

Optimization of ARTM regularization coefficients via stochastic variational inference

Michael Khalman, Dmitry Vetrov, Michael Figurnov

November 13, 2015

Outline

1. Topic modeling
 - 1.1 ARTM probabilistic model
2. Evidence maximization
 - 2.1 Variational inference
 - 2.2 Reparametrization trick for posterior
 - 2.3 Logistic-normal distribution
 - 2.4 How to optimize with normalizing constant of prior unknown?
3. Review

Topic modeling

Topic modeling an algorithmic tool that help us organize, search and understand vast amount of information.

Applications

- ▶ Exploratory search
- ▶ Clustering
- ▶ Creating annotations
- ▶ Recommendation systems
- ▶ Looking for similar documents in a huge collection
- ▶ ...

Probabilistic Topic Modeling

Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

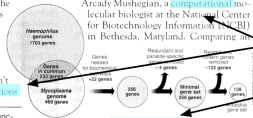
Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions** "are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson of the University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** numbers game. Some, particularly **more** and **more genes** are **sequenced**, **more** and **more genes** are **sequenced**, **more** and **more genes** are **sequenced**, **more** and **more genes** are **sequenced**.

Some, particularly **more** and **more genes** are **sequenced**, **more** and **more genes** are **sequenced**, **more** and **more genes** are **sequenced**, **more** and **more genes** are **sequenced**.

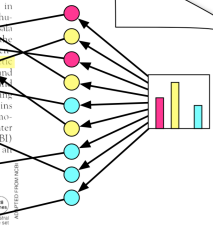


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments



Probabilistic Topic Modeling

Consider a set of documents D , dictionary W . Each document d consists of some words $w \in W$.

- ▶ There's a set of latent topics T . Each $t \in T$ is a distribution over W :

$$\varphi_{wt} = p(w|t)$$

- ▶ Each document $d \in D$ has specific distribution over T : $\theta_{td} = p(t|d)$.

- ▶ Bag of words assumption: order of words in d doesn't matter.

Words in d are i.i.d. from $p(w|d)$

- ▶ Conditional independence assumption: $p(w|d, t) = p(w|t)$.

- ▶ Law of total probability: $p(w|d) = \sum_t \underbrace{p(w|t)}_{\varphi_{wt}} \underbrace{p(t|d)}_{\theta_{td}}$

ARTM

- ▶ $N = (n_{dw})_{W \times D}$ — matrix of counts of words in documents
- ▶ $F = (n_{dw}/n_d)_{W \times D}$ — normalized matrix of counts; $n_d \equiv \sum_w n_{dw}$.
- ▶ $(\varphi_{wt})_{W \times T} \equiv \Phi \in \mathbb{R}^{W \times T}$ — topic-specific distribution over words.
- ▶ $(\theta_{td})_{T \times D} \equiv \Theta \in \mathbb{R}^{T \times D}$ — document-specific distribution over topics.

$$p(w|d) = \sum_t \varphi_{wt} \theta_{td}$$

Maximization of log-likelihood:

$$p(N|\Phi, \Theta) = \prod_d \prod_w \left(\sum_t \varphi_{wt} \theta_{td} \right)^{n_{dw}} \rightarrow \max_{\Phi, \Theta}$$

ARTM

- ▶ $N = (n_{dw})_{W \times D}$ — matrix of counts of words in documents
- ▶ $F = (n_{dw}/n_d)_{W \times D}$ — normalized matrix of counts; $n_d \equiv \sum_w n_{dw}$.
- ▶ $(\varphi_{wt})_{W \times T} \equiv \Phi \in \mathbb{R}^{W \times T}$ — topic-specific distribution over words.
- ▶ $(\theta_{td})_{T \times D} \equiv \Theta \in \mathbb{R}^{T \times D}$ — document-specific distribution over topics.

$$p(w|d) = \sum_t \varphi_{wt} \theta_{td}$$

Maximization of log-likelihood:

$$p(N|\Phi, \Theta) = \prod_d \prod_w (\sum_t \varphi_{wt} \theta_{td})^{n_{dw}} \rightarrow \max_{\Phi, \Theta}$$

$$\log p(N|\Phi, \Theta) = \sum_d \sum_w n_{dw} \log \sum_t \varphi_{wt} \theta_{td} \rightarrow \max_{\Phi, \Theta}$$

ARTM

- ▶ $N = (n_{dw})_{W \times D}$ — matrix of counts of words in documents
- ▶ $F = (n_{dw}/n_d)_{W \times D}$ — normalized matrix of counts; $n_d \equiv \sum_w n_{dw}$.
- ▶ $(\varphi_{wt})_{W \times T} \equiv \Phi \in \mathbb{R}^{W \times T}$ — topic-specific distribution over words.
- ▶ $(\theta_{td})_{T \times D} \equiv \Theta \in \mathbb{R}^{T \times D}$ — document-specific distribution over topics.

$$p(w|d) = \sum_t \varphi_{wt} \theta_{td}$$

Maximization of log-likelihood:

$$p(N|\Phi, \Theta) = \prod_d \prod_w (\sum_t \varphi_{wt} \theta_{td})^{n_{dw}} \rightarrow \max_{\Phi, \Theta}$$

$$\log p(N|\Phi, \Theta) = \sum_d \sum_w n_{dw} \log \sum_t \varphi_{wt} \theta_{td} \rightarrow \max_{\Phi, \Theta}$$

$$\log p(N|\Phi, \Theta) = - \sum_d n_d \text{KL}(F_d || (\Phi\Theta)_d) \rightarrow \max_{\Phi, \Theta}$$

Non-negative matrix factorization: $F \approx \Phi\Theta$:

$$\mathcal{L}(\Phi, \Theta) \equiv \log p(N|\Phi, \Theta) \rightarrow \max_{\Phi, \Theta}$$

$$\mathcal{L}(\Phi, \Theta) = - \sum_d n_d \text{KL}(F_d || (\Phi\Theta)_d)$$

Ill-posed problem: $F = \Phi\Theta = (\Phi S)(S^{-1})\Theta$

Prior on Φ, Θ :

$$\hat{p}(\Phi, \Theta|\tau) \propto \exp\left(\sum_i \tau_i R_i(\Phi, \Theta)\right)$$

$$p(N, \Phi, \Theta|\tau) = p(N|\Phi, \Theta)\hat{p}(\Phi, \Theta|\tau) \rightarrow \max_{\Phi, \Theta}$$

Regularized non-negative matrix factorization: $F \approx \Phi\Theta$:

$$\underbrace{\mathcal{L}(\Phi, \Theta)}_{\text{log-likelihood}} + \underbrace{\mathbf{R}(\Phi, \Theta)}_{\text{regularization}} \rightarrow \max_{\Phi, \Theta}$$

$$\mathcal{L}(\Phi, \Theta) = - \sum_d n_d \text{KL}(F_d || (\Phi\Theta)_d)$$

$$\mathbf{R}(\Phi, \Theta) = \sum_i \tau_i R_i(\Phi, \Theta)$$

$R_i(\Phi, \Theta)$ — any differentiable regularization function

ARTM — regularizers

We can use $R(\Phi, \Theta)$ to:

- ▶ Make model better (easy to interpret, coherent topics)

- ▶ Smoothing/sparsing (LDA)

$$R = \sum_{t \in T} \sum_{w \in W} \beta_{wt} \log \varphi_{wt} + \sum_{d \in D} \sum_{t \in T} \alpha_{td} \log \theta_{td}$$

- ▶ Decorellation: $R = -\tau \sum_{t \in T} \sum_{s \in T \setminus t} \sum_{w \in W} \varphi_{ws} \varphi_{wt}$

- ▶ Add soft restrictions to the model

- ▶ Semi-supervised regularizer

$$R = \sum_{t \in T_0} \sum_{w \in W_t} \beta_{wt} \log \varphi_{wt} + \sum_{d \in D_0} \sum_{t \in T_d} \alpha_{td} \log \theta_{td}$$

- ▶ Add new information to the model

- ▶ Author-topic models

- ▶ Make the model specified for solving specific problem

- ▶ Topic models with time
 - ▶ Topic models for classification

- ▶ ...

ARTM — advantages

- ▶ **Really** easy to infer new models
- ▶ There is an iterative algorithm for any differentiable $R(\Phi, \Theta)$ based on fixed-point iteration method for stationarity conditions
- ▶ You can combine any amount of regularizers in one model

ARTM — problems

How to fit hyperparameters (regularization coefficients τ)?

Current approach: heuristic, manual fitting.

How to fit them **automatically** is an **open problem**

- ▶ Too many parameters for grid-search
- ▶ We don't know what to optimize

Maximum-evidence approach

$$p(N|\tau) = \int p(N|\Phi, \Theta)p(\Phi, \Theta|\tau)d\Phi d\Theta \longrightarrow \max_{\tau}$$

Or, if we denote (Φ, Θ) as \mathbf{z} , and N as \mathbf{x}

$$p(\mathbf{x}|\tau) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z}|\tau)d\mathbf{z} = \mathbb{E}_{\mathbf{z}|\tau}p(\mathbf{x}|\mathbf{z}) \longrightarrow \max_{\tau}$$

We'll use variational inference to optimize it.

Variational inference

Recall our probabilistic model:

$$p(\mathbf{x}|\mathbf{z}) \equiv p(N|\Phi, \Theta) = \prod_d \prod_w \left(\sum_t \varphi_{wt} \theta_{td} \right)^{n_{dw}}$$
$$p(\mathbf{z}|\boldsymbol{\tau}) \equiv p(\Phi, \Theta|\boldsymbol{\tau}) \propto \exp\left(\sum_i \tau_i R_i(\Phi, \Theta) \right)$$

We introduce some family of distributions $q(\mathbf{z}|\boldsymbol{\lambda}) \approx p(\mathbf{z}|\mathbf{x})$

$$\log p(\mathbf{x}|\boldsymbol{\tau}) \geq \mathcal{F}(\boldsymbol{\lambda}, \boldsymbol{\tau}) \equiv \int q(\mathbf{z}|\boldsymbol{\lambda}) \log \frac{p(\mathbf{x}, \mathbf{z}|\boldsymbol{\tau})}{q(\mathbf{z}|\boldsymbol{\lambda})} d\mathbf{z} =$$
$$= \mathbb{E}_{q_{\boldsymbol{\lambda}}} \log p(\mathbf{x}|\mathbf{z}) - \mathbb{E}_{q_{\boldsymbol{\lambda}}} \log \frac{q(\mathbf{z}|\boldsymbol{\lambda})}{\hat{p}(\mathbf{z}|\boldsymbol{\tau})} - \log Z(\boldsymbol{\tau}) \rightarrow \max_{\boldsymbol{\tau}, \boldsymbol{\lambda}}$$

$\mathcal{F}(\boldsymbol{\lambda}, \boldsymbol{\tau})$ — ELBO (evidence lower-bound)

Variational inference can be used to optimize evidence over hyperparameters of the model.

Reparametrization trick

Stochastic optimization problem:

$$\begin{aligned} \mathbf{z} &\sim p(\mathbf{z}) \\ \mathbb{E}_{\mathbf{z}} f(\mathbf{z}, \boldsymbol{\lambda}) &\rightarrow \max_{\boldsymbol{\lambda}} \\ \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{\mathbf{z}} f(\mathbf{z}, \boldsymbol{\lambda}) &=? \end{aligned}$$

$$\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{\mathbf{z}} f(\mathbf{z}, \boldsymbol{\lambda}) = \mathbb{E}_{\mathbf{z}} \nabla_{\boldsymbol{\lambda}} f(\mathbf{z}, \boldsymbol{\lambda}) \approx \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\lambda}} f(\mathbf{z}_i, \boldsymbol{\lambda})$$

We can use stochastic gradient descent

Reparametrization trick

What if probability $p(\mathbf{z})$ depends on λ ?

$$\mathbf{z} \sim p(\mathbf{z}|\lambda) \equiv p_\lambda(\mathbf{z})$$

$$\mathbb{E}_{p_\lambda} f(\mathbf{z}) \rightarrow \max_{\lambda}$$

$$\nabla_{\lambda} \mathbb{E}_{\mathbf{z}} f(\mathbf{z}) \text{ — ?}$$

Reparametrization trick

Let's assume that if $\mathbf{z} \sim q(\mathbf{z}|\boldsymbol{\lambda}) \equiv q_{\boldsymbol{\lambda}}$,
 $\mathbf{z} = g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda})$, where $\boldsymbol{\varepsilon} \sim q_0$

$$\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q_{\boldsymbol{\lambda}}} f(\mathbf{z}) = \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} f(g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda})) \approx \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\lambda}} f(g(\boldsymbol{\varepsilon}_i, \boldsymbol{\lambda}))$$

Example:

$$\begin{aligned} \mathbf{z} &\sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}); & \boldsymbol{\lambda} &\equiv (\boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda}) &= \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \boldsymbol{\varepsilon}; & \boldsymbol{\varepsilon} &\sim \mathcal{N}(0, \mathbf{I}) \\ \nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q_{\boldsymbol{\lambda}}} f(\mathbf{z}) &= \mathbb{E}_{q_0} f(\boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \boldsymbol{\varepsilon}) \end{aligned}$$

Estimation has low variance.

Gradients of ELBO

Reparametrization trick for $q(\mathbf{z}|\boldsymbol{\lambda})$:

$$\mathbf{z} = g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda}); \quad \boldsymbol{\varepsilon} \sim q_0$$

$$\mathcal{F}(\boldsymbol{\lambda}, \boldsymbol{\tau}) = \mathbb{E}_{q_{\boldsymbol{\lambda}}} \log p(\mathbf{x}|\mathbf{z}) - \mathbb{E}_{q_{\boldsymbol{\lambda}}} \log \frac{q(\mathbf{z}|\boldsymbol{\lambda})}{\hat{p}(\mathbf{z}|\boldsymbol{\tau})} - \log Z(\boldsymbol{\tau})$$

$$\begin{aligned} \frac{\partial \mathcal{F}}{\partial \boldsymbol{\lambda}} &= \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \log p(\mathbf{x}|g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda})) - \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \log q(g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda})|\boldsymbol{\lambda}) \\ &\quad + \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \log \hat{p}(g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda})|\boldsymbol{\tau}) \end{aligned} \quad (1)$$

$$\frac{\partial \mathcal{F}}{\partial \boldsymbol{\tau}} = \mathbb{E}_{q_0} \nabla_{\boldsymbol{\tau}} \log \hat{p}(g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda})|\boldsymbol{\tau}) - \mathbb{E}_{q_0} \nabla_{\boldsymbol{\tau}} \log Z(\boldsymbol{\tau})$$

Algorithm

Algorithm 1 Sketch

Data: Flow of data \mathbf{x} ; learning rate κ

Result: λ , τ

initialize λ and τ

repeat

$\mathbf{x} \leftarrow \{\text{Get mini-batch}\}$

Sample $\varepsilon \leftarrow q_0$

Compute stochastic estimations of $\frac{\partial \mathcal{F}}{\partial \tau}$ and $\frac{\partial \mathcal{F}}{\partial \lambda}$ according to (1)

$\tau = \tau - \kappa \frac{\partial \mathcal{F}}{\partial \tau}$ // or any other stochastic-gradient-kind

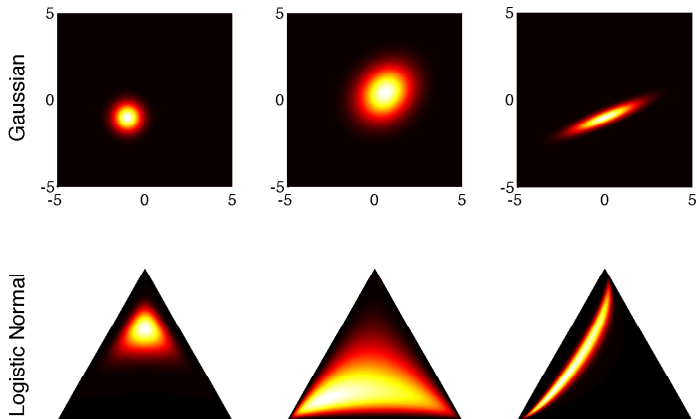
$\lambda = \lambda - \kappa \frac{\partial \mathcal{F}}{\partial \lambda}$ // optimization algorithm, e.g. AdaGrad

until convergence

Logistic-normal distribution

$$\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}); \quad \mathbf{z} = \mathcal{P}(\boldsymbol{\varepsilon}) \implies \mathbf{z} \sim \mathcal{P}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})),$$

where $\mathcal{P}(\mathbf{x}) = \mathbf{z} \implies z_i = \frac{\exp(x_i)}{\sum_j \exp(x_j)}$ — softmax function



Logistic-normal distribution — Properties

- ▶ Distribution over simplex $\mathcal{S}_x = \{x : x_i \geq 0; \sum_i x_i = 1\}$;
- ▶ Similar to Dirichlet, but **never** exactly the same;
- ▶ Although can be approximated with Dirichlet with large α
- ▶ No analytical solution for mode/mean/variance.
- ▶ We can write down component-wise median: $\mathcal{P}(\boldsymbol{\mu})$
- ▶ Easy to sample from

Reparametrization of q

$$\mathbf{z} = g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda}) = \mathcal{P}(\text{diag}(\boldsymbol{\sigma})\boldsymbol{\varepsilon} + \boldsymbol{\mu}); \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$$
$$\mathbf{z} \sim \mathcal{P}(\mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma})))$$

where $\boldsymbol{\lambda}$ denotes $(\boldsymbol{\mu}, \boldsymbol{\sigma})$.

Individual vector of parameters $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}$ for each column of Φ and Θ

$$\log q(g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda})) = \log q_0(\boldsymbol{\varepsilon}) + \log \left| \frac{\partial g}{\partial \boldsymbol{\varepsilon}} \right| = \log q_0(\boldsymbol{\varepsilon}) + \sum_k (\log z_k + \log \sigma_k)$$

Recall gradients of ELBO

$$\mathbf{z} = g(\boldsymbol{\varepsilon}, \boldsymbol{\lambda}); \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I})$$

$$\frac{\partial \mathcal{F}}{\partial \boldsymbol{\lambda}} = \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \log p(\mathbf{x}|\mathbf{z}) - \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \log q(\mathbf{z}|\boldsymbol{\lambda}) + \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \log \hat{p}(\mathbf{z}|\boldsymbol{\tau}) \quad (2)$$

$$\frac{\partial \mathcal{F}}{\partial \boldsymbol{\tau}} = \mathbb{E}_{q_0} \nabla_{\boldsymbol{\tau}} \log \hat{p}(\mathbf{z}|\boldsymbol{\tau}) - \mathbb{E}_{q_0} \nabla_{\boldsymbol{\tau}} \log Z(\boldsymbol{\tau})$$

- ▶ $\log p(\mathbf{x}|\mathbf{z})$ — log-likelihood of our model
- ▶ $\log \hat{p}(\mathbf{z}|\boldsymbol{\tau}) = \sum_i \tau_i R_i(\mathbf{z})$ — regularization term.
- ▶ $\mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \log q(\mathbf{z}|\boldsymbol{\lambda}) = \mathbb{E}_{q_0} \nabla_{\boldsymbol{\lambda}} \sum_k (\log z_k + \log \sigma_k) + \text{const}$
- ▶ $\nabla_{\boldsymbol{\lambda}} \log \hat{p}(\mathbf{z}|\boldsymbol{\tau}) = \sum_i \tau_i \nabla_{\boldsymbol{\lambda}} R_i(\mathbf{z})$
- ▶ $\mathbb{E}_{q_0} \nabla_{\boldsymbol{\tau}} \log \hat{p}(\mathbf{z}|\boldsymbol{\tau}) = \mathbb{E}_{q_0} \mathbf{R}(\mathbf{z})$
- ▶ $\nabla_{\boldsymbol{\tau}} \log Z(\boldsymbol{\tau}) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\tau})} \mathbf{R}(\mathbf{z})$

How to optimize with $Z(\boldsymbol{\tau})$ unknown?

$$\nabla_{\boldsymbol{\tau}} \log Z(\boldsymbol{\tau}) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\tau})} R(\mathbf{z})$$

We need samples from $p(\mathbf{z}|\boldsymbol{\tau})$

- ▶ Contrastive divergence

- ▶ After sampling $\mathbf{z}_i \sim q_{\lambda}$ do **a few** iterations of MCMC (starting from \mathbf{z}_i) to obtain samples from $p(\mathbf{z}|\boldsymbol{\tau})$.

- ▶ Importance sampling

- ▶ $Z(\boldsymbol{\tau}) = \mathbb{E}_{q_0} \frac{\hat{p}(\mathbf{z}|\boldsymbol{\tau})}{q_0(\mathbf{z})}$
- ▶ $\mathbb{E}_{p(\mathbf{z}|\boldsymbol{\tau})} R(\mathbf{z}) = \frac{1}{Z(\boldsymbol{\tau})} \mathbb{E}_{q_0} \frac{R(\mathbf{z}) \hat{p}(\mathbf{z}|\boldsymbol{\tau})}{q_0(\mathbf{z})}$
- ▶ We use the same samples $\mathbf{z}_i \sim q_{\lambda}$ to estimate $Z(\boldsymbol{\tau})$ and $\mathbb{E}_{q_0} \frac{R(\mathbf{z}) \hat{p}(\mathbf{z}|\boldsymbol{\tau})}{q_0(\mathbf{z})}$

Review

Our algorithm

- ▶ Our algorithm is scalable on data size
- ▶ Time of each iteration on one batch of D docs with N total words is $O(N + TD + TW)$ (no DW term)
- ▶ ...or $O(N + T^2D + TW)$ with decorrelation
- ▶ Returns both fitted regularization coefficients τ and matrices (Φ, Θ) whole distribution over (Φ, Θ) .
- ▶ Coefficients τ and distributions (Φ, Θ) are fitted **simultaneously** as long as algorithms iterates over data (over batches)

Questions still not answered

- ▶ What to do with that distribution (what is the final answer)
 - ▶ Sample from it (probably multiple times)
 - ▶ Return median $\mathcal{P}(\mu)$
 - ▶ Compute mode $\operatorname{argmax}_{\Phi, \Theta} q(\Phi, \Theta | \lambda)$
- ▶ Is maximum-evidence approach really going to work in this scenario?
- ▶ How to organize online-learning process carefully?

Thank you for your attention!