

PCFGs, Topic Models, Adaptor Grammars and Learning Topical Collocations and the Structure of Proper Names

Mark Johnson
Department of Computing
Macquarie University
Mark.Johnson@mq.edu.au

July 4, 2010

Outline

LDA topic models as PCFGs

Adaptor grammars

Finding topic-specific collocations

Learning the structure of proper nouns

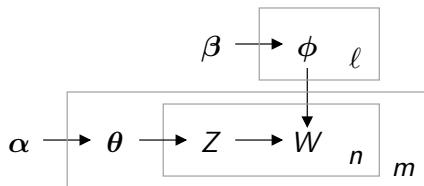
Conclusion

LDA topic models

- LDA topic models are generative models for documents
 - ▶ identifies documents about similar topics
 - ▶ identifies words characteristic of topics
- Each topic i is a distribution over words ϕ_i
- Each document j has a *distribution* θ_j over topics
- To generate document j :
 - ▶ for each word position in document:
 - choose a topic z according to θ_j , and then
 - choose a word belonging to that topic according to ϕ_z
- “Sparse priors” on ϕ and θ
 - ⇒ most documents have few topics
 - ⇒ most topics have few words
- Bayesian inference (Gibbs sampling, Variational Bayes)
See: Blei, Ng and Jordan (2002), Griffiths and Steyvers (2004)

LDA topic models: formal description

$$\begin{aligned}\phi_i &\sim \text{Dir}(\beta) & i = 1, \dots, \ell = \text{number of topics} \\ \theta_j &\sim \text{Dir}(\alpha) & j = 1, \dots, m = \text{number of documents} \\ z_{j,k} &\sim \theta_j & j = 1, \dots, m \\ & & k = 1, \dots, n = \text{number of words in a document} \\ w_{j,k} &\sim \phi_{z_{j,k}} & j = 1, \dots, m \\ & & k = 1, \dots, n\end{aligned}$$



Context-Free Grammars

- A CFG (N, W, R, S) defines *sets of trees* \mathcal{T}_X for each $X \in N \cup W$:
 - ▶ if $X \in W$ then $\mathcal{T}_X = \{X\}$ (the 1-node tree labelled X)
 - ▶ if $X \in N$ then:

$$\mathcal{T}_X = \bigcup_{X \rightarrow B_1 \dots B_n \in R_X} \text{TREE}_X(\mathcal{T}_{B_1}, \dots, \mathcal{T}_{B_n})$$

where $R_A = \{A \rightarrow \beta : A \rightarrow \beta \in R\}$ for each $A \in N$, and

$$\text{TREE}_X(\mathcal{T}_{B_1}, \dots, \mathcal{T}_{B_n}) = \left\{ \begin{array}{c} X \\ \wedge \\ t_1 \dots t_n \end{array} : \begin{array}{l} t_i \in \mathcal{T}_{B_i}, \\ i = 1, \dots, n \end{array} \right\}$$

That is, $\text{TREE}_X(\mathcal{T}_{B_1}, \dots, \mathcal{T}_{B_n})$ consists of the set of trees with whose root node is labelled X and whose i th child is a member of \mathcal{T}_{B_i} .

Probabilistic Context-Free Grammars

- A PCFG is a CFG (N, W, R, S) and multinomials θ_X over R_X for each $X \in N$
 - ▶ $\theta_{X \rightarrow \beta}$ is the probability of X expanding to β
- A PCFG associates each $X \in N \cup W$ with a *distribution* G_X over trees \mathcal{T}_X
 - ▶ if $X \in W$ then $G_X(X) = 1$
 - ▶ if $X \in N$ then:

$$G_X(t) = \sum_{X \rightarrow B_1 \dots B_n \in R_X} \theta_{X \rightarrow B_1 \dots B_n} \text{TD}_X(G_{B_1}, \dots, G_{B_n})(t) \quad (1)$$

where:

$$\text{TD}_A(G_1, \dots, G_n) \left(\begin{array}{c} X \\ \wedge \\ t_1 \dots t_n \end{array} \right) = \prod_{i=1}^n G_i(t_i).$$

That is, $\text{TD}_A(G_1, \dots, G_n)$ is a distribution over \mathcal{T}_A where each subtree t_i is generated independently from G_i .

Bayesian PCFGs

- Place Dirichlet priors $\text{Dir}(\alpha_X)$ on each rule probability multinomial θ_X for each $X \in N$

$$\theta_X \sim \text{Dir}(\alpha_X) \quad X \in N$$

- “Sparse priors” \Rightarrow prefer to use as few rules as possible
- Unsupervised Bayesian inference for PCFGs from strings:
 - ▶ MCMC sampling
 - ▶ Variational Bayes

See: Kurihara and Sato (2006), Johnson et al (2007)

LDA topic models as PCFGs (1)

- Prefix strings from document j with a *document identifier* “ $-j$ ”

Sentence \rightarrow Doc $'_j$ $j \in 1, \dots, m$

Doc $'_j \rightarrow -j$ $j \in 1, \dots, m$

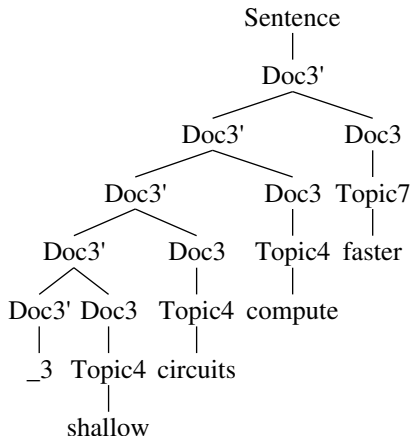
Doc $'_j \rightarrow$ Doc $'_j$ Doc $_j$ $j \in 1, \dots, m$

Doc $_j \rightarrow$ Topic $_i$ $i \in 1, \dots, \ell$

Topic $_i \rightarrow w$ $j \in 1, \dots, m$

Topic $_i \rightarrow w$ $i \in 1, \dots, \ell$

$w \in \mathcal{V}$



LDA topic models as PCFGs (2)

- Spine *propagates document id up through tree*

Sentence \rightarrow Doc' $_j$ $j \in 1, \dots, m$

Doc' $_j \rightarrow$ $-j$ $j \in 1, \dots, m$

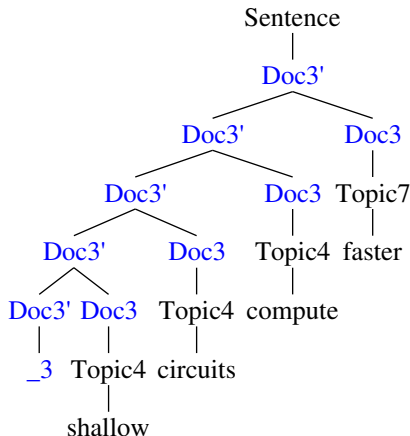
Doc' $_j \rightarrow$ Doc' $_j$ Doc $_j$ $j \in 1, \dots, m$

Doc $_j \rightarrow$ Topic $_i$ $i \in 1, \dots, \ell$

Topic $_i \rightarrow$ w $j \in 1, \dots, m$

Topic $_i \rightarrow$ w $i \in 1, \dots, \ell$

$w \in \mathcal{V}$



LDA topic models as PCFGs (3)

- $\text{Doc}_j \rightarrow \text{Topic}_i$; rules map *documents to topics*

Sentence $\rightarrow \text{Doc}'_j \quad j \in 1, \dots, m$

$\text{Doc}'_j \rightarrow _j \quad j \in 1, \dots, m$

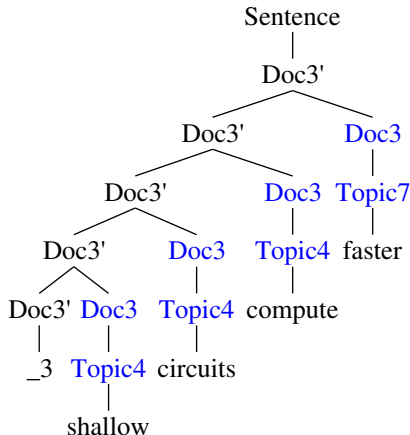
$\text{Doc}'_j \rightarrow \text{Doc}'_j \text{Doc}_j \quad j \in 1, \dots, m$

$\text{Doc}_j \rightarrow \text{Topic}_i \quad i \in 1, \dots, \ell$

$j \in 1, \dots, m$

$\text{Topic}_i \rightarrow w \quad i \in 1, \dots, \ell$

$w \in \mathcal{V}$



LDA topic models as PCFGs (4)

- $\text{Topic}_i \rightarrow w$ rules map *topics to words*

$\text{Sentence} \rightarrow \text{Doc}'_j \quad j \in 1, \dots, m$

$\text{Doc}'_j \rightarrow _j \quad j \in 1, \dots, m$

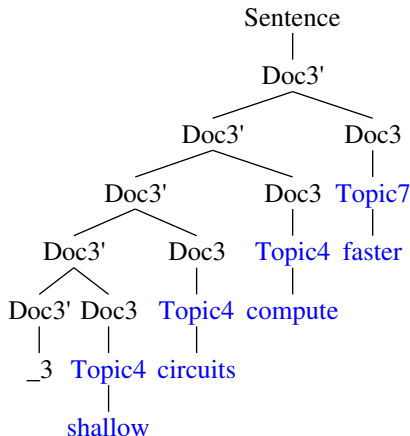
$\text{Doc}'_j \rightarrow \text{Doc}'_j \text{ Doc}_j \quad j \in 1, \dots, m$

$\text{Doc}_j \rightarrow \text{Topic}_i \quad i \in 1, \dots, \ell$

$\quad \quad \quad j \in 1, \dots, m$

$\text{Topic}_i \rightarrow w \quad i \in 1, \dots, \ell$

$w \in \mathcal{V}$

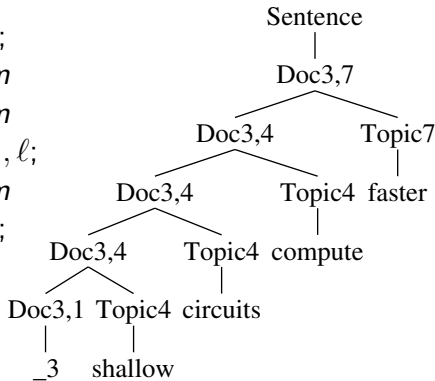


LDA topic models as PCFGs (5)

- *Not* suggesting blind use of PCFG inference for topic models
 - ▶ One iteration of LDA inference is *linear* in document length
 - ▶ One iteration of PCFG inference is *cubic* in document length
- Reduction of LDA topic models to PCFGs suggests ways of extending both kinds of models

“Sticky” topic models

$\text{Sentence} \rightarrow \text{Doc}_{j,i} \quad i \in 1, \dots, \ell;$
 $\quad \quad \quad \quad \quad \quad \quad j \in 1, \dots, m$
 $\text{Doc}_{j,1} \rightarrow _j \quad j \in 1, \dots, m$
 $\text{Doc}_{j,i} \rightarrow \text{Doc}_{j,i'} \text{ Topic}_i \quad i, i' \in 1, \dots, \ell;$
 $\quad \quad \quad \quad \quad \quad \quad j \in 1, \dots, m$
 $\text{Topic}_i \rightarrow w \quad i \in 1, \dots, \ell;$
 $\quad \quad \quad \quad \quad \quad \quad w \in \mathcal{V}$



- Prefer *adjacent words to have same topic*
- $\text{Doc}_{j,i}$ means “document j , topic i ”
- *Non-uniform Dirichlet prior* disprefers topic shift
 - ▶ $\alpha_{\text{Doc}_{j,i} \rightarrow \text{Doc}_{j,i} \text{ Topic}_i} \gg \alpha_{\text{Doc}_{j,i} \rightarrow \text{Doc}_{j,i'} \text{ Topic}_i}$ for $i' \neq i$

Outline

LDA topic models as PCFGs

Adaptor grammars

Finding topic-specific collocations

Learning the structure of proper nouns

Conclusion

From Multinomials to Dirichlet Processes

- Dirichlet Processes (DPs) are the *infinite-dimensional generalisation* of Dirichlet-Multinomials
- *Predictive distribution*: predict z_{n+1} given observations $\mathbf{z} = (z_1, \dots, z_n)$
 - ▶ *Finite set* of outcomes $(1, \dots, m)$:
Dirichlet-multinomial with prior $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_m)$

$$P(Z_{n+1} = k \mid \mathbf{z}) \propto n_k(\mathbf{z}) + \alpha_k$$

where $n_k(\mathbf{z})$ is the number of times k appears in $\mathbf{z} = (z_1, \dots, z_n)$

- ▶ *Infinite set* of outcomes Ω :
Dirichlet process $DP(\alpha, P_0)$ with *base distribution* $P_0(Z)$ and *concentration parameter* α

$$P(Z_{n+1} = z' \mid \mathbf{z}) \propto n_{z'}(\mathbf{z}) + \alpha P_0(z')$$

Dirichlet Processes as Adaptors

- DPs generalise Dirichlet-multinomials

$$P(Z_{n+1} = z' | \mathbf{z}) \propto n_{z'}(\mathbf{z}) + \alpha P_0(z')$$

- DPs follow a “*rich get richer*” law
 - ▶ frequent outcomes are increasingly likely to be predicted
- The DP is *stochastic*:
in general, every sample $\mathbf{z} = (z_1, z_2, \dots)$ is different
 - ⇒ DPs map a *base distribution* P_0 to a *distribution over distributions* $\text{DP}(\alpha, P_0)$
- *Pitman-Yor Processes* (PYPs) generalise Dirichlet Processes
- An *adaptor* is a function that *maps a base distribution* P_0 to a *distribution over distributions* with the same support as P_0
 - ▶ Dirichlet Processes and Pitman-Yor Processes are adaptors

Adaptor grammars as generalised PCFGs

- An *adaptor grammar* is a PCFG with a set $A \subseteq N$ of *adapted nonterminals*, and *adaptors* C_X for each $X \in A$
- *Dirichlet Process Adaptor Grammar*:
 - ▶ If $X \in W$ then $G_X(X) = 1$ (all mass on singleton tree X)
 - ▶ If $X \in N \setminus A$ is *not adapted* then X expands as in PCFG, i.e.,:

$$G_X = \sum_{X \rightarrow Y_1 \dots Y_m \in R_X} \theta_{X \rightarrow Y_1 \dots Y_m} \text{TD}_X(G_{Y_1}, \dots, G_{Y_m})$$

- ▶ If $X \in A$ is *adapted*, then PCFG distribution is adapted:

$$G_X \sim \text{DP}(\alpha, H_X)$$

$$H_X = \sum_{X \rightarrow Y_1 \dots Y_m \in R_X} \theta_{X \rightarrow Y_1 \dots Y_m} \text{TD}_X(G_{Y_1}, \dots, G_{Y_m})$$

- Other kinds of adaptor grammars use different adaptors
 - ▶ *Pitman-Yor adaptor grammars* use Pitman-Yor Processes as

Predictive distribution of DP adaptor grammars

- Predictive distribution: predict next tree t_{n+1} given previously generated trees $t = (t_1, \dots, t_n)$
- Predictive model “caches” adapted subtrees:
 - ▶ An *unadapted nonterminal* B expands using $B \rightarrow \beta$ with probability $\theta_{B \rightarrow \beta}$
 - ▶ Each adapted nonterminal B is associated with a DP that caches previously generated subtrees in \mathcal{T}_B
 - ▶ An *adapted nonterminal* B expands:
 - to a subtree $t' \in \mathcal{T}_B$ probability proportional to the number of times t' was previously generated
 - using $B \rightarrow \beta$ with probability proportional to $\alpha \theta_{B \rightarrow \beta}$

Adaptor grammars for word segmentation

- Input: phoneme sequences with *sentence boundaries* (Brent)
- Task: identify *word boundaries*, and hence words

y _Δ u _▲ w _Δ a _Δ n _Δ t _▲ t _Δ u _▲ s _Δ i _▲ D _Δ 6 _▲ b _Δ U _Δ

Words → Word

Words → Word Words

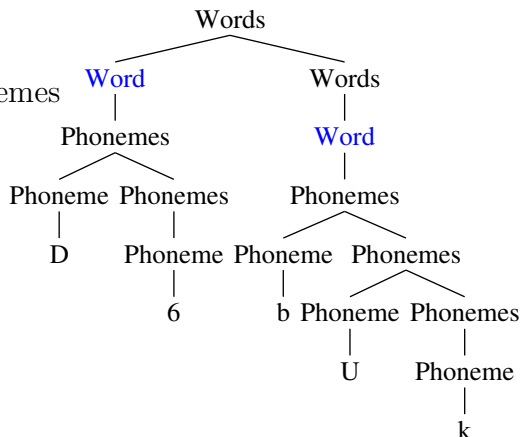
Word → Phonemes

Phonemes → Phoneme Phonemes

Phonemes → Phoneme

Phoneme → a | ... | z

- Adapted nonterminals (e.g., Word) highlighted and underlined



Outline

LDA topic models as PCFGs

Adaptor grammars

Finding topic-specific collocations

Learning the structure of proper nouns

Conclusion

Topic model with collocations

- Combines *PCFG topic model* and *segmentation adaptor grammar*

Sentence \rightarrow Doc_{*j*} $j \in 1, \dots, m$

Doc_{*j*} \rightarrow -*j* $j \in 1, \dots, m$

Doc_{*j*} \rightarrow Doc_{*j*} Topic_{*i*} $i \in 1, \dots, \ell;$

$j \in 1, \dots, m$

Topic_{*i*} \rightarrow Words

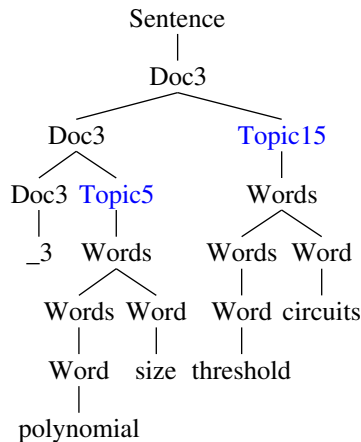
$i \in 1, \dots, \ell$

Words \rightarrow Word

Words \rightarrow Words Word

Word \rightarrow *w*

$w \in \mathcal{V}$



Finding topical collocations in NIPS abstracts

- Run topical collocation adaptor grammar on NIPS corpus
- Run with $\ell = 20$ topics (i.e., 20 distinct Topic_i nonterminals)
- Corpus is segmented by punctuation
 - ▶ terminal strings are fairly short
 - ⇒ inference is fairly efficient
- Used Pitman-Yor adaptors
 - ▶ sampled Pitman-Yor a and b parameters
 - ▶ flat and “vague Gamma” priors on Pitman-Yor a and b parameters

See: Griffiths et al (2007), Johnson and Goldwater (2009)

Sample output on NIPS corpus, 20 topics

- Multiword subtrees learned by adaptor grammar:

T_0 → gradient descent	T_1 → associative memory
T_0 → cost function	T_1 → standard deviation
T_0 → fixed point	T_1 → randomly chosen
T_0 → learning rates	T_1 → hamming distance
T_3 → membrane potential	T_10 → ocular dominance
T_3 → action potentials	T_10 → visual field
T_3 → visual system	T_10 → nervous system
T_3 → primary visual cortex	T_10 → action potential
- Sample skeletal parses:
 - _3 (T_5 polynomial size) (T_15 threshold circuits)
 - _4 (T_11 studied) (T_19 pattern recognition algorithms)
 - _4 (T_2 feedforward neural network) (T_1 implements)
 - _5 (T_11 single) (T_10 ocular dominance stripe) (T_12 low)
(T_3 ocularity) (T_12 drift rate)

Outline

LDA topic models as PCFGs

Adaptor grammars

Finding topic-specific collocations

Learning the structure of proper nouns

Conclusion

Learning the structure of proper nouns

- Grammars offer *structural* and *positional sensitivity* not captured in topic models: can we use this somehow?
- The Penn WSJ assigns flat structures to names and other base NPs
- Identifying structure within names can be useful
 - ▶ *Bill Clinton* and *Hillary Clinton* are unlikely to corefer because *Bill* and *Hillary* are both first names
 - ▶ *Secretary Clinton* and *Hillary Clinton* can corefer because *Secretary* is an honorific
- There are many different types of names (e.g., company names, person names)
- Some components of a name can be filled by multi-word sequences
 - ▶ In *Jean-Claude van Damme*, *van Damme* is the surname

An adaptor grammar for names

NP \rightarrow (A0) (A1) ... (A6) NP \rightarrow (B0) (B1) ... (B6)

A0 \rightarrow Word⁺

B0 \rightarrow Word⁺

...

...

A6 \rightarrow Word⁺

B6 \rightarrow Word⁺

NP \rightarrow Unordered⁺

Unordered \rightarrow Word⁺

- *Sample parses:*

(A0 barrett) (A3 smith)

(A0 albert) (A2 j.) (A3 smith) (A4 jr.)

(A0 robert) (A2 b.) (A3 van dover)

(B0 aim) (B1 prime rate) (B2 plus) (B5 fund) (B6 inc.)

(B0 balfour) (B1 maclaine) (B5 international) (B6 ltd.)

(B0 american express) (B1 information services) (B6 co)

(U abc) (U sports)

(U sports illustrated)

(U sports unlimited)

See: Elsner, Charniak and Johnson (2009)

Outline

LDA topic models as PCFGs

Adaptor grammars

Finding topic-specific collocations

Learning the structure of proper nouns

Conclusion

Conclusion

- LDA topic models can be expressed as Bayesian PCFGs
 - ▶ makes it easier to combine grammars and topic models
 - ▶ may help us to design new topic models that incorporate configurational sensitivity that is easy to express with grammars
- Adaptor grammars are a non-parametric extension of PCFGs which associate probabilities with entire subtrees
- Adaptor grammars can be used to express generalised topic models
 - ▶ learning topical collocations
 - ▶ learning the structure of names

Interested in Bayesian Inference and Language?

We're recruiting *PhD students* and *post-docs*.

Contact Mark.Johnson@mq.edu.au for more information.

