

Complexity of Inference in Latent Dirichlet Allocation

David Sontag, Daniel Roy
(NYU, Cambridge)

W66

Topic models are powerful tools for exploring large data sets and for making inferences about the content of documents

Documents



Topics

<u>politics</u> .0100
president .0095
obama .0090
washington .0085
religion .0060
...

<u>religion</u> .0500
hindu .0092
judiasm .0080
ethics .0075
buddhism .0016
...

<u>sports</u> .0105
baseball .0100
soccer .0055
basketball .0050
football .0045
...

$$\beta_t = \{ p(w | z = t) \}$$

Almost all uses of topic models (e.g., for unsupervised learning, information retrieval, classification) require **probabilistic inference**:

New document



Words w_1, \dots, w_N



What is this document about?

weather .50
finance .49
sports .01

Distribution of topics θ

Complexity of Inference in Latent Dirichlet Allocation

David Sontag, Daniel Roy
(NYU, Cambridge)

W66

Main Results

Maximize $p(z_{1:N} | w_{1:N})$

For any α

	# topics in MAP assignment	Complexity	Intuition
Most common setting →	Small	Easy	First choose topic sizes, then match words to topics
	Large	NP-hard	Reduction from set packing

Maximize $p(\theta | w_{1:N})$

	Dirichlet hyper-parameters	Complexity	Intuition
Most common setting →	$\alpha_t \geq 1$	Easy	Maximizing concave function
	$\alpha_t < 1$	NP-hard	Reduction from set cover

Sample from $p(\theta | w_{1:N})$

	Dirichlet hyper-parameters	Complexity	Intuition
	$\alpha_t \geq 1$	Easy	Log-concave distribution
	$\alpha_t \approx 0$	NP-hard	Reduction from set cover