

MMDS 2012
Stanford University

Large Scale Machine Learning for Query Document Matching in Web Search*

Hang Li

Huawei Technologies

* Work was done at Microsoft Research, with former colleagues and interns

Talk Outline

- Motivation
- Regularized Latent Semantic Indexing
- Group Matrix Factorization
- Matching in Latent Space
- Conclusion

Same Search Intent, Different Query Representations

Example = “Distance between Sun and Earth”

- "how far" earth sun
- "how far" sun
- "how far" sun earth
- average distance earth sun
- average distance from earth to sun
- average distance from the earth to the sun
- distance between earth & sun
- distance between earth and sun
- distance between earth and the sun
- distance from earth to the sun
- distance from sun to earth
- distance from sun to the earth
- distance from the earth to the sun
- distance from the sun to earth
- distance from the sun to the earth
- distance of earth from sun
- distance between earth sun
- how far away is the sun from earth
- how far away is the sun from the earth
- how far earth from sun
- how far earth is from the sun
- how far from earth is the sun
- how far from earth to sun
- how far from the earth to the sun
- distance between sun and earth

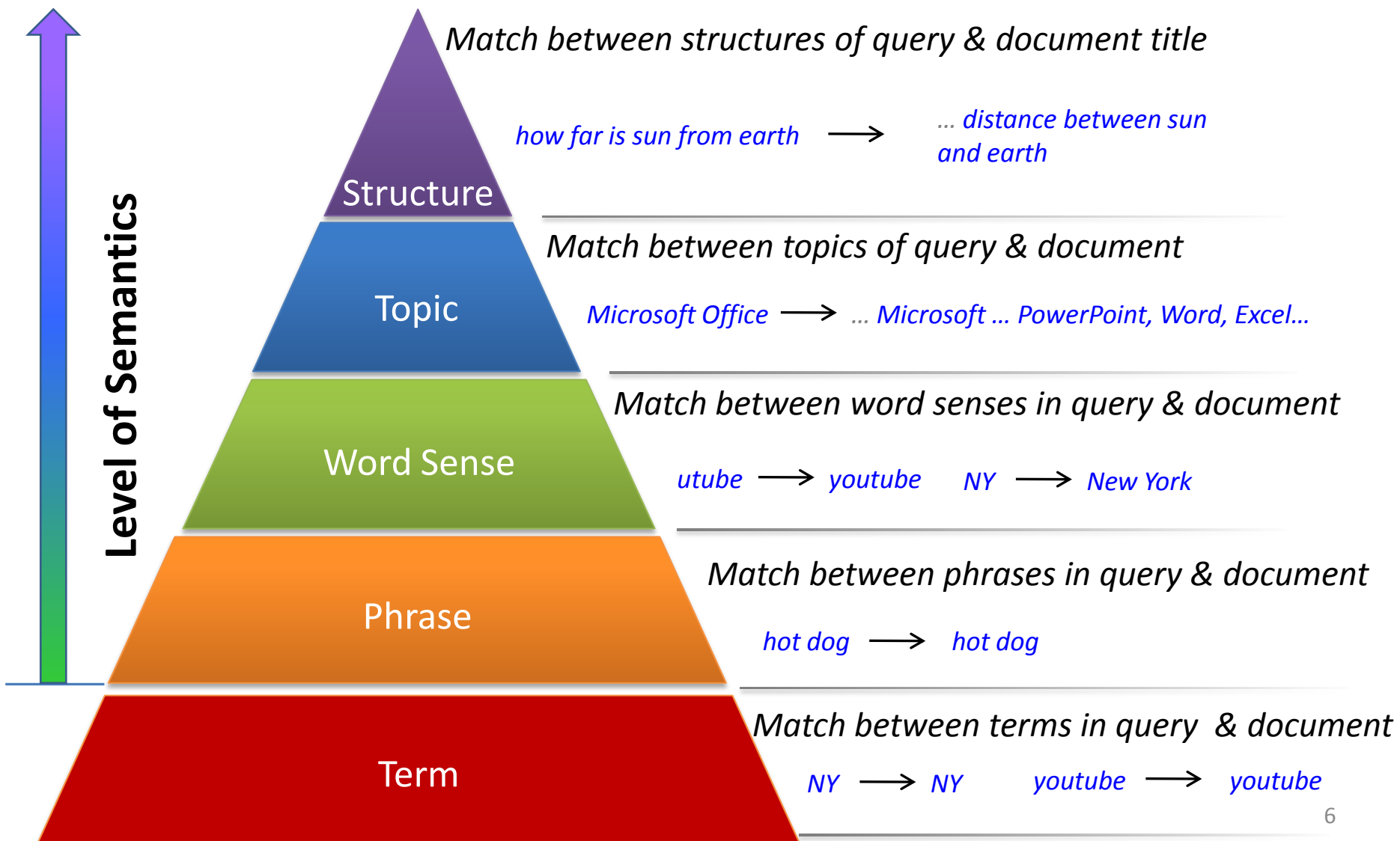
Same Search Intent, Different Query Representations

Example = “Youtube”

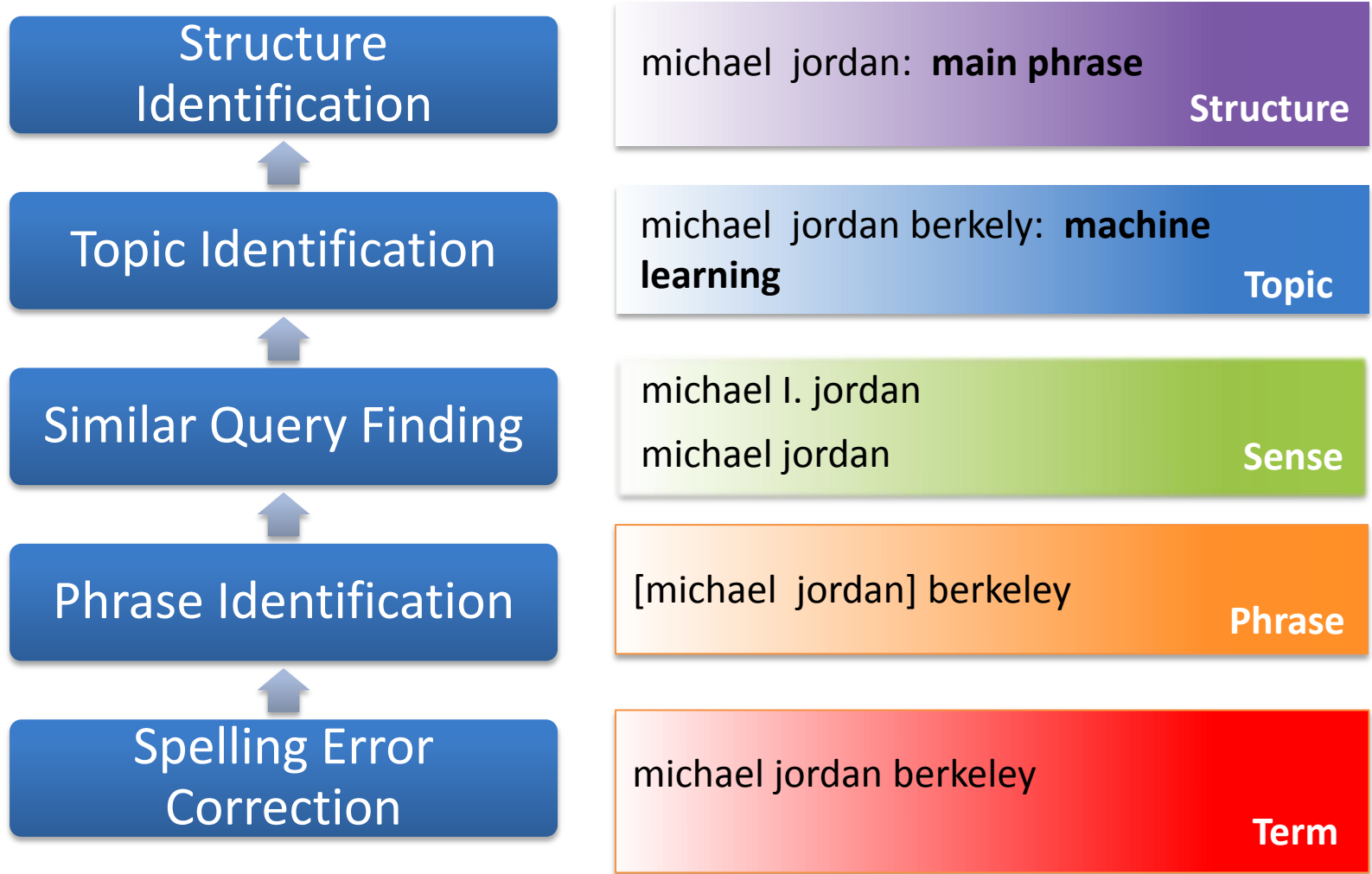
- | | | |
|-------------------|-----------------------|----------------------|
| • youtube | yuotube | yuo tube |
| • ytube | youtubr | yu tube |
| • youtubo | youtuber | youtubecom |
| • youtube om | youtube music videos | youtube videos |
| • youtube | youtube com | youtube co |
| • youtub com | you tube music videos | yout tube |
| • youtub | you tube com yourtube | your tube |
| • you tube | you tub | you tube video clips |
| • you tube videos | www you tube com | www youtube com |
| • www youtube | www youtube com | www youtube co |
| • yotube | www you tube | www utube com |
| • ww youtube com | www utube | www u tube |
| • utube videos | utube com | utube |
| • u tube com | utub | u tube videos |
| • u tube | my tube | toutube |
| • outube | our tube | toutube |

Semantic Matching Project: Solving Document Mismatch in Web Search

Matching at Different Levels



Query Understanding



michael jordan berkele

Document Understanding

Title Structure
Identification



Topic Identification



Key Phrase
Identification



Phrase
Identification

Michael Jordan: *main phrase in Title* **Structure**

Michael Jordan is Professor in the
Department of Electrical Engineering: *machine learning* **Topic**

[Michael Jordan], [Professor]
[Electrical Engineering]: *keyphrase* **Phrase**

[Michael Jordan] is [Professor] in the
[Department] of [Electrical Engineering] **Phrase**

Homepage of Michael Jordan

Michael Jordan is Professor in the
Department of Electrical Engineering

Online Matching

[Michael I. Jordan's Home Page](#)

Models of visuomotor and other learning (Univ. of California, Berkeley, USA)
www.cs.berkeley.edu/~jordan · [Cached page](#) · [Mark as spam](#)

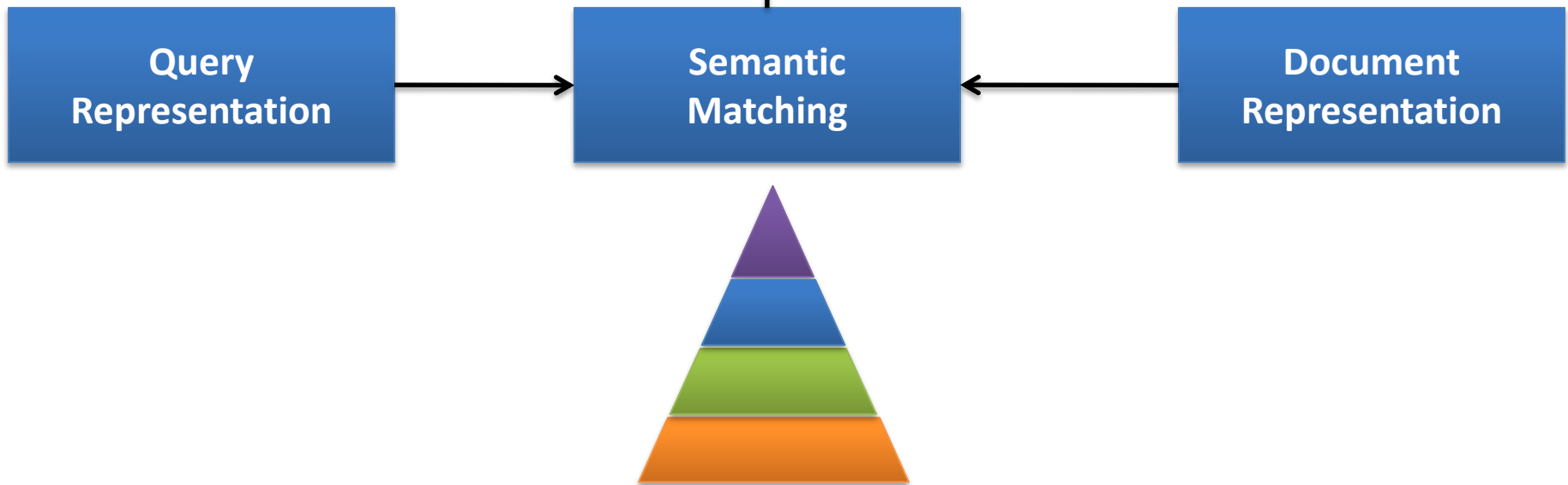
[Michael Jordan | EECS at UC Berkeley](#)

Michael Jordan Professor Research Areas Artificial Intelligence (AI) Biosystems & Computational Biology (BIO) Control, Intelligent Systems, and Robotics (CIR)
www.eecs.berkeley.edu/Faculty/Homepages/jordan.html · [Cached page](#) · [Mark as spam](#)

[Publications](#)

Jordan. In M.-H. Chen, D. Dey, P. Mueller, D. Sun, and K. Ye (Eds.), *Frontiers of ...*
Technical Report 661, Department of Statistics, University of California, Berkeley, 2004.
www.cs.berkeley.edu/~jordan/publications.html · [Cached page](#) · [Mark as spam](#)

Ranking



Matching can be conducted at different levels

Related Work

- Studied in long history of IR
- Query expansion, pseudo relevance feedback
- Latent Semantic Indexing, Probabilistic Latent Semantic Indexing, Latent Dirichlet Allocation
-
- New problem setting
 - Large amount of data available
 - New machine learning techniques

Matching vs Ranking

In search, first matching and then ranking

	Matching	Ranking
Prediction	Matching degree between query and document	Ranking list of documents
Model	$f(q, d)$	$f(q, d_1), f(q, d_2), \dots, f(q, d_n)$
Challenge	Mismatch	Correct ranking on top

Matching between Heterogeneous Data is Everywhere

- Matching between user and product (collaborative filtering)
- Matching between text and image (image annotation)
- Matching between people (dating)
- Matching between languages (machine translation)
- Matching between receptor and ligand (drug design)

Regularized Latent Semantic Indexing

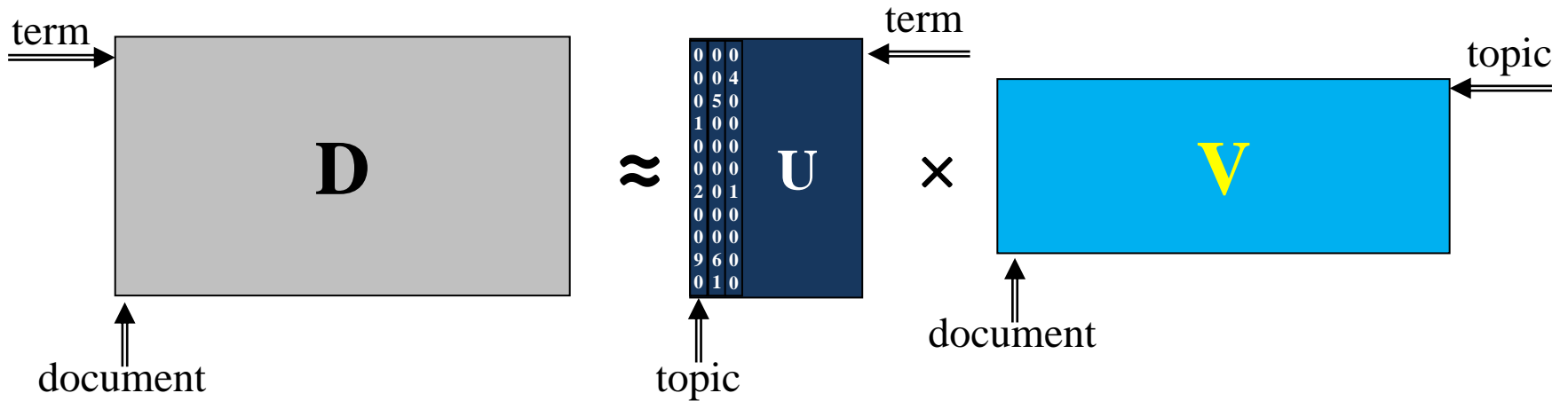
Joint Work with Quan Wang, Jun Xu,
and Nick Craswell

SIGIR 2011

Regularized Latent Semantic Indexing

- Motivation
 - Matching between query and document at topic level
 - Scale up to large datasets (vs. existing methods)
- Approach
 - Matrix Factorization
 - Regularization on topics and documents (vs. Sparse Coding)
 - Learning problem can be easily decomposed
- Results
 - l_1 on topics leads to sparse topics and l_2 on documents leads to accurate matching
 - Comparable with existing methods in topic discovery and search relevance
 - But can easily scale up to large document sets

Regularized Latent Semantic Indexing

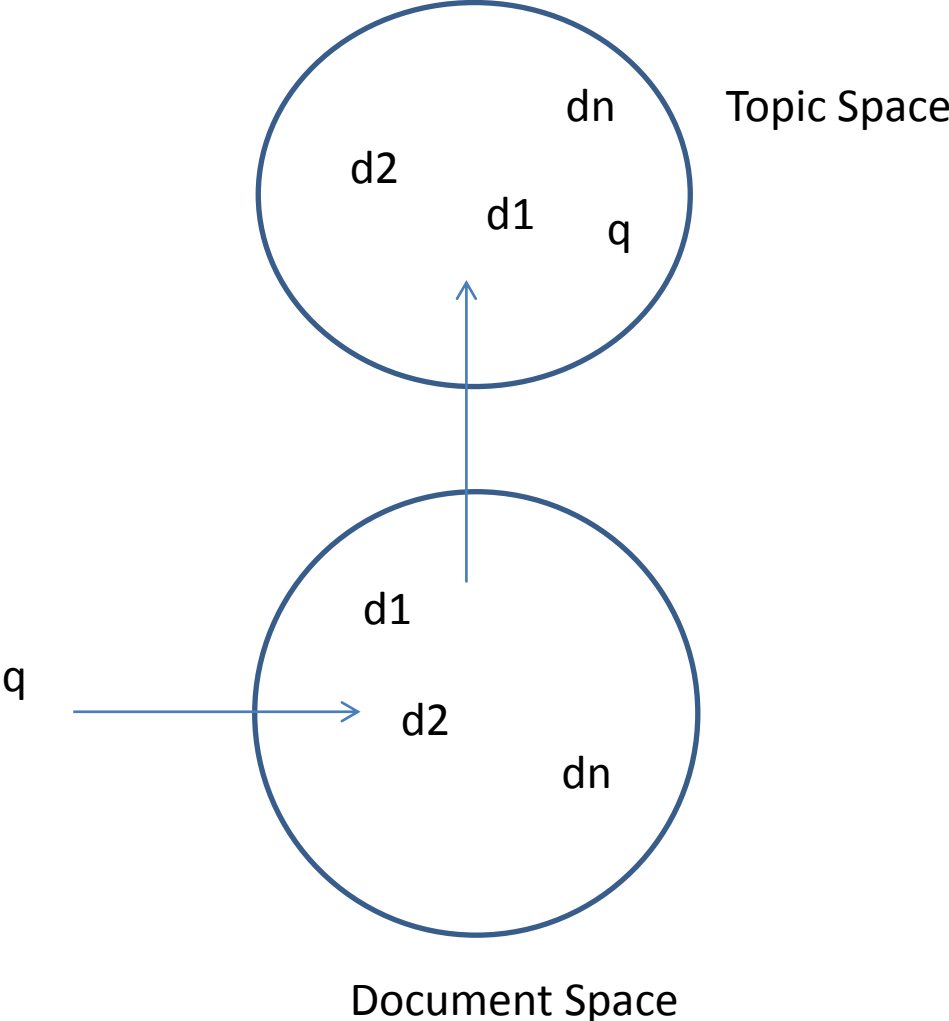


term representation of doc n topics topic representation of doc n documents are smooth

$$\min_{\mathbf{U}, \{\mathbf{v}_n\}} \sum_{n=1}^N \|\mathbf{d}_n - \mathbf{U}\mathbf{v}_n\|_2^2 + \lambda_1 \sum_{k=1}^K \|\mathbf{u}_k\|_1 + \lambda_2 \sum_{n=1}^N \|\mathbf{v}_n\|_2^2$$

topics are sparse

Query and Document Matching in Topic Space



Optimization Strategy

$$\min_{\mathbf{U}, \{\mathbf{v}_n\}} \sum_{n=1}^N \|\mathbf{d}_n - \mathbf{U}\mathbf{v}_n\|_2^2 + \lambda_1 \sum_{k=1}^K \|\mathbf{u}_k\|_1 + \lambda_2 \sum_{n=1}^N \|\mathbf{v}_n\|_2^2$$

Coordinate Decent

Update \mathbf{U}

Update \mathbf{V}

$$\min_{\{\bar{\mathbf{u}}_m\}} \sum_{m=1}^M \|\bar{\mathbf{d}}_m - \mathbf{V}^T \bar{\mathbf{u}}_m\|_2^2 + \lambda_1 \sum_{m=1}^M \|\bar{\mathbf{u}}_m\|_1$$

$$\min_{\{\mathbf{v}_n\}} \sum_{n=1}^N \|\mathbf{d}_n - \mathbf{U}\mathbf{v}_n\|_2^2 + \lambda_2 \sum_{n=1}^N \|\mathbf{v}_n\|_2^2$$

$$\min_{\bar{\mathbf{u}}_m} \|\bar{\mathbf{d}}_m - \mathbf{V}^T \bar{\mathbf{u}}_m\|_2^2 + \lambda_1 \|\bar{\mathbf{u}}_m\|_1$$

for $m = 1, \dots, M$

$$\min_{\mathbf{v}_n} \|\mathbf{d}_n - \mathbf{U}\mathbf{v}_n\|_2^2 + \lambda_2 \|\mathbf{v}_n\|_2^2$$

for $n = 1, \dots, N$

$$u_{mk} = \begin{cases} \frac{(r_{mk} - \sum_{l \neq k} s_{kl} u_{ml}) - \frac{1}{2} \lambda_1}{s_{kk}}, & \text{if } u_{mk} > 0 \\ \frac{(r_{mk} - \sum_{l \neq k} s_{kl} u_{ml}) + \frac{1}{2} \lambda_1}{s_{kk}}, & \text{if } u_{mk} < 0 \end{cases}$$

$$\mathbf{v}_n^* = (\mathbf{U}^T \mathbf{U} + \lambda_2 \mathbf{I})^{-1} \mathbf{U}^T \mathbf{d}_n$$

Analytic Solution

RLSI Algorithm

- Single machine multi core version
- Multiple machine version (MapReduce and MPI)

Require: $\mathbf{D} \in \mathbb{R}^{M \times N}$

- 1: $\mathbf{V}^{(0)} \in \mathbb{R}^{K \times N} \leftarrow$ random matrix
- 2: **for** $t = 1 : T$ **do**
- 3: $\mathbf{U}^{(t)} \leftarrow$ Update $\mathbf{U}(\mathbf{D}, \mathbf{V}^{(t-1)})$
- 4: $\mathbf{V}^{(t)} \leftarrow$ Update $\mathbf{V}(\mathbf{D}, \mathbf{U}^{(t)})$
- 5: **end for**
- 6: **return** $\mathbf{U}^{(T)}, \mathbf{V}^{(T)}$

terms
processed
in parallel

Algorithm 2 Update \mathbf{U}

Require: $\mathbf{D} \in \mathbb{R}^{M \times N}, \mathbf{V} \in \mathbb{R}^{K \times N}$

- 1: $\mathbf{S} \leftarrow \mathbf{V}\mathbf{V}^T$
- 2: $\mathbf{R} \leftarrow \mathbf{D}\mathbf{V}^T$
- 3: **for** $m = 1 : M$ **do**
- 4: $\bar{\mathbf{u}}_m \leftarrow \mathbf{0}$
- 5: **repeat**
- 6: **for** $k = 1 : K$ **do**
- 7: $w_{mk} \leftarrow r_{mk} - \sum_{l \neq k} s_{kl} u_{ml}$
- 8: $u_{mk} \leftarrow \frac{(|w_{mk}| - \frac{1}{2} \lambda N)_+ \text{sign}(w_{mk})}{s_{kk}}$
- 9: **end for**
- 10: **until** convergence
- 11: **end for**
- 12: **return** \mathbf{U}

Algorithm 3 Update \mathbf{V}

Require: $\mathbf{D} \in \mathbb{R}^{M \times N}, \mathbf{U} \in \mathbb{R}^{M \times K}$

- 1: $\Sigma \leftarrow (\mathbf{U}^T \mathbf{U} + \theta \mathbf{I})^{-1}$
- 2: $\Phi \leftarrow \mathbf{U}^T \mathbf{D}$
- 3: **for** $n = 1 : N$ **do**
- 4: $\mathbf{v}_n \leftarrow \Sigma \phi_n$, where ϕ_n is the n^{th} column
- 5: **end for**
- 6: **return** \mathbf{V}

docs
processed
in parallel

Scalability Comparison

algorithm	max dataset applied (#docs; #words)	# topics	# processors used
PLDA and PLDA+ (by Google)	Wiki-200T(2,112,618; 200,000)	1000	2, 048
AD-LDA	NY Times (300,000; 102,660) PubMed (8,200,000; 141,043)	200	16
RLSI	B01 (1,562,807; 7,014,881) Bing News (940,702; 500,033) Wiki-All (3,239,884; 6,043,069) MSWeb Data (2,635,158; 2,371,146)	500 ~ 1000	single machine, 24 cores

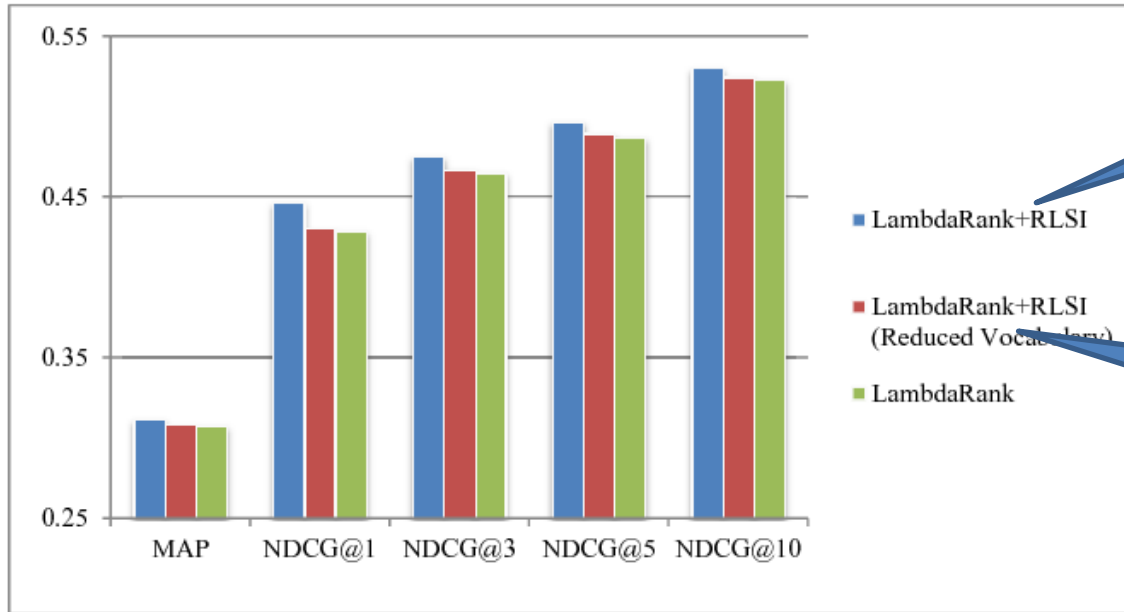
Experimental Results on Topic Discovery

Topics discovered by RLSI are equally readable compared with LDA, PLSI, LSI

Table 8: Topics discovered by RLSI, LDA, PLSI, and LSI from AP dataset.

	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6	Topic 7	Topic 8	Topic 9	Topic 10
RLSI AvgComp = 0.0075	opec oil cent barrel price	africa south african angola apartheid	aid virus infect test patient	school student teacher educate college	noriega panama panamanian delval canal	percent billion rate 0 trade	plane crash flight air airline	israeli palestinian israel arab plo	nuclear soviet treaty missile weapon	bush dukakis campaign quayle bentsen
LDA AvgComp = 1	soviet nuclear union state treaty	school student year educate university	dukakis democrat campaign bush jackson	party govern minister elect nation	year new time television film	water year fish animal 0	price year market trade percent	court charge case judge attorney	air plane flight crash airline	iran iranian ship iraq navy
PLSI AvgComp = 0.9534	company million share billion stock	israeli iran israel palestinian arab	year state new nation govern	year state new nation 0	bush dukakis democrat campaign republican	court charge attorney judge trial	soviet treaty missile nuclear gorbachev	year state new nation govern	plane flight airline crash air	year state new people nation
LSI AvgComp = 1	soviet percent police govern state	567 234 0 percent 12	0 yen dollar percent tokyo	earthquake quake richter scale damage	drug school test court dukakis	0 dukakis bush jackson dem	israel israeli student palestinian africa	yen dukakis bush dollar jackson	urgent oil opec dukakis cent	student school noriega panama teacher

Experimental Results on Web Search



RLSI can help improve search relevance

Reducing vocabulary hurts ranking accuracy

Group Matrix Factorization

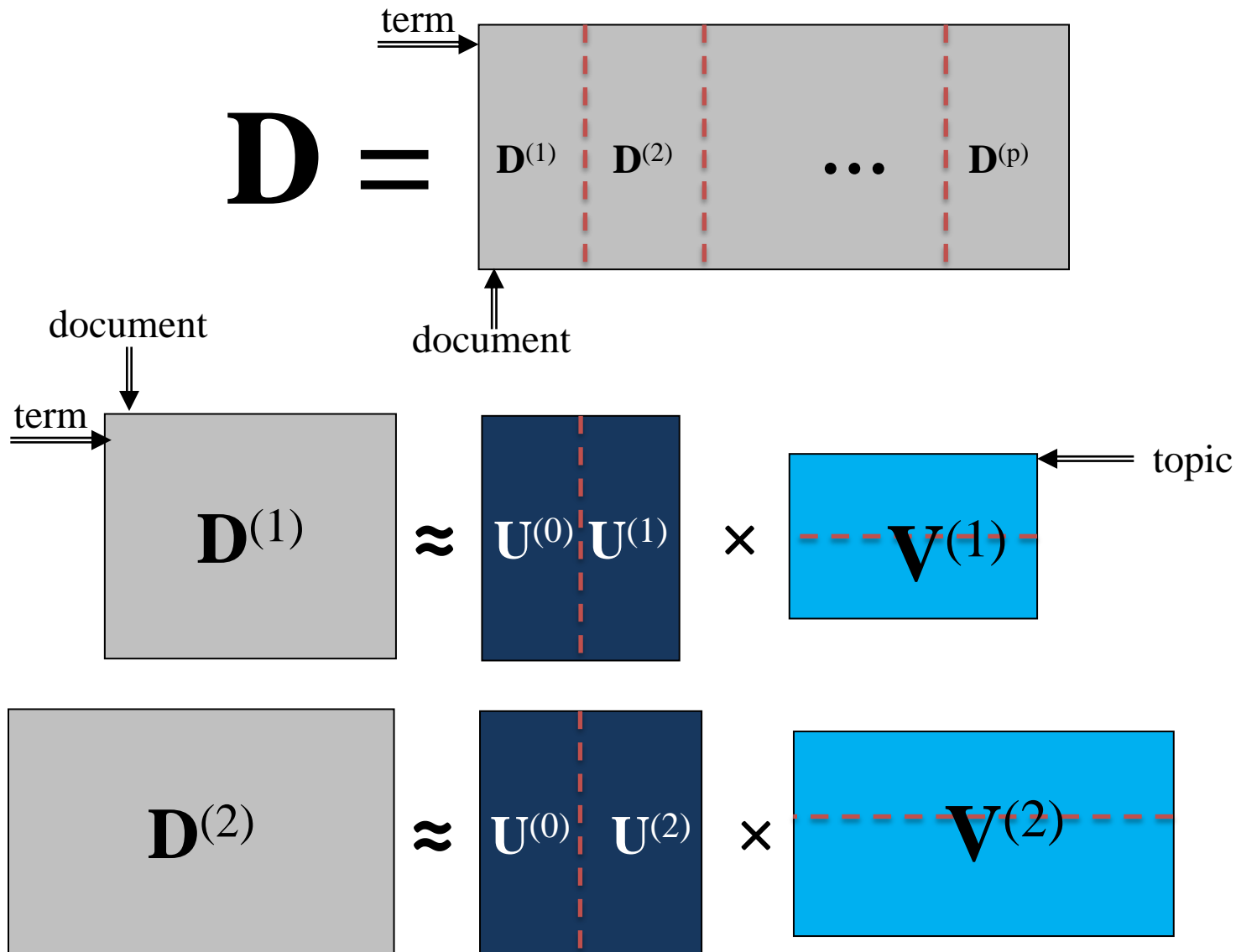
Joint Work with Quan Wang, Zheng
Cao, and Jun Xu

SIGIR 2012

Group Matrix Factorization

- Motivation
 - Matching between query and document at topic level
 - Having even better scalability
- Approach
 - Matrix Factorization (RLSI and NMF)
 - Assuming documents have been classified into classes
 - Class specific topics and shared topics
- Results
 - Comparable with existing methods in topic discovery and search relevance
 - Can scale up to even large document sets

Group Matrix Factorization



Group Matrix Factorization (cont')

$$\min_{\{\mathbf{u}_k^{(0)}\}, \{\mathbf{u}_k^{(p)}\}, \{\mathbf{v}_n^{(p)}\}} \sum_{p=1}^P \sum_{n=1}^{N_p} \mathcal{L}(\mathbf{d}_n^{(p)} \| \tilde{\mathbf{U}}_p \mathbf{v}_n^{(p)}) + \theta_1 \sum_{k=1}^{K_s} \mathcal{R}_1(\mathbf{u}_k^{(0)})$$
$$+ \theta_2 \sum_{p=1}^P \sum_{k=1}^{K_c} \mathcal{R}_2(\mathbf{u}_k^{(p)}) + \theta_3 \sum_{p=1}^P \sum_{n=1}^{N_p} \mathcal{R}_3(\mathbf{v}_n^{(p)})$$

$$s.t. \quad \mathbf{u}_k^{(0)} \in \mathcal{C}_1, \quad k = 1, \dots, K_s,$$

$$\mathbf{u}_k^{(p)} \in \mathcal{C}_2, \quad k = 1, \dots, K_c, p = 1, \dots, P,$$

$$\mathbf{v}_n^{(p)} \in \mathcal{C}_3, \quad n = 1, \dots, N_p, p = 1, \dots, P,$$

Group Regularized Latent Semantic Indexing

$$\min_{\{\mathbf{u}_k^{(0)}\}, \{\mathbf{u}_k^{(p)}\}, \{\mathbf{v}_n^{(p)}\}} \sum_{p=1}^P \sum_{n=1}^{N_p} \|\mathbf{d}_n^{(p)} - \tilde{\mathbf{U}}_p \mathbf{v}_n^{(p)}\|_2^2 + \lambda_1 \sum_{k=1}^{K_s} \|\mathbf{u}_k^{(0)}\|_1$$
$$+ \lambda_1 \sum_{p=1}^P \sum_{k=1}^{K_c} \|\mathbf{u}_k^{(p)}\|_1 + \lambda_2 \sum_{p=1}^P \sum_{n=1}^{N_p} \|\mathbf{v}_n^{(p)}\|_2^2$$

G-RLSI Algorithm

- Single machine multi core version
- Multiple machine version (MapReduce and MPI)

Algorithm 1 Group RLSI

Require: $\mathbf{D}_1, \dots, \mathbf{D}_P$

- 1: **for** $p = 1 : P$ **do**
- 2: $\mathbf{U}_p \leftarrow$ zero matrix
- 3: $\mathbf{V}_p \leftarrow$ random matrix
- 4: **end for**
- 5: **repeat**
- 6: $\mathbf{U}_0 \leftarrow$ Update $U_0(\{\mathbf{D}_p\}, \{\mathbf{U}_p\}, \{\mathbf{V}_p\})$
- 7: **for** $p = 1 : P$ **do**
- 8: $\mathbf{U}_p \leftarrow$ Update $U_p(\mathbf{D}_p, \mathbf{U}_0, \mathbf{V}_p)$
- 9: $\mathbf{V}_p \leftarrow$ Update $V_p(\mathbf{D}_p, \mathbf{U}_0, \mathbf{U}_p)$
- 10: **end for**
- 11: **until** convergence
- 12: **return** $\mathbf{U}_0, \mathbf{U}_1, \dots, \mathbf{U}_P, \mathbf{V}_1, \dots, \mathbf{V}_P$

Algorithm 4 Update V_p

Require: $\mathbf{D}_p, \mathbf{U}_0, \mathbf{U}_p$

- 1: $\Sigma_p \leftarrow (\tilde{\mathbf{U}}_p^T \tilde{\mathbf{U}}_p + \lambda_2 \mathbf{I})^{-1}$
- 2: $\Phi_p \leftarrow \tilde{\mathbf{U}}_p^T \mathbf{D}_p$
- 3: **for** $n = 1 : N_p$ **do**
- 4: $v_n^{(p)} \leftarrow \Sigma_p \phi_n^{(p)}$
- 5: **end for**
- 6: **return** \mathbf{V}_p

Algorithm 2 Update U_0

Require: $\mathbf{D}_1, \dots, \mathbf{D}_P, \mathbf{U}_1, \dots, \mathbf{U}_P, \mathbf{V}_1, \dots, \mathbf{V}_P$

- 1: $\mathbf{S}_0 \leftarrow \sum_{p=1}^P \mathbf{H}_p \mathbf{H}_p^T$
- 2: $\mathbf{R}_0 \leftarrow \sum_{p=1}^P \mathbf{D}_p \mathbf{H}_p^T - \sum_{p=1}^P \mathbf{U}_p \mathbf{W}_p \mathbf{H}_p^T$
- 3: **for** $m = 1 : M$ **do**
- 4: $\bar{\mathbf{u}}_m^{(0)} \leftarrow \mathbf{0}$
- 5: **repeat**
- 6: **for** $k = 1 : K_s$ **do**
- 7: $x_{mk} \leftarrow r_{mk}^{(0)} - \sum_{l \neq k} s_{kl}^{(0)} u_{ml}^{(0)}$
- 8: $u_{mk}^{(0)} \leftarrow \frac{(|x_{mk}| - \frac{1}{2} \lambda_1)_+ \text{sign}(x_{mk})}{s_{kk}^{(0)}}$
- 9: **end for**
- 10: **until** convergence
- 11: **end for**
- 12: **return** \mathbf{U}_0

Algorithm 3 Update U_p

Require: $\mathbf{D}_p, \mathbf{U}_0, \mathbf{V}_p$

- 1: $\mathbf{S}_p \leftarrow \mathbf{W}_p \mathbf{W}_p^T$
- 2: $\mathbf{R}_p \leftarrow \mathbf{D}_p \mathbf{W}_p^T - \mathbf{U}_0 \mathbf{H}_p \mathbf{W}_p^T$
- 3: **for** $m = 1 : M$ **do**
- 4: $\bar{\mathbf{u}}_m^{(p)} \leftarrow \mathbf{0}$
- 5: **repeat**
- 6: **for** $k = 1 : K_c$ **do**
- 7: $x_{mk} \leftarrow r_{mk}^{(p)} - \sum_{l \neq k} s_{kl}^{(p)} u_{ml}^{(p)}$
- 8: $u_{mk}^{(p)} \leftarrow \frac{(|x_{mk}| - \frac{1}{2} \lambda_1)_+ \text{sign}(x_{mk})}{s_{kk}^{(p)}}$
- 9: **end for**
- 10: **until** convergence
- 11: **end for**
- 12: **return** \mathbf{U}_p

Efficiency Comparison

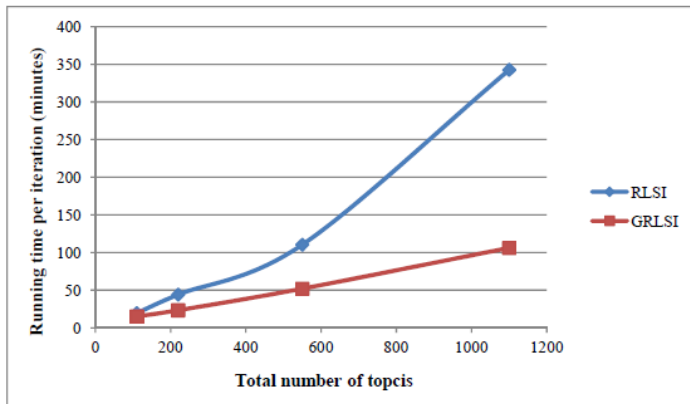


Figure 2: Execution time of RLSI and GRLSI on Wikipedia.

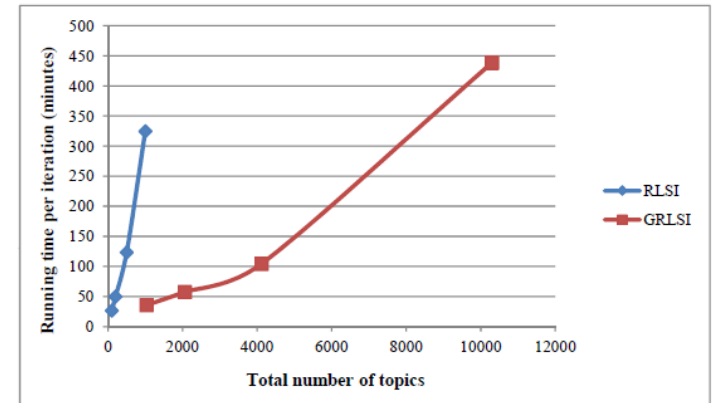


Figure 4: Execution time of RLSI and GRLSI on Web-I.

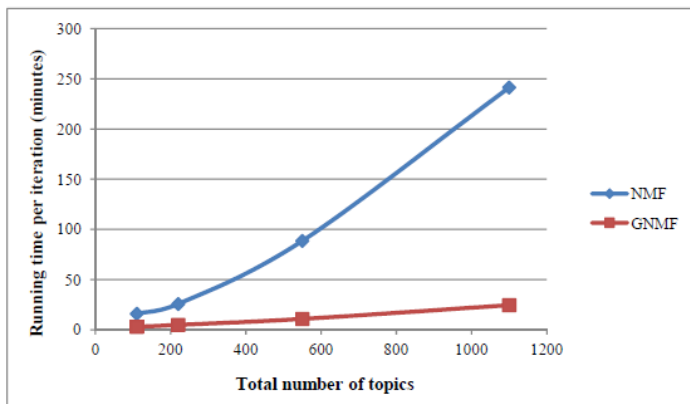


Figure 3: Execution time of NMF and GNMF on Wikipedia.

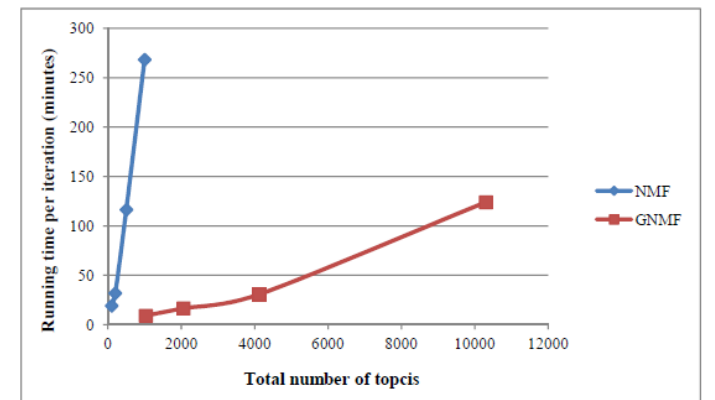


Figure 5: Execution time of NMF and GNMF on Web-I.

Experimental Results on Topic Discovery

Table 10: Topics discovered by GRLSI (top) and GNMF (bottom) on Wikipedia.

	Shared topics			Arts			Geography			Politics		
GRLSI	commune	state	political	album	rock	groups	province	municipality	communes	elections	states	kingdom
	communes	highways	party	albums	american	musical	state	municipalities	commune	election	congressional	political
	department	route	colour	singers	musicians	music	village	gmina	department	weapon	delegations	parties
	places	highway	india	musicians	singers	rappers	villages	voivodeship	france	party	elections	country
	france	india	canada	track	country	metal	highways	population	departments	parties	united	party
	populated	brazil	australia	listing	english	heavy	united	germany	places	political	senate	fascism
	municipality	oregan	parties	pop	pop	created	states	spain	county	results	tennessee	submarine
GNMF	places	new	language	album	groups	rappers	village	district	department	elections	war	military
	populated	york	japanese	albums	rock	musicians	villages	germany	commune	election	world	country
	village	city	films	track	american	american	england	districts	communes	results	poland	units
	azerbaijan	zealand	cast	listing	musical	singers	india	town	france	members	weapons	formations
	population	jersey	chinese	released	metal	singles	population	administrative	departments	parties	conflict	army
	municipality	routes	english	band	musicians	wiley	central	towns	home	held	union	infantry
	census	south	spanish	records	pop	blues	within	cities	western	general	force	established

Table 11: Topics discovered by GRLSI (top) and GNMF (bottom) on Web-I.

	Shared topics			Arts/literature			Business/healthcare			Computers/internet		
GRLSI	video	business	games	poems	harry	book	dental	healthcare	care	chat	facebook	web
	phone	services	game	poetry	potter	chapter	dentist	practice	medical	teen	people	hosting
	mobile	company	cheats	poem	books	summary	care	test	health	online	connect	design
	tv	service	xbox	poets	rowling	books	dentistry	management	equipment	people	sign	website
	cell	products	ign	love	series	analysis	dentists	exam	ppo	friends	web	domain
	phones	management	pc	poet	children	author	health	patient	supplies	join	password	internet
	iphone	jobs	updates	american	deathly	study	teeth	jobs	company	video	friends	services
GNMF	www	products	day	poems	harry	books	dentist	healthcare	medical	google	facebook	design
	http	product	october	quotes	potter	children	dentists	management	equipment	maps	people	web
	org	quality	september	shakespear	rowling	read	dentistry	patient	supplies	blog	connect	website
	website	buy	july	william	series	reading	dr	hospital	surgical	gmail	sign	development
	net	accessories	june	poetry	deathly	list	dental	solutions	patient	map	friends	marketing
	html	store	august	poets	hallows	readers	cosmetic	nursing	hospital	engine	password	graphic
	web	supplies	april	poem	stone	fiction	teeth	hospitals	device	web	share	logo

Experimental Results on Web Search

Table 12: Relevance performance of RLSI families on Web-II.

Method	MAP	NDCG@1	NDCG@3	NDCG@5	NDCG@10
BM25	0.3006	0.3043	0.3490	0.3910	0.4805
BM25+RLSI	0.3050	0.3076	0.3539	0.3943	0.4858
BM25+CRLSI	0.3027	0.3051	0.3509	0.3927	0.4840
BM25+GRLSI	0.3039	0.3066	0.3520	0.3934	0.4855

Table 13: Relevance performance of NMF families on Web-II.

Method	MAP	NDCG@1	NDCG@3	NDCG@5	NDCG@10
BM25	0.3006	0.3043	0.3490	0.3910	0.4805
BM25+NMF	0.3057	0.3091	0.3546	0.3960	0.4895
BM25+CNMF	0.3033	0.3055	0.3512	0.3934	0.4869
BM25+GNMF	0.3046	0.3080	0.3530	0.3955	0.4887

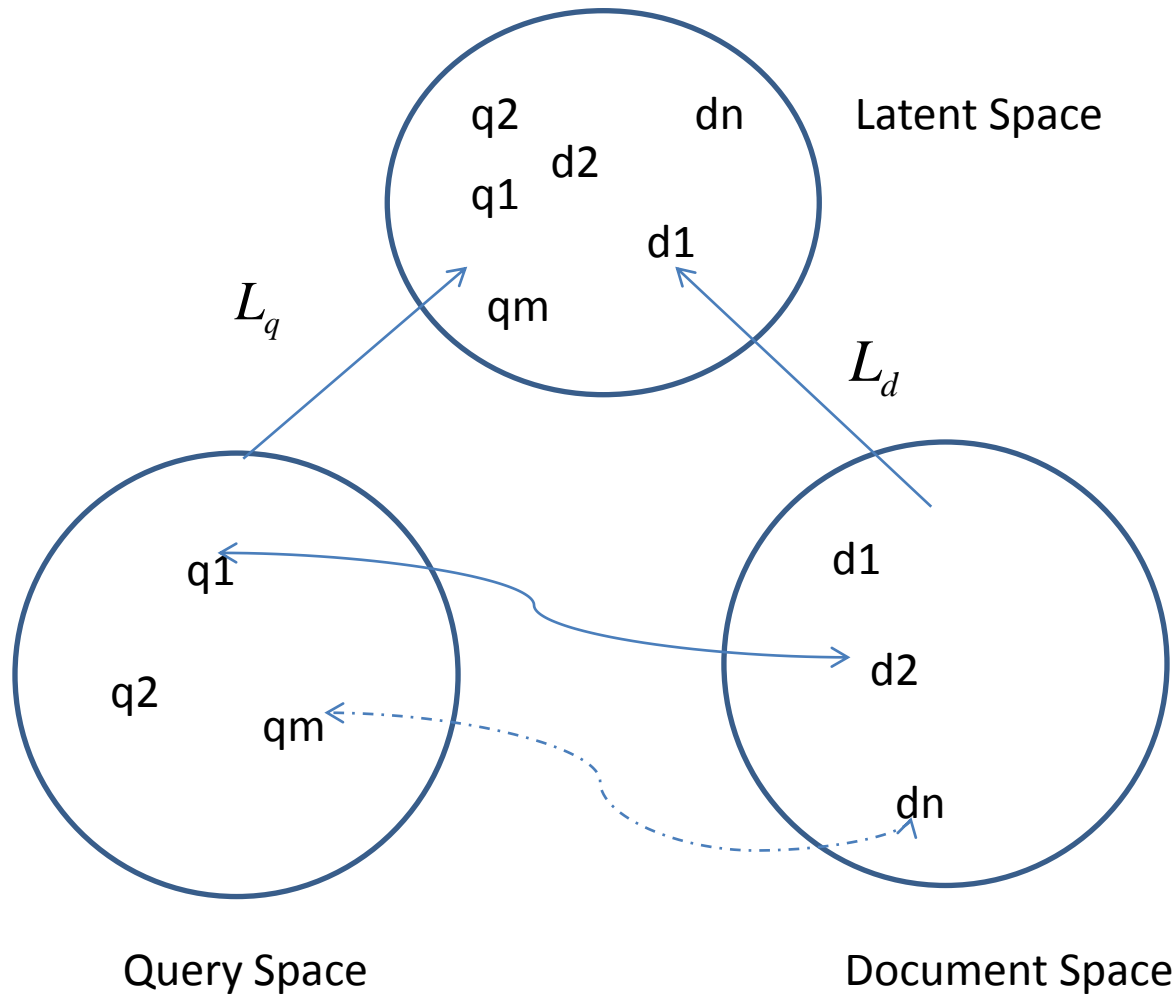
Matching in Latent Space

Joint Work with Wei Wu, Zhengdong Lv
Under review

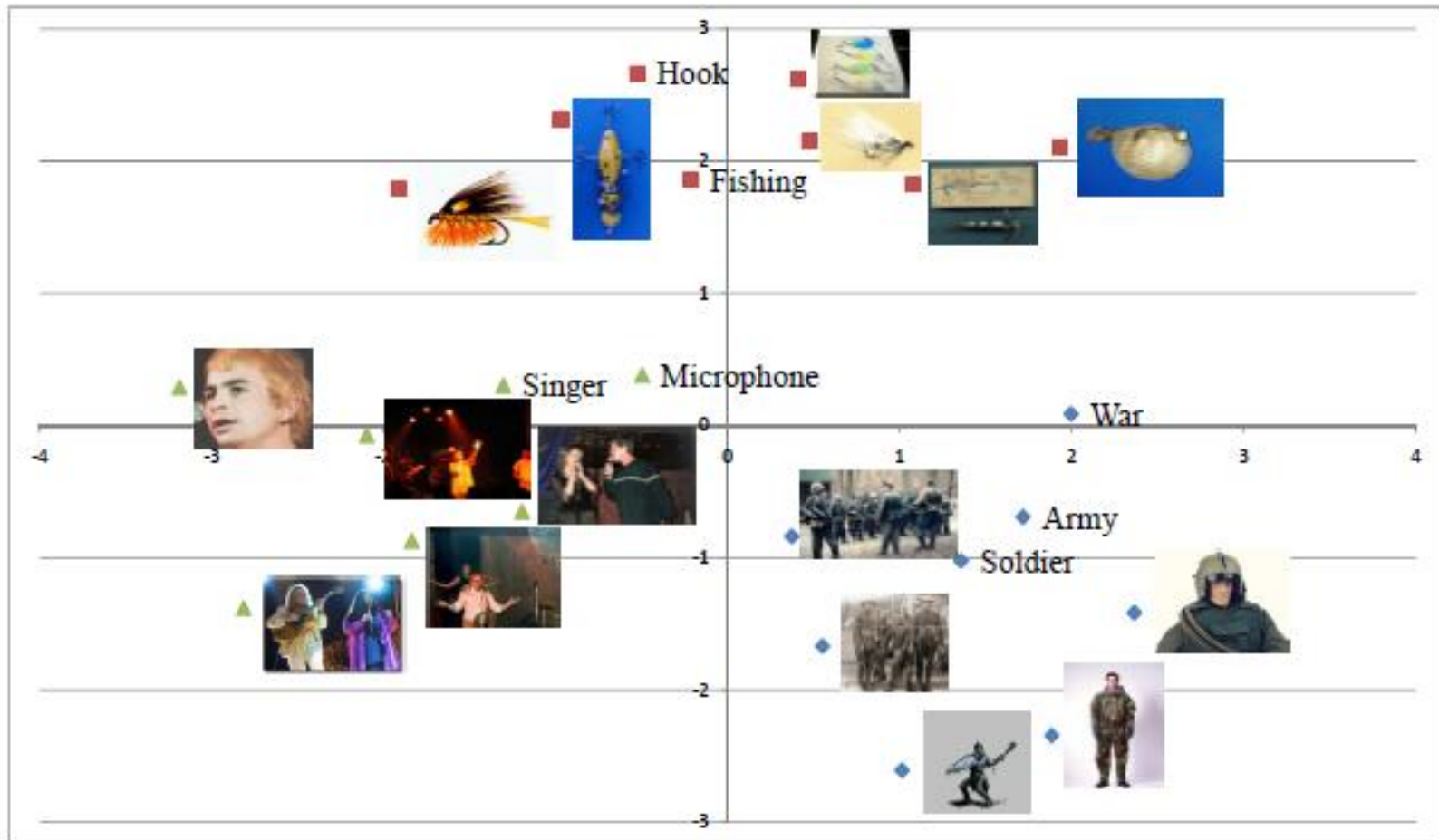
Matching in Latent Space

- Motivation
 - Matching between query and document in latent space
- Assumption
 - Queries have similarity
 - Document have similarity
 - Click-through data represent “similarity” relations between queries and documents
- Approach
 - Projection to latent space
 - Regularization or constraints
- Results
 - Significantly enhance accuracy of query document matching

Matching in Latent Space



Example: Projecting Keywords and Images into Latent Space



Partial Least Square (PLS)

- Setting
 - Two spaces: $\mathcal{X} \subset \mathbb{R}^m$ and $\mathcal{Y} \subset \mathbb{R}^n$.
- Input
 - Training data: $\{(x_i, y_i, r_i)\}_{1 \leq i \leq N}$, $r_i \in \{+1, -1\}$
- Output
 - Similarity function $f(x, y)$
- Assumption
 - Two linear (and orthonormal) transformations L_x and L_y
 - Dot product as similarity function $\langle L_x^T x, L_y^T y \rangle = x^T L_x L_y^T y$

- Optimization

$$\operatorname{argmax}_{L_x, L_y} \sum_{r_i=+1} x_i^T L_x L_y^T y_i - \sum_{r_i=-1} x_i^T L_x L_y^T y_i$$
$$\text{subject to } L_x^T L_x = I_{k \times k}, L_y^T L_y = I_{k \times k}$$

Solution of Partial Least Square

- Non-convex optimization
- Can prove that global optimal solution exists
- Global optimal can be found by solving SVD (Singular Value Decomposition)
- SVD of Matrix $M_S - M_D = U\Sigma V^T$

Regularized Mapping to Latent Space (RMLS)

- Setting
 - Two spaces: $\mathcal{X} \subset \mathbb{R}^m$ and $\mathcal{Y} \subset \mathbb{R}^n$.
- Input
 - Training data: $\{(x_i, y_i, r_i)\}_{1 \leq i \leq N}$, $r_i \in \{+1, -1\}$
- Output
 - Similarity function $f(x, y)$
- Assumption
 - L1 and L2 regularization on L_x and L_y (sparse transformations)
 - Dot product as similarity function $\langle L_x^T x, L_y^T y \rangle = x^T L_x L_y^T y$
- Optimization

$$\begin{aligned} \operatorname{argmax}_{L_x, L_y} & \sum_{r_i=+1} x_i^T L_x L_y^T y_i - \sum_{r_i=-1} x_i^T L_x L_y^T y_i \\ \text{subject to} & \|L_x\| \leq \vartheta_x, \|L_y\| \leq \vartheta_y, \|L_x\| \leq \lambda_x, \|L_y\| \leq \lambda_y, \end{aligned}$$

Solution of Regularized Mapping to Latent Space

- Coordinate Descent
- Repeat
 - Fix Lx , update Ly
 - Fix Ly , update Lx
- Update can be parallelized by rows

Comparison

	PLS	RMLS
Assumption	Orthogonal	L1 and L2 Regularization
Optimization Method	Singular Value Decomposition	Coordinate Descent
Optimality	Global optimum	Local optimum
Efficiency	Low	High
Scalability	Low	High

Experimental Results

Enterprise Search				Web Search			
	NDCG@1	NDCG@3	NDCG@5		NDCG@1	NDCG@3	NDCG@5
MPLS _{Com}	0.715	0.733	0.747	MPLS _{Com}	0.681	0.731	0.739
MPLS _{Conca}	0.700	0.728	0.742	MPLS _{Conca}	0.676	0.728	0.736
MPLS _{Word}	0.688	0.718	0.739	MPLS _{Word}	0.674	0.726	0.732
MPLS _{Bipar}	0.659	0.684	0.705	MPLS _{Bipar}	0.612	0.680	0.693
BM25	0.653	0.657	0.663	BM25	0.637	0.690	0.690
RW	0.654	0.683	0.700	RW	0.655	0.704	0.704
RW+BM25	0.664	0.688	0.705	RW+BM25	0.671	0.718	0.716
LSI	0.656	0.676	0.695	LSI	0.588	0.665	0.676
LSI+BM25	0.692	0.701	0.712	LSI+BM25	0.649	0.705	0.706

- RMLS and PLS work better than BM25, Random Walk, Latent Semantic Indexing
- RMLS works equally well as PLS, with higher learning efficiency and scalability

Conclusion

Conclusion

- Large scale topic modeling techniques
 - Regularized Latent Semantic Indexing
 - Group Matrix Factorization
- Large scale matching techniques
 - Matching in Latent Space
- Comparable with existing methods in terms of accuracy, much better in terms of efficiency and scalability
- Useful for web search

Publications of the Project

- **Quan Wang, Zheng Cao, Jun Xu, Hang Li, Group Matrix Factorization for Scalable Topic Modeling, In Proceedings of the 35th Annual International ACM SIGIR Conference (SIGIR'12), to appear, 2012.**
- Xiaobing Xue, Yu Tao, Daxin Jiang and Hang Li, Automatically Mining Question Reformulation Patterns from Search Log Data, In Proceedings of the 50th Annual Meeting of Association for Computational Linguistics (ACL'12), to appear, 2012.
- Fan Bu, Hang Li, Xiaoyan Zhu, String Re-Writing Kernel, In Proceedings of the 50th Annual Meeting of Association for Computational Linguistics (ACL'12), to appear, 2012.
- Chen Wang, Keping Bi, Yunhua Hu, Hang Li, and Guihong Cao. Extracting Search-Focused Key N-Grams for Relevance Ranking in Web Search. In Proceedings of the 3rd ACM International Conference on Web Search and Data Mining (WSDM'12), 343-352, 2012.
- Wei Wu, Jun Xu, Hang Li, and Satoshi Oyama, Learning A Robust Relevance Model for Search Using Kernel Methods, Journal of Machine Learning Research, 12, 1429-1458. 2011.
- **Quan Wang, Jun Xu, Hang Li, Nick Craswell, Regularized Latent Semantic Indexing, In Proceedings of the 34th Annual International ACM SIGIR Conference (SIGIR'11), 685-694, 2011.**
- Ziqi Wang, Gu Xu, Hang Li and Ming Zhang, A Fast and Accurate Method for Approximate String Search, In Proceedings of the 49th Annual Meeting of Association for Computational Linguistics: Human Language Technologies (ACL-HLT'11), 52-61, 2011.
- Jun Xu, Hang Li, Chaoliang Zhong, Relevance Ranking Using Kernels, In Proceedings of the 6th Asian Information Retrieval Societies Symposium (AIRS'10), Best Paper Award, 1-12, 2010.
- Jiafeng Guo, Gu Xu, Hang Li, Xueqi Cheng. A Unified and Discriminative Model for Query Refinement. In Proceedings of the 31st Annual International ACM SIGIR Conference (SIGIR'08), 379-386, 2008.
- **Wei Wu, Zhengdong Lv, Hang Li, Regularized Mapping to Latent Structures and Its Application to Web Search, under review.**

Thank You!

Contact: hangli.hl@huawei.com