RECOMMENDATION MODELS FOR WEB USERS

Dr. Şule Gündüz Öğüdücü sgunduz@itu.edu.tr

What is Web Mining?

- The use of data mining techniques to automatically discover and extract information from Web documents and services (Etzioni, 1996)
- Web mining research integrate research from several research communities (Kosala and Blockeel, 2000) such as:

- Database (DB)
- Information Retrieval (IR)
- The sub-areas of machine learning (ML)
- Natural language processing (NLP)







Challenges on WWW Interactions

- Searching for usage patterns, Web structures, regularities and dynamics of Web contents
- Finding relevant information
 99% of info of no interest to 99% of people
- Creating knowledge from information available
 - Limited query interface based on keyword oriented search
- Personalization of the information
 - Limited customization to individual users

Web Mining Taxonomy







What information can be included in User **Representation?**

13

- The order of visited Web pages
- The visiting page time
- The content of the visited Web page
- The change of user behavior over time
- The difference in usage and behavior from different geographic areas
- User profile



ddress	Use	r ID	Timestamp	Method/ URL	Status		Size
S	Source of Request		Date and Time of Request	Meth	id,URL, Protocoli	Status Code	Num. o Bytes
216	216.35.116.28		[11/Jan/2002:00:58:43 -05	00] "GET/	HTTP/L0"	200	6557
216	216.35.116.28		[11/Jan/2002:00:58:53 -05	00] "GET a.gi	"GET a.gif HTTP/1.0"		5478
216	216.35.116.28		[11/Jan/2002:00:59:25-05	00] "GET b.gi	"GET b.gif HTTP/1.0"		6057
216	216.35.116.28		[11/Jan/2002:00:59:54 -05	00] "GET B.htr	"GET B.html HTTP/1.1"		59825
216	35.116.28		[11/Jan/2002:00:59:54 -05	00] "GET B.g	f HITP/L1*	200	2050
24.	102.227.6		[11/Jan/2002:00:59:55 -05	00] *GET index.h	stml HTTP/1.1*	200	6557
216	35.116.28		[11/Jan/2002:00:59:55 -05	00] "GET C.htr	nl HTTP/1.1"	200	2560
24.	102.227.6		[11/Jan/2002:00:59:56 -05	00] "GET a.gi	CHTTP/1.1*	200	5478
24.	102.227.6	•	[11/Jan/2002:00:59:56 -05	00] "GET b.gi	CHTTP/L1*	200	6057
24.	102.227.6		[11/Jan/2002:00:59:57 -05	00] "GET D.HT	ML HTTP/L1"	200	12800
24.	102.227.6		[11/Jan/2002:00:59:58 -05	00] "GET G.g	f HTTP/1.1*	200	1500
24.	102.227.6		[11/Jan/2002:00:59:58-05	00] "GET e.gi	CHTTP/L1*	200	1230
24.	102.227.6		[11/Jan/2002:00:59:59 -05	00] "GET e.jp	g HTTP/1.1*	200	3345
216	.35.116.28		[11/Jan/2002:00:59:59 -05	00] "GET c.jp	g HTTP/1.1*	200	2247
216	35.116.28		[11/Jan/2002:01:00:00 -05	00] *GET E.jp	g HTTP/1.1"	200	2247
216	35.116.28		[11/Jan/2002:01:00:00 -05	00] "GET D.htt	nl HTTP/1.1°	200	32768
216	35.116.28		[11/Jan/2002:01:00:01 -05	00] "GET D.gi	fHTTP/1.1*	200	7977
216	35.116.28		[11/Jan/2002:01:00:01 -05	00] "GET d.jp	8 HTTP/1.1*	200	6121
216	35.116.28		[11/Jan/2002:01:00:02 -05	00] "GET e.jp	g HTTP/1.1*	200	3567
24.	102.227.6	-	[11/Jan/2002:01:00:02 -05	00] "GET C.htr	nl HTTP/1.1"	200	32768



Data Preprocessing (1)

- Data cleaning
 - Remove log entries with filename suffixes such as gif, jpeg, GIF, JPEG, jpg, JPG
 - Remove the page requests made by the automated agents and spider programs
 - Remove the log entries that have a status code of 400 and 500 series
- Normalize the URLs: Determine URLs that correspond to the same Web page Data integration
- Merge data from multiple server logs Integrate semantics (meta-data)
- Integrate registration data



Discovery of Usage Patterns

- Pattern Discovery is the key component of the Web mining, which converges the algorithms and techniques from data mining, machine learning, statistics and pattern recognition etc research categories
- Separate subsections:
 - Statistical analysis
 - Association rules
 - Clustering
 - Classification
 - Sequential pattern





23

Classification

- The technique to map data item into one of several predefined classes
- Application
 - Developing a usage profile belonging to a particular class or category
- Examples:
 - WebLogMiner (Zaiane et al., 1998)

Sequential Pattern

- Discovers frequent subsequences as patterns
- Applications:
 - The analysis of customer purchase behavior
 - Optimization of Web site structure
- Examples:
 - WUM (Spiliopoulou and Faulstich, 1998)
 - Longest Repeated Subsequence (Pitkow and Pirolli, 1999)

Online Module: Recommendation

- The discovered patterns are used by the online component of the model to provide dynamic recommendations to users based on their current navigational activity
- The produced recommendation set is added to the last requested page as a set of links before the page is sent to the client browser

25

27

29

Probabilistic models of browsing behavior

- Useful to build models that describe the browsing behavior of users
- Can generate insight into how we use Web
- Provide mechanism for making predictions
- Can help in pre-fetching and personalization

26

Markov models for page prediction

General approach is to use a finite-state Markov chain

- Each state can be a specific Web page or a category of Web pages
- If only interested in the order of visits (and not in time), each new request can be modeled as a transition of states
- Issues
 - Self-transition
 - Time-independence

Markov models for page prediction For simplicity, consider order-dependent, time-independent finite-state Markov chain with M states tes se a sequence of observed states of length L. e.g. s = ABBCAABBCCBBAA with three states A, B and C. st is state at position t (1<=t<=L). In general, \$\$\mathcal{D}_{\mathcal{D}}(\beta_{\mathcal{D}}(s_{\mathcal

Markov models for page prediction

- If we denote $T_{ij} = P(s_t = j|s_{t-1} = i)$, we can define a M x M transition matrix
- Properties
 - Strong first-order assumption
- Simple way to capture sequential dependence
- If each page is a state and if W pages, $O(W^2),\,W$ can be of the order 10^5 to 10^6 for a CS dept. of a university
- To alleviate, we can cluster W pages into M clusters, each assigned a state in the Markov model
- Clustering can be done manually, based on directory structure on the Web server, or automatic clustering using clustering techniques





- First-order Markov model assumes that the next state is based only on the current state
- Limitations
 - Doesn't consider `long-term memory'
- We can try to capture more memory with *k*th-order Markov chain

31

33

35

 $P(s_t \mid s_{t-1}, ..., s_1) = P(s_t \mid s_{t-1}, ..., s_{t-k})$

- Limitations
 - Inordinate amount of training data O(M^{k+1})

Fitting Markov models to observed pagerequest data

- Assume that we collected data in the form of N sessions from server-side logs, where *l*^h session s_i, 1<= *i* <= N, consists of a sequence of L_i page requests, categorized into M 1 states and terminating in E. Therefore, data D = {s₁, ..., s_N}
- Let θ denote the set of parameters of the Markov model, θ consists of M² -1 entries in T
- Let θ_i denote the estimated probability of transitioning from state *i* to *j*.

32



where n_{ij} is the number of times we see a transition from state i to state j in the observed data D.



Bayesian parameter estimation for Markov models

- In practice, M is large (~10²-10³), we end up estimating M² probabilities
- D may contain potentially millions of sequences, so some n_{ij} = 0
 Better way would be to incorporate prior knowledge prior probability distribution P(0) and then maximize, the posterior
- distribution on $P(\theta|D)$ given the data (rather than $P(D|\theta)$) Prior distribution reflects our prior belief about the parameter set θ
- The posterior reflects our posterior belief in the parameter set now informed by the data D



Bayesian parameter estimation for Markov models

The MP posterior parameter estimates are

$$\theta_{ij}^{MP} = \frac{n_{ij} + \alpha q_{ij}}{n_i + \alpha}$$

- If n_{ij} = 0 for some transition (*i*, *j*) then instead of having a parameter estimate of 0 (ML), we will have αq_i /(n_i + α) allowing prior knowledge to be incorporated
- If n_{jj} > 0, we get a smooth combination of the datadriven information (n_{jj}) and the prior



38

40



41

37

Predicting page requests with Markov models with a mixture of first-order Markov chains

- $P(s_t \mid s_{t-1}, ..., s_1) = \sum_{k=1}^{K} P(s_t \mid s_{t-1}, c = k) P(c = k)$
- where c is a discrete-value hidden variable taking K values $Sum_k P(c=k) = 1 and P(s_\ell \mid s_{\ell\cdot l'} c=k)$ is the transition matrix for the *k*th mixture component
- One interpretation of this is user behavior consists of K different navigation behaviors described by the K Markov chains
- Cadez *et al.* use this model to cluster sequences of page requests into K groups, parameters are learned using the EM algorithm

















53

Email Product Recommendation

- What was Hotmail's primary form of advertising?
 - Small link to the sign up page at the bottom of every email sent by a subscriber
 - 'Spreading Activation'
 - Implicit recommendation

Spreading Activation

Network effects

- Even if a small number of people who receive the message subscribe (~0.1%), the service will spread rapidly
- This can be contrasted with the current practice of SPAM
 - SPAM is not sent by friends, family, co-workers
 - No implicit recommendation
 - SPAM is often viewed as not providing a good service

Page Prediction

- Click-straem Tree Model (Gunduz et. al. 2003)
- Pair-wise similarities of user sessions are calculated
 - the order of pages
 - the distance between identical pages the time spent on these pages
- A graph is constructed
 - Vertices: User sessions
 - Edges: Pair-wise similarities
- A graph-based clustering approach is applied
- Each cluster is then represented by a click-stream tree whose nodes are pages of user sessions of that cluster
- When a request is received from an active user, a recommendation set consisting of three different pages that the user has not yet visited, is produced using the best matching user session

Problems

- Evaluation of recommender systems
- Poor data
- Lack of data

55

57

Privacy control

References

- http://ibook.ics.uci.edu/slides.html
- http://www-users.cs.umn.edu/~kumar/dmbook/
- O. Etzioni. The world wide web: Quagmire or gold mine. Communications of the ACM, 39(11):65-68, 1996.
- R. Kosala and H. Blockeel. Web Mining Research: A summary. SIGKDD Explorations, 2(1):1-15, 2000
- L. Page, S. Brin, R. Motwani, T. Winograd, The PageRank Citation Ranking: Bringing Order to the Web. In Proc. of the 7th Intenational World Widw Web Conference, 161-172, 1998

References

- L. V.S. Lakshmanan, F. Sadri, and I. N. Subramanian. A Declarative Language for Querying and Restructing the World Wide web. Post-ICDE IEEE Workshop on Research Issues in Data Engineering (RIDE-NDS'96). New Orleans, February 1996. G. O. Arocena and A. O. Mendelzon. WebOQL: Restructuring Documents, Databases and Webs. In Proc. ICDE'98, 1998 M. Perkowitz and O. Etzioni, Adaptive Web Sites: Automatically Synthesizing Web Pages. In Proc. of Fifteenth National Conference on Artificial Intelligence, 1998 L. Larsen, L. K. Hansen, A. Szymkowiak, T. Christiansen and T.

- J. Larsen, L. K.Hansen, A. Szymkowiak, T. Christiansen and T. Kolenda, Webmining: Learning from the World Wide Web.
 Computational Statistics and Data Analysis, 2001
 Q. Yang, H. H. Zhang and I. T. Yi Li, Mining web logs for prediction models in WWW caching and prefetching. Knowledge Discovery and Data Mining, 473-478, 2001

References

- Y. Yan, M. Jacobsen and H. G. Molina and U. Dayal, From User Access Patterns to Dynamic Hypertext Linking. In Proc. 5th Int. World Wide Web Conference, 1007-1014, 1996 C. Shahabi, A. Zarkesh, J. Adibi and V. Shah, Knowledge Discovery from Users Web-Page Navigation. In Proc. 7th Int. Workshop on Research Issues in Data Engineering, 20-29, 1997
- R. R. Sarukkai, Link Prediction and Path Analysis Using Markov Chains. In Proc. 9th Int. World Wide Web Conference, 377-386, 2000.
- O. R. Zaiane, M. Xin and J. Han, *Discovering Web Access Patterns and ternds by applying OLAP and data mining technology on Web logs.* Advances in Digital Libraries, 19-29, 1998
- M. Spiliopoulou and L. C. Faulstich, WUM: A Tool for Web Utilization Analysis, extended version of Proc. EDBT Workshop WebDB'98, 184-203, 1998

59

References

- J.Pitkow and P.Pirolli, Mining Longest Repeating Subsequences to Predict World Wide Web Surfing. In Proc. USENIX Symp. on Internet Technologies and Systems (USITS'99), 1999.
- B. Mobasher and H. Dai and T. Luo and M. Nakagawa, *Discovery* of Aggregate Usage Profiles for Web Personalization. In roc. International {WEBKDD} Workshop -- Web Mining for E-Commerce: Challenges and Opportunities, 2000.
- B. Mobasher and H. Dai and T. Luo and M. Nakagawa, Effective Personalization Based on Association Rule Discovery from Web Usage Data. In Proc. 3rd ACM Workhop on Web Information and Data Management, 2001.

56