



Welcome to the class of
Advanced Topics in
Information Retrieval !

Min ZHANG (张敏)

z-m@tsinghua.edu.cn



Course basic information

Instructor

Instructor: Min ZHANG (张敏)



- Associate Professor
- IR group, DCST, Tsinghua Uni.
- Office: FIT 1-506A Tel: 62798279
- Email: z-m@tsinghua.edu.cn
- <http://www.thuir.cn/group/~mzhang>

TA: ZHANG Fan (张帆)

- Email: frankyzf94@gmail.com
- Tel: 18671829106
- Office: FIT 1-506

3

A little bit about Min...



■ Experiences

- Dept. of CST, Tsinghua University
 - 1995~1999 (Bachelor), 1999~2003 (Ph.D.)
- Associate Professor in THUIR group, CST Dept.
- Visiting scholar/researcher in DFKI (Germany), Kyoto University, City University of HongKong, MSRA and NUS.

■ Research Interests

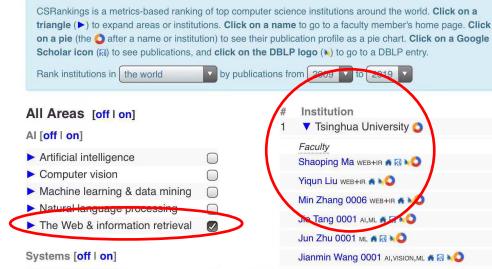
- Information Retrieval & Recommendation,
- User modeling and behavior analysis,
- Machine learning applications.

4

A little bit about Min...

■ Selected Achievements and Awards

- Published 100+ papers on important conferences & journals
 - ~3500 Citations(on Feb. 2019, Google Scholar) , H-index 32
- Multiple Top Perform. in TREC and NTCIR since 2002.
- One of the top ranked World-wide IR researchers
- Awards:
 - Beijing Science & Technology Progress Award, 1st prize, 2016
 - State of Minist. of Radio, Film & Television, Guangxi Science &Technology, etc.
 - Excellent Young Faculty Teaching Award, Tsinghua Uni.



5

A little bit about Min...

■ Academic Activities and Industry Connections

- Vice Director of the state key lab of intelligent technology and systems (AI institute)
- Vice Director of MOE MSRA Key Lab on Web and multimedia, Tsinghua University
- Associate Editor of Transaction on Information Systems (TOIS) (CCF A journal) and reviewers for multiple top journals
- Program Chairs: WSDM2019 workshop, SIGIR 2019 Tutorials, SIGIR2018 short papers, WSDM 2017, AIRS 2016
- Area Chair, Senior PC or PC for
 - IJCAI, SIGIR, WWW, KDD, WSDM, CIKM, ACL.....
- Committee members of China Associations
 - Chinese Information Processing Committee (CCF)
 - Information Retrieval Committee (CIPS)
 - Machine Learning Committee (CAAI)

6

Course Arrangements

- Friday afternoon (13:30-14:15, 14:20-15:05), VI-6B206
- The first 13~14 weeks
 - Lectures by the teacher on selected topics
 - With “Tea Time” discussions on the news/progresses on Web and IR industry/research (by one student)
 - Send me the email by Wed. noon 11:59am to apply a tea time show on Thursday's course.
 - [WolframAlpha](#), [Adidas miCoach Dark net SE Man vs. Machine Facebook Scandal](#)
- The last 2~3 weeks -- A workshop
 - Seminar presented by the students
 - 10~20 mins' talk + 5~10 mins' QA (tentative, depending on #students)
 - Awards and celebration to “[the Best Project](#)”

7

Part I-1. Lectures on Web IR Fundamentals (subject to modifications)

Week	Date	Topic	Content
1	3. 1	Introduction	Course intro, IR history
2	3. 8	IR key tech (I)	General procedure. Key tech: crawler , Index
3	3. 15	IR key tech (II)	Content-based Ranking Models: Term weighting; Boolean, VSM, probability(classical), LM
4	3. 22	IR key tech (III)	Evaluation: methodology; metrics (pre, recall, MAP, MRR, NDCG); Consistency analysis
5	3. 29	IR key tech (IV): Web Link Analysis	Link Analysis: Counting degree, HITS, PageRank, TrustRank, Spreading Activation, Anchor
6	4. 5	Holiday	
7	4. 12	Recommendation	Recommendation

Part I-2. Lectures on Advanced Topics (subject to modifications)

Week	Date	Topic	Content
8	4. 19	Advanced Topics (I)	User Behavior Analyses
9	4. 26	Advanced Topics (II)	Challenges: anti-spam
10	5. 3→5. 5 Sunday	Advanced Topics (III)	Challenges: Multi-source fusion, UI
11	5. 10	Advanced Topics (IV)	Challenges: Scale
12	5. 17	Evaluation Rethinking	What to evaluate, How to evaluate
13	5. 24	Social IR Or Invited Industry Talk	Social network & Human computation, Crowd computing
14	5. 31	Visual IR Or Invited Industry Talk	Visual IR

Part III. Project and Course workshop

- III-1. Design& implement a [prototype Search Engine](#)
 - Two students as a team
 - You [can use general IR toolkits](#) to build the prototype SE.
 - Both [general and domain-specific SEs](#) are acceptable.
 - A specific one with some interesting functions might be more attractive to your audiences
 - A [demo system](#) should be released [online](#) (accessible in THU net) [at least 1 week before the presentation.](#)
 - Final score including both pre-test and the live show.
- III-2. Submit a paper on this project (1 paper per student)

Week	Date	Topic	Content
15	6.7→6.10 Mon	Course Workshop	SE prototype design & implementation
16	6.14	Course Workshop	SE prototype design & implementation



11

cafe

Paradiso Coffee - TsingHua University

Coffee Shop

Address:

100084
北京
Building 21 7iina

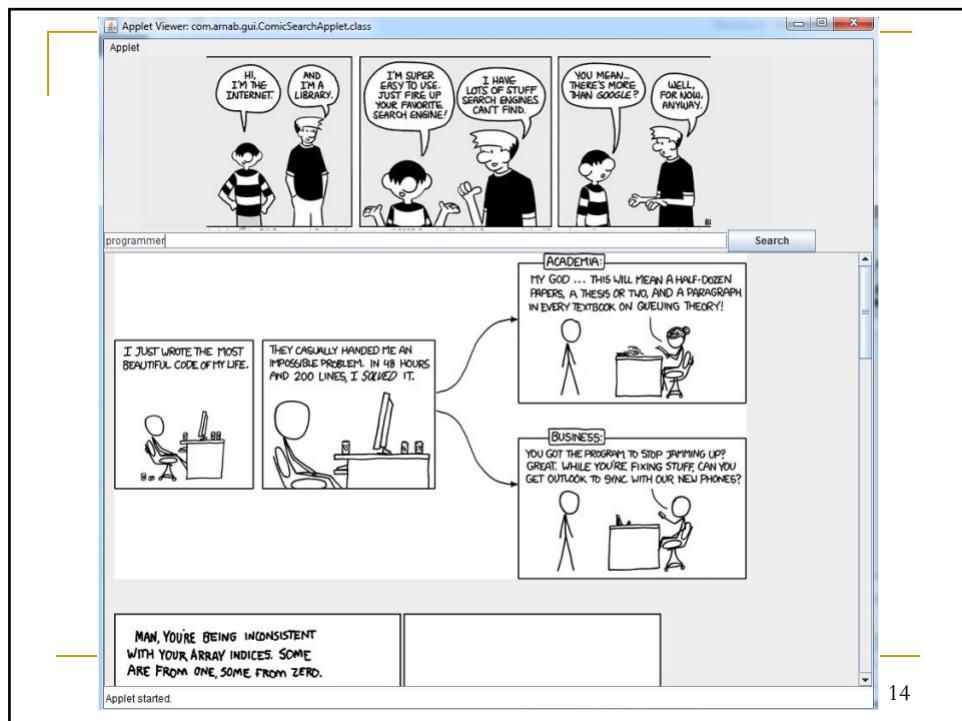
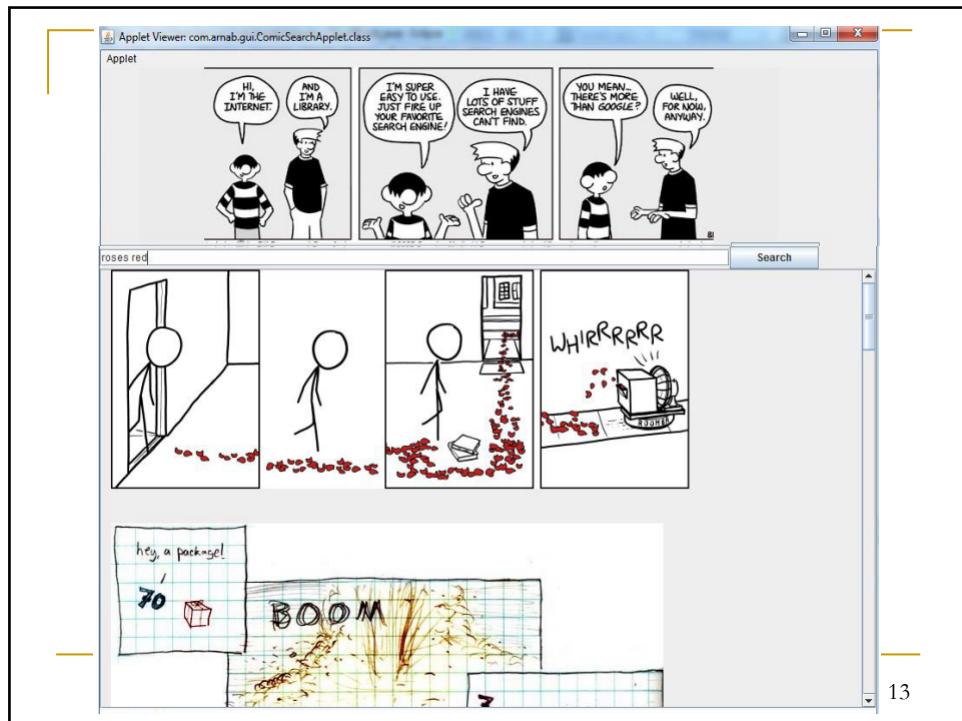
To finish one...

To search one...

To search all...

Comic Searcher

12



Supports fuzzy search (maximum edit distance of 3)

LATeX Symbol Search

esplon

Mouse over equations to display their TeX commands

Symbol in TeX	TeX Command	Can search content	
E, ε and ε	E, \epsilon and \varepsilon	LATEX Symbol Search	

the set of all real numbers

Mouse over equations to display their TeX commands

Symbol in TeX	TeX Command	Name	Explanation
		Read as	
		Category	
R	\mathbb{R}	real numbers	R; the (set of) real numbers; the reals
R	\mathbf{R}		R means the set of real numbers.
		numbers	

15

teamazingemojisearcher.herokuapp.com/?q=happy

teamazingemojisearcher.herokuapp.com/emoji/3

The Amazing Emoji Searcher

Home / 😂 face with tears of joy

Face with tears of joy

happy



 smi

 Description

查看详情

A laughing emoji which at small sizes is often mistaken for being tears of sadness. In fact, this emoji is laughing so much that it is crying tears of joy. Tears are coming from both eyes, not due to sadness, but overwhelmed with happiness or laughter.

Also known as

- 😂 Laughing Emoji
- 😂 Laughing Crying Emoji
- 😂 Happy Tears Emoji
- 😂 Laughing Tears Emoji
- 😂 LOL Emoji

Different versions

iPhone  Android  Twitter  Gmail  Windows  Black and white 

Made with love by Richard Luong Available on Github

素问 关于素问 联系我们 [用微博登录](#)

Tradition Chinese Medicine Doctors Search

素問

专业、安全、免费的医疗信息服务平台

Article Generator

搜索病名, 例如: 胃

Keywords

- learning
- machine
- data
- machine learning**
- method
- algorithm
- model
- program

Useful Links

- en.wikipedia.org
- azure.microsoft.com
- online.stanford.edu
- www.springer.com
- whatis.techtarget.com
- research.microsoft.com
- ocw.mit.edu
- azure.microsoft.com

Machine learning is a subfield of computer science that evolved from the study of pattern recognition and computational learning theory in artificial intelligence. Machine learning explores the construction and study of algorithms that can learn from and make predictions on data. Machine learning is closely related to and often overlaps with computational statistics; a discipline that also specializes in prediction-making. Machine learning is employed in a range of computing tasks where designing and programming explicit, rule-based algorithms is infeasible.

Machine learning is sometimes conflated with data mining, although that focuses more on exploratory data analysis. When employed in industrial contexts, machine learning methods may be referred to as predictive analytics or predictive modelling. In 1959, Arthur Samuel defined machine learning as a "Field of study that gives computers the ability to learn without being explicitly programmed". Already in the early days of AI as an academic discipline, some researchers were interested in having machines learn from data. Machine learning and data mining often employ the same methods and overlap significantly. The two areas overlap in many ways: data mining uses many machine learning methods, but often with a slightly different goal in mind. On the other hand, machine learning also employs data mining methods as "unsupervised learning".

Back PDF Clear Highlights

Web Art Search

清华大学
Tsinghua University

WEB ART SEARCH

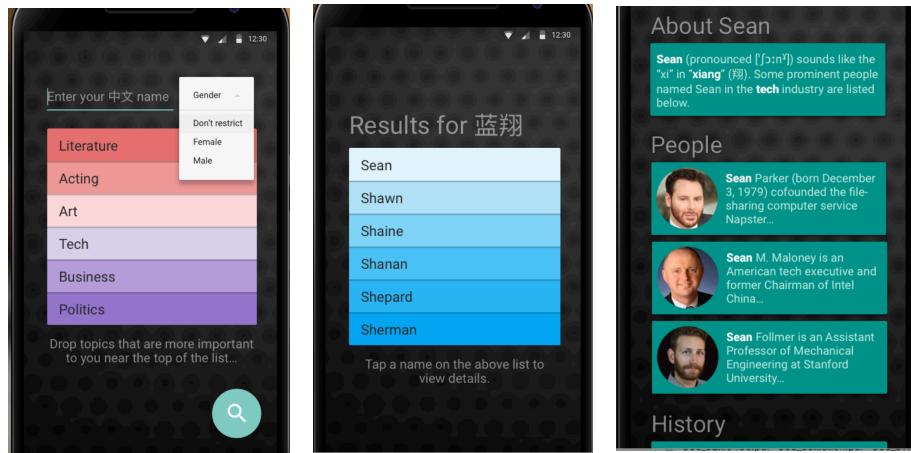
DEPARTMENT OF COMPUTER SCIENCE

Art Search

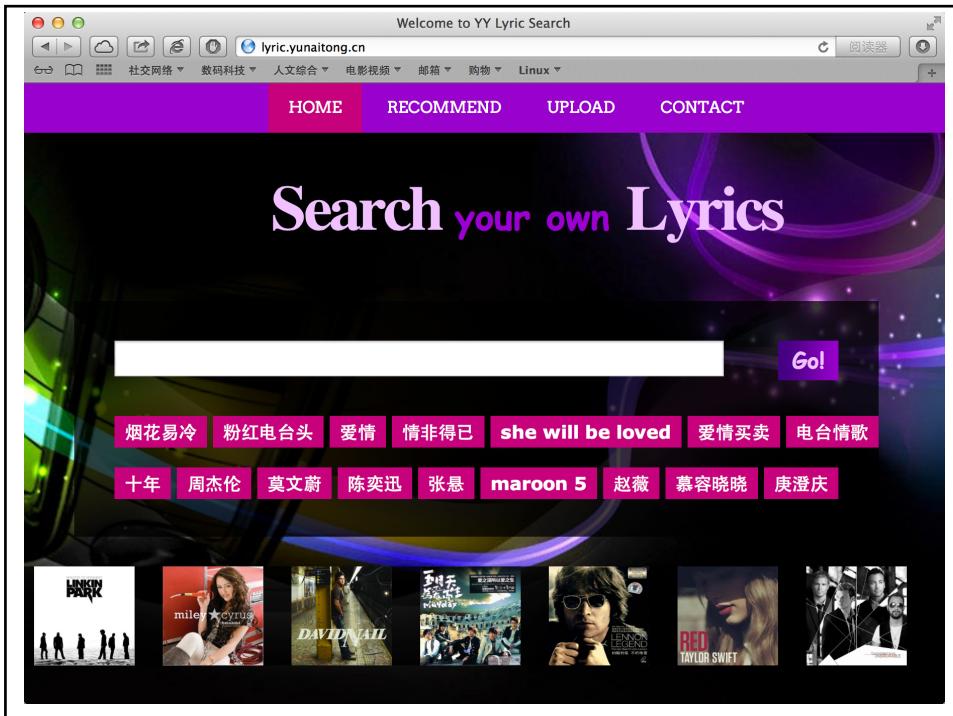
SEARCH

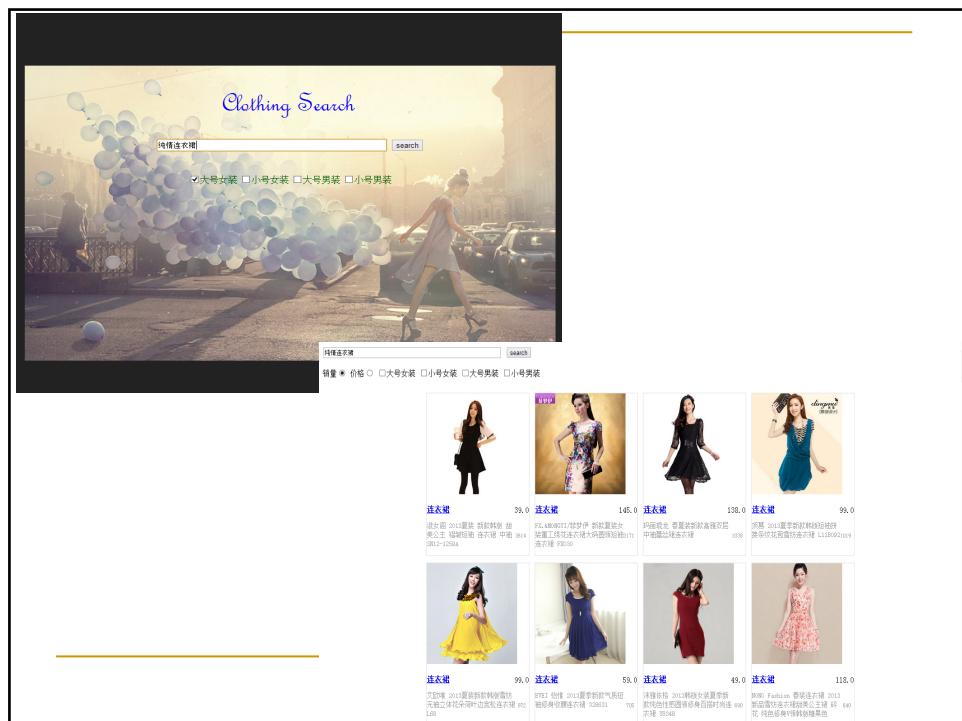
18

Bilingual Name Search Engine



19





CWePS

Search

Tag ClusterVisualization

All Persons

王伟_0

曾国藩

京东 (6)

王伟_1

美国 (4)

王伟_2

中国科学院 (3)

王伟_3

昆山

财政部 (3)

王伟_4

鹿鼎记 (2)

王伟_5

海南

中华网 (2)

王伟_6

新安

中国共产党 (2)

This cluster with 4 documents

1王伟 坠机前美国侦察机上拍摄的录像[视频]_历史频道_凤凰网

>2012年4月1日我们可以看到，王伟 驾驶着低速性能很差的歼-8II，接近飞行速度很慢的美军EP-3。

http://news.ifeng.com/history/zhongguoxiandaishi/detail_2012_04/01/1

2王伟_明星库_娱乐频道_凤凰网

>凤凰网娱乐频道艺人资料库为您提供最新、最全的 王伟 资料、新闻、照片、视频。

<http://app.ent.ifeng.com/star/3730>

9海空卫士 王伟：尽职爱国的飞行员-搜狐视频

> 1 Apr 2013

<http://tv.sohu.com/20130402/n371449684.shtml>

25[视频]海空卫士 王伟：“9时05分”至“9时07分”_新闻频道_央视网

>2013年4月1日他的父母回忆说，王伟 想当飞行员”简直到了痴迷的程度”。他的老家湖州有个空军 机场，只要飞机从城市上空飞过，哪怕是吃饭，王伟 也要丢下碗筷，。

<http://news.cntv.cn/2013/04/01/VIDE1364831560361674.shtml>

Evaluation (Subject to modifications)

- [Homework](#) (0%~15%)
 - [Workshop](#) (55%~70%), evaluated by
 - the other students (30%)
 - + the teacher and TA (30%)
 - [Paper](#) (~20%)
 - [QA activities](#) (~10%)
 - Activities in QA sessions of seminars
 - Question answering in classes
 - [Bonus:](#)
 - **The best project** on the workshop
 - **Tea time presentations** on weekly classes
- Active thinking and discussions are highly encouraged !***

25

What you will learn

- Overview and In-sight of Search (and Recommendation)Technologies
 - Beat 90% general SE users (but still not an expert, to be honest).
- Obtaining knowledge for (web) [information processing](#) and [big data analysis](#) (not only SE and IR), Such as
 - Evaluation strategies
 - Consistencies, correlation, and metrics (precision, recall, AUC, NDCG,)
 - Important classical and state-of-arts models
 - Better understanding on challenges
- Getting to know about [new techniques](#) in Web information services
 - Search, Recommender Systems, User modeling, Social analysis, Visual processing ...
- Be used to always keep an eye on the [new trends of Internet and AI](#)

26

References

- We're not having an official textbook
 - There isn't one with good coverage of all & only the topics we'll discuss
 - A changing field, advanced topics
- A list of references:
 - Books
 - W. Bruce Croft, Donald Metzler, Trevor Strohman, **Search Engine: information retrieval in practice**
 - Christopher D. Manning, Prabhakar Raghavan ,Hinrich Schütze, **Introduction to information retrieval**
 - I. Witten, A. Moffat, and T. Bell, **Managing Gigabytes**
 - Proceedings of Conferences
 - SIGIR, WWW, IJCAI, WSDM, CIKM, TREC, NTCIR ...
 - Very important: Web resources, Search engines

27



What's IR?



Figure Copyright by TREC

What is Information Retrieval (IR)?

■ Narrow-sense:

- IR = Search Engine Technologies (i.e. IR =
 - Google, Yahoo, Bing, Ask, Baidu, Sogou, ...
 - Library info search, enterprise search, in-site search, desktop search...
 - PicSearch, Greplin, Blekko, SkyScanner, KooXoo, Qunar, ...



29

What's IR? (cont.)

■ Broad-sense: IR ~ **Information Management**

- General problem: **how to manage information?**
- How to **find** useful information? (**retrieval & recommendation**)
 - Beyond search engine:
 - e.g. in news feed, movie, travel, e-commerce, financial... scenarios
 - e.g. in social media platform, e.g. Twitter, Facebook, YouTube, WeChat, Weibo, Zhihu,
- How to **organize** information? (**classification & filtering**)
 - e.g., automatically assign email to different folders
- How to **discover** information (or even knowledge) from the data? (**mining**)
 - e.g., discover correlation of events

30

What's IR? (cont.)

- Goal:

- Find documents *relevant* to **an information need** from a large **document set**

- And now:

- Beyond relevance
 - Multi-modal documents
 - **Users' (implicit) information need**
 - Heterogeneous environment



Figure Copyright by TREC

31

IR is Hard!

- Under/over-specified query

- Ambiguous: "buying CDs" (certificate deposit? or compact disc?)
 - Incomplete: what kind of CDs?
 - What if "CD" is never mentioned in document?

- Vague semantics of documents

- Ambiguity: word-sense, structural
 - e.g. "bank"
 - Incomplete: Inferences required
 - E.g. "windows" "apple"

- A difficult task **even for human beings!**

- Only 80% agreement in human judgments

32

IR is “Easy”!

- IR CAN be easy in a particular case
 - Ambiguity in query/document is RELATIVE to the database
 - So, if the query is SPECIFIC enough, just one keyword may get all the relevant documents
- PERCEIVED IR performance is usually better than the actual performance
 - Users can NOT judge the completeness of an answer
 - E.g. Web Search vs. Machine Translation

33