

# Web Mining – Data Mining im Internet

Johannes Fürnkranz

`fuernkranz@informatik.tu-darmstadt.de`

# General Information

- Web-page:
  - <http://www.ke.informatik.tu-darmstadt.de/lehre/ss06/web-mining/>
- Text:
  - Soumen Chakrabarti: *Mining the Web – Discovering Knowledge from Hypertext Data*, Morgan Kaufmann Publishers 2003.
  - Johannes Fürnkranz: *Web Mining*. Draft book chapter with many pointers to the literature
  - Various other articles available from the Web-page
- Slides:
  - available from course page

# Motivation

- The Web is now over 10 years old
  - ca. 1990, Tim Berners-Lee, CERN developed the first graphical hypertext browser
- The information on the Web has grown exponentially
  - on probably every topic you can think of, there is some information available on some Web page
- However, it is still very hard to find relevant information
  - The query interface to search engines has not changed since the early days of the Web!
    - Users have adapted to the interface instead of the other way around



Web [Bilder](#) [Groups](#) [Verzeichnis](#) [News](#) [Froogle](#) <sup>Neu!</sup> [Desktop](#)

wer unterrichtet web mining in Darmstadt

Suche

[Erweiterte Suche](#)  
[Einstellungen](#)

Suche:  Das Web  Seiten auf Deutsch  Seiten aus Deutschland

Die folgenden Wörter kommen sehr häufig vor und wurden daher in Ihrer Suchanfrage ignoriert: **wer in**. [[Einzelheiten](#)]

Web

Ergebnisse **1 - 10** von ungefähr 35 für **wer unterrichtet web mining in Darmstadt**. (0,18 Sekunden)

### [GULP - GULP Profildatenbank: Mitarbeiter von Dienstleistern](#)

... über die Innovationen des gesamten Softwarebereichs **unterrichtet**. ... in Lindau (Bodensee) und Wien, sowie Projektbüros in Stuttgart und **Darmstadt**. ...

[www.gulp.de/itex2/hotlist/itexprofile.html](http://www.gulp.de/itex2/hotlist/itexprofile.html) - 101k - 19. Apr. 2005 - [Im Cache](#) - [Ähnliche Seiten](#)

### [\[PDF\] Informatik Info Oktober 2000](#)

Dateiformat: PDF/Adobe Acrobat - [HTML-Version](#)

... **Darmstadt**, **Web**-Security, 29. Juni 2000, **Darmstadt**, Deutschland. Kryptologie: Von der Geheimwissen- ... **Web Mining**, Electronic Negotiation und ...

[www.ifi.uni-klu.ac.at/Friends/iinfo/info-10-2000/Info-10-00.pdf](http://www.ifi.uni-klu.ac.at/Friends/iinfo/info-10-2000/Info-10-00.pdf) - [Ähnliche Seiten](#)

### [Seminare Ruhrgebiet Oracle Java PHP XML C # sharp C++ VBA ...](#)

... Cuxhaven, Dannstadt-Schauernheim, **Darmstadt**, Dassel, Dattenberg, Deesen, ... Dieser Kurs **unterrichtet** in die Datenzentrierte Anwendung und in **Web** ...

[www.kurse-nrw.de/](http://www.kurse-nrw.de/) - 60k - [Im Cache](#) - [Ähnliche Seiten](#)

### [\[PDF\] Übersicht 7. Semester \(Stand: 24.06.04\) Änderungen möglich!](#)

Dateiformat: PDF/Adobe Acrobat - [HTML-Version](#)

... Einzelthemen und organisatorische Details werden auf der Kursseite im **Web** ... dieses Kurses wird Portfoliomanagement **unterrichtet**, welches eine im ...

[www.ebs.de/uploads/media/Vorlesungsverzeichnis\\_Semester\\_7\\_02.pdf](http://www.ebs.de/uploads/media/Vorlesungsverzeichnis_Semester_7_02.pdf) - [Ähnliche Seiten](#)

### [Telefonmarketing Klatte](#)

... wird in Deutschland beispielsweise von der technischen Hochschule **Darmstadt** betrieben. ... So wird die Empfangsstation von dem Datenstau **unterrichtet**. ...

[www.octokom.de/glossar/glossar.htm](http://www.octokom.de/glossar/glossar.htm) - 520k - [Im Cache](#) - [Ähnliche Seiten](#)

### [\[PDF\] INFORMATIONEN ZUR FORSCHUNGSFÖRDERUNG](#)

Dateiformat: PDF/Adobe Acrobat

... fahren des Information **Mining** sowie zur Integration, Ex- ... Antragstellung **unterrichtet** werden. Seine Kontaktadresse lautet: Im Neuenheimer Feld 366 ...

[www.zuv.uni-heidelberg.de/d6/foerderung/Infor0204.pdf](http://www.zuv.uni-heidelberg.de/d6/foerderung/Infor0204.pdf) - [Ähnliche Seiten](#)

Anzeigen

### [Web Data Extraction](#)

Extract data from target websites  
Save content to your Access MDB  
[www.knowledsys.com](http://www.knowledsys.com)



Web [Bilder](#) [Groups](#) [Verzeichnis](#) [News](#) [Froogle](#) <sup>Neu!</sup> [Desktop](#)

who teaches web mining in darmstadt

Suche

[Erweiterte Suche](#)

[Einstellungen](#)

Suche:  Das Web  Seiten auf Deutsch  Seiten aus Deutschland

Die folgenden Wörter kommen sehr häufig vor und wurden daher in Ihrer Suchanfrage ignoriert: **who in**. [\[Einzelheiten\]](#)

Web

Ergebnisse **1 - 10** von ungefähr 355 für **who teaches web mining in darmstadt**. (0,43 Sekunden)

### [\[PDF\] Chapter # WEB MINING](#)

Dateiformat: PDF/Adobe Acrobat - [HTML-Version](#)

... **WEB MINING**. Johannes Fürnkranz. TU **Darmstadt**, Knowledge Engineering Group ... them to answer queries like "Who **teaches** course X at university Y? " or ...

[www.ke.informatik.tu-darmstadt.de/~juffi/publications/web-mining-chapter.pdf](http://www.ke.informatik.tu-darmstadt.de/~juffi/publications/web-mining-chapter.pdf) -

[Ähnliche Seiten](#)

### [\[PDF\] Knowledge Engineering Group](#)

Dateiformat: PDF/Adobe Acrobat - [HTML-Version](#)

... problems to industrial applications in the areas of data or **web mining** . ... The tutoring system DaMIT **teaches** basics and applications of data **mining** . ...

[www.ke.informatik.tu-darmstadt.de/research/Leaflet.pdf](http://www.ke.informatik.tu-darmstadt.de/research/Leaflet.pdf) - [Ähnliche Seiten](#)

### [Paolo Buono - Home Page](#) - [\[Diese Seite übersetzen\]](#)

... Data **Mining**, Information Visualization, Human-Computer Interaction, **Web**-based ... scientist for several short periods at Fraunhofer IPSI (**Darmstadt**). ...

[lacam.di.uniba.it/8000/people/buono.htm](http://lacam.di.uniba.it/8000/people/buono.htm) - 16k - [Im Cache](#) - [Ähnliche Seiten](#)

### [Langtech 2003 - Knowledge Management & Semantic Web Session](#) - [\[Diese Seite übersetzen\]](#)

... He also **teaches** primary school level at the Dutch School in Oslo (NTC) and is a ... information analysis, document **mining**, information retrieval and ...

[www.lang-tech.org/Speakers%20and%20Presentations/knowledge](http://www.lang-tech.org/Speakers%20and%20Presentations/knowledge) - 27k -

[Im Cache](#) - [Ähnliche Seiten](#)

### [DB: Browsing Object-Oriented Databases over the Web](#) - [\[Diese Seite übersetzen\]](#)

... intelligent agents, data **mining** applications and countless others. ... Conference on the World-Wide **Web**, April 10-14, 1995, **Darmstadt**, Germany. ...

[www.w3.org/Conferences/MWWW4/Papers2/282/](http://www.w3.org/Conferences/MWWW4/Papers2/282/) - 41k - [Im Cache](#) - [Ähnliche Seiten](#)

Anzeigen

### [Web Data Extraction](#)

Extract data from target websites

Save content to your Access MDB

[www.knowledsys.com](http://www.knowledsys.com)

# Hard queries

- For many queries, the information that is needed to answer the query is readily available on the Web:
  - What are the cheapest hotels in Vienna's first district?
- The problems are
  - finding the pages that contain relevant information
    - pages of hotels in Vienna
  - extracting the relevant pieces of information from these pages
    - finding the prices, names, address of these hotels
  - connecting the information that is extracted from the pages
    - comparing the prices, sorting the hotels
  - apply common-sense reasoning in all phases
    - e.g., look for pages of bed & breakfast (Pension) as well
    - know about different currencies and conversions, etc.

# Example Application: Citeseer

- Citeseer is a very popular search engine for publications in Computer Science
  - <http://citeseer.ist.psu.edu/>
- It provides
  - keyword search for articles
  - on-line access to the articles
  - pointers to articles that the articles cites
  - pointers to articles that cite an article
  - pointers to related articles
  - identification of important papers (citation analysis)
  - identification of important publication media
- All of that is generated automatically!

Searching for **PHRASE** web mining.

Restrict to: [Author](#) [Title](#) Order by: [Expected citations](#) [Date](#) Hits: [100](#) Try: [Google \(CiteSeer\)](#) [Google \(Web\)](#) [Yahoo!](#) [MSN](#) [CSB](#) [DBLP](#)  
596 citations found. Retrieving citations...

[Context](#) [Doc](#) **12** (9): Robert Cooley, Bamshad Mobasher, and Jaideep Srivastava. **Web mining: Information and pattern discovery on the world wide web.** In ICTAIS'97, Dec. 1997.

Looking for an author? You may be seeing only a fraction of all citations. Try: [web w/2 mining or w w/2 mining](#) (w/2 means within 2 words)

[Context](#) [Doc](#) **34** (7): B. Mobasher, N. Jain, E-H. Han, and J. Srivastava "**Web mining: Pattern discovery from world wide web transactions,**" Technical Report 96-050, University of Minnesota, Sep. 1996.

[Context](#) [Doc](#) **34** (1): R. Kosala and H. Blockeel, "**Web Mining Research: A Survey,**" in SIGKDD: SIGKDD Explorations: Newsletter of the Special Interest Group (SIG) on Knowledge Discovery & Data **Mining**, ACM, ACM Press, 2000, pp. 1--15.

[Context](#) [Doc](#) **14** (11): O. Nasraoui, H. Frigui, R. Krishnapuram, and A. Joshi. **Mining web access logs using relational competitive fuzzy clustering.** In Eighth International Fuzzy Systems Association Congress, Hsinchu, Taiwan, Aug. 1999.

[Context](#) [Doc](#) **14** (8): M. Craven, S. Slattery, and K. Nigam. **First-order learning for Web mining.** In C. Ndellec and C. Rouveirol, editors, Proceedings of the 10th European Conference on Machine Learning (ECML-98), pages 250--255, Chemnitz, Germany, 1998. Springer-Verlag.

[Context](#) [Doc](#) **12** (0): Karuna P. Joshi, Anupam Joshi, Yelena Yesha, Raghu Krishnapuram, "**Warehousing and Mining Web Logs**", Proc. of 2nd Workshop on **Web** Information and Data Management (WIDM99) (in conj. with CIKM '99), Kansas City (November 1999).

[Context](#) [Doc](#) **10** (3): A. Banerjee and J. Ghosh. **Clickstream Clustering Using Weighted Longest Common Subsequences.** In Proceedings of the **Web Mining** Workshop at the 1st SIAM Conf. on Data **Mining**, pages 34--40, Chicago, IL, April 2001.

[Context](#) [Doc](#) **10** (0): B. Sarwar, G. Karypis, J.A. Konstan, and J.T. Riedl. **Application of Dimensionality Reduction in Recommender System -- A Case Study.** In ACM WebKDD 2000 **Web Mining** for E-Commerce Workshop.

[Context](#) [Doc](#) **10** (0): M. Mulvenna, S. Anand, and A. Buchner. **Personalization on the net using web mining.** CACM, 43(8):122--125, 2000.

[Context](#) [Doc](#) **9** (2): Myra Spiliopoulou. **The laborious way from data mining to web mining.** submitted, June 1998.

citation counts



# Web Mining: Information and Pattern Discovery on the World Wide Web (1997) [\(Make](#)

[Corrections\)](#) [\(82 citations\)](#)

R. Cooley, B. Mobasher, J. Srivastava

**CiteSeer**  
Scientific Literature Digital Library

[Home](#) [Search](#) [Context](#) [Related](#)

View or download:

[depaul.edu/~mobashe...webminertai97.ps](http://depaul.edu/~mobashe...webminertai97.ps)  
[depaul.edu/~mobasher/WebKI...cmstai.ps](http://depaul.edu/~mobasher/WebKI...cmstai.ps)  
[depaul.edu/~mobasher/clas...cmstai.pdf](http://depaul.edu/~mobasher/clas...cmstai.pdf)  
Cached: [PS.gz](#) [PS](#) [PDF](#) [Image](#) [Update](#) [Help](#)

From: [depaul.edu/~mobasher/pubs](http://depaul.edu/~mobasher/pubs) [\(more\)](#)  
Homepages: [R.Cooley](#) [HPSearch](#) [\(Update Links\)](#)

Rate this article: 1 2 3 4 5 (best)  
[Comment on this article](#)

[\(Enter summary\)](#)

**Abstract:** Application of data mining techniques to the World Wide Web has been the focus of several recent research projects and papers. The term Web mining has been used in two distinct ways. The first, called Web content mining, is the process of information discovery from sources across the World Wide Web. The second, called Web usage mining, is the process of mining for user browsing and access patterns.

**Cited by:** [More](#)

WUM: A Tool for Web Utilization Analysis - Myra Spillo  
P-Jigsaw: Extending Jigsaw with Rules Assisted Cache  
Combining Web Usage Mining and Fuzzy Inference for

**Similar documents (at the sentence level):**

8.5%: [Mir: A Tool For Visual Presentation Of Web Acc](#)  
5.5%: [Web Mining: Pattern Discovery from World Wide](#)

**Active bibliography (related documents):** [More](#) [All](#)

0.7: [Grouping Web Page References into Transactions](#)  
0.5: [Document Categorization and Query Generation on](#)  
0.5: [Software Environments in Support of Wide-Area Di](#)

**Similar documents based on text:** [More](#) [All](#)

0.8: [Some Experiences on Large Scale Web Mining - I](#)  
0.7: [Blockmodeling Techniques for Web Mining - Schoi](#)  
0.6: [Usage Mining for and on the Semantic Web - Stun](#)

**Related documents from co-citation:** [More](#) [All](#)

25: [Data preparation for mining world wide web browsi](#)  
24: [Fast Algorithms for Mining Association Rules - Agr](#)  
20: [From user access patterns to dynamic hypertext li](#)

**BibTeX entry:** [\(Update\)](#)

## Web Mining: Information and Pattern Discovery on the World Wide Web \*

R. Cooley, B. Mobasher, and J. Srivastava

Department of Computer Science and Engineering  
University of Minnesota  
Minneapolis, MN 55455, USA

### Abstract

Application of data mining techniques to the World Wide Web, referred to as Web mining, has been the focus of several recent research projects and papers. However, there is no established vocabulary, leading to confusion when comparing research efforts. The term Web mining has been used in two distinct ways. The first, called Web content mining in this paper, is the process of information discovery from sources across the World Wide Web. The second, called Web usage mining, is the process of mining for user browsing and access patterns. In this paper we define Web mining and present an overview of the various research issues, techniques, and development efforts. We briefly describe WEBMINER, a system for Web usage mining, and conclude this paper by listing research issues.

## 1 Introduction

With the explosive growth of information sources available on the World Wide Web, it has become increasingly necessary for users to utilize automated tools in find the desired information resources, and to track and analyze their usage patterns. These factors

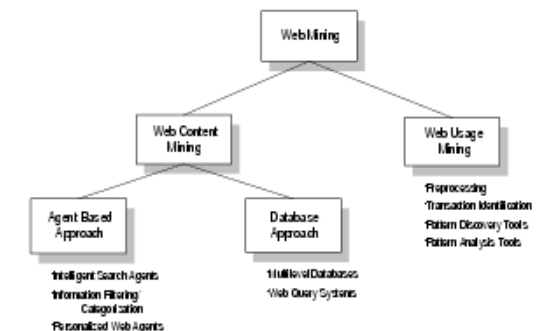


Figure 1: Taxonomy of Web Mining

context. There are several important issues, unique to the Web paradigm, that come into play if sophisticated types of analyses are to be done on server side data collections. These include integrating various data sources such as server access logs, referrer logs, user registration or profile information; resolving difficulties in the identification of users due to missing unique key attributes in collected data; and the importance of identifying user sessions or transactions

82 citations found. Retrieving documents...

Robert Cooley, Bamshad Mobasher, and Jaideep Srivastava. *Web mining: Information and Dec.* 1997.

**CiteSeer** [Home/Search](#) [Document Details and Download](#) [Summary](#) [Related Article](#)

This paper is cited in the following contexts:

[First 50 documents](#) [Next 50](#)

[Low-Complexity Fuzzy Relational Clustering - Algorithms For Web](#) (Correct)

...In particular, Han et al. [36] create a MOI AP based warehouse from Web logs and allow users to **time dependent patterns in the access logs** [9] [10]. However, both these approaches are used and the clients are willing to release

However, it is not clear how the similar clusters. There is also a recent body structured, database-like entity. In particular, Web logs, and allow users to perform patterns in the access logs [53]. Similar [1], have been proposed in [9], [10].

ids, which is not true in the real world. the clients are willing to release the momentum is the idea that we can let their *clickstreams*, which is of great interest

An important component of personalization is the extraction of structure from unlabeled information. The logs kept by Web servers can be viewed as a special case of the mining. It can be said to have three operations

## Low-Complexity Fuzzy Relational Clustering Algorithms for Web Mining

Raghu Krishnapuram  
IBM India Research Lab  
Indian Institute of Technology, Hauz Khas, New Delhi 110016  
kraghura@in.ibm.edu

On leave from Dept of Mathematical and Computer Sciences, Colorado School of Mines, Golden, CO 80401

Anupam Joshi  
Department of Computer Science and Electrical Engineering  
University of Maryland Baltimore County, Baltimore, MD 21250  
joshi@cs.umbc.edu

Olfa Nasraoui  
Department of Electrical Engineering  
University of Memphis, Memphis, TN 38152  
Liyu Yi

### REFERENCES

- [1] R. Agrawal and R. Srikant. Fast algorithms for mining association rules. In *Proceedings of the 20th VLDB Conference*, pages 487–499, Santiago, Chile, 1994.
- [2] R. Armstrong, T. Joachims D. Freitag, and T. Mitchell. Webwatcher: A learning apprentice for the World Wide Web. In *Proceedings of the AAAI Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*, pages 6–13, Stanford, CA, March 1995.
- [3] G. Arocena and A. Mendelz. Webowl: Restructuring documents, databases, and web. In *Proc. IEEE Intl. Conf. Data Engineering '98*, pages 24–33. IEEE Press, 1998.
- [4] P. Bajcsy and N. Ahuja. Location- and density-based hierarchical clustering using similarity analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20:1011–1015, 1998.
- [5] G. Beni and X. Liu. A least biased fuzzy clustering method. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 16:954–960, September 1994.
- [6] J. C. Bezdek. *Pattern Recognition with Fuzzy Objective Function Algorithms*. Plenum Press, New York, 1981.
- [7] J. Abidi C. Shahabi, A.M. Zarkesh and V. Shah. Knowledge discovery from users web-page navigation. In *Proceedings of the Seventh IEEE Intl. Workshop on Research Issues in Data Engineering (RIDE)*, pages 20–29, Birmingham, UK, 1997.
- [8] J. Chen, A. Mikulcic, and D. H. Kraft. An integrated approach to information retrieval with fuzzy clustering and fuzzy inferencing. In O. Pons, M. Ampara Vila, and J. Kacprzyk, editors, *Knowledge Management in Fuzzy Databases*, volume 163. Physica Verlag, Heidelberg, Germany, 2000.
- [9] M.S. Chen, J.-S. Park, and P. S. Yu. Efficient data mining for path traversal patterns. *IEEE Trans. Knowledge and Data Engineering*, 10(2):209–221, April 1998.
- [10] R. Cooley, B. Mobasher, and J. Srivastav. Web Mining: Information and pattern discovery on the World Wide Web. In *Proc. IEEE Intl. Conf. Tools with AI*, pages 558–567, Newport Beach, CA, 1997.
- [11] R. N. Davé and R. Krishnapuram. Robust clustering methods: A unified view. *IEEE Transactions on Fuzzy Systems*, 5(2):270–293, 1997.
- [12] E. Diday. La methode des nuées dynamiques. *Rev. Stat. Appliquee*, XIX(2):19–34, 1975.
- [13] D. Riecken: Guest Editor. Special issue on personalization. *Communications of the ACM*, 43(9), Sept. 2000.
- [14] J. Fink, A. Kobsa, and J. Schreck. Personalized hypermedia information provision through adaptive and adaptable system features. <http://zeus.gmd.de/hci/projects/avanti/publications/ISandN97/ISandN97.html>, 1997.
- [15] K. S. Fu. *Syntactic Pattern Recognition and Applications*. Academic Press, San Diego, CA, 1982.
- [16] K. C. Gowda and E. Diday. Symbolic clustering using a new similarity measure. *IEEE Transactions on Systems, Man, and Cybernetics*, 20:368–377, 1992.
- [17] S. Guha, R. Rastogi, and K. Shim. CURE: An efficient algorithm for large databases. In *Proceedings of SIGMOD '98*, pages 73–84, Seattle, June 1998.
- [18] R. J. Hathaway and J. C. Bezdek. Switching regression models and fuzzy clustering. *IEEE Transactions on Fuzzy Systems*, 1(3):195–204, 1993.

### Citations (may not include all citations):

- 866 Fast algorithms for mining association rules - Agrawal, Srikant - 1994
- 359 Data cube: A relational aggregation operator generalizing gr. - Gray, Bosworth et al. - 1991
- 321 A query language and optimization techniques for unstructure.. - Buneman, Davidson et al
- 262 Finding Groups in Data: an Introduction to Cluster Analysis (context) - Kaufman, Rousseeur
- 239 Efficient and effective clustering method for spatial data m.. - Ng, Han - 1994
- 236 Implementing data cubes efficiently - Harinarayan, Rajaraman et al. - 1996
- 235 Information Retrieval Data Structures and Algorithms (context) - Frakes, Baeza-Yates - 19!
- 198 Webwatcher: A learning apprentice for the world wide web - Armstrong, Freitag et al. - 199
- 183 Discovering frequent episodes in sequences (context) - Mannila, Toivonen et al. - 1995
- 174 word of mouth (context) - Shardanand, Maes et al. - 1995
- 169 A scalable comparison shopping agent for the world wide web - Doorenbos, Etzioni et al. -
- 164 the computation of multidimensional aggregates - Agrawal, Agrawal et al. - 1996
- 162 An efficient algorithm for mining association rules in large.. (context) - Savasere, Omiecins
- 154 Mining sequential patterns: Generalizations and performance .. - Srikant, Agrawal - 1996
- 144 Wq query system world wide web - Shmueli, system et al. - 1995
- 116 A declarative language for querying and restructuring the we.. - Lakshmanan, Sadri et al. -
- 114 Syntactic clustering of the web (context) - Broder, Glassman et al. - 1997
- 114 Data-driven discovery of quantitative rules in relational da.. (context) - Han, Cai et al. - 199!
- 113 webert: Identifying interesting web sites (context) - Pazzani, Muramatsu et al. - 1996
- 107 Silk from a sow's ear: Extracting usable structures from the.. - Pirolli, Pitkow et al. - 1996
- 100 Querying semistructured heterogeneous information - Quass, Rajaraman et al. - 1995
- 99 Computer Systems that Learn: Classification and Prediction M.. (context) - Weiss, Kulikow
- 89 Planning to gather information - Kwok, Weld - 1996
- 87 The information manifold - Kirk, Levy et al. - 1995
- 82 Web mining: Information and pattern discovery on the world w.. - Cooley, Mobasher et al. -
- 71 Parasite: mining structural information on the web (context) - Spertus - 1997
- 64 Dmql: A data mining query language for relational databases - Han, Fu et al. - 1996
- 53 Category translation: learning to understand information on .. - Perkowitz, Etzioni - 1995
- 53 Storage estimation for multidimensional aggregates in the pr.. - Shukla, Deshpande et al. -
- 50 Hypursuit: a hierarchical network search engine that exploit.. (context) - Weiss, Velez et al.
- 45 The tsimmis project: Integration of heterogenous information.. (context) - Chawathe, Garcia-
- 42 Semistructured and structured data in the web: Going back an.. - Merialdo, Atzeni et al. - 1!
- 42 Aliweb - archie-like indexing in the web (context) - Koster - 1994
- 41 Web mining: Pattern discovery from world wide web transactio.. - Mobasher, Jain et al. - 19!
- 36 Data mining for path traversal patterns in a web environment - Chen, Park et al. - 1996
- 28 An adaptive agent for automated web browsing - Balabanovic, Shoham et al. - 1995
- 22 Finding salient features for personal web page categorizatio.. - Wulfekuhler, Punch - 1997
- 22 Faq-finder: A case-based approach to knowledge navigation (context) - Hammond, Burke et
- 21 Automatically organizing bookmarks per content (context) - Maarek, Shaul - 1996

### References

- [1] R. Agrawal and R. Srikant. Fast algorithms for mining association rules. In *Proc. of the 20th VLDB Conference*, pages 487–499, Santiago, Chile, 1994.
- [2] S. Agrawal, R. Agrawal, P.M. Deshpande, A. Gupta, J. Naughton, R. Ramakrishna, and S. Sarawagi. On the computation of multidimensional aggregates. In *Proc. of the 22nd VLDB Conference*, pages 506–521, Mumbai, India, 1996.
- [3] R. Armstrong, D. Freitag, T. Joachims, and T. Mitchell. Webwatcher: A learning apprentice for the world wide web. In *Proc. AAAI Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*. 1995.
- [4] M. Balabanovic, Yoav Shoham, and Y. Yun. An adaptive agent for automated web browsing. *Journal of Visual Communication and Image Representation*, 6(4), 1995.
- [5] A. Z. Broder, S. C. Glassman, M. S. Manasse, and G. Zweig. Syntactic clustering of the web. In *Proc. of 6th International World Wide Web Conference*, 1997.
- [6] C. M. Brown, B. B. Danzig, D. Hardy, U. Manber, and M. F. Schwartz. The harvest information discovery and access system. In *Proc. 2nd International World Wide Web Conference*, 1994.
- [7] P. Buneman, S. Davidson, G. Hillebrand, and D. Suciu. A query language and optimization techniques for unstructured data. In *Proc. of 1996 ACM-SIGMOD Int. Conf. on Management of Data*, 1996.
- [8] P. Buneman, S. Davidson, and D. Suciu. Programming constructs for unstructured data. In *Proceedings of ICDT'95, Gubbio, Italy*, 1995.
- [9] C. Chang and C. Hsu. Customizable multi-engine search tool with clustering. In *Proc. of 6th International World Wide Web Conference*, 1997.

[CiteSeer.IST Home](#) **Check:** The following citations are predicted to all refer to the same paper. [Details](#)

COOLEY, R., SRIVASTAVA, J., MOBASHER, B., *Web Mining: Information and Pattern Discovery on the World Wide Web*, Proceedings of the 9th IEEE International Conference on Tools with Artificial Intelligence (ICTAI97), November 1997.

Robert Cooley, Bamshad Mobasher, and Jaideep Srivastava. *Web mining: Information and pattern discovery on the world wide web*. In ICTAI97, Dec. 1997.

R. Cooley, B. Mobasher, and J. Srivastava. *Web mining: Information and pattern discovery on the world wide web*. Technical Report TR 97-027, University of Minnesota, Dept. of Computer Science, Minneapolis, 1997.

Robert Cooley, Bamshad Mobasher, and Jaideep Srivastava. *Web mining: Information and pattern discovery on the world wide web*. In ICTAI97, Dec. 1997.

Cooley, R., Mobasher, R. & Srivastava, J. (1997) *Web Mining: Information and Pattern Discovery on the World Wide Web*, Proc. 9 th IEEE Int'l Conf. on Tools with Artificial Intelligence.

Cooley, R., Mobasher, B., and Srivastava, J. (1997b). *Web mining: Information and pattern discovery on the world wide web*. In ICTAI97.

R. Cooley, B. Mobasher, and J. Srivastava, "*Web mining: Information and Pattern discovery on the World Wide Web*," Proc. IEEE Intl. Conf. Tools with AI, Dec, 1997.

R. Cooley, B. Mobasher, and J. Srivastava. *Web mining: Information and patterns discovery on the world wide web*. In Proc. of the 9th IEEE Int. Conf. on Tools with Artificial Intelligence, pages 558-567, 1997.

R. Cooley, B. Mobasher and J. Srivastava. *Web Mining: Information and Pattern Discovery on the Word Wide Web*. Technical Report TR 97-027, University of Minnesota, Dept. of Computer Science, Minneapolis, 1997.

R. Cooley, B. Mobasher, and J. Srivastava. *Web mining: Information and pattern discovery on the world wide web*. In International Conference on Tools with Artificial Intelligence, pages 558-567, Newport Beach, 1997. IEEE.

Robert Cooley, Bamshad Mobasher, and Jaideep Srivastava. *Web mining: Information and patterns discovery on the world wide web*. In Proc. of the ninth IEEE International Conference on Tools with Artificial Intelligence (ICTAI97), November 1997.

R. Cooley, B. Mobasher, and J. Srivastava. *Web mining: Information and pattern discovery on the world wide web*. In International Conference on Tools for Artificial Intelligence, Newport Beach, CA, November 1997.

Cooley, R., Mobasher, B., and Srivastava, J. (1997). *Web mining: Information and pattern discovery on the world wide web*. In International Conference on Tools for Artificial Intelligence, Newport Beach, CA.



# Task that need to be solved

- Information Retrieval
  - search for research papers on the Web
- Information Extraction
  - extract relevant information (title, author, journal/conference, publication year,...) from the research papers
  - extract citations from the research papers
- Information Integration
  - match extracted citations with the text where they are cited
  - match extracted citations with other extracted citations
  - identify similar documents
- Citation analysis
  - build and analyze a graph of citations of papers
  - build and analyze a co-authorship graph
- and many more...

# Web Mining

Web Mining is Data Mining for Data on the World-Wide Web

- Text Mining:
  - Application of Data Mining techniques to unstructured (free-format) text
- Structure Mining:
  - taking into account the structure of (semi-)structured hypertext (HTML tags, hyperlinks)
- Usage Mining:
  - taking into account user interactions with the text data (click-streams, collaborative filtering, ...)

# Web Mining Tasks

- Message Filter or Message Sorter
- Intelligent Browsing Assistants
- Formation or Update of Web Catalogues
- Ranking or Clustering of Search Results
- Building the Semantic Web / World-Wide Knowledge Base
- Click-stream Analysis
- Product Recommendations
- Digital libraries and Citation Analysis
- ...

# The Web

- The Web is a unique kind of hypertext document
  - a large number of pages
  - on a wide variety of topics
  - originating by a large variety of authors
  - speaking many different languages
  - annotated via hyperlinks
  - accessible to everybody
- Main Problem:
  - How can I find the information I am looking for?
- Web Mining:
  - finding and extracting relevant information from the Web



# A Brief History of Hypertext

- On Paper

- Annotated books (e.g., the Talmud)
- Dictionaries and encyclopedias
  - cross-references are hyperlinks
- Scientific literature
  - citations of other works is another form of hyperlinks

- Electronic

- Memex (Vannevar Bush, 1945)
  - design for a photo-electrical, mechanical storage device that could link documents
  - On-line Demo <http://www.dynamicdiagrams.com/demos/memex1a.zip>
- Xanadu (Engelbart & Nelson 1965) <http://xanadu.com/>
  - first conventional hypertext system, also pioneered wikis
  - too complex to be realized, first use of word „hypertext“
- Many successor systems



# A Brief History of the Web

- Tim Berners-Lee (CERN)
  - first proposals around 1980
  - 1990: work on the „World Wide Web“
  - first graphical interfaces
- 1993:
  - Mosaic (Mark Andressen, NCSA): intuitive hypertext GUI for UNIX
  - HTML: hypertext markup language
  - HTTP: hypertext transport protocol
- 1994:
  - Netscape was founded
  - 1<sup>st</sup> World Wide Web Conference <http://www.w3.org/>
  - World Wide Web Consortium founded by CERN and MIT

# HTTP (hypertext transport protocol)

- Built on top of the Transport Control Protocol (TCP)
- Steps(from client end) <http://www.w3.org/Protocols>
  - resolve the server host name to an Internet address (IP)
    - Use Domain Name Server (DNS)
    - DNS is a distributed database of name-to-IP mappings maintained at a set of known servers
  - contact the server using TCP
    - connect to default HTTP port (80) on the server.
    - Enter the HTTP requests header (E.g.: GET)
    - Fetch the response header
      - MIME (Multipurpose Internet Mail Extensions)
        - A meta-data standard for email and Web content transfer
    - Fetch the HTML page

## Host Port

```
% telnet www.cse.iitb.ac.in 80
Trying 144.16.111.14...
Connected to www.cse.iitb.ac.in.
Escape character is '^]'.
GET / Http/1.0
```

GET / Http/1.0

↑  
Pfad

Header

```
Http/1.1 200 OK
Date: Sat, 13 Jan 2001 09:01:02 GMT
Server: Apache/1.3.0 (Unix) PHP/3.0.4
Last-Modified: Wed, 20 Dec 2000 13:18:38 GMT
ETag: "5c248-153d-3a40b1ae"
Accept-Ranges: bytes
Content-Length: 5437
Connection: close
Content-Type: text/html
X-Pad: avoid browser bug
```

HTML  
of Web  
page

```
<html>
<head><title>IIT Bombay CSE Department Home Page</title></head>
<body>...<a href="http://www.iitb.ac.in">IIT Bombay</a>...
</body></html>
Connection closed by foreign host.
```

# HTML

<http://www.w3.org/MarkUp/>

- HyperText Markup Language
- Lets the author
  - specify document structure
    - browser converts structure to layout
    - direct specification of layout and typeface possible
  - embed diagrams
  - create hyperlinks.
    - expressed as an anchor tag with a HREF attribute
    - HREF names another page using a Uniform Resource Locator (URL),
  - URL =
    - protocol field (“HTTP”) +
    - a server hostname (“www.cse.iitb.ac.in”) +
    - file path (/, the `root' of the published file system).

# DOM Tree

- DOM = Document Object Model      <http://www.w3.org/DOM/>
- An HTML document can be viewed as a tree
  - markup items are interior nodes
  - text are leafs
  - Xpath: language for denoting the path from the root to a tree  
<http://www.zvon.org/xxl/XPathTutorial/General/examples.html>
- document structure can be exploited
  - sectioning of documents
  - recognition of important text parts (e.g., anchor text)
  - structural patterns (XPath) may identify important information on the page
- Firefox->Tools/Web Development/DOM Inspector

# Web: A populist, participatory medium

- number of writers =(approx) number of readers.
- the evolution of MEMES
  - ideas, theories etc that spread from person to person by imitation.
  - good memes survive, bad memes die out
- but the Web archives them all

# Abundance and authority crisis

- liberal and informal culture of content generation and dissemination.
  - despite a few commercial niches we still have anarchy
- Very little uniform civil code.
- redundancy and non-standard form and content.
- millions of qualifying pages for most broad queries
  - Example: java or kayaking
- no authoritative information about the reliability of a site



# Problems due to Uniform accessibility

- little support for adapting to the background of specific users.
- commercial interests routinely influence the operation of Web search
  - “Search Engine Optimization“ !!

# Data Mining - Motivation

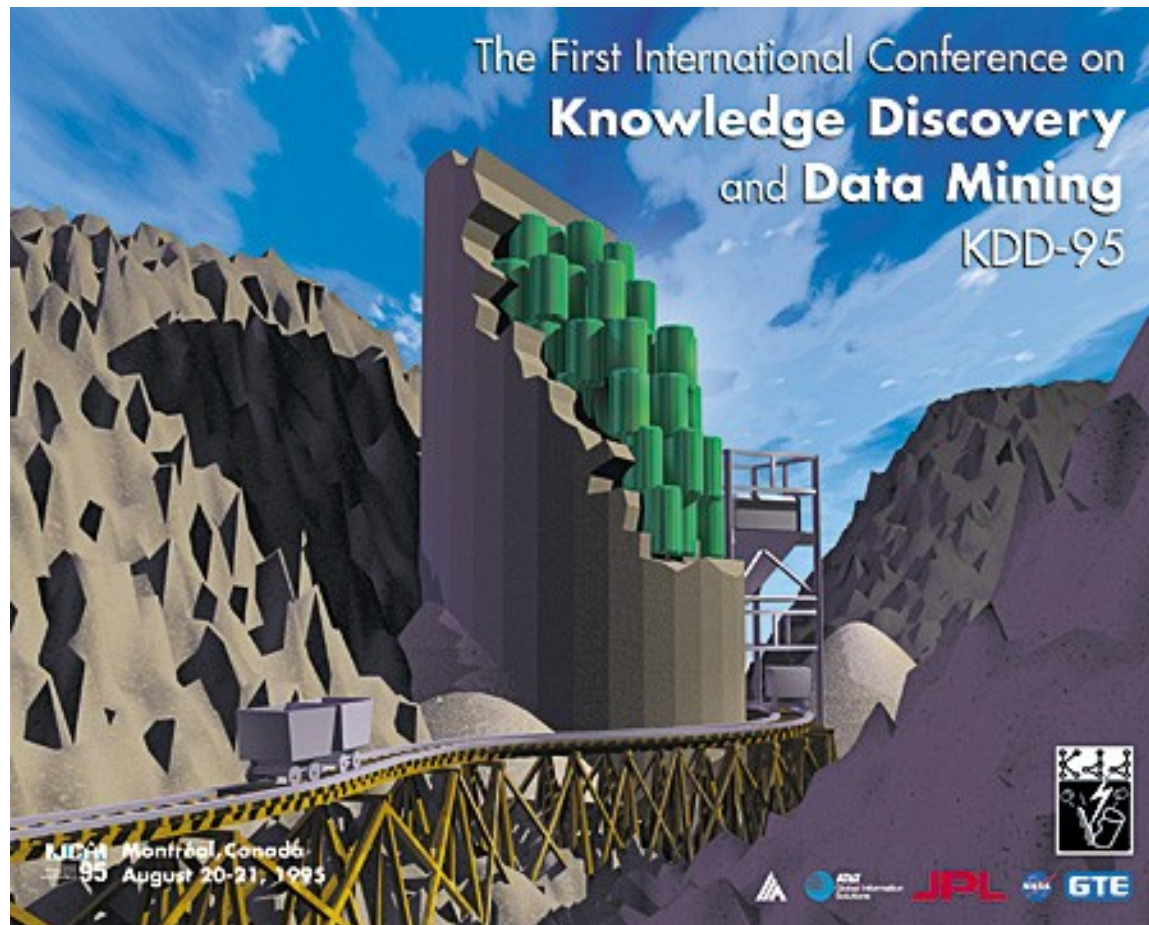
"Computers have promised us a fountain of wisdom but delivered a flood of data."

"It has been estimated that the amount of information in the world doubles every 20 months."

*(Frawley, Piatetsky-Shapiro, Matheus, 1992)*

# Data Mining

Mining for nuggets of knowledge in mountains of Data.



# Definition

Data Mining is a non-trivial *process* of identifying

- valid
- novel
- potentially useful
- ultimately understandable

patterns in data.

*(Fayyad et al. 1996)*

It employs techniques from

- machine learning
- statistics
- databases

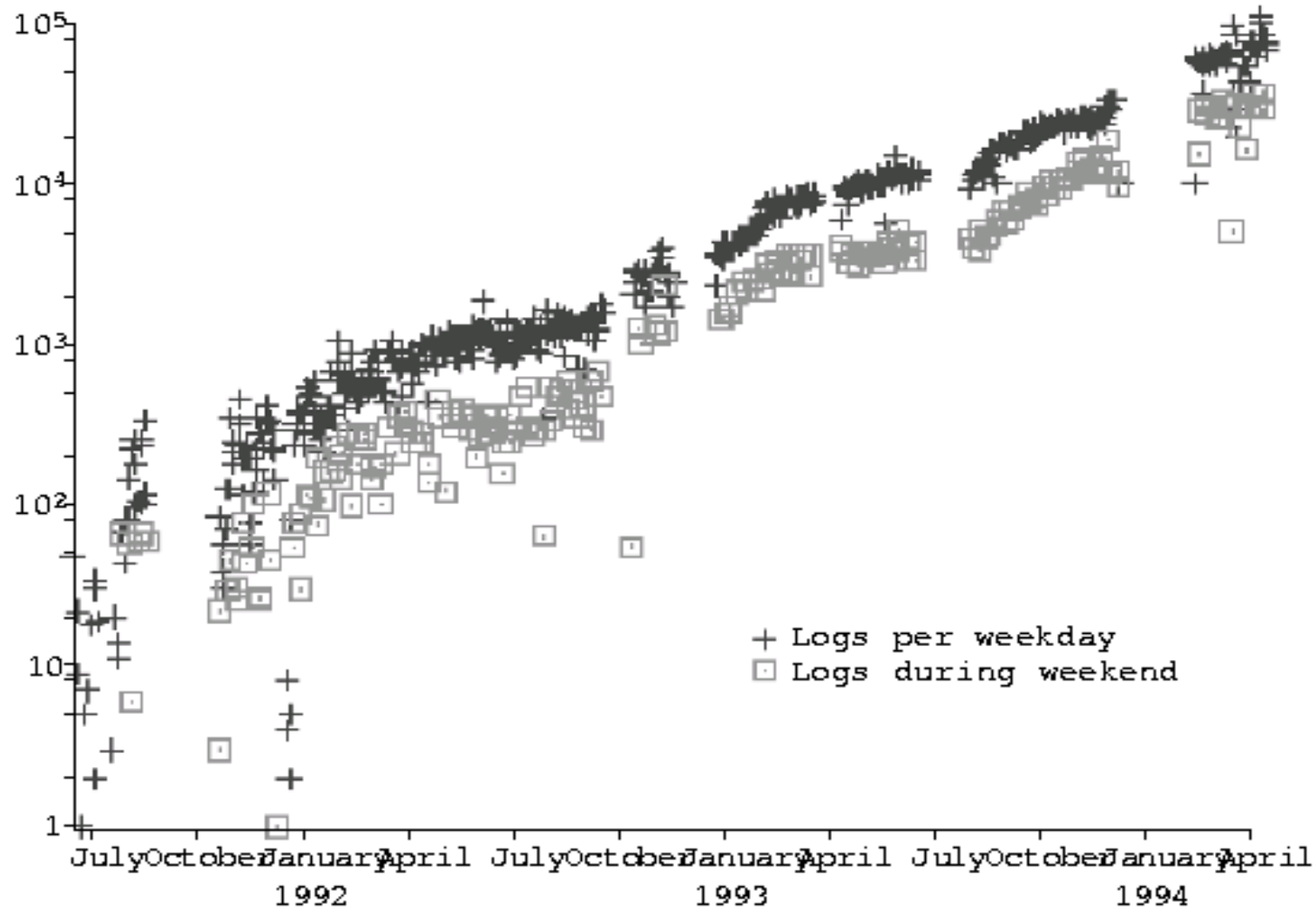
Or maybe:

- Data Mining is torturing your database until it confesses.

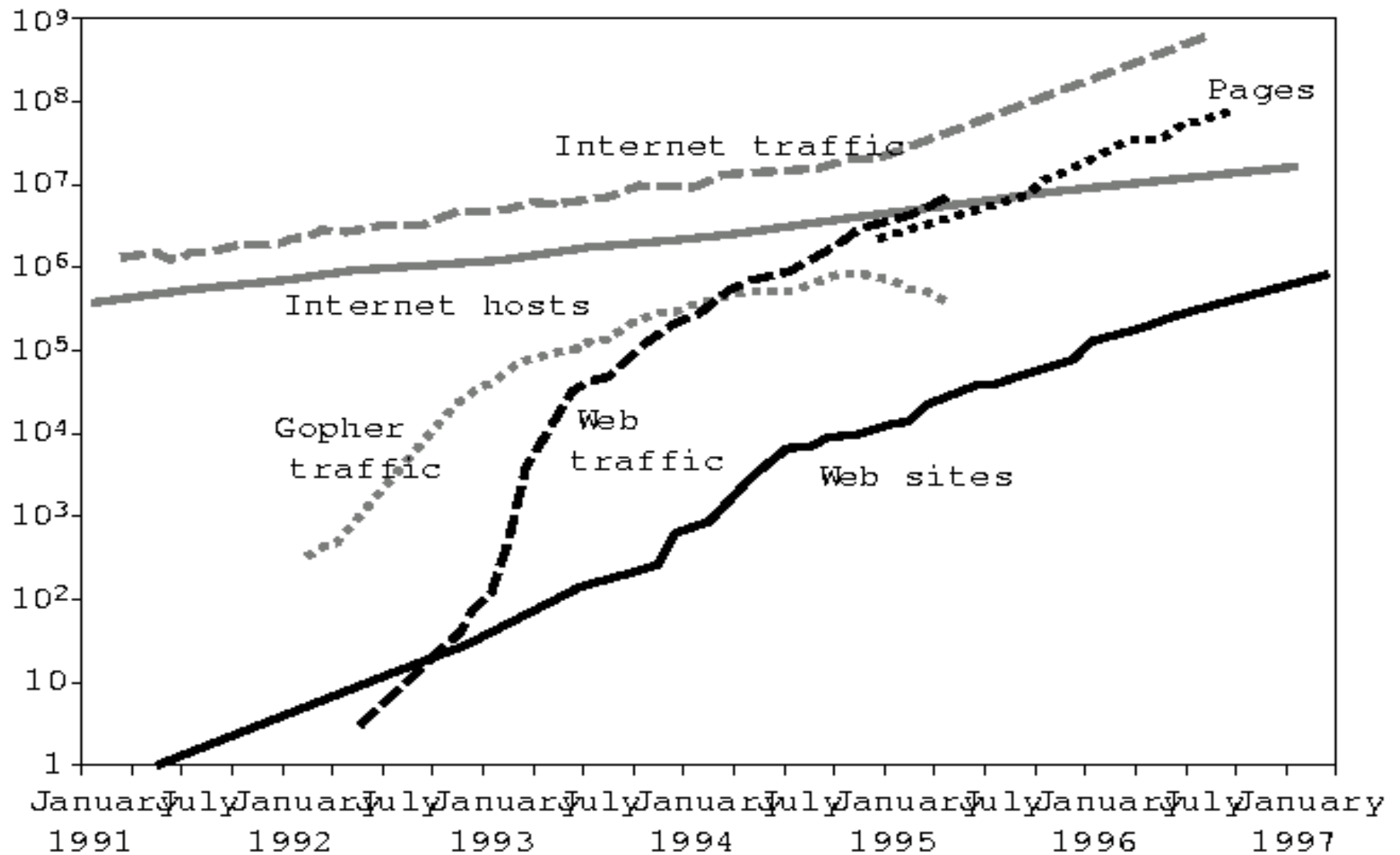
*(Mannila (?))*

# World-Wide Data Growth

- Science
  - satellite monitoring
  - human genome
- Business
  - OLTP (on-line transaction processing)
  - data warehouses
  - e-commerce
- Industry
  - process data
- World-Wide Web



The early days of the Web : CERN HTTP traffic grows by 1000 between 1991–1994 (image courtesy W3C)



The early days of the Web: The number of servers grows from a few hundred to a million between 1991 and 1997 (image courtesy Nielsen)

# How Big is the Web?

- Google:
  - early 2001: 1,346,966,000 web pages
  - 11.2.2002: 2,073,418,204
  - 2004: 4,285,199,774
  - 28.4.2005: 8,058,044,651
- Size of the Web
  - Results from 1998 estimate that the best search engines index about 30% of the Web.
- Gulli & Signorini (2005)
  - estimate the size of the Web to 11.5 billion pages,
  - Coverage of search engines
    - Google=76.16%, Msn Beta=61.90%, Ask/Teoma=57.62%, Yahoo!=69.32%



# Structured vs. Web data mining

- traditional data mining
  - data is structured and relational
  - well-defined tables, columns, rows, keys, and constraints.
- Web data
  - semi-structured and unstructured
  - readily available
  - rich in features and patterns
  - spontaneous formation and evolution of
    - topic-induced graph clusters
    - hyperlink-induced communities

# Structured Data

- Attribute-Value data:
  - Each example is described with values for a fixed number of attributes
    - **Nominal Attributes:**
      - store an unordered list of symbols (e.g., *color*)
    - **Numeric Attributes:**
      - store a number (e.g., *income*)
    - **Other Types:**
      - hierarchical attributes
      - set-valued attributes
  - the data corresponds to a single relation (spreadsheet)
- Multi-Relational data:
  - The relevant information is distributed over multiple relations
  - Inductive Logic Programming

# Structured Data

<i>Day</i>	<i>Temperature</i>	<i>Outlook</i>	<i>Humidity</i>	<i>Windy</i>	<i>Play Golf?</i>
07-05	hot	sunny	high	false	no
07-06	hot	sunny	high	true	no
07-07	hot	overcast	high	false	yes
07-09	cool	rain	normal	false	yes
07-10	cool	overcast	normal	true	yes
07-12	mild	sunny	high	false	no
07-14	cool	sunny	normal	false	yes
07-15	mild	rain	normal	false	yes
07-20	mild	sunny	normal	true	yes
07-21	mild	overcast	high	true	yes
07-22	hot	overcast	normal	false	yes
07-23	mild	rain	high	true	no
07-26	cool	rain	normal	true	no
07-30	mild	rain	high	false	yes

today	cool	sunny	normal	false	?
tomorrow	mild	sunny	normal	false	?

# Semi-Structured and Unstructured Data

- Semi-structured Data
  - no clear tables
    - it may be hard to identify the attributes for each example
    - it may also be hard to identify the examples themselves
  - some structure implicit in the data
    - e.g., formatting via HTML
  - large parts without structure
    - free text
  - <http://weather.yahoo.com/forecast/GMXX0020.html>

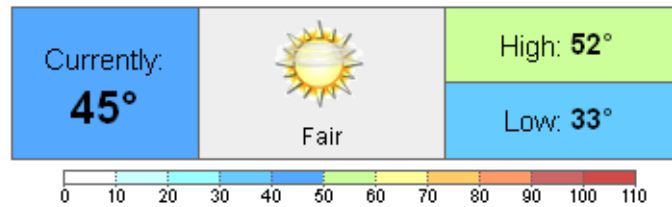
# Semi-Structured

## Darmstadt Weather

at 9:50 am CEST

F° | C°

[Text Forecast](#)



### 5 Day Forecast

Today	Tomorrow	Sat	Sun	Mon	6-10 Day
					<a href="#">Extended Forecast</a>
Sunny	Sunny	PM Showers	Light Rain	Light Rain	
High: <b>52°</b> Low: <b>33°</b>	High: <b>57°</b> Low: <b>38°</b>	High: <b>63°</b> Low: <b>38°</b>	High: <b>61°</b> Low: <b>47°</b>	High: <b>56°</b> Low: <b>45°</b>	

Featured Forecasts at weather.com:

[Allergies](#) | [Golf](#) | [Driving Conditions](#)

### More Current Conditions

<b>Feels Like:</b>	45°	<b>Dewpoint:</b>	28°
<b>Barometer:</b>	30.09 in and steady	<b>Wind:</b>	NNE 9 mph
<b>Humidity:</b>	53%	<b>Sunrise:</b>	6:21 am
<b>Visibility:</b>	9.99 mi	<b>Sunset:</b>	8:28 pm

### Local Forecast - ([How to Read This](#))

**Today:** Abundant sunshine. High 52F. Winds NE at 5 to 10 mph.

**Tonight:** Mainly clear. Cold. Low 33F. Winds ENE at 5 to 10 mph.

**Tomorrow:** Mainly sunny. High 57F. Winds ESE at 5 to 10 mph.

**Tomorrow night:** A few clouds from time to time. Low 38F. Winds light and variable.

**Saturday:** Showers possible in the afternoon. Highs in the low 60s and lows in the upper 30s.

**Sunday:** Light rain. Highs in the low 60s and lows in the upper 40s.

### Sponsored Links

[Darmstadt, Germany](#)

Pioneer Military Loans, offering loans up to \$10,000, 24 hours, 7 days a week worldwide for active and retired military and Federal GS employees.

[www.themilitaryzone.com](http://www.themilitaryzone.com)

[Darmstadt Germany Tourism Information](#)

Visit our site for information on German Cities, Hotels, Restaurants, Tours, Airports, Activities and everything German.

[www.cometogermanynow.com](http://www.cometogermanynow.com)

([What's this?](#))

- Semi-structured Data
  - no clear tables
    - it may be hard to identify
    - it may also be hard to identify
  - some structure implicit in
    - e.g., formatting via HTML
  - large parts without structure
    - free text
  - <http://weather.yahoo.com/>

# Semi-Structured and Unstructured Data

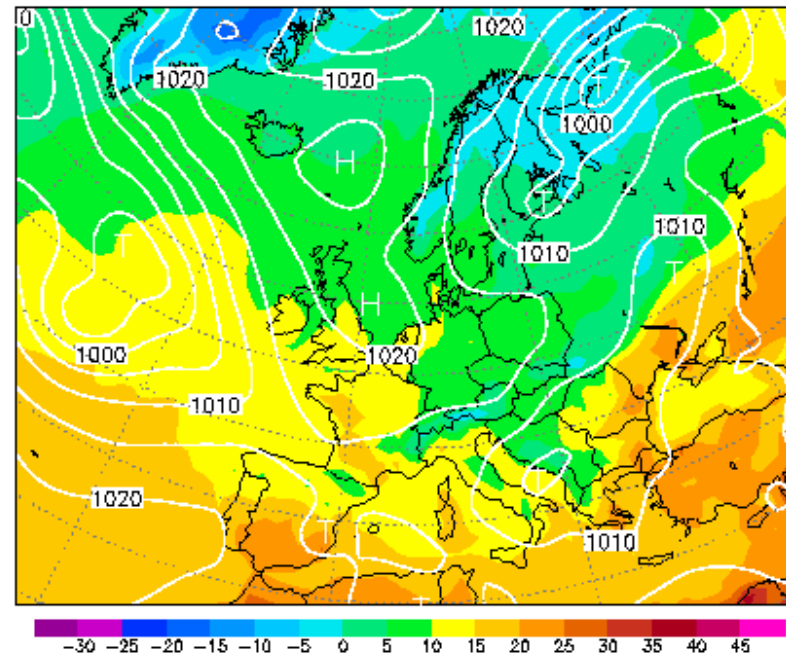
- Semi-structured Data
  - no clear tables
    - it may be hard to identify the attributes for each example
    - it may also be hard to identify the examples themselves
  - some structure implicit in the data
    - e.g., formatting via HTML
  - large parts without structure
    - free text
  - <http://weather.yahoo.com/forecast/GMXX0020.html>
- Unstructured Data
  - free text
  - <http://www.wetterzentrale.de/wzwb.html>

# Der Wetterzentrale Wetterbericht ausgegeben am 21. April 2005, 8:09 MESZ

## Lage:

Die aus Nordosten eingeflossene Kaltluft gelangt rasch unter schwachen Hochdruckeinfluss. Bereits am Samstag greifen die Ausläufer westeuropäischer Tiefs auf den Südwesten über und führen mildere und feuchte Luft heran.

Temperatur und Druckverteilung in Europa Thu,21APR2005 12Z



## Vorhersage für Deutschland:

Heute nach Auflösung örtlichen Nebels meist heiter bis wolkig und trocken. Am Alpenrand anfangs noch stark bewölkt, aber kaum noch Regen. Im Norddeutschen Tiefland ab dem Mittag einige Wolkenfelder. Höchsttemperaturen 8 bis 13 Grad. Dabei am Rhein am mildesten. Schwacher bis mäßiger Wind, im Norden auf West drehend, sonst aus Nordost bis Nord. In der kommenden Nacht im Norden wolkig. Sonst klar. Tiefstwerte zwischen 3 Grad im Norden und bis -3 Grad im Süden.

Morgen östlich der Elbe wolkig, es bleibt aber trocken. Sonst sonnig und trocken. Höchsttemperaturen zwischen 10 Grad an der Oder und bis 16 Grad am Rhein.

## Tendenz für die Folgetage:

Am Samstag im Südwesten bereits am Vormittag zunehmende Bewölkung und ab dem Mittag einsetzender Regen. In der Mitte freundlich und mild. Im Nordosten wolkig und immer noch kühl.

Am Sonntag im Norddeutschen Tiefland heiter bis wolkig und trocken. Bei kräftigem Ostwind recht kühl. In der Mitte und im Süden wolkig bis stark bewölkt mit gebietsweisem Regen oder einzelnen Schauern und mild.

Am Wochenbeginn auch im Norden unbeständiger.

Ab der Wochenmitte deutet sich trockenes und wärmeres Wetter an.

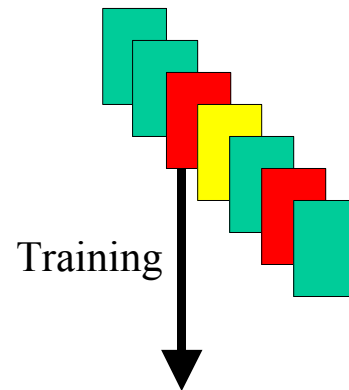
# Web Tasks for ML/DM Techniques

- Classifiers:
  - assigning categories to documents (E-mail/newsgroup sorting and filtering, building a Web catalogue, user modelling,...)
- Regression:
  - predict numerical values (ratings, GUI settings,...)
- Clustering:
  - grouping documents (structuring search results, ...)
- Association Rule Discovery:
  - finding events and event sequences that co-occur frequently (click stream analysis,...)
- Reinforcement Learning:
  - learning to improve agents (crawlers, relevance feedback, ...)



# Induction of Classifiers

*Inductive Machine Learning* algorithms induce a classifier from *labeled training examples*. The classifier *generalizes* the training examples, i.e. it is able to assign labels to new cases.



An inductive learning algorithm searches in a given family of hypotheses (e.g., *decision trees*, *neural networks*) for a member that optimizes given *quality criteria* (e.g., estimated predictive accuracy or misclassification costs).

