

Advanced Topics in Information Retrieval Learning to Rank

Vinay Setty

Jannik Strötgen

vsetty@mpi-inf.mpg.de

jannik.stroetgen@mpi-inf.mpg.de

ATIR - July 14, 2016



Before we start

oral exams

- July 28, the full day
- if you have any temporal constraints, let us know

Q-A sessions – suggestion

- Thursday, July 21: Vinay and "his topics"
- Monday, July 25: Jannik and "his topics"





Advanced Topics in Information Retrieval Learning to Rank

Vinay Setty

Jannik Strötgen

vsetty@mpi-inf.mpg.de

jannik.stroetgen@mpi-inf.mpg.de

ATIR - July 14, 2016

The Beginning of LeToR

learning to rank (LeToR)

- builds on established methods from Machine Learning
- allows different targets derived from different kinds of user input
- active area of research for past 10 15 years
- early work already (end of) 1980s (e.g., Fuhr 1989)



The Beginning of LeToR

why wasn't LeToR successful earlier?

- IR and ML communities were not very connected
- sometimes ideas take time
- Iimited training data
 - it was hard to gather (real-world) test collection queries and relevance judgments that are representative of real user needs and judgments on returned documents
 - this changed in academia and industry
- poor machine learning techniques
- insufficient customization to IR problem
- not enough features for ML to show value



The Beginning of LeToR

standard ranking based on

- term frequency / inverse document frequency
- Okapi BM25
- Ianguage models
- •

traditional ranking functions in IR exploit very few features

standard approach to combine different features

- normalize features (zero mean, unit standard deviation)
- feature combination function (typically: weighted sum)
- tune weights (either manually or exhaustively via grid search)

traditional ranking functions easy to tune



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

Why learning to rank nowadays?



Why learning to rank?

modern systems use a **huge number of features** (especially Web search engines)

- textual relevance (e.g., using LM, Okapi BM25)
- proximity of query keywords in document content
- Iink-based importance (e.g., determined using PageRank)
- depth of URL (top-level page vs. leaf page)
- spamminess (e.g., determine using SpamRank)
- host importance (e.g., determined using host-level PageRank)
- readability of content
- Iocation and time of the user
- Iocation and time of documents

Why learning to rank?

high creativity in feature engineering task

- query word in color on page?
- number of images on page?
- URL contains ~?
- number of (out) links on a page?
- page edit recency
- page length

learning to rank makes combining features more systematic



Outline I



- 2 Modeling Approaches
- 3 Gathering User Feedback
- Evaluating Learning to Rank
- 5 Learning-to-Rank for Temporal IR
- Learning-to-Rank Beyond Search



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

Outline



- 2 Modeling Approaches
- 3 Gathering User Feedback
- 4 Evaluating Learning to Rank
- 5 Learning-to-Rank for Temporal IR
- 6 Learning-to-Rank Beyond Search





open issues

- how do we model the problem?
- is it a regression or classification problem?
- what about our prediction target?





scoring as function

- different input signals (features) x_i
- with weights α_i

score(d,q) =
$$f(x_1, ..., x_m, \alpha_1, ..., \alpha_m)$$

where

- weights are learned
- features derived from d, q, and context

planck institut

formatik

LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
Outline					

1 LeToR Framework

- 2 Modeling Approaches
- 3 Gathering User Feedback
- 4 Evaluating Learning to Rank
- 5 Learning-to-Rank for Temporal IR
- 6 Learning-to-Rank Beyond Search



Classification – Regression

classification example

- dataset of $\langle q, d, r \rangle$ triples
- r: relevance (binary or multiclass)
- d: document represented by feature vector
- train ML model to predict class r of a d-q pair
- decide relevant if score is above threshold

classification problems

result in an unordered set of classes



Classification – Regression

classification problems

result in an unordered set of classes

regression problems

map to real values

ordinal regression problems

result in ordered set of classes



LeToR Modeling

LeToR can be modeled in three ways:

- pointwise: predict goodness of individual documents
- pairwise: predict users' relative preference for pairs of documents
- Istwise: predict goodness of entire query results
- each has advantages and disadvantages
- for each concrete approaches exist

in-depth discussion of concrete approaches by Liu 2009





pointwise approaches

- predict for every document based on its feature vector x
- document goodness y (e.g., label or measure of engagement)
- training determines the parameter θ based on a loss function (e.g., root-mean-square error)

main disadvantage

 as input is single document, relative order between documents cannot be naturally considered in the learning process





pairwise approaches

- predict for every pair of documents based on feature vector x
- users' relative preference regarding the documents (+1 shows preference for document 1; -1 for document 2)
- training determines the parameter θ based on a loss function (e.g., the number of inverted pairs)
- advantage: models relative order main disadvantages:

planck institut

- no distinction between excellent-bad and fair-bad
- sensitive to noisy labels (1 wrong label, many mislabeled pairs)



listwise approaches

- predict for ranked list of documents based on feature vector x
- effectiveness of ranked list y (e.g., MAP or nDCG)
- training determines the parameter θ based on a loss function advantage: positional information visible to loss function disadvantage: high training complexity, ...



Typical Learning-to-Rank Pipeline

learning to rank

 is typically deployed as a re-ranking step (infeasible to apply it to entire document collection)

step 1

 Determine a top-K result (K ~1,000) using a proven baseline retrieval method (e.g., Okapi BM25 + PageRank)

step 2

 Re-rank documents from top-K using learning to rank approach, then return top-k (k ~100) to user



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

Outline

- 1 LeToR Framework
- 2 Modeling Approaches
- 3 Gathering User Feedback
- 4 Evaluating Learning to Rank
- 5 Learning-to-Rank for Temporal IR
- 6 Learning-to-Rank Beyond Search



Gathering User Feedback

independent of pointwise, pairwise, or listwise modeling

some **input from the user is required** to determine prediction target **y**

two types of user input

- explicit user input (e.g., relevance assessments)
- implicit user input (e.g., by analyzing their behavior)



Relevance Assessments

procedure

- construct a collection of (difficult) queries
- pool results from different baselines
- gather graded relevance assessments from human assessors

problems

- hard to represent query workload within 50, 500, 5K queries
- difficult for queries that require personalization or localization
- expensive, time-consuming, and subject to Web dynamics



Clicks

track user behavior and measure their engagement with results

- click-through rate of document when shown for query
- dwell time, i.e., how much time user spent on document

problems

- position bias (consider only first result shown)
- spurious clicks (consider only clicks with dwell time above threshold)
- feedback loop (add some randomness to results)

reliability of click data

Joachims et al. (2007) and Radlinski & Joachims (2005)



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
Skips					

user behavior tells us more:

skips in addition to clicks as a source of implicit feedback

top 5:
$$d_7$$
 d_1 d_3 d_9 d_8 click
no click

- **skip previous**: $d_1 > d_7$ and $d_9 > d_3$ (user prefers d_1 over d_7)
- skip above: *d*₁ > *d*₇ and *d*₉ > *d*₃, *d*₉ > *d*₇

user study (Joachims et al., 2007): derived relative preferences

- are less biased than measures merely based on clicks
- show moderate agreement with explicit relevance assessments



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

Outline

- 1 LeToR Framework
- 2 Modeling Approaches
- 3 Gathering User Feedback
- 4 Evaluating Learning to Rank
- 5 Learning-to-Rank for Temporal IR
- 6 Learning-to-Rank Beyond Search



Learning to Rank – Evaluation

Several **benchmark datasets** have been released to allow for a comparison of different learning-to-rank methods

LETOR 2.0, 3.0, 4.0 (2007–2009) by Microsoft Research Asia

- based on publicly available document collections
- come with precomputed low-level features and relevance assessments

Yahoo! Learning to Rank Challenge (2010) by Yahoo! Labs

 comes with precomputed low-level features and relevance assessments

Microsoft Learning to Rank Datasets by Microsoft Research U.S.

 comes with precomputed low-level features and relevance assessments



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
Feature	es				
		Yahool Fea	atures		

queries, ulrs and features descriptions are not given, only the feature values!

- feature engineering critical for any commercial search engine
- releasing queries, URLs leads to risk of reverse engineering
- reasonable consideration, but prevent IR researchers from studying what feature are most effective

LETOR / Microsoft Features

each query-url pair is represented by a 136-dimensional vector



_eToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

feature id	feature description	stream	comments
1	covered query term number	body	
2		anchor	
3		title	
4		url	
5		whole document	
6	covered query term ratio	body	
7		anchor	
8		title	
9		url	
10		whole document	
11	stream length	body	
12		anchor	
13		title	
14		url	
15		whole document	



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
	F	_			

16	IDF(Inverse document	body
17	frequency)	anchor
18		title
19		url
20		whole document
21	sum of term frequency	body
22		anchor
23		title
24		url
25		whole document
26	min of term frequency	body
27		anchor
28		title
29		url
30		whole document



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

31	max of term frequency	body
32		anchor
33		title
34		url
35		whole document
36	mean of term frequency	body
37		anchor
38		title
39		url
40		whole document
41	variance of term frequency	body
42		anchor
43		title
44		url
45		whole document



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
LETOR	Feature	S			

46	sum of stream length	body
47	normalized term frequency	anchor
48		title
49		url
50		whole document
51	min of stream length	body
52	normalized term frequency	anchor
53		title
54		url
55		whole document
56	max of stream length	body
57	normalized term frequency	anchor
58		title
59		url
60		whole document



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
-					

61	mean of stream length	body
62	normalized term frequency	anchor
63		title
64		url
65		whole document
66	variance of stream length	body
67	normalized term frequency	anchor
68		title
69		url
70		whole document
71	sum of tf*idf	body
72		anchor
73		title
74		url
75		whole document



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

76	min of tf*idf	body
77		anchor
78		title
79		url
80		whole document
81	max of tf*idf	body
82		anchor
83		title
84		url
85		whole document
86	mean of tf*idf	body
87		anchor
88		title
89		url
90		whole document



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
		-			

91	variance of tf*idf	body
92		anchor
93		title
94		url
95		whole document
96	boolean model	body
97		anchor
98		title
99		url
100		whole document
101	vector space model	body
102		anchor
103		title
104		url
105		whole document


LeTo	R Framework	Modeling	User Feedback	Evaluation	Time	Beyond Searcl
	LETC	R Features				
	106	BM25	body			
	107		anchor			
	108		title			
	109		url			
	110		whole d	ocument		
	111	LMIR.ABS	body		Language model	
	112		anchor		approach for info	rmation
	113		title		retrieval (IR) with	n nting
	114		url		absolute discount	
	115		whole d	locument	smoothing	
	116	LMIR.DIR	body		Language model	
	117		anchor		approach for IR w	ith
	118		title		Bayesian smoothin	ng using
	110		url		Dirichlet priors	

 119
 url
 Diric

 120
 whole document



LeTo	R Framework	Modeling	User Feed	back	Evaluation	Time	Beyond Search
	LETO	R Features					
	121	LMIRJM		body		Language mo	del
	122			anchor		approach for	IR with
	123			title		Jelinek-Merce	r
	124			url		smoothing	
	125			whole do	cument		
	126	Number of slash in URL					
	127	Length of URL					
	128	Inlink number					
	129	Outlink number					
	130	PageRank					
	131	SiteRank				Site level Page	eRank
	132	QualityScore				The quality sc web page. The outputted by quality classifi	ore of a e score is a web page er.



еTo	oR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
	LETC	R Features				
	133	QualityScore2			The quality score web page. The s outputted by a v quality classifier, measures the ba a web page.	e of a core is web page which adness of
	134	Query-url click count			The click count of query-url pair at engine in a period	of a : a search od
	135	url click count			The click count of aggregated from	of a url n user

136 url dwell time

The average dwell time of a url aggregated from user browsing data in a period

browsing data in a period



all details

http://research.microsoft.com/en-us/um/beijing/projects/letor/

https://www.microsoft.com/en-us/research/project/mslr/

- datasets
- dataset descriptions
- partitioned in subsets (for cross-validation)
- evaluation scripts, significance test scripts
- feature list

everything required to get started is available



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

Outline

- 1 LeToR Framework
- 2 Modeling Approaches
- 3 Gathering User Feedback
- 4 Evaluating Learning to Rank
- 5 Learning-to-Rank for Temporal IR
- 6 Learning-to-Rank Beyond Search



Learning-to-Rank for Temporal IR

Kanhabua & Nørvåg (2012)

learning-to-rank approach for time-sensitive queries

standard temporal IR approaches

- mixture model linearly combining textual similarity and temporal similarity
- probabilistic model generating a query from the textual and temporal part of document independent

learning-to-rank approach

- two classes of features: entity features and time features
- both derived from annotations (NER, temporal tagging)



Learning-to-Rank for Temporal IR

document model

- a document collection over time
- document is composed of a bag of words and time
 - publication date
 - temporal expressions mentioned in document
- annotated document composed of
 - set of named entities
 - set of temporal expressions
 - set of annotated sentences

temporal query model

- $q = \{q_{text}, q_{time}\}$
- *q*_{time} might be explicit or implicit

Learning-to-Rank for Temporal IR

learning-to-rank

- a wide range of **temporal features**
- a wide range of entity features
- models trained using labeled query/document pairs
- documents ranked according to weighted sum of its feature scores

experiments

- show improvement of baselines and other time-aware models
- (many queries also contained entities, news corpus)



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

Outline

- 1 LeToR Framework
- 2 Modeling Approaches
- 3 Gathering User Feedback
- Evaluating Learning to Rank
- 5 Learning-to-Rank for Temporal IR
- 6 Learning-to-Rank Beyond Search



Learning-to-Rank – Beyond Search

learning to rank is applicable beyond web search

Example: matching in eharmony.com

```
Slides by Vaclav Petricek:
```

http:

//www.slideshare.net/VaclavPetricek/data-science-of-love

basic idea:

- standard approach is search-based, filter out non-matches
- eharmony approach is learning to rank: suggest potential matches



starting point in the 1990s

distinguish marriages that work well and those that don't

step 1: compatibility matching

- based on 150 questions: personality, values, attitudes, beliefs important attributes for the long term
- predict marital satisfaction

even if people are compatible,

they might not be interested to talk to each other





step 2: affinity matching

- based on other features: distance, height difference, zoom level of photo
- predict probability of message exchange

however

who should be introduced to whom and when?

match distribution

based on graph optimization problem (constrained max flow)







LeToR Framewo	rk Modeling	User Feedback	Evaluation	Time	Beyond Search
C	compatibility	Matching	Romantic		
Ť			ro·man·tic /roˈmantik/ +> Adjective Inclined toward or sugges mystery associated with to Noun A person with romantic be Synonyms romanticist	tive of the feeling of e vve.	excitement and
() ()	max planck ins informatik	titut	© Jannik Strötgen – ATIF	R-10	51 / 72

even if people are compatible,

they might not be interested to talk to each other



LeToR	Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
	Affinity	Matching				
		1	<i></i>	Ŕ		
		Å	i.	1		
		Ŕ				
	mp	max planck institut informatik	© Ja	nnik Strötgen – ATIR-10)	53 / 72













LeToR	Framework	Modeling	User Feedback	Evaluation	Time Be	eyond Search
	Affinit	v Matchi		zoom s	size matters	
		y watorin		only face: c no face: so ratio: fa	loesn't tell mu meone is hidi ace / pic size	ch ng
					٢	
(G		max planck institut informatik		© Jannik Strötgen – ATIF	₹-10	59 / 72

LeToR Framework	Modeling	User Feedback	Evalu	ation	Time	Beyond Search
Affini	ty Matchin	g > Foo	d preference	e		
	1		•		-	
	25%	-1%	-24%	20%	13%	
	9%	-5%	-27%	7%	0%	
Ý 🤞	-12%	-21%	-42%	-19%	-23%	
1 Alexandre	19%	0%	-28%	28%	10%	
-	9%	-11%	-35%	11%	44%	
	max planck institut informatik		© Jannik Ströte	gen – ATIR-10		60 / 72

Affinity Matching Food preference Image: Second preference Image: Second preference Image: Second preference Image: Seco	ToR Framework	Modeling	User Feedback	Eval	uation	Time	Beyond Search
25% -1% 20% 13% 9% -1% 20% 0% 10% 02% 23% 13%	Affinit	y Matchin	g > Foo	od preferenc	æ		
25% -1% -24% 20% 13% 9% -5% -27% 7% 0% 12% -21% -42% -19% -23% 19% 0% -28% 28% 10%				†			
9% -5% -27% 7% 0% 12% -21% -42% -19% -23% 19% 0% -28% 28% 10%		25%	-1%	-24%	20%	13%	
-12% -21% -42% -19% -23% 19% 0% -28% 28% 10%	. ·	9%	-5%	-27%	7%	0%	
19% 0% -28% 28% 10%		-12%	-21%	-42%	-19%	-23%	
		19%	0%	-28%	28%	10%	
9% -11% -35% 11% 44%	╡ ┥	9%	-11%	-35% © Jannik Strö	11% tgen – ATIR-10	44%	61 / 72

LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search

however

who should be introduced to whom and when?



LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
Matc	h Distribut	ion > Grap	h optimization		
	4		Ŵ		
	4		Ť		
	1		Ť		
	1		Ť		
() mp	max planck institut informatik	© J	lannik Strötgen – ATIR	-10	63 / 72













LeToR Framework	Modeling	User Feedback	Evaluation	Time	Beyond Search
Summary	/				

Learning to Rank

provides systematic ways to combine features

modeling

- pointwise: predict goodness of individual document
- pairwise: predict relative preference for document pairs
- Istwise: predict effectiveness of ranked list of documents

explicit and implicit user inputs

include relevance assessments, clicks, and skips

Thank you for your attention!





- Fuhr (1989): Optimum Polynomial Retrieval Functions based on the the Probability Ranking Principle, ACM TOIS 7(3).
- Liu (2009): *Learning to Rank for Information Retrieval*, Foundations and Trends in Information Retrieval 3(3):225â331.
- Joachims et al. (2007): Evaluating the Accuracy of Implicit Feedback from Clicks and Query Reformulations in Web Search, ACM TOIS 25(2).
- Radlinski & Joachims (2005): Query Chains: Learning to Rank from Implicit Feedback, KDD.
- Kanhabua & Nørvåg (2012): Learning to Rank Search Results for Time-Sensitive Queries, CIKM.



Thanks

some slides / examples are taken from / similar to those of:

- Klaus Berberich, Saarland University, previous ATIR lecture
- Manning, Raghavan, Schütze: Introduction to Information Retrieval (including slides to the book)
- and the eharmony.com slides by Vaclav Petricek: http://www.slideshare.net/VaclavPetricek/data-science-of-love

