TOPIC MODELING WITH STRUCTURED PRIORS FOR TEXT-DRIVEN SCIENCE

MICHAEL J. PAUL JOHNS HOPKINS UNIVERSITY

















EXAMPLES

• Disease monitoring in Twitter



Lamb, Paul, Dredze (2013) Separating fact from fear: Tracking flu infections on Twitter. NAACL.

Broniatowski, Paul, Dredze (2013) National and local influenza surveillance through Twitter: An analysis of the 2012-2013 influenza epidemic. *PLOS ONE* 8(12): e83672.

• Air pollution in Chinese social media



Wang, Paul, Dredze (2015) Social media as a sensor of air quality and public response in China. *Journal of Medical Internet Research*.

• Measuring healthcare quality from online reviews

Wallace, Paul, Sarkar, Trikalinos, Dredze (2014) A large-scale quantitative analysis of latent factors and sentiment in online doctor reviews. *Journal of the American Medical Informatics Association* 21(6), 1098-1103. 7

MAKING SENSE







A topic model is a probabilistic model of text

We pretend that our data (text) are the output of a probabilistic process that generates data

sick sore throat feel fever flu ...

allergies
nose
eyes
allergy
allergic
sneezing

watch
watching
tv
killing
movie
seen

class school read test doing finish

. . .









	Michael Paul @mjp39 · Jan 24						Ŷ	
	•		*	•••				



	Michael Paul @mjp39 · Jan 24					
	 ← ti ★ 	•••				



Michael Paul @mjp39 · Jan 24	Ŷ
◆ t3 ★ ···	



	👧 Mic	hael Paul	@mjp39 ·	Jan 24	\$
	(Cherny)			fever	
	•		*	•••	



	De M	lichae	el Paul @	≷mjp39 · J	lan 24		•
	s ann an				fever		
	•	h		*	•••		



	Micha	ael Paul @	∮mjp39 · 、	Jan 24	Q
	steen v			fever	
	•		*	•••	



	Michael Paul @mjp39 · Jan 24				0	
	S. Sector			fever		
	•		*	•••		



	Michael Paul @mjp39 · Jan 24	Q				
_	fever					
	watching					
	◆ t3 ★ ···					



Michael Paul @mjp39 · Jan 24	•
fever	
◆ t3 ★ ···	









Our imaginary process also needs to generate all these distributions





Our imaginary process also needs to generate all these distributions

- We need a distribution over distributions
 - Called a prior distribution





TOPIC M	ODELING	PRI	ORS	

























Latent Dirichlet Allocation (LDA) Blei, Ng, Jordan 2003

The topic and word distributions have Dirichlet priors: $\phi_t \sim \text{Dirichlet}(\tilde{\phi})$ $\theta_m \sim \text{Dirichlet}(\tilde{\theta})$





Paul, Dredze (2011) You are what you tweet: Analyzing Twitter for public health. 5th International Conference on Weblogs and Social Media (ICWSM).

Paul, Dredze (2014) **Discovering** health topics in social media using topic models. *PLOS ONE* 9(8).



Topics can be organized in ways that are more interpretable to users

ADDING STRUCTURE

Topics in online doctor reviews:



Both have positive sentiment



Both about staff/office issues

TOPIC MODELING ADDING STRUCTURE

Topics in online doctor reviews:



	Staff/Office	Personality	Surgery
Positive	time	best	surgery
	staff	years	first
	great	caring	son
	helpful	care	life
	feel	patients	surgeon
	questions	patient	daughter
	office	recommend	recommend
	friendly	family	thank
Negative	office	care	pain
Ŭ	time	medical	told
	appointment	patients	went
	rude	doesn't	said
	staff	help	surgery
	room	know	later
	didn't	don't	didn't
	wait	problem	months
A multi-dimensional topic model

Word distributions are grouped into different concepts

• e.g. sentiment and aspect

Paul and Dredze (2012) Factorial LDA: Sparse multi-dimensional text models. Proceedings of *Advances in Neural Information Processing Systems (NIPS)*.

DRUG DISCUSSIONS

Analyzing online drug forums:





Paul, Dredze (2013) Summarizing drug experiences with multi-dimensional topic models. North American ACL (NAACL). Paul, Chisolm, Johnson, Vandrey, Dredze (in preparation) Who participates in online drug communities? A largescale analysis of demographic and temporal trends³⁸.

DRUG DISCUSSIONS

3-dimensional model:

Drugs-Forum

- Route of administration (i.e. method of intake)
- Aspect

Drug type

Drug (22 total)	Route	Aspect		
 Alcohol Amphetamine Cannabis Cocaine 	 Injection Oral Smoking Snorting 	 Chemistry Culture Effects Health 		
SalviaTobacco		• Usage		



Suppose we want to model: (Marijuana, Oral, Chemistry)



DRUG DISCUSSIONS

Marijuana

weed cannabis thc marijuana stoned bowl bud joint blunt herb bong pot sativa blaze indica smoking blunts

. . .

Oral

capsules consumes toast stomach chewing ambien digestion juice absorbed ingestion meal tiredness chew juices gelatin yogurt fruit

. . .

solvent extraction evaporate evaporated solvents evaporation yield chloride alkaloids tek compounds evaporating atom aromatic non-polar purified jar

. . . .

Chemistry

DRUG DISCUSSIONS



DRUG DISCUSSIONS



L LDA DRUG DISCUSSIONS

FACTORIAL LDA

Dirichlet(

thc method extraction plant material cannabis simple coffee oil contains jar dried process dry water extract results . . .

DRUG DISCUSSIONS

word distribution for the triple:



oil water butter thc weed hash cannabis alcohol make milk high marijuana add . . . mixture hours trv brownies

~ Dirichlet(

thc method extraction plant material cannabis simple coffee oil contains jar dried process dry water extract results



. . .

DRUG DISCUSSIONS



oil water butter thc weed hash cannabis alcohol make milk high marijuana add mixture hours try brownies

~ Dirichlet(

thc method extraction plant material cannabis simple coffee oil contains jar dried process dry water extract results

. . .

DRUG DISCUSSIONS

word distribution for the triple:



oil water butter thc weed hash cannabis alcohol make milk high marijuana add . . . mixture hours trv brownies

~ Dirichlet(

thc method extraction plant material cannabis simple coffee oil contains jar dried process dry water extract results

. . .

DEFINITION

Prior for triple (*i*,*j*,*k*):

$$\tilde{\phi}_{(i,j,k)v} = \exp(\omega_{iv}^{(\text{drug})} + \omega_{jv}^{(\text{route})} + \omega_{kv}^{(\text{aspect})})$$

$$\phi_{(i,j,k)} \sim \text{Dirichlet}(\tilde{\phi}_{(i,j,k)})$$

distribution over words for this triple

In general, prior for tuple *t*:

$$\tilde{\phi}_{\vec{t}v} = \exp(\sum_{k=1}^{K} \omega_{\vec{t}_k v}^{(k)})$$

$$\phi_{\vec{t}} \sim \text{Dirichlet}(\tilde{\phi}_{\vec{t}})$$

DEFINITION

Prior for triple (*i*,*j*,*k*):

$$\tilde{\phi}_{(i,j,k)v} = \exp(\omega_{iv}^{(\text{drug})} + \omega_{jv}^{(\text{route})} + \omega_{kv}^{(\text{aspect})})$$

$$\phi_{(i,j,k)} \sim \text{Dirichlet}(\tilde{\phi}_{(i,j,k)})$$

distribution over words for this triple

In general, prior for tuple *t*:

$$\tilde{\phi}_{\vec{t}v} = \exp(\sum_{k=1}^{K} \omega_{\vec{t}_k v}^{(k)})$$

$$\phi_{\vec{t}} \sim \text{Dirichlet}(\tilde{\phi}_{\vec{t}})$$

Document priors: $\tilde{\theta}_{m\bar{t}} = \exp(\sum_{k=1}^{K} \alpha_{m\bar{t}_{k}}^{(k)})$ $\theta_{m} \sim \text{Dirichlet}(\tilde{\theta}_{m})$

SEMI-SUPERVISION

Marijuana

weed cannabis thc marijuana stoned bowl bud joint blunt herb bong pot sativa blaze indica smoking blunts

. . .

Oral

capsules consumes toast stomach chewing ambien digestion juice absorbed ingestion meal tiredness chew juices gelatin yogurt fruit

. . .

solvent extraction evaporate evaporated solvents evaporation yield chloride alkaloids tek compounds evaporating atom aromatic non-polar purified jar

. . . .

Chemistry

Where did these vectors come from?

SEMI-SUPERVISION

Weights learned from a supervised model are then used to create a **Gaussian prior** over the FLDA weights:

Health





We can use this model to extract specific information about new drugs

• e.g. dosage, desired effects, negative effects

Drugs-Forum

"What is the dosage when taking mephedrone orally?"



We can use this model to extract specific information about new drugs

• e.g. dosage, desired effects, negative effects

"What is the dosage when taking mephedrone orally?"



Drugs-Forum



We can use this model to extract specific information about new drugs

• e.g. dosage, desired effects, negative effects

Drugs-Forum

"What is the dosage when taking mephedrone orally?"



If it is [someone who isn't you]'s first time using Mephedrone [someone who isn't me] recommends a 100mg oral dose on an empty stomach.



We can use this model to extract specific information about new drugs

• e.g. dosage, desired effects, negative effects

Drugs-Forum

"What is the dosage when taking mephedrone orally?"



If it is [someone who isn't you]'s first time using Mephedrone [someone who isn't me] recommends a 100mg oral dose on an empty stomach.

Reference text:

It is recommended by users that Mephedrone be taken on an empty stomach. Doses usually vary between 100mg – 1g.



We can use this model to extract specific information about new drugs



Components:





SPARSITY



This Cartesian product can be huge!

• And not all triples make sense...



SPARSITY





• Proposed solution: learn a sparsity pattern

$$\tilde{\theta}_{m\bar{t}} = b_{\bar{t}} \exp(\sum_{k=1}^{K} \alpha_{m\bar{t}_{k}}^{(k)}) \qquad b_{\bar{t}} \in$$

$$b_{\vec{t}} \in (0,1); \ b_{\vec{t}} \sim \text{Beta}(\rho < 1)$$

SPARSITY



61

SPARSITY

"Topic"				"Approach"		"Focus"		
"SPEECH"	"I.R."	"М.Т."		"EMPIRICAL"	"THEORETICAL"	"METHODS"	"APPLICATIONS"	
speech	document	translation		task	theory	word	user	
spoken	retrieval	machine		tasks	description	algorithm	research	
recognition	documents	source		performance	formal	method	project	
state	question	mt		improve	forms	accuracy	technology	
vocabulary	web	parallel		accuracy	treatment	best	processing	
recognizer	answering	french		learning	linguistics	sentence	science	
utterances	query	bilingual		demonstrate	syntax	statistical	natural	
synthesis	answer	transfer		using	ed	previously	development	

To	Copic SPEECH		DATA		MODELING		GRAMMAR		
Focus		METHODS	APPL.	METHODS	APPL.	METHODS	APPL.	METHODS	APPL.
Approach		(b=0.20)	(b=1.00)	(b=1.00)	(b=1.00)	(b=1.00)	(b=0.50)	(b=1.00)	(b=0.57)
	AL		dialogue	corpus	data	models		parsing	grammar
			spoken	data	corpus	model		parser	parsing
	SIC		speech	training	annotation	approach		syntactic	based
	EMPIH		dialogues	model	annotated	shown		tree	robust
			understanding	tagging	corpora	error		parse	component
			task	annotated	collection	errors		dependency	processing
			recognition	test	xml	statistical		treebank	linguistic
		(b=0.99)	(b=0.00)	(b=0.07)	(b=0.02)	(b=1.00)	(b=0.01)	(b=1.00)	(b=1.00)
	ICAL	speech				rules		grammar	grammar
		words				rule		parsing	grammars
	ETJ	recognition				model		grammars	formalism
	ORI	prosodic				shown		structures	parsing
	IEC	written				models		paper	based
	TE	phonological				right		formalism	efficient
		spoken				left		based	unifi <u>ga</u> tion

FACTORIZATION













(SPARSE) DAG



Structured-prior topic models

A family of topic models in which the Dirichlet priors are functions of underlying components

Paul and Dredze (2015) SPRITE: Generalizing topic models with structured priors. *Transactions of the Association for Computational Linguistics (TACL)* 3: 43-57.

DEFINITION

The priors over distributions are weighted combinations of components:

$$\tilde{\phi}_{tv} = \exp(\sum_{c=1}^{C(\phi)} \beta_{tc} \omega_{cv})$$

$$\phi_t \sim \text{Dirichlet}(\tilde{\phi}_t)$$

distribution over words in *i*th topic



DEFINITION

The priors over distributions are weighted combinations of components:

$$\tilde{\theta}_{mt} = \exp(\sum_{c=1}^{C(\theta)} \alpha_{mc} \delta_{ct})$$

$$\theta_{m} \sim \text{Dirichlet}(\tilde{\theta}_{m})$$

$$\theta_{m} \sim \text{Dirichlet}(\tilde{\theta}_{m})$$

$$\theta_{m} = \frac{\theta_{m}}{\theta_{0}} \frac{\theta_{0}}{\theta_{0}} \frac$$

distribution over topics in *m*th document

DEFINITION

The priors over distributions are weighted combinations of components:



DEFINITION

The priors over distributions are weighted combinations of components:



DEFINITION

The priors over distributions are weighted combinations of components:



We can induce different structures by constraining the values of β
DEFINITION

The priors over distributions are weighted combinations of components:



DEFINITION

The priors over distributions are weighted combinations of components:



DEFINITION

The priors over distributions are weighted combinations of components:



DEFINITION

The priors over distributions are weighted combinations of components:



Tree: Each topic's β vector is zero in all but one component

DEFINITION

The priors over distributions are weighted combinations of components:



Factorization: Like a tree, but a nonzero component in each factor









• Initial solution: learn a sparsity pattern















SPRITE

SPECIAL CASES

SPRITE generalizes many existing topic models:

Model	Document priors	Topic priors
LDA	Single component	Single component
SCTM	Single component	Sparse binary β
SAGE	Single component	Sparse ω
FLDA	Binary δ is transpose of β	Factored binary β
PAM	α are supertopic weights	Single component
DMR	α are feature values	Single component

PARAMETER ESTIMATION

Need to estimate values for the parameters:

- Word and topic distributions
 - Collapsed Gibbs sampling
- Component parameters (i.e. β , ω)
 - Gradient ascent



PARAMETER ESTIMATION

What if β has constraints?

$$\beta_{tc} \in \{0,1\}, \forall c$$
$$\sum_{c} \beta_{tc} = 1$$









EXAMPLE 1



Modeling perspective in text:

"What opinions are people tweeting about gun control?"

Benton, Paul, Hancock, Dredze. A structured model of topic and perspective in social media. In preparation.

89

EXAMPLE 1

Suppose we want to model how **perspective** influences topics

• e.g. certain topics are "pro" or "anti" gun control

A single-component SPRITE model:

 $\tilde{\phi}_{tv} = \exp(r_t \omega_v)$ The *v*th word's perspective association
The *t*th topic's perspective association

EXAMPLE 1

Suppose we want to model how **perspective** influences topics

• e.g. certain topics are "pro" or "anti" gun control

A single-component SPRITE model:

 $\tilde{\phi}_{tv} = \exp(r_t \omega_v)$ The *v*th word's perspective association $\tilde{\theta}_{mt} = \exp(\alpha_m r_t)$ The *m*th tweet's perspective association

EXAMPLE 1

Suppose we want to model how **perspective** influences topics

• e.g. certain topics are "pro" or "anti" gun control

A single-component SPRITE model:

 $\widetilde{\phi}_{tv} = \exp(r_t \omega_v)$ The *v*th word's perspective association $\widetilde{\theta}_{mt} = \exp(\alpha_m r_t)$ Prior is a function of:
• Hashtags (#GunControlNow vs #NoGunControl)
• Survey data (% gun ownership in each state)



EXAMPLE 1



• Mean error: 8.4



EXAMPLE 2

Suppose we want to model **perspective and** we want to organize topics in a **hierarchy**

EXAMPLE 2

Suppose we want to model **perspective and** we want to organize topics in a **hierarchy**

$$\tilde{\phi}_{tv} = \exp(r_t \omega_{0v} + \sum_{c=1}^{C(\phi)} \beta_{tc} \omega_{cv})$$

with (soft) constraints that β_i is an indicator vector



EXAMPLE 2



SUMMARY

- Organizes topics in a variety of useful ways
 - Can be tailored toward different applications
- Generalizes many topic models
 - While opening up new possibilities
- Allows practitioners to make sense of big text data
 - Can drive new scientific research



CONCLUSION

There's no end to exciting questions we can ask of big, open data

We need methods to understand what people are saying on the web and learn meaningful trends

This requires models that can discover patterns automatically, while accommodating user expectations



THANK YOU

WITH HELP FROM:

Committee: Mark Dredze, Jason Eisner, Eric Horvitz, Hanna Wallach **Funding:** Microsoft Research, NSF, JHU Dean's office

Public opinion:

- Adrian Benton
- Braden Hancock

Doctor reviews:

- Byron Wallace
- Urmimala Sarkar
- Thomas Trikalinos

Drug forums:

- Meg Chisolm
- Matthew Johnson
- Ryan Vandrey





Microsoft Research



center of excellenge



QUESTIONS?