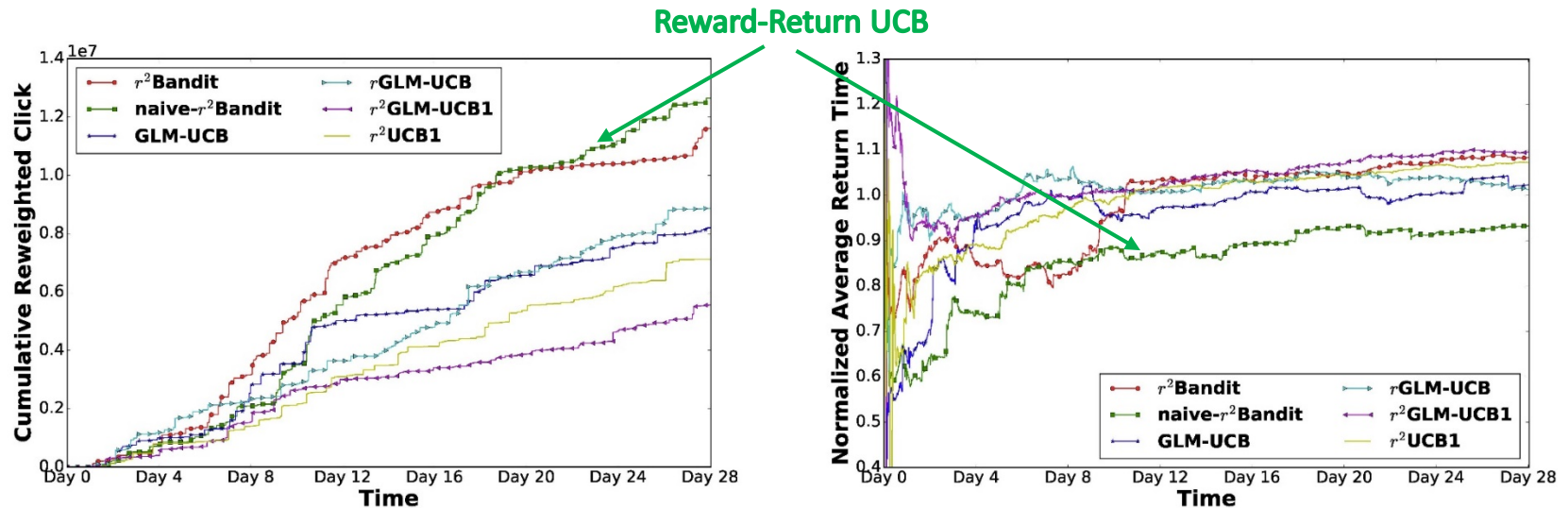
**Strategy?**

**Goal:**

**Challenge:**
1) Unknown preference;
2) Feedback is acquired on the fly, and it is not free!

# What Am I Doing at UVa?

- Yahoo frontpage news recommendation



18,882 users, 188,384 articles, and 9,984,879 logged events segmented into 1,123,583 sessions.

# What Am I Doing at UVa?

- Personalization techniques raise serious public concerns about privacy infringement

*No means for users to opt-out data collection!*

# What Am I Doing at UVa?

- Privacy-preserving personalization



Stronger privacy guarantee than *k*-anonymity

True search intent:
"polycystic ovary syndrome treatment"

latent group 1 — apple, ipad, iphone
latent group 2 — plan, health, insurance
latent group 3 — hotel, flight, rental
latent group *k* — soccer, basketball, tennis

Step *t*: genuine & cover-up queries
Step *t+1*: search results
Step *t+2*: genuine & pseudo clicks

Inferred search intent:
29% health, 8% sports, 5% vacation….

User

Client: balance privacy-preservation and personalization

Service provider

# What are about you?

- Why do you choose this course?

- Anything specific you want me to know?

- What type of text data do you often encounter in your projects?

- What kind of knowledge do you want to extract from it?

# What is "Text Mining"?

- *"Text mining, also referred to as **text data mining**, roughly equivalent to text analytics, refers to the process of deriving high-quality information from text."* - wikipedia

- *"Another way to view text data mining is as a process of **exploratory** data analysis that leads to **heretofore unknown** information, or to answers for questions for which the answer is not currently known."* - Hearst, 1999

# Two different definitions of mining

- Goal-oriented (effectiveness driven)
  - Any process that generates useful results that are non-obvious is called "mining".
  - Keywords: "**useful**" + "**non-obvious**"
  - Data isn't necessarily massive
- Method-oriented (efficiency driven)
  - Any process that involves extracting information from massive data is called "mining"
  - Keywords: "**massive**" + "**pattern**"
  - Patterns aren't necessarily useful

# Knowledge discovery from text data

- IBM's Watson wins at Jeopardy! - 2011

# An overview of Watson



- On questions, at the start of question analysis
- On primary search results, before candidate answer generation
- On supporting evidence, before deep evidence scoring

# What is inside Watson?

- *"Watson had access to <u>200 million pages</u> of structured and unstructured content consuming four terabytes of disk storage including the full text of Wikipedia" – PC World*

- *"The sources of information for Watson include encyclopedias, dictionaries, thesauri, newswire articles, and literary works. Watson also used databases, taxonomies, and ontologies. Specifically, DBPedia, WordNet, and Yago were used." – AI Magazine*

# What is inside Watson?

- DeepQA system
  - *"Watson's main innovation was not in the creation of a new algorithm for this operation but rather its ability to **quickly** execute hundreds of proven language analysis algorithms simultaneously to find the correct answer."* – New York Times
  - The DeepQA Research Team

# Text mining around us

- Sentiment analysis

# Text mining around us

- Sentiment analysis

# Text mining around us

- Document summarization

# Text mining around us

- Document summarization

# Text mining around us

- Movie recommendation

# Text mining around us

- Restaurant/hotel recommendation

# Text mining around us

- News recommendation

# Text mining around us

- Text analytics in financial services

# Text mining around us

- Text analytics in healthcare

# How to perform text mining?

- As computer scientists, we view it as
  - Text Mining = Data Mining + Text Data

Applied machine learning

Natural language processing

Information retrieval

Blogs

Web pages

Software documentations

Emails

Tweets

News articles

Scientific literature

# Text mining v.s. NLP, IR, DM…

- How does it relate to data mining in general?
- How does it relate to computational linguistics?
- How does it relate to information retrieval?

| | Finding Patterns | Finding "Nuggets" | |
|---|---|---|---|
| | | Novel | Non-Novel |
| Non-textual data | General data-mining | Exploratory data analysis | Database queries |
| Textual data | Comp Ling  **Text Mining** | | Information retrieval |

# Text mining in general



**Access**  Serve for IR applications  Sub-area of DM research  **Mining**

Filter information

Discover knowledge

**Organization**

Based on NLP/ML techniques

Add Structure/Annotations

# Challenges in text mining

- Data collection is "free text"
  - Data is not well-organized
    - Semi-structured or unstructured
  - Natural language text contains ambiguities on many levels
    - Lexical, syntactic, semantic, and pragmatic
  - Learning techniques for processing text typically need annotated training examples
    - Expensive to acquire at scale
- What to mine?

# Text mining problems we will solve

- Lexical semantics and word senses
  - Identifying which sense of a word (i.e. meaning) is used in a sentence, when the word has multiple meanings



Bass: fish

???

Bass: instrument

# Text mining problems we will solve

- Document categorization
  - Adding structures to the text corpus

# Text mining problems we will solve

- Text clustering
  - Identifying structures in the text corpus

# Text mining problems we will solve

- Topic modeling
  - Identifying structures in the text corpus

# Text mining problems we will solve

- Social media and social network analysis
  - Exploring additional structure in the text corpus

# We will also briefly cover

- Natural language processing pipeline
  - Tokenization
    - "Studying text mining is fun!" -> "studying" + "text" + "mining" + "is" + "fun" + "!"
  - Part-of-speech tagging
    - "Studying text mining is fun!" ->



  - Dependency parsing
    - "Studying text mining is fun!" ->

# We will also briefly cover

- Machine learning techniques
  - Supervised methods
    - Naïve Bayes, k Nearest Neighbors, Logistic Regression
  - Unsupervised methods
    - K-Means, hierarchical clustering, topic models
  - Semi-supervised methods
    - Expectation Maximization

# Text mining in the era of Big Data

- Huge in size
  - Google processes 5.13B queries/day (2013)
  - Twitter receives 340M tweets/day (2012)
  - Facebook has 2.5 PB of user data + 15 TB/day (4/2009)
  - eBay has 6.5 PB of user data + 50 TB/day

- 80% data is unstructured (IBM, 2010)

**640K** ought to be enough for anybody.

# Scalability is crucial

- Large scale text processing techniques
  - MapReduce framework

# State-of-the-art solutions

- Apache Spark ([spark.apache.org](spark.apache.org))
  - In-memory MapReduce
    - Specialized for machine learning algorithms
  - Speed
    - 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.



Logistic regression in Hadoop and Spark

# State-of-the-art solutions

- Apache Spark ([spark.apache.org](spark.apache.org))
  - In-memory MapReduce
    - Specialized for machine learning algorithms
  - Generality
    - Combine SQL, streaming, and complex analytics

# State-of-the-art solutions

- GraphLab ([graphlab.com](graphlab.com))
  - Graph-based, high performance, distributed computation framework

# State-of-the-art solutions

- GraphLab ([graphlab.com](graphlab.com))
  - Specialized for sparse data with local dependencies for iterative algorithms

LOGISTIC REGRESSION PERFORMANCE (ACCURACY AND SPEED)

# Text mining in the era of Big Data



**Knowledge Discovery**

**Human-generated data**

*Text data*          *Behavior data*

*Knowledge service system*

**Decision Support**
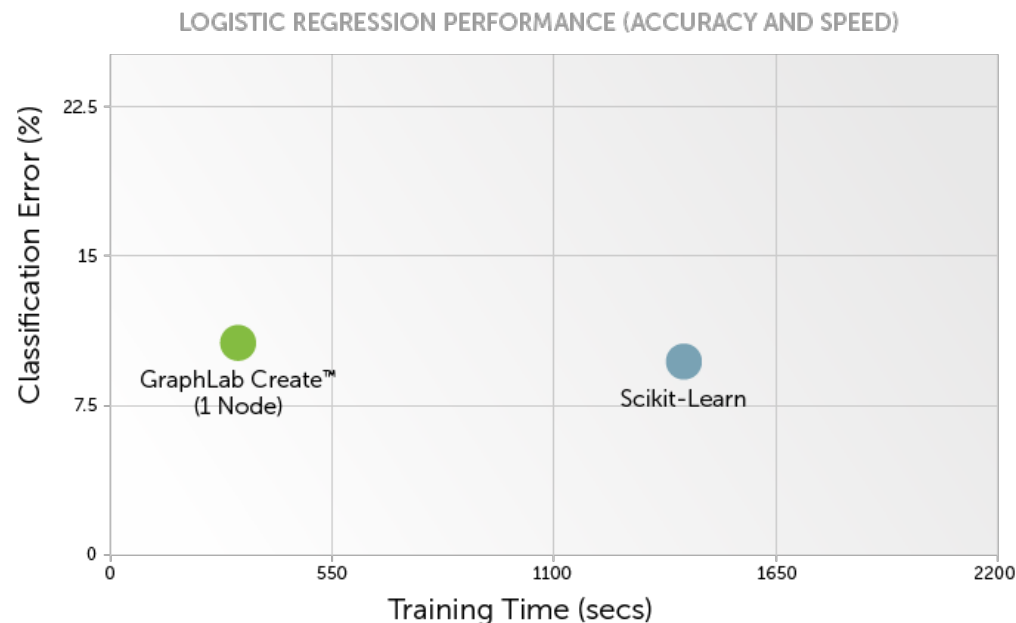
**Data Generation Modeling**

*As knowledge consumer*

**Challenges:**
1. **Implicit feedback**
2. **Diverse and dynamic**

*Human: big data producer and consumer*

*As data producer*

**Challenges:**
1. **Unstructured data**
2. **Rich semantic**

# Text books

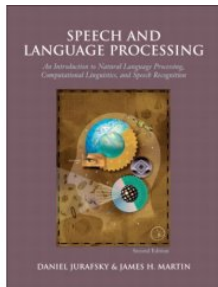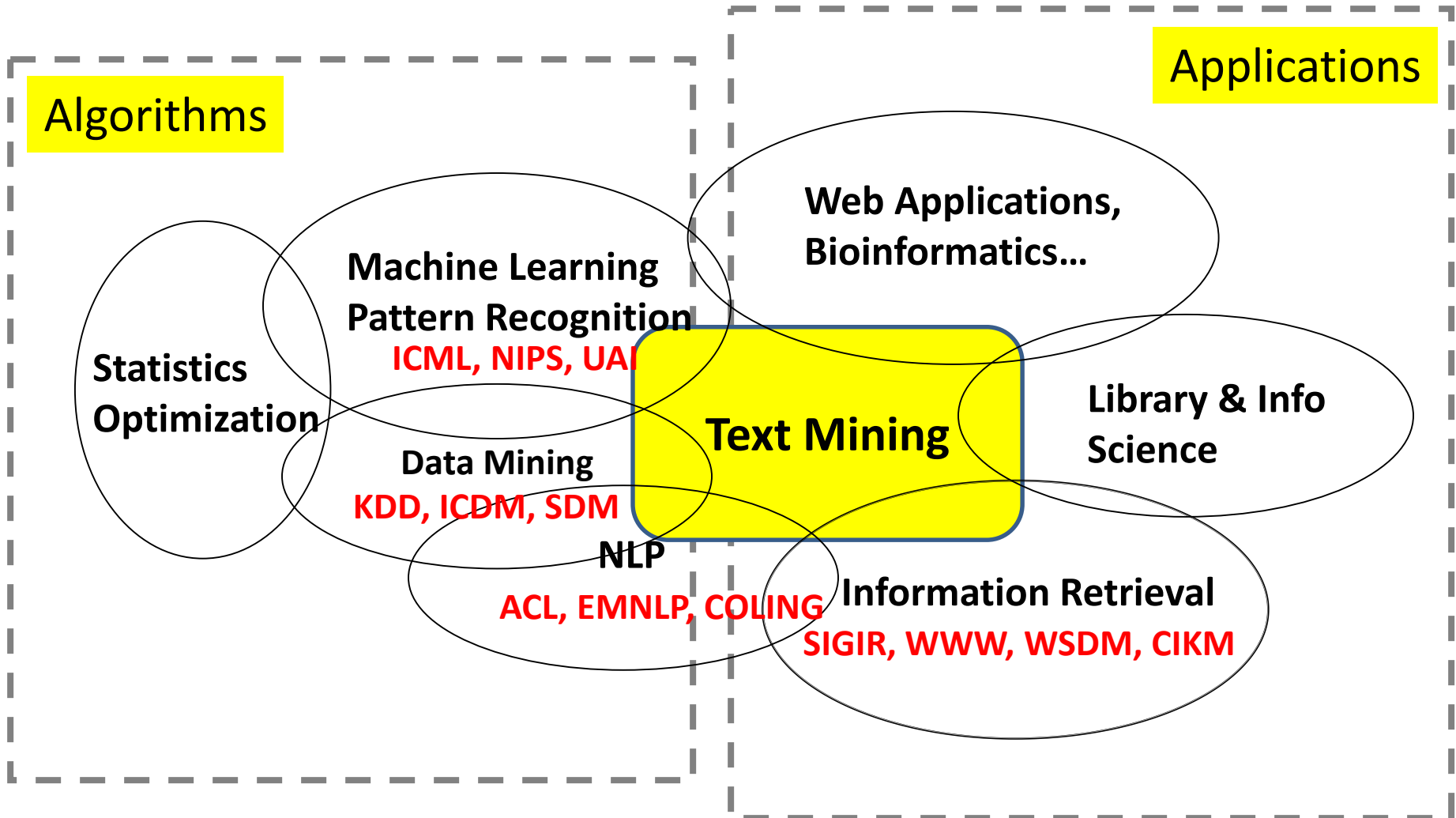- ***Introduction to Information Retrieval***. Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schuetze, Cambridge University Press, 2007.

- ***Speech and Language Processing***. Daniel Jurafsky and James H. Martin, Pearson Education, 2000.

- ***Mining Text Data***. Charu C. Aggarwal and ChengXiang Zhai, Springer, 2012.

# What to read?



Algorithms

Applications

Statistics Optimization

Machine Learning Pattern Recognition
**ICML, NIPS, UAI**

Data Mining
**KDD, ICDM, SDM**

NLP
**ACL, EMNLP, COLING**

Text Mining

Web Applications, Bioinformatics…

Library & Info Science

Information Retrieval
**SIGIR, WWW, WSDM, CIKM**

- Find more on course website for resource

# Welcome to the class of "Text Mining"!