Adding Domain Knowledge to Latent Topic Models

Xiaojin Zhu

Department of Computer Sciences University of Wisconsin-Madison

2011

- 4 目 ト - 4 日 ト - 4 日 ト

Acknowledgments



David M. Andrzejewski did most work; now postdoc at Livermore National Laboratory



Mark Craven Ben Liblit Ben Recht University of Wisconsin-Madison

The Wisconsin Alumni Research Foundation, NSF CAREER Award IIS-0953219, AFOSR FA9550-09-1-0313, and NIH/NLM R01 LM07050

< 3 > < 3 >

New Year's Eve, Times Square



Knowledge \rightarrow LDA

3

・ロト ・聞 ト ・ ヨト ・ ヨト …

The Wish Corpus

[Goldberg et al., NAACL 2009]

- Peace on earth
- own a brewery
- I hope I get into Univ. of Penn graduate school.
- The safe return of my friends in Iraq
- find a cure for cancer
- To lose weight and get a boyfriend
- I Hope Barack Obama Wins the Presidency
- To win the lottery!

Corpus-wide word frequencies



Knowledge \rightarrow LDA

くほと くほと くほと

Some Topics by Latent Dirichlet Allocation (LDA) [Blei *et al.*, JMLR 2003]



- 4 3 6 4 3 6

Some Topics by Latent Dirichlet Allocation (LDA) [Blei *et al.*, JMLR 2003]

$p(word \mid topic)$





"election"



"love"

4 3 > 4 3 >

"troops"

Major applications: exploratory data analysis

- Research trends [Wang & McCallum, 2006]
- Scientific influence [Gerrish & Blei, 2009]
- Matching papers to reviewers [Mimno & McCallum, 2007]

Dirichlet [20, 5, 5]

2

通 ト イヨ ト イヨト



7 / 42

B ▶ < B ▶



Observed counts [3, 1, 2]

A, A, B, C, A, C

B ▶ < B ▶

Knowledge \rightarrow LDA



7 / 42



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\begin{array}{l} \phi_t \sim {\rm Dirichlet}(\beta) \\ \mbox{For each document } d \\ \theta \sim {\rm Dirichlet}(\alpha) \\ \mbox{For each word position in } d \\ \mbox{topic } z \sim {\rm Multinomial}(\theta) \\ \mbox{word } w \sim {\rm Multinomial}(\phi_z) \\ \mbox{Inference goals: } p(z \mid w, \alpha, \beta), {\rm argmax}_{\phi, \theta} \, p(\phi, \theta \mid w, \alpha, \beta) \\ \mbox{(reminder on top)} \end{array}
```



```
A generative model for p(\phi, \theta, z, w \mid \alpha, \beta):
For each topic t
           \phi_t \sim \mathsf{Dirichlet}(\beta)
For each document d
          \theta \sim \text{Dirichlet}(\alpha)
          For each word position in d
                    topic z \sim \text{Multinomial}(\theta)
                    word w \sim \text{Multinomial}(\phi_z)
Inference goals: p(z \mid w, \alpha, \beta), \operatorname{argmax}_{\phi, \theta} p(\phi, \theta \mid w, \alpha, \beta)
(reminder on top)
```



```
A generative model for p(\phi, \theta, z, w \mid \alpha, \beta):
For each topic t
          \phi_t \sim \text{Dirichlet}(\beta)
For each document d
          \theta \sim \text{Dirichlet}(\alpha)
          For each word position in d
                    topic z \sim \text{Multinomial}(\theta)
                    word w \sim \text{Multinomial}(\phi_z)
Inference goals: p(z \mid w, \alpha, \beta), \operatorname{argmax}_{\phi, \theta} p(\phi, \theta \mid w, \alpha, \beta)
(reminder on top)
```



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t $\phi_t \sim \mathsf{Dirichlet}(\beta)$ For each document d $\theta \sim \text{Dirichlet}(\alpha)$ For each word position in dtopic $z \sim \text{Multinomial}(\theta)$ word $w \sim \text{Multinomial}(\phi_z)$ Inference goals: $p(z \mid w, \alpha, \beta)$, $\operatorname{argmax}_{\phi, \theta} p(\phi, \theta \mid w, \alpha, \beta)$ (reminder on top)



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\begin{array}{l} \phi_t \sim \mathsf{Dirichlet}(\beta) \\ \text{For each document } d \\ \theta \sim \mathsf{Dirichlet}(\alpha) \\ \text{For each word position in } d \\ \text{topic } z \sim \mathsf{Multinomial}(\theta) \\ \text{word } w \sim \mathsf{Multinomial}(\phi_z) \\ \text{Inference goals: } p(z \mid w, \alpha, \beta), \operatorname{argmax}_{\phi, \theta} p(\phi, \theta \mid w, \alpha, \beta) \\ (\text{reminder on top}) \end{array}
```



```
A generative model for p(\phi, \theta, z, w \mid \alpha, \beta):
For each topic t
           \phi_t \sim \mathsf{Dirichlet}(\beta)
For each document d
          \theta \sim \text{Dirichlet}(\alpha)
          For each word position in d
                    topic z \sim \text{Multinomial}(\theta)
                    word w \sim \text{Multinomial}(\phi_z)
Inference goals: p(z \mid w, \alpha, \beta), \operatorname{argmax}_{\phi, \theta} p(\phi, \theta \mid w, \alpha, \beta)
(reminder on top)
```

3 K K 3 K



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
\begin{array}{l} \phi_t \sim {\sf Dirichlet}(\beta) \\ {\sf For \ each \ document \ d} \\ \theta \sim {\sf Dirichlet}(\alpha) \\ {\sf For \ each \ word \ position \ in \ d} \\ {\sf topic \ z \sim {\sf Multinomial}(\theta)} \\ {\sf word \ w \sim {\sf Multinomial}(\phi_z)} \\ {\sf Inference \ goals: \ p(z \mid w, \alpha, \beta), {\rm argmax}_{\phi, \theta} \ p(\phi, \theta \mid w, \alpha, \beta)} \\ ({\sf reminder \ on \ top}) \end{array}
```



A generative model for $p(\phi, \theta, z, w \mid \alpha, \beta)$: For each topic t

```
 \begin{array}{l} \phi_t \sim {\sf Dirichlet}(\beta) \\ {\sf For \ each \ document \ } d \\ \theta \sim {\sf Dirichlet}(\alpha) \\ {\sf For \ each \ word \ position \ in \ } d \\ {\sf topic \ } z \sim {\sf Multinomial}(\theta) \\ {\sf word \ } w \sim {\sf Multinomial}(\phi_z) \\ {\sf Inference \ goals: \ } p(z \mid w, \alpha, \beta), {\rm argmax}_{\phi, \theta} \, p(\phi, \theta \mid w, \alpha, \beta) \\ {\sf (reminder \ on \ top)} \end{array}
```

When LDA Alone is not Enough

- LDA is unsupervised
- \bullet Often domain experts have knowledge in addition to data, want better topics ϕ

When LDA Alone is not Enough

- LDA is unsupervised
- \bullet Often domain experts have knowledge in addition to data, want better topics ϕ
- There has been many specialized LDA variants
- This talk: how to do (general) "knowledge + LDA" with 3 models:
 - topic-in-set
 - 2 Dirichlet forest
 - Fold.all

- E > - E >

Model 1: Topic-in-Set

2

イロト イヨト イヨト イヨト

topic dice $\phi \sim \text{Dir}(\beta)$, doc dice $\theta \sim \text{Dir}(\alpha)$, topic z Model 1: Topic-in-Set

Example Application: Statistical Software Debugging [Andrzejewski *et al.*, ECML 2007], [Zheng *et al.*, ICML 2006]

• Insert predicates into a software:

```
int x = my_func()
if (x > 5) {
    branch_42_true++
    '''
else {
    branch_42_false++
    '''
}
```

A B M A B M

topic dice $\phi \sim \text{Dir}(\beta)$, doc dice $\theta \sim \text{Dir}(\alpha)$, topic z Model 1: Topic-in-Set

Example Application: Statistical Software Debugging [Andrzejewski *et al.*, ECML 2007], [Zheng *et al.*, ICML 2006]

• Insert predicates into a software:

```
int x = my_func()
if (x > 5) {
    branch_42_true++
    ...
}
else {
    branch_42_false++
    ...
}
```

- $\bullet \ {\rm predicates} \to {\rm words} \ w$
- a software run $\rightarrow \operatorname{doc} d$
- we know which runs crashed and which didn't (extra knowledge)
- the hope: run LDA on crashed runs, some topics ϕ will correspond to "buggy behaviors"

- - E + - E +

Normal software usage topics dominate. No "bug" topic.

3

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Normal software usage topics dominate. No "bug" topic. Toy example:

• Actual usage (left) and bug (right) topics. Each pixel is a predicate.



・ 何 ト ・ ヨ ト ・ ヨ ト

Normal software usage topics dominate. No "bug" topic. Toy example:

• Actual usage (left) and bug (right) topics. Each pixel is a predicate.



• Synthetic success (left) and crashed (right) runs





· · · · · · · · ·

Normal software usage topics dominate. No "bug" topic. Toy example:

Actual usage (left) and bug (right) topics. Each pixel is a predicate.



• Synthetic success (left) and crashed (right) runs







・ 何 ト ・ ヨ ト ・ ヨ ト

LDA topics on crashed runs



∆LDA [Andrzejewski *et al.*, ECML 2007]

Model success and crashed runs jointly with T topics:

- fix t < T
- for all words in success runs, $z \in \{1 \dots t\}$ (restricted)
- for all words in crashed runs, $z \in \{1 \dots T\}$

New hope: $\phi_1 \dots \phi_t$ usage topics, $\phi_{t+1} \dots \phi_T$ bug topics

くほと くほと くほと

New Hope Succeeds

- Actual usage (left) and bug (right) topics. Each pixel is a predicate.
- Synthetic success (left) and crashed (right) runs





イロト イポト イヨト イヨト

• LDA topics on crashed runs



• Δ LDA topics on success and crashed runs







Generalize \triangle LDA to Topic-in-Set

[Andrzejewski & Zhu, NAACL'09 WS]

- The domain knowledge:
 - ▶ For each word position *i* in corpus, we are given a set $C_i \subset \{1 \dots T\}$, such that $z_i \in C_i$.

> < 프 > < 프 >

Generalize \triangle LDA to Topic-in-Set

[Andrzejewski & Zhu, NAACL'09 WS]

- The domain knowledge:
 - ▶ For each word position i in corpus, we are given a set $C_i \subset \{1 \dots T\}$, such that $z_i \in C_i$.
- Very easy to implement in collapsed Gibbs sampling:

$$P(z_i = v | \mathbf{z}_{-i}, \mathbf{w}, \alpha, \beta) \propto$$
$$\left(\frac{n_{-i,v}^{(d)} + \alpha}{\sum_{u}^{T} (n_{-i,u}^{(d)} + \alpha)}\right) \left(\frac{n_{-i,v}^{(w_i)} + \beta}{\sum_{w'}^{W} (\beta + n_{-i,v}^{(w')})}\right) \delta(v \in C_i)$$

- ▶ $n_{-i,v}^{(d)}$ is the number of times topic v is used in document d
- $n_{-i,v}^{(w_i)}$ is the number of times word w_i is generated by topic v
- both excluding position i
- Easy to relax the hard constraints

イロト イポト イヨト イヨト

Model 2: Dirichlet Forest

2

イロン イヨン イヨン イヨン

Dirichlet Forest Enables Interactive Topic Modeling [Andrzejewski *et al.* ICML 2009] LDA Topics on Wish Corpus:

Topic	Top words sorted by $\phi = p(word topic)$
0	love i you me and will forever that with hope
1	and health for happiness family good my friends
2	year new happy a this have and everyone years
3	that is it you we be t are as not s will can
4	my to get job a for school husband s that into
5	to more of be and no money stop live people
6	to our the home for of from end safe all come
7	to my be i find want with love life meet man
8	a and healthy my for happy to be have baby
9	a 2008 in for better be to great job president
10	i wish that would for could will my lose can
11	peace and for love all on world earth happiness
12	may god in all your the you s of bless 2008
13	the in to of world best win 2008 go lottery
14	me a com this please at you call 4 if 2 www

イロト イヨト イヨト

Dirichlet Forest Enables Interactive Topic Modeling [Andrzejewski *et al.* ICML 2009] LDA Topics on Wish Corpus:

Горіс	Top words sorted by $\phi = p(word topic)$
0	love i you me and will forever that with hope
1	and health for happiness family good my friends
2	year new happy a this have and everyone years
3	that is it you we be t are as not s will can
4	my to get job a for school husband s that into
5	to more of be and no money stop live people
6	to our the home for of from end safe all come
7	to my be i find want with love life meet man
8	a and healthy my for happy to be have baby
9	a 2008 in for better be to great job president
10	i wish that would for could will my lose can
11	peace and for love all on world earth happiness
12	may god in all your <mark>the</mark> you s of bless 2008
13	the in to of world best win 2008 go lottery
14	me a com this please at you call 4 if 2 www

(日) (周) (三) (三)
isolate(50 stopwords from existing topics)

Topic	Top words sorted by $\phi = p(word topic)$
0	love forever marry happy together mom back
1	health happiness good family friends prosperity
2	life best live happy long great time ever wonderful
3	out not up do as so what work don was like
4	go school cancer into well free cure college
5	no people stop less day every each take children
6	home safe end troops iraq bring war husband house
7	love peace true happiness hope joy everyone dreams
8	happy healthy family baby safe prosperous everyone
9	better job hope president paul great ron than person
10	make money lose weight meet finally by lots hope married
Isolate	and to for a the year in new all my 2008
12	god bless jesus loved know everyone love who loves
13	peace world earth win lottery around save
14	com call if 4 2 www u visit 1 3 email yahoo
Isolate	i to wish my for and a be that the in

3

イロト イヨト イヨト イヨト

isolate(50 stopwords from existing topics)

Topic	Top words sorted by $\phi = p(word topic)$
0	love forever marry happy together mom back
1	health happiness good family friends prosperity
2	life best live happy long great time ever wonderful
3	out not up do as so what work don was like
MIXED	go school cancer into well free cure college
5	no people stop less day every each take children
6	home safe end troops iraq bring war husband house
7	love peace true happiness hope joy everyone dreams
8	happy healthy family baby safe prosperous everyone
9	better job hope president paul great ron than person
10	make money lose weight meet finally by lots hope married
Isolate	and to for a the year in new all my 2008
12	god bless jesus loved know everyone love who loves
13	peace world earth win lottery around save
14	com call if 4 2 www u visit 1 3 email yahoo
Isolate	i to wish my for and a be that the in

3

(日) (周) (三) (三)

split([cancer free cure well],[go school into college])

0	love forever happy together marry fall
1	health happiness family good friends
2	life happy best live love long time
3	as not do so what like much don was
4	out make money house up work grow able
5	people no stop less day every each take
6	home safe end troops iraq bring war husband
7	love peace happiness true everyone joy
8	happy healthy family baby safe prosperous
9	better president hope paul ron than person
10	lose meet man hope boyfriend weight finally
Isolate	and to for a the year in new all my 2008
12	god bless jesus loved everyone know loves
13	peace world earth win lottery around save
14	com call if 4 www 2 u visit 1 email yahoo 3
Isolate	i to wish my for and a be that the in me get
Split	job go school great into good college
Split	mom husband cancer hope free son well

イロト イ理ト イヨト イヨト

split([cancer free cure well],[go school into college])

LOVE	love forever happy together marry fall
1	health happiness family good friends
2	life happy best live love long time
3	as not do so what like much don was
4	out make money house up work grow able
5	people no stop less day every each take
6	home safe end troops iraq bring war husband
7	love peace happiness true everyone joy
8	happy healthy family baby safe prosperous
9	better president hope paul ron than person
LOVE	lose meet man hope boyfriend weight finally
Isolate	and to for a the year in new all my 2008
12	god bless jesus loved everyone know loves
13	peace world earth win lottery around save
14	com call if 4 www 2 u visit 1 email yahoo 3
Isolate	i to wish my for and a be that the in me get
Split	job go school great into good college
Split	mom husband cancer hope free son well

(日) (周) (三) (三)

merge([love ... marry...],[meet ... married...]) (10 words total)

Topic	Top words sorted by $\phi = p(word topic)$
Merge	love lose weight together forever marry meet
success	health happiness family good friends prosperity
life	life happy best live time long wishes ever years
-	as do not what someone so like don much he
money	out make money up house work able pay own lots
people	no people stop less day every each other another
iraq	home safe end troops iraq bring war return
јоу	love true peace happiness dreams joy everyone
family	happy healthy family baby safe prosperous
vote	better hope president paul ron than person bush
Isolate	and to for a the year in new all my
god	god bless jesus everyone loved know heart christ
peace	peace world earth win lottery around save
spam	com call if u 4 www 2 3 visit 1
Isolate	i to wish my for and a be that the
Split	job go great school into good college hope move
Split	mom hope cancer free husband son well dad cure

æ

▲日 ▶ ▲圖 ▶ ▲ 国 ▶ ▲ 国 ▶ →

From LDA to Dirichlet Forest



```
For each topic t

\phi_t \sim \text{Dirichlet}(\beta) \ \phi_t \sim \text{DirichletForest}(\beta, \eta)

For each doc d

\theta \sim \text{Dirichlet}(\alpha)

For each word w

z \sim \text{Multinomial}(\theta)

w \sim \text{Multinomial}(\phi_z)
```

< 3 > < 3 >

Richer Knowledge Enabled by Dirichlet Forest

Two pairwise relational primitives:

- Must-Link(u, v)
 - ▶ e.g., "school" and "college" should be in the same topic
 - encoded as $\phi_t(u) \approx \phi_t(v)$ for $t = 1 \dots T$
 - can be both large or both small in a given topic

· · · · · · · · ·

Richer Knowledge Enabled by Dirichlet Forest

Two pairwise relational primitives:

- Must-Link(u, v)
 - ▶ e.g., "school" and "college" should be in the same topic
 - encoded as $\phi_t(u) \approx \phi_t(v)$ for $t = 1 \dots T$
 - can be both large or both small in a given topic
- Cannot-Link(u, v)
 - e.g., "college" and "cure" should not be in the same topic
 - no topic has $\phi_t(u), \phi_t(v)$ both high

通 ト イヨ ト イヨト

Richer Knowledge Enabled by Dirichlet Forest

Two pairwise relational primitives:

- Must-Link(u, v)
 - ▶ e.g., "school" and "college" should be in the same topic
 - encoded as $\phi_t(u) \approx \phi_t(v)$ for $t = 1 \dots T$
 - can be both large or both small in a given topic
- Cannot-Link(u, v)
 - e.g., "college" and "cure" should not be in the same topic
 - ▶ no topic has $\phi_t(u), \phi_t(v)$ both high

Complex knowledge:

- split = CL between groups, ML within group
- merge = ML between groups
- isolate = CL to all high prob words in all topics

・ 同 ト ・ ヨ ト ・ ヨ ト …

The Dirichlet Distribution is Insufficient for Must-Link

- Must-Link(school, college) with vocabulary {school, college, lottery}
- \bullet desired topics ϕ distribution on 3-simplex



< 3 > < 3 >

The Dirichlet Distribution is Insufficient for Must-Link

- Must-Link(school, college) with vocabulary {school, college, lottery}
- \bullet desired topics ϕ distribution on 3-simplex



通 ト イヨ ト イヨト

[Dennis III 1991], [Minka 1999]

• Dirichlet variance $V(i) = \frac{E(i)(1-E(i))}{1+\sum_j \beta_j}$ tied to mean $E(i) = \frac{\beta_i}{\sum_j \beta_j}$

• More flexible: $\phi \sim \text{DirichletTree}(\gamma)$

- Draw multinomial at each internal node
- Multiply probabilities from leaf to root





[Dennis III 1991], [Minka 1999]

- Dirichlet variance $V(i) = \frac{E(i)(1-E(i))}{1+\sum_j \beta_j}$ tied to mean $E(i) = \frac{\beta_i}{\sum_j \beta_j}$
- More flexible: $\phi \sim \text{DirichletTree}(\gamma)$
 - Draw multinomial at each internal node
 - Multiply probabilities from leaf to root





イロト 不得下 イヨト イヨト

[Dennis III 1991], [Minka 1999]

• Dirichlet variance $V(i) = \frac{E(i)(1-E(i))}{1+\sum_j \beta_j}$ tied to mean $E(i) = \frac{\beta_i}{\sum_j \beta_j}$

• More flexible: $\phi \sim \text{DirichletTree}(\gamma)$

- Draw multinomial at each internal node
- Multiply probabilities from leaf to root





イロト 不得下 イヨト イヨト

$$p(\phi|\gamma) = \left(\prod_{k}^{L} \phi^{(k)\gamma^{(k)}-1}\right) \left(\prod_{s}^{I} \frac{\Gamma\left(\sum_{k}^{C(s)} \gamma^{(k)}\right)}{\prod_{k}^{C(s)} \Gamma\left(\gamma^{(k)}\right)} \left(\sum_{k}^{L(s)} \phi^{(k)}\right)^{\Delta(s)}\right)$$

• C(s) = children of s, L(k) = leaves of subtree k

• $\Delta(s) = \text{InDegree}(s) - \text{OutDegree}(s)$ being zero recovers Dirichlet

< 回 ト < 三 ト < 三 ト

$$p(\phi|\gamma) = \left(\prod_{k}^{L} \phi^{(k)\gamma^{(k)}-1}\right) \left(\prod_{s}^{I} \frac{\Gamma\left(\sum_{k}^{C(s)} \gamma^{(k)}\right)}{\prod_{k}^{C(s)} \Gamma\left(\gamma^{(k)}\right)} \left(\sum_{k}^{L(s)} \phi^{(k)}\right)^{\Delta(s)}\right)$$

• C(s) =children of s, L(k) =leaves of subtree k

• $\Delta(s) = \text{InDegree}(s) - \text{OutDegree}(s)$ being zero recovers Dirichlet Dirichlet tree is conjugate to multinomial

$$p(\mathbf{w}|\gamma) = \prod_{s}^{I} \left(\frac{\Gamma\left(\sum_{k}^{C(s)} \gamma^{(k)}\right)}{\Gamma\left(\sum_{k}^{C(s)} \left(\gamma^{(k)} + n^{(k)}\right)\right)} \prod_{k}^{C(s)} \frac{\Gamma\left(\gamma^{(k)} + n^{(k)}\right)}{\Gamma(\gamma^{(k)})} \right)$$

Encode Must-Links with a Dirichlet Tree [Andrzejewski *et al.* ICML 2009]

- Compute Must-Link transitive closures
- ② Each transitive closures is a subtree with edges ηeta
- The root connects to
 - each transitive closure with edge $|closure|\beta|$
 - each individual word not in any closure with edge β



The Dirichlet Distribution is Insufficient for Cannot-Link

- Cannot-Link(school, cancer) with vocabulary {school, cancer, cure}
- desired topics ϕ distribution on 3-simplex



- cannot be represented by a Dirichlet ($\beta < 1$ not robust)
- cannot be represented by a Dirichlet Tree
- requires mixture of Dirichlet Trees (Dirichlet Forest)

Dirichlet Forest Example

- A toy example:
 - Vocabulary: {A,B,C,D,E,F,G}
 - Must-Link(A,B)
 - Cannot-Link(A,D), Cannot-Link(C,D), Cannot-Link(E,F)
- Cannot-Link graph on Must-Link closures and other words



Identify Cannot-Link connected components (they are independent)



Zhu (Wisconsin)

Knowledge \rightarrow LDA

Flip edges: "Can"-Links within components



"Can"-maximal-cliques (no more words can have high prob together)



Zhu (Wisconsin)

Knowledge \rightarrow LDA

Picking a Dirichlet Tree from the Dirichlet Forest Pick a Can-clique (subtree) $q^{(1)}$ for the first connected component



(**F**)

Picking a Dirichlet Tree from the Dirichlet Forest We picked the first Can-clique $q^{(1)} = 1$: give *ABC* most prob mass





D has little mass, will not occur with any of A, B, C: satisfying Cannot-Links



Zhu (Wisconsin)

Knowledge \rightarrow LDA

Mass distribution between (AB) and C flexible; A, B under Must-Link subtree



Picking a Dirichlet Tree from the Dirichlet Forest Pick a Can-clique $q^{(2)}$ for the second connected component



Picking a Dirichlet Tree from the Dirichlet Forest We picked the 2nd Can-clique $q^{(2)} = 2$



Zhu (Wisconsin)

Here is the chosen subtree: ${\cal F}$ will get prob mass, not ${\cal E}$



Inference with Dirichlet Forest

- In theory the number of Can-cliques is exponential, in practice the largest we've seen is 3
- Collapsed Gibbs sampling over topics $z_1 \dots z_N$ and subtree indices $q_j^{(r)}$ for topic $j = 1 \dots T$ and Cannot-Link connected components r

$$p(z_i = v | \mathbf{z}_{-i}, \mathbf{q}_{1:T}, \mathbf{w}) \propto (n_{-i,v}^{(d)} + \alpha) \prod_{s}^{I_v(\uparrow i)} \frac{\gamma_v^{(C_v(s \downarrow i))} + n_{-i,v}^{(C_v(s \downarrow i))}}{\sum_k^{C_v(s)} \left(\gamma_v^{(k)} + n_{-i,v}^{(k)}\right)}$$

$$p(q_j^{(r)} = q' | \mathbf{z}, \mathbf{q}_{-j}, \mathbf{q}_j^{(-r)}, \mathbf{w}) \propto \left(\sum_{k}^{M_{rq'}} \beta_k \right) \prod_{s}^{I_{j,r=q'}} \left(\frac{\Gamma\left(\sum_{k}^{C_j(s)} \gamma_j^{(k)}\right)}{\Gamma\left(\sum_{k}^{C_j(s)} (\gamma_j^{(k)} + n_j^{(k)})\right)} \prod_{k}^{C_j(s)} \frac{\Gamma(\gamma_j^{(k)} + n_j^{(k)})}{\Gamma(\gamma_j^{(k)})} \right)$$

通 ト イヨ ト イヨト

Model 3: Fold.all

æ

<ロ> (日) (日) (日) (日) (日)

Logic can Encode Very General Knowledge

• Topic-in-set and Dirichlet Forest (ML, CL) not general enough

· · · · · · · · ·

Logic can Encode Very General Knowledge

- Topic-in-set and Dirichlet Forest (ML, CL) not general enough
- Fold.all = First-Order Logic latent Dirichlet ALLocation
 - easy for domain experts to write rules
 - can describe very general domain knowledge
 - ★ can encode many existing LDA variants
 - efficient inference

Logic can Encode Very General Knowledge

- Topic-in-set and Dirichlet Forest (ML, CL) not general enough
- Fold.all = First-Order Logic latent Dirichlet ALLocation
 - easy for domain experts to write rules
 - can describe very general domain knowledge
 - ★ can encode many existing LDA variants
 - efficient inference
- A hybrid Markov Logic Network (MLN) [Wang & Domingos 2008] [Richardson & Domingos 2006], but with fast stochastic optimization

・ 何 ト ・ ヨ ト ・ ヨ ト

Domain Knowledge in Logic

- Key hidden predicate: Z(i,t) TRUE if topic $z_i = t$
- Observed predicates (anything goes):
 - W(i, v) TRUE if word $w_i = v$
 - D(i,j) TRUE if word position i is in document j
 - ▶ HasLabel(j, l) TRUE if document j has label l
 - S(i,k) TRUE if word position i is in document k

▶ ...

・ 何 ト ・ ヨ ト ・ ヨ ト

Domain Knowledge in Logic

- Key hidden predicate: Z(i,t) TRUE if topic $z_i = t$
- Observed predicates (anything goes):
 - W(i, v) TRUE if word $w_i = v$
 - D(i,j) TRUE if word position i is in document j
 - ▶ HasLabel(j, l) TRUE if document j has label l
 - S(i,k) TRUE if word position i is in document k
 - ▶ ...
- Domain knowledge-base $(\lambda_1, \psi_1) \dots (\lambda_L, \psi_L)$
 - $\blacktriangleright \ {\rm rules} \ \psi$
 - \blacktriangleright positive weights λ indicate strength of rule

Example:

$$\lambda_1 = 1, \psi_1 = "\forall i : W(i, \mathsf{embryo}) \Rightarrow Z(i, 3)''$$

 $\lambda_2 = 100, \psi_2 = \text{``}\forall i, j, t : W(i, \mathsf{movie}) \land W(j, \mathsf{film}) \Rightarrow \neg (Z(i, t) \land Z(j, t))''$

・ロト ・四ト ・ヨト ・ヨト ・ヨ
Propositionalization

• Let $G(\psi)$ be all ground formulas of ψ .

- $\blacktriangleright \ \psi = ``\forall i, j, t : W(i, \mathsf{movie}) \land W(j, \mathsf{film}) \Rightarrow \neg (Z(i, t) \land Z(j, t))''$
- One ground formula $g \in G(\psi)$ is $W(123, \text{movie}) \land W(456, \text{film}) \Rightarrow \neg(Z(123, 9) \land Z(456, 9))$

(本部)と 本語 と 本語を

Propositionalization

• Let $G(\psi)$ be all ground formulas of ψ .

- $\blacktriangleright \ \psi = ``\forall i, j, t : W(i, \mathsf{movie}) \land W(j, \mathsf{film}) \Rightarrow \neg (Z(i, t) \land Z(j, t))''$
- One ground formula $g \in G(\psi)$ is $W(123, \text{movie}) \land W(456, \text{film}) \Rightarrow \neg(Z(123, 9) \land Z(456, 9))$

• $|G(\psi)|$ combinatorial.

Let

$$\mathbb{1}_{g}(\mathbf{z}) = \left\{ \begin{array}{ll} 1, & \text{if } g \text{ is TRUE under } \mathbf{z} \\ 0, & \text{otherwise.} \end{array} \right.$$

3

topic dice $\phi \sim \text{Dir}(\beta)$, doc dice $\theta \sim \text{Dir}(\alpha)$, topic z Mo

$\mathsf{Fold.all} = \mathsf{LDA} + \mathsf{MLN}$

[Andrzejewski et al. IJCAI 2011]



$$p(\mathbf{z}, \phi, \theta \mid \mathbf{w}, \alpha, \beta) \\ \propto \left(\prod_{t}^{T} p(\phi_t \mid \beta)\right) \left(\prod_{j}^{D} p(\theta_j \mid \alpha)\right) \left(\prod_{i}^{N} \phi_{z_i}(w_i) \theta_{d_i}(z_i)\right)$$

Zhu (Wisconsin)

æ

イロト イヨト イヨト

topic dice $\phi \sim \text{Dir}(\beta)$, doc dice $\theta \sim \text{Dir}(\alpha)$, topic z Mo

$\mathsf{Fold.all} = \mathsf{LDA} + \mathsf{MLN}$

[Andrzejewski et al. IJCAI 2011]



$$p(\mathbf{z}, \phi, \theta \mid \mathbf{w}, \alpha, \beta)$$

$$\propto \left(\prod_{t}^{T} p(\phi_{t} \mid \beta)\right) \left(\prod_{j}^{D} p(\theta_{j} \mid \alpha)\right) \left(\prod_{i}^{N} \phi_{z_{i}}(w_{i})\theta_{d_{i}}(z_{i})\right)$$

$$\times \exp\left[\sum_{l}^{L} \sum_{g \in G(\psi_{l})} \lambda_{l} \mathbb{1}_{g}(\mathbf{z})\right]$$

Zhu (Wisconsin)

Fold.all Inference

MAP estimate, non-convex objective

$$Q(\mathbf{z}, \phi, \theta) \equiv \sum_{t}^{T} \log p(\phi_t | \beta) + \sum_{j}^{D} \log p(\theta_j | \alpha)$$
$$+ \sum_{i}^{N} \log \phi_{z_i}(w_i) \theta_{d_i}(z_i) + \sum_{l}^{L} \sum_{g \in G(\psi_l)} \lambda_l \mathbb{1}_g(\mathbf{z})$$

Alternating optimization. Repeat:

- fixing z, let $(\phi^*, \theta^*) \leftarrow \operatorname{argmax}_{\phi, \theta} Q(\mathbf{z}, \phi, \theta)$ (easy)
- fixing ϕ, θ , let $\mathbf{z}^* \leftarrow \operatorname{argmax}_z Q(\mathbf{z}, \phi, \theta)$ (integer)

・ 同 ト ・ ヨ ト ・ ヨ ト …

$$g = \mathsf{Z}(i,1) \lor \neg \mathsf{Z}(j,2)$$
, and $t \in \{1,2,3\}$

1 Take complement $\neg g$



(日) (周) (三) (三)

$$g = \mathtt{Z}(i,1) \vee \neg \mathtt{Z}(j,2) \text{, and } t \in \{1,2,3\}$$

- **1** Take complement $\neg g$
- **2** Remove negations $(\neg g)_+$

 $\neg \mathsf{Z}(i,1) \land \mathsf{Z}(j,2)$ $(\mathsf{Z}(i,2) \lor \mathsf{Z}(i,3)) \land \mathsf{Z}(j,2)$

(日) (周) (三) (三)

$$g = \mathsf{Z}(i,1) \lor \neg \mathsf{Z}(j,2)$$
, and $t \in \{1,2,3\}$

- **1** Take complement $\neg g$
- **2** Remove negations $(\neg g)_+$
- $I Numeric z_{it} \in \{0,1\}$

$$\neg \mathbf{Z}(i,1) \land \mathbf{Z}(j,2)$$
$$(\mathbf{Z}(i,2) \lor \mathbf{Z}(i,3)) \land \mathbf{Z}(j,2)$$
$$(z_{i2}+z_{i3})z_{j2}$$

(日) (周) (三) (三)

$$g = \mathtt{Z}(i,1) \vee \neg \mathtt{Z}(j,2) \text{, and } t \in \{1,2,3\}$$

- **1** Take complement $\neg g$
- **2** Remove negations $(\neg g)_+$
- I With Model States States
- Polynomial $1_g(\mathbf{z})$

$$\neg Z(i,1) \land Z(j,2) (Z(i,2) \lor Z(i,3)) \land Z(j,2) (z_{i2} + z_{i3})z_{j2} 1 - (z_{i2} + z_{i3})z_{j2}$$

- 4 週 ト - 4 三 ト - 4 三 ト

$$g = \mathtt{Z}(i,1) \vee \neg \mathtt{Z}(j,2) \text{, and } t \in \{1,2,3\}$$

- **1** Take complement $\neg g$
- 2 Remove negations $(\neg g)_+$
- I With Model States States
- Polynomial $\mathbb{1}_g(\mathbf{z})$
- **(5)** Relax discrete z_{it}

 $\neg Z(i, 1) \land Z(j, 2)$ $(Z(i, 2) \lor Z(i, 3)) \land Z(j, 2)$ $(z_{i2} + z_{i3})z_{j2}$ $1 - (z_{i2} + z_{i3})z_{j2}$ $z_{it} \in \{0, 1\} \rightarrow z_{it} \in [0, 1]$

$$g = \mathtt{Z}(i,1) \vee \neg \mathtt{Z}(j,2) \text{, and } t \in \{1,2,3\}$$

- **1** Take complement $\neg g$
- 2 Remove negations $(\neg g)_+$
- I With Model States States
- Polynomial $\mathbb{1}_g(\mathbf{z})$
- **(5)** Relax discrete z_{it}

 $\neg Z(i, 1) \land Z(j, 2)$ $(Z(i, 2) \lor Z(i, 3)) \land Z(j, 2)$ $(z_{i2} + z_{i3})z_{j2}$ $1 - (z_{i2} + z_{i3})z_{j2}$ $z_{it} \in \{0, 1\} \rightarrow z_{it} \in [0, 1]$

$$g = \mathtt{Z}(i,1) \vee \neg \mathtt{Z}(j,2) \text{, and } t \in \{1,2,3\}$$

1 Take complement $\neg g$ $\neg Z(i,1) \land Z(j,2)$ **2** Remove negations $(\neg g)_+$ $(Z(i,2) \lor Z(i,3)) \land Z(j,2)$ **3** Numeric $z_{it} \in \{0,1\}$ $(z_{i2} + z_{i3})z_{j2}$ **4** Polynomial $\mathbb{1}_g(\mathbf{z})$ $1 - (z_{i2} + z_{i3})z_{j2}$ **5** Relax discrete z_{it} $z_{it} \in \{0,1\} \rightarrow z_{it} \in [0,1]$

$$\mathbb{1}_{g}(\mathbf{z}) = 1 - \prod_{g_i \neq \emptyset} \left(\sum_{\mathbf{Z}(i,t) \in (\neg g_i)_+} z_{it} \right)$$

Optimizing z Step 2: Stochastic Optimization

• $\sum_{l}^{L} |G(\psi_l)| + NT$ terms in Q related to z:

$$\begin{split} \max_{\mathbf{z}} & \sum_{l}^{L} \sum_{g \in G(\psi_l)} \lambda_l \mathbb{1}_g(\mathbf{z}) + \sum_{i}^{N} \sum_{t}^{T} z_{it} \log \phi_t(w_i) \theta_{d_i}(t) \\ \text{s.t.} & z_{it} \geq 0, \quad \sum_{t}^{T} z_{it} = 1. \end{split}$$

- Entropic Mirror Descent [Beck & Teboulle, 2003]. Repeat:
 - select a term f at random
 - descent with decreasing step size η

$$z_{it} \leftarrow \frac{z_{it} \exp\left(\eta \nabla_{z_{it}} f\right)}{\sum_{t'} z_{it'} \exp\left(\eta \nabla_{z_{it'}} f\right)}$$

Example: Movie Reviews

 $\forall i, j, t: W(i, \mathsf{movie}) \land W(j, \mathsf{film}) \Rightarrow \neg(Z(i, t) \land Z(j, t))$





- 4 @ ▶ 4 @ ▶ 4 @ ▶

Generalization and Scalability

Cross-validation

- Training: do Fold.all MAP inference to estimate $\hat{\phi}$
- Testing: use trainset $\hat{\phi}$ to infer testset \hat{z} (no logic rules)
- Evaluation: testset objective Q
- "-": runs more than 24 hours

Data	Fold.all	LDA	Alchemy	$\sum_{l} G(\psi_l) $
Synth	9.86	-2.18	-1.73	10^{5}
Comp	2.40	1.19	-	10^{4}
Con	2.51	1.09	-	10^{3}
Pol	5.67	5.67	-	10^{9}
HDG	10.66	3.59	-	10^{8}

Summary

- "Knowledge + data" for latent topic models
- Increasingly more general models
 - Topic-in-set
 - Dirichlet forest (Must & Cannot links)
 - Fold.all (logic)
- Easy for users
- Scalable inference

(B)