

Variational Inference: Foundations and Modern Methods

David Blei, Rajesh Ranganath, Shakir Mohamed

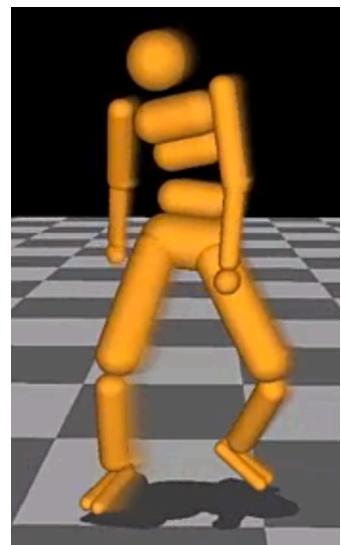
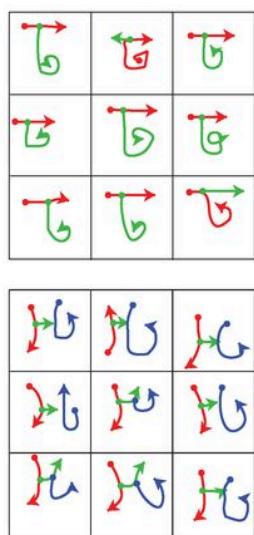
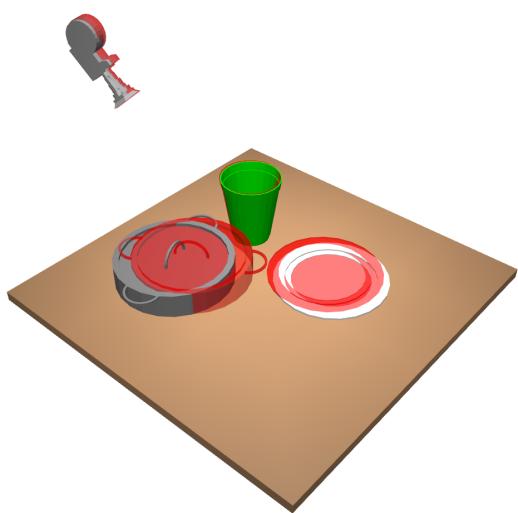
NIPS 2016 Tutorial · December 5, 2016





Topics found in 1.8M articles from the New York Times

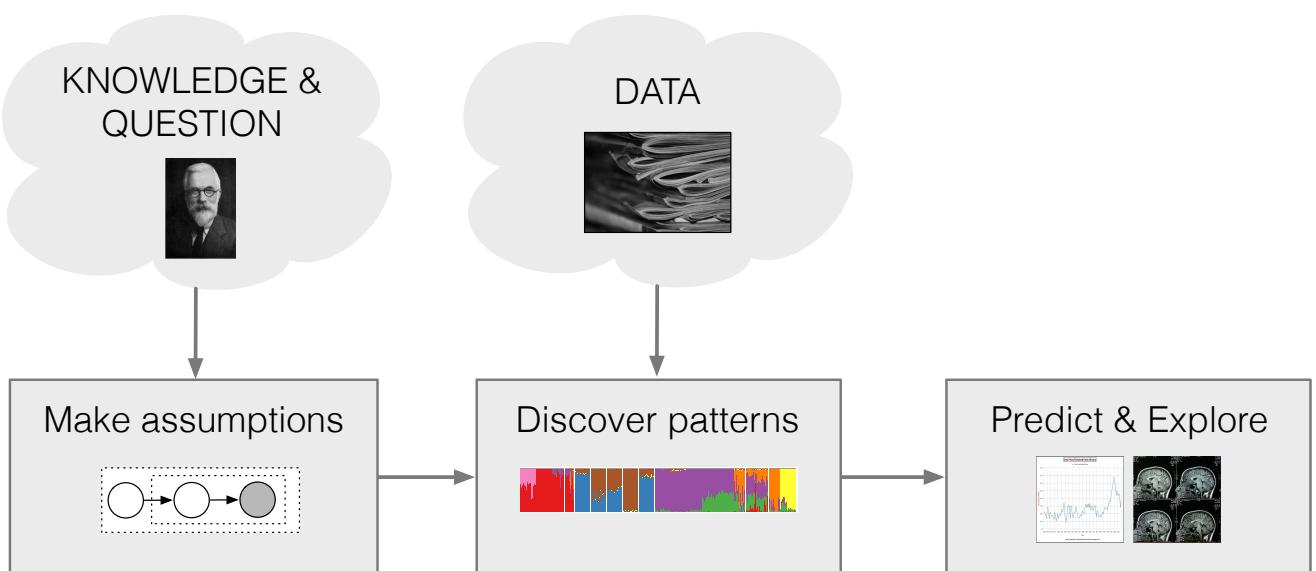
[Hoffman, Blei, Wang, Paisley, JMLR 2013]



Scenes, concepts and control.

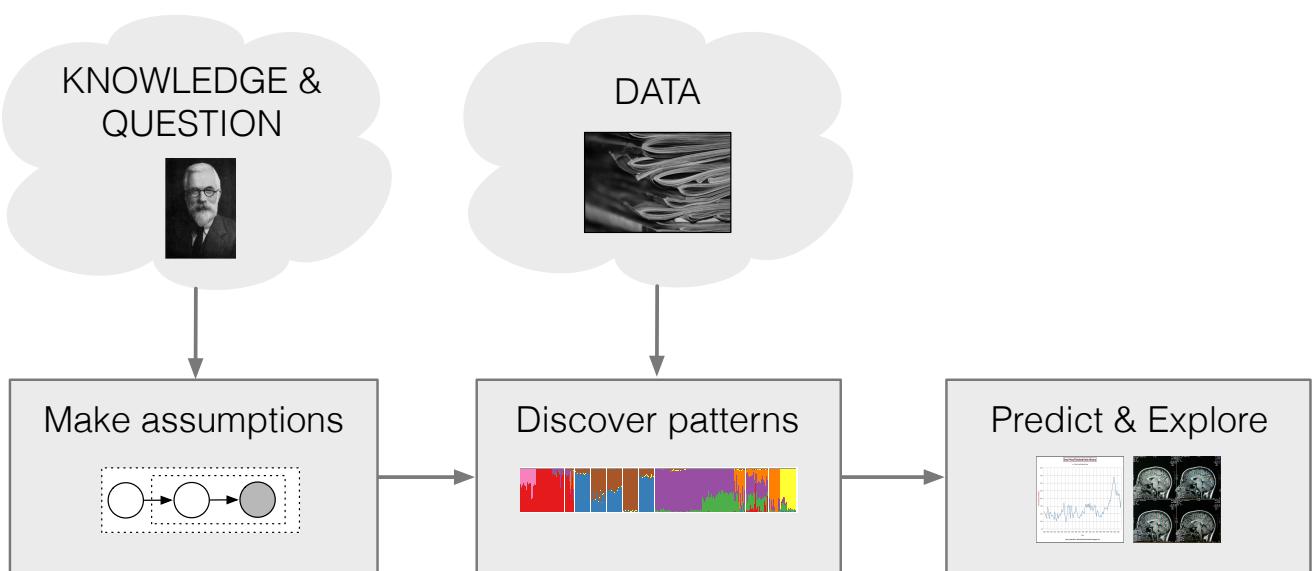
[Eslami et al., 2016, Lake et al. 2015]

The probabilistic pipeline

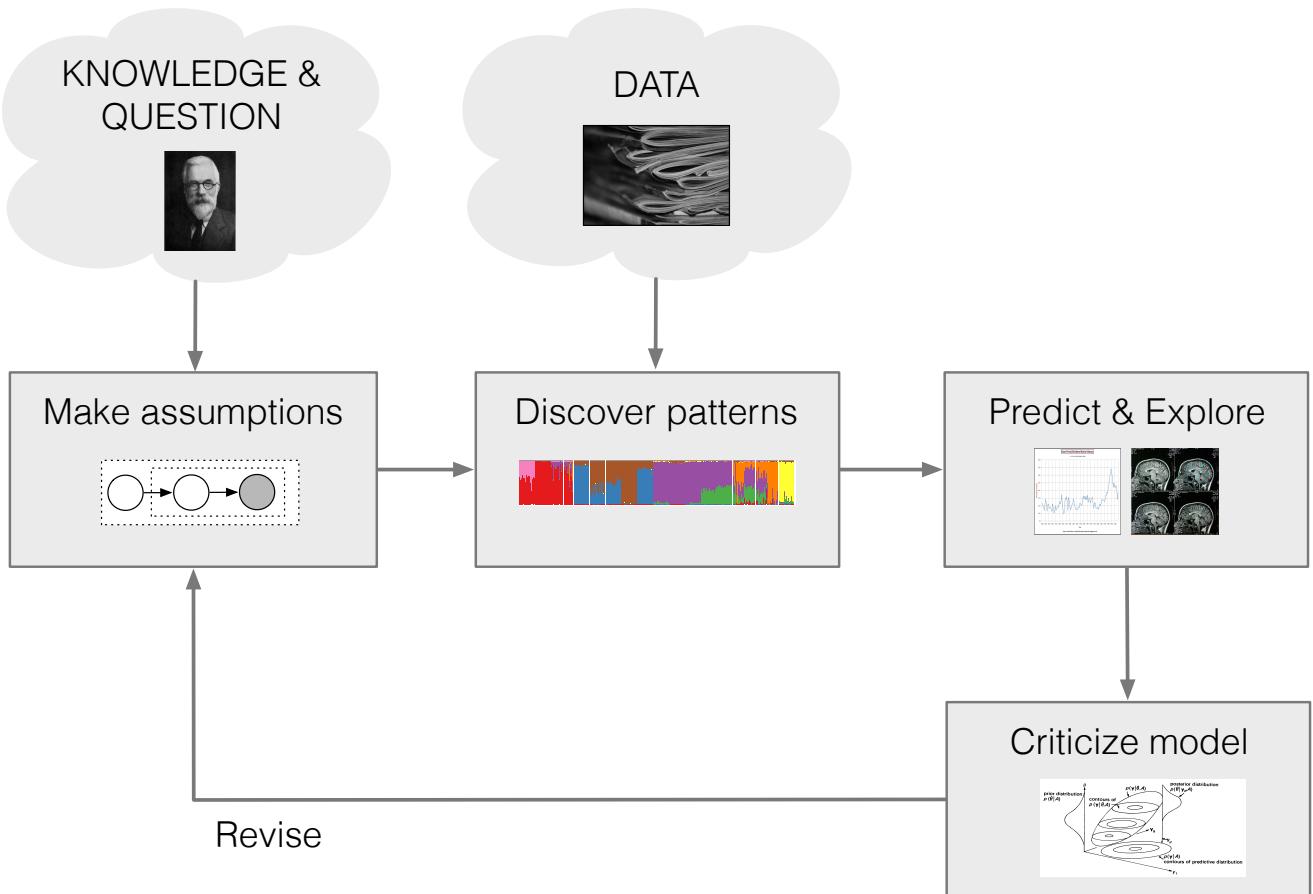


- Customized data analysis is important to many fields.
- Pipeline separates **assumptions, computation, application**
- Eases collaborative solutions to statistics problems

The probabilistic pipeline



- **Inference** is the key algorithmic problem.
- Answers the question: What does this model say about this data?
- Our goal: **General** and **scalable** approaches to inference



[Box, 1980; Rubin, 1984; Gelman et al., 1996; Blei, 2014]

PART I

Main ideas and historical context

Probabilistic Machine Learning

- A probabilistic model is a joint distribution of hidden variables \mathbf{z} and observed variables \mathbf{x} ,

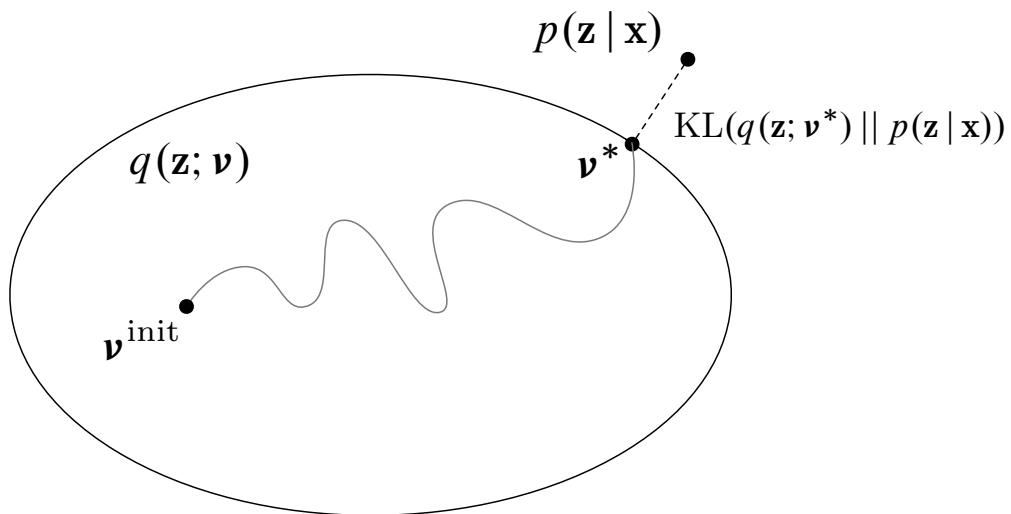
$$p(\mathbf{z}, \mathbf{x}).$$

- Inference about the unknowns is through the **posterior**, the conditional distribution of the hidden variables given the observations

$$p(\mathbf{z} | \mathbf{x}) = \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})}.$$

- For most interesting models, the denominator is not tractable. We appeal to **approximate posterior inference**.

Variational Inference



- VI turns **inference into optimization**.
- Posit a **variational family** of distributions over the latent variables,

$$q(\mathbf{z}; \boldsymbol{\nu})$$

- Fit the **variational parameters** $\boldsymbol{\nu}$ to be close (in KL) to the exact posterior.
(There are alternative divergences, which connect to algorithms like EP, BP, and others.)

Variational Inference: Foundations and Modern Methods

Part II: Mean-field VI and stochastic VI

Jordan+, *Introduction to Variational Methods for Graphical Models*, 1999

Ghahramani and Beal, *Propagation Algorithms for Variational Bayesian Learning*, 2001

Hoffman+, *Stochastic Variational Inference*, 2013

Part III: Stochastic gradients of the ELBO

Kingma and Welling, *Auto-Encoding Variational Bayes*, 2014

Ranganath+, *Black Box Variational Inference*, 2014

Rezende+, *Stochastic Backpropagation and Approximate Inference in Deep Generative Models*, 2014

Part IV: Beyond the mean field

Agakov and Barber, *An Auxiliary Variational Method*, 2004

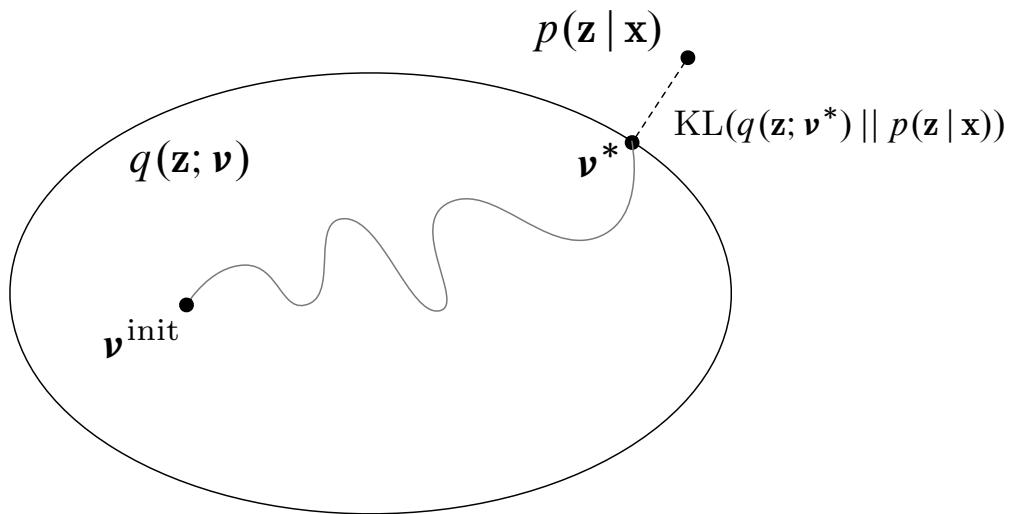
Gregor+, *DRAW: A recurrent neural network for image generation*, 2015

Rezende+, *Variational Inference with Normalizing Flows*, 2015

Ranganath+, *Hierarchical Variational Models*, 2015

Maaløe+, *Auxiliary Deep Generative Models*, 2016

Variational Inference: Foundations and Modern Methods



VI approximates difficult quantities from complex models.

With **stochastic optimization** we can

- scale up VI to massive data
- enable VI on a wide class of difficult models
- enable VI with elaborate and flexible families of approximations

PART II

Mean-field variational inference and stochastic variational inference

Motivation: Topic Modeling



Topic models use posterior inference to discover the hidden thematic structure in a large collection of documents.

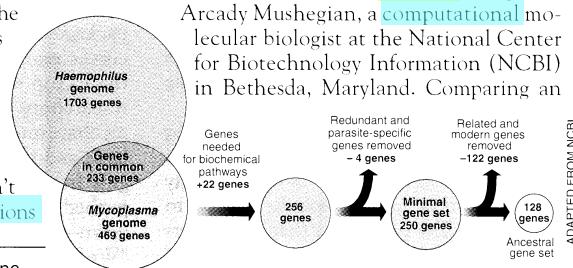
Example: Latent Dirichlet Allocation (LDA)

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 Uppsala 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, particularly as more and more genomes are completely mapped and sequenced. “It may be a way of organizing any newly sequenced genome,” explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an

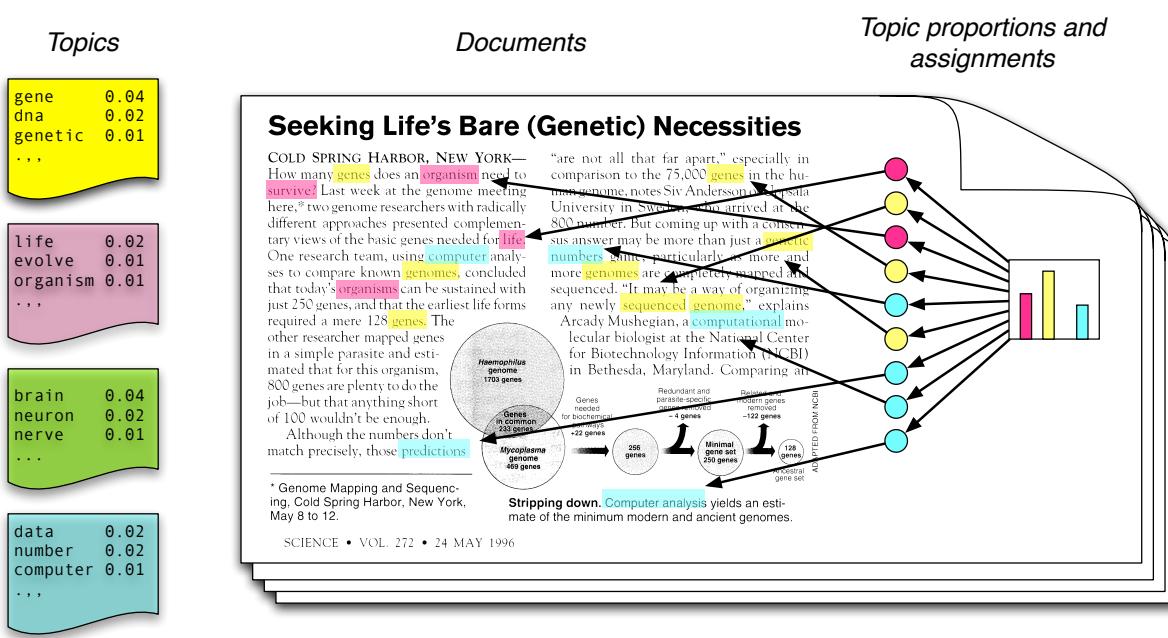


Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

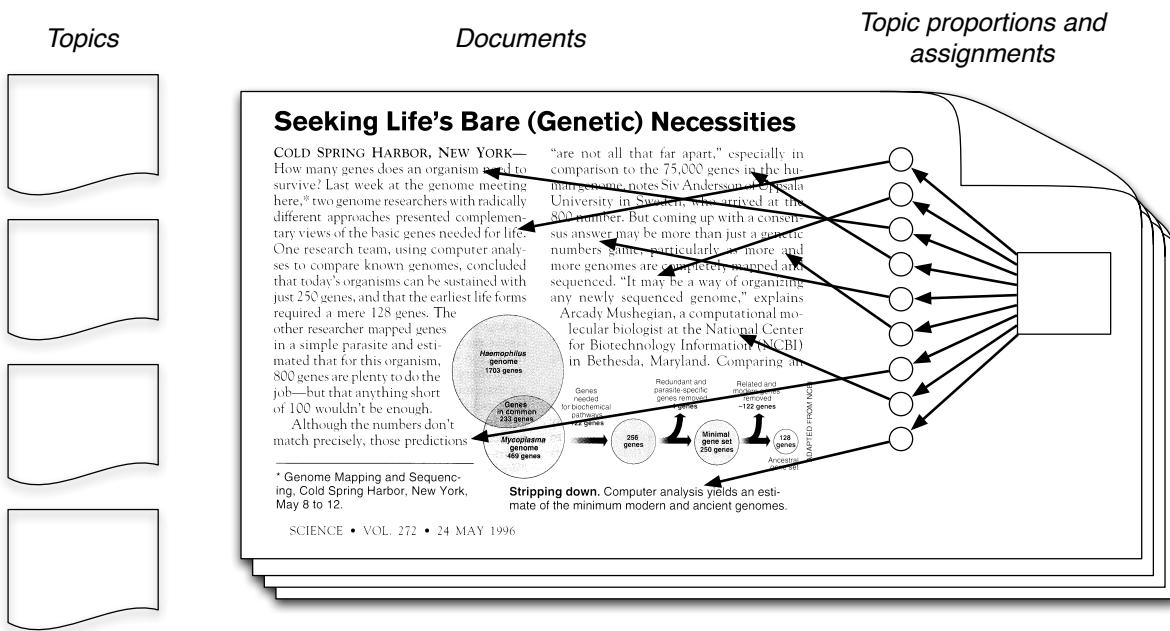
Documents exhibit multiple topics.

Example: Latent Dirichlet Allocation (LDA)



- Each **topic** is a distribution over words
- Each **document** is a mixture of corpus-wide topics
- Each **word** is drawn from one of those topics

Example: Latent Dirichlet Allocation (LDA)

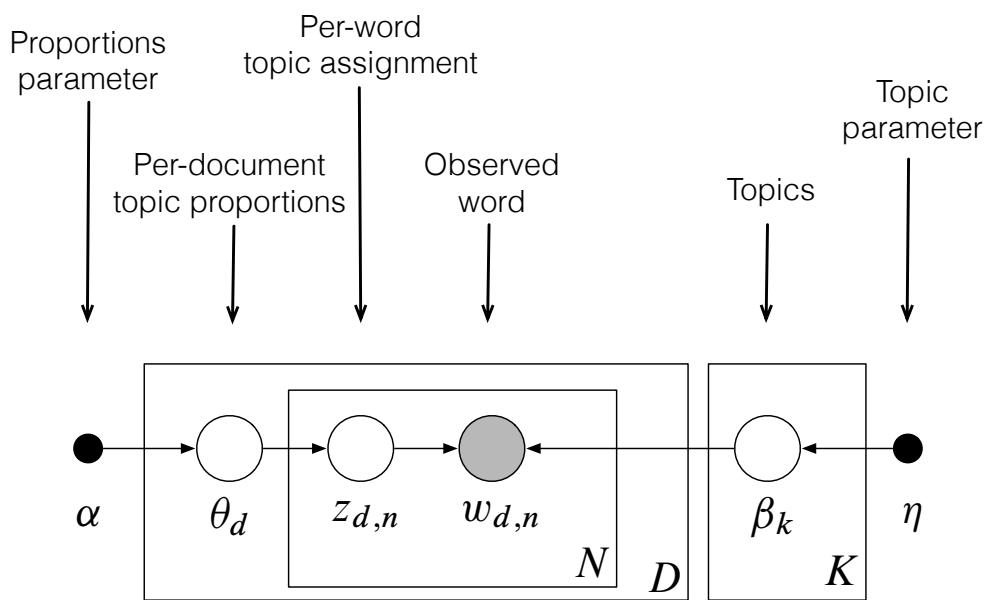


- But we only observe the documents; everything else is hidden.
- So we want to calculate the posterior

$$p(\text{topics, proportions, assignments} \mid \text{documents})$$

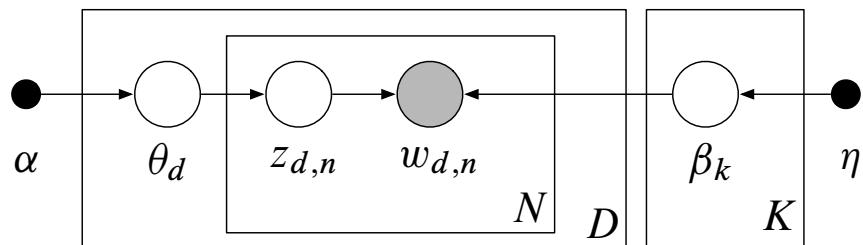
(Note: millions of documents; billions of latent variables)

LDA as a Graphical Model



- Encodes **assumptions** about data with a factorization of the joint
- Connects assumptions to **algorithms** for computing with data
- Defines the **posterior** (through the joint)

Posterior Inference



- The posterior of the latent variables given the documents is

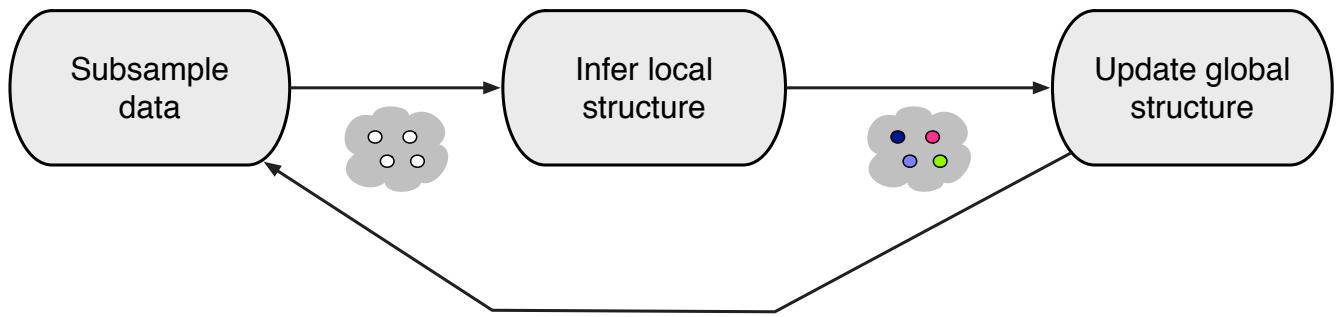
$$p(\beta, \theta, z | w) = \frac{p(\beta, \theta, z, w)}{\int_{\beta} \int_{\theta} \sum_z p(\beta, \theta, z, w)}.$$

- We can't compute the denominator, the marginal $p(w)$.
- We use approximate inference.



Topics found in 1.8M articles from the New York Times

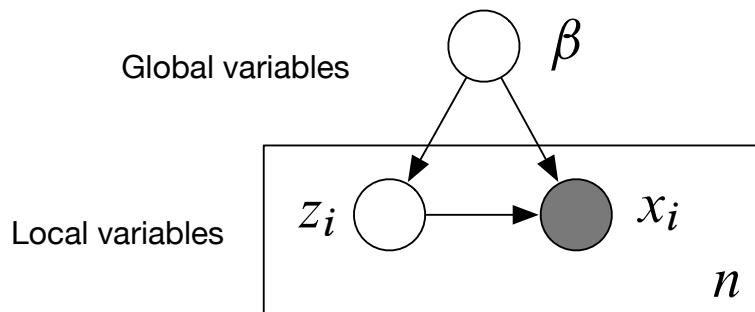
Mean-field VI and Stochastic VI



Road map:

- Define the generic class of conditionally conjugate models
- Derive classical mean-field VI
- Derive stochastic VI, which scales to massive data

A Generic Class of Models

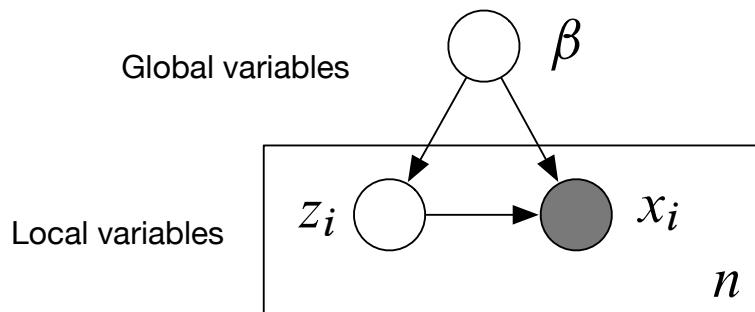


$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

- The observations are $\mathbf{x} = x_{1:n}$.
- The **local** variables are $\mathbf{z} = z_{1:n}$.
- The **global** variables are β .
- The i th data point x_i only depends on z_i and β .

Compute $p(\beta, \mathbf{z} | \mathbf{x})$.

A Generic Class of Models

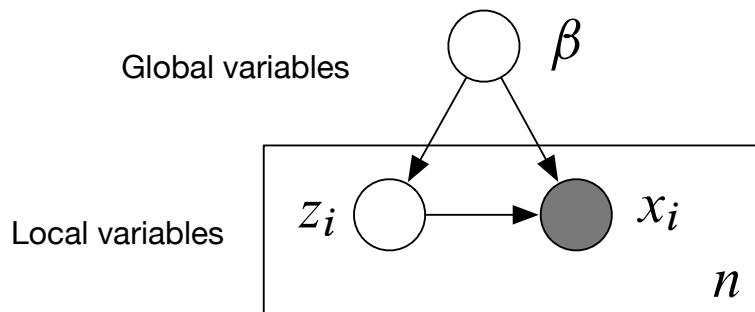


$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

- A **complete conditional** is the conditional of a latent variable given the observations and other latent variables.
- Assume each complete conditional is in the exponential family,

$$\begin{aligned} p(z_i | \beta, x_i) &= h(z_i) \exp\{\eta_\ell(\beta, x_i)^\top z_i - a(\eta_\ell(\beta, x_i))\} \\ p(\beta | \mathbf{z}, \mathbf{x}) &= h(\beta) \exp\{\eta_g(\mathbf{z}, \mathbf{x})^\top \beta - a(\eta_g(\mathbf{z}, \mathbf{x}))\}. \end{aligned}$$

A Generic Class of Models



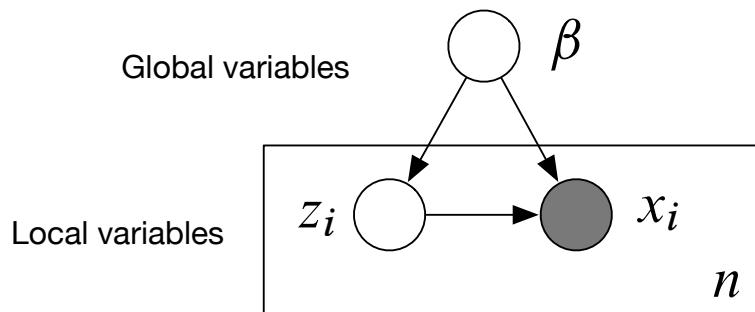
$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

- A **complete conditional** is the conditional of a latent variable given the observations and other latent variable.
- The global parameter comes from conjugacy [Bernardo and Smith, 1994]

$$\eta_g(\mathbf{z}, \mathbf{x}) = \alpha + \sum_{i=1}^n t(z_i, x_i),$$

where α is a hyperparameter and $t(\cdot)$ are sufficient statistics for $[z_i, x_i]$.

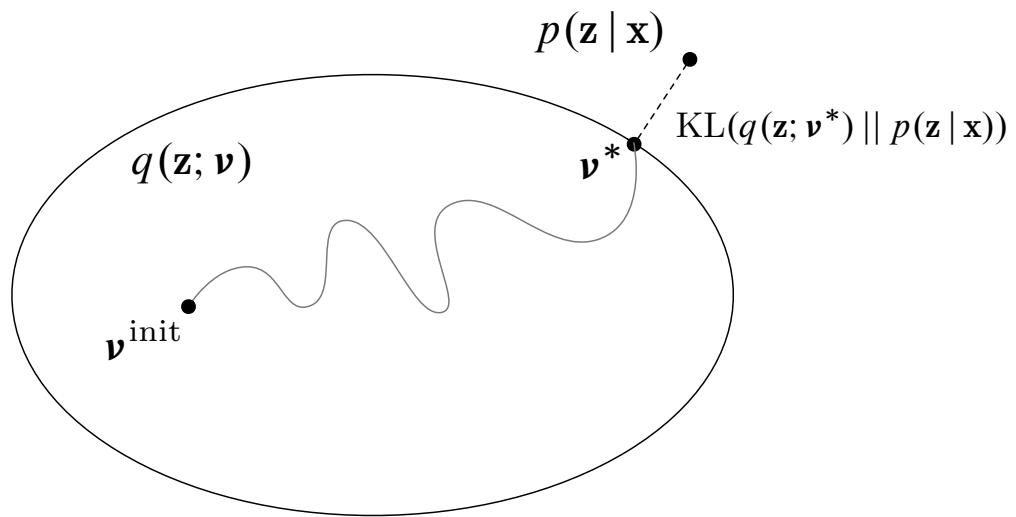
A Generic Class of Models



$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

- Bayesian mixture models
- Time series models
(HMMs, linear dynamic systems)
- Factorial models
- Matrix factorization
(factor analysis, PCA, CCA)
- Dirichlet process mixtures, HDPs
- Multilevel regression
(linear, probit, Poisson)
- Stochastic block models
- Mixed-membership models
(LDA and some variants)

Variational Inference



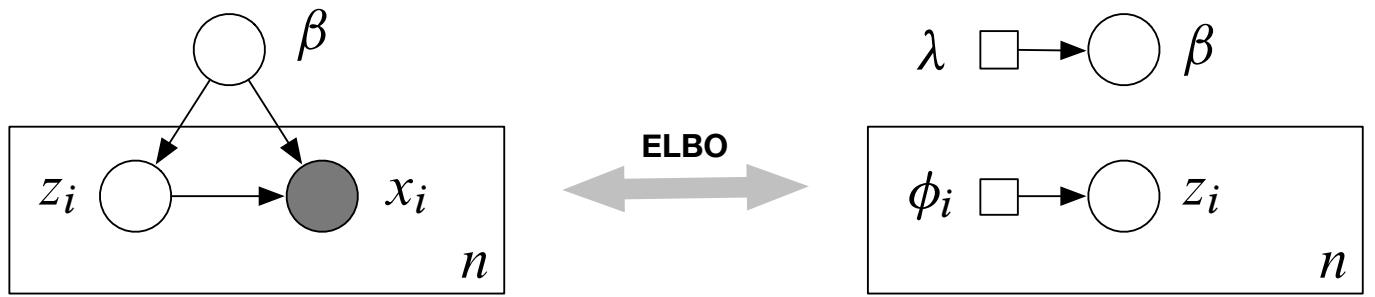
Minimize KL between $q(\beta, \mathbf{z}; \boldsymbol{\nu})$ and the posterior $p(\beta, \mathbf{z} | \mathbf{x})$.

The Evidence Lower Bound

$$\mathcal{L}(\nu) = \mathbb{E}_q [\log p(\beta, \mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\beta, \mathbf{z}; \nu)]$$

- KL is intractable; VI optimizes the **evidence lower bound** (ELBO) instead.
 - It is a lower bound on $\log p(\mathbf{x})$.
 - Maximizing the ELBO is equivalent to minimizing the KL.
- The ELBO trades off two terms.
 - The first term prefers $q(\cdot)$ to place its mass on the MAP estimate.
 - The second term encourages $q(\cdot)$ to be diffuse.
- Caveat: The ELBO is not convex.

Mean-field Variational Inference



- We need to specify the form of $q(\beta, \mathbf{z})$.
- The **mean-field family** is fully factorized,

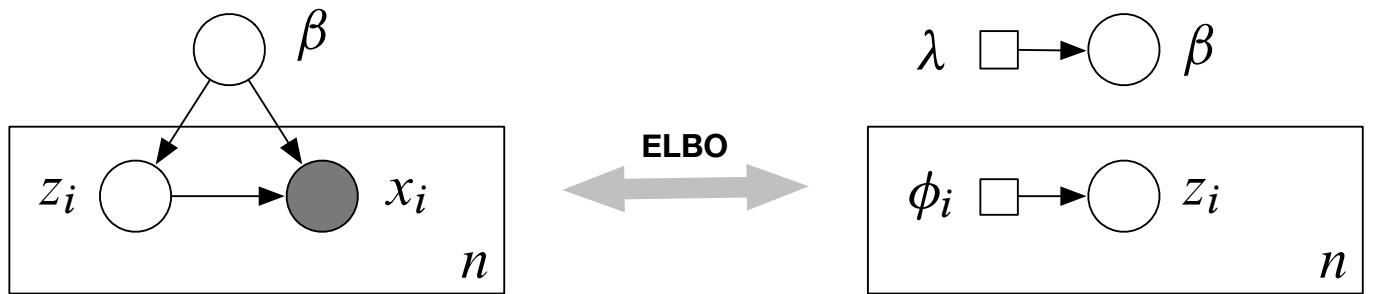
$$q(\beta, \mathbf{z}; \lambda, \boldsymbol{\phi}) = q(\beta; \lambda) \prod_{i=1}^n q(z_i; \phi_i).$$

- Each factor is the same family as the model's complete conditional,

$$p(\beta | \mathbf{z}, \mathbf{x}) = h(\beta) \exp\{\eta_g(\mathbf{z}, \mathbf{x})^\top \beta - a(\eta_g(\mathbf{z}, \mathbf{x}))\}$$

$$q(\beta; \lambda) = h(\beta) \exp\{\lambda^\top \beta - a(\lambda)\}.$$

Mean-field Variational Inference



- Optimize the ELBO,

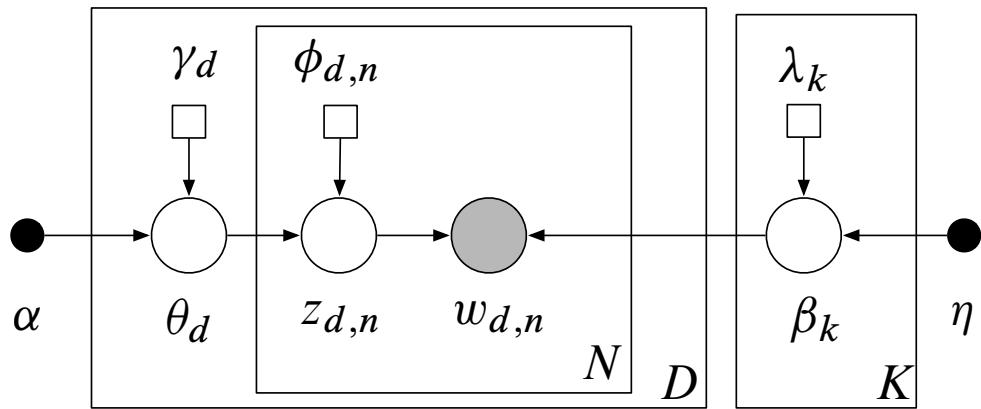
$$\mathcal{L}(\lambda, \phi) = \mathbb{E}_q [\log p(\beta, \mathbf{z}, \mathbf{x})] - \mathbb{E}_q [\log q(\beta, \mathbf{z})].$$

- Traditional VI uses coordinate ascent [Ghahramani and Beal, 2001]

$$\lambda^* = \mathbb{E}_\phi [\eta_g(\mathbf{z}, \mathbf{x})]; \phi_i^* = \mathbb{E}_\lambda [\eta_\ell(\beta, x_i)]$$

- Iteratively update each parameter, holding others fixed.
 - Notice the relationship to Gibbs sampling [Gelfand and Smith, 1990].
 - Caveat: The ELBO is not convex.

Mean-field Variational Inference for LDA



- The local variables are the per-document variables θ_d and \mathbf{z}_d .
- The global variables are the topics β_1, \dots, β_K .
- The variational distribution is

$$q(\beta, \theta, \mathbf{z}) = \prod_{k=1}^K q(\beta_k; \lambda_k) \prod_{d=1}^D q(\theta_d; \gamma_d) \prod_{n=1}^N q(z_{d,n}; \phi_{d,n})$$

Mean-field Variational Inference for LDA

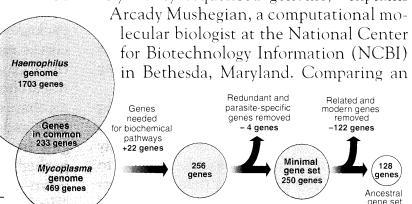
Seeking Life's Bare (Genetic) Necessities

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 Uppsala 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, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains

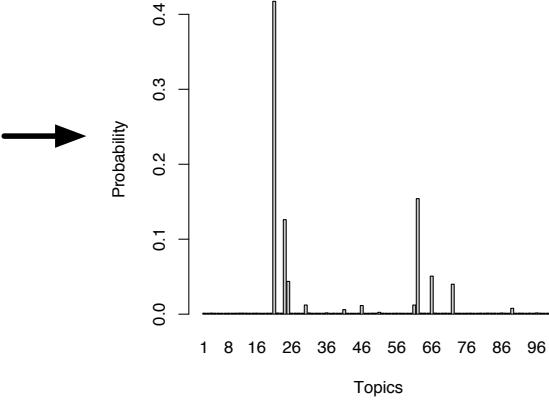
Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



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

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

SCIENCE • VOL. 272 • 24 MAY 1996



Mean-field Variational Inference for LDA

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

Classical Variational Inference

Input: data \mathbf{x} , model $p(\beta, \mathbf{z}, \mathbf{x})$.

Initialize λ randomly.

repeat

for each data point i **do**

 | Set local parameter $\phi_i \leftarrow \mathbb{E}_\lambda [\eta_\ell(\beta, x_i)]$.

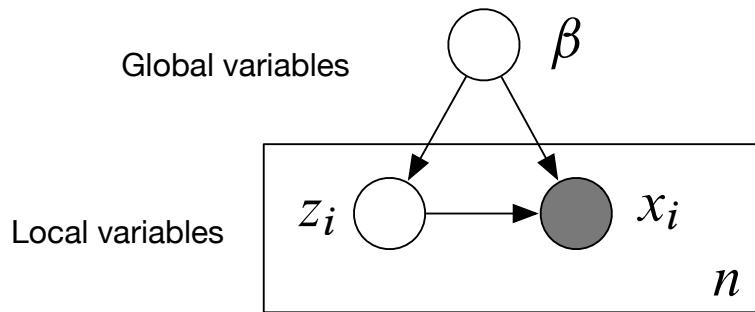
end

 Set global parameter

$$\lambda \leftarrow \alpha + \sum_{i=1}^n \mathbb{E}_{\phi_i} [t(Z_i, x_i)].$$

until the ELBO has converged

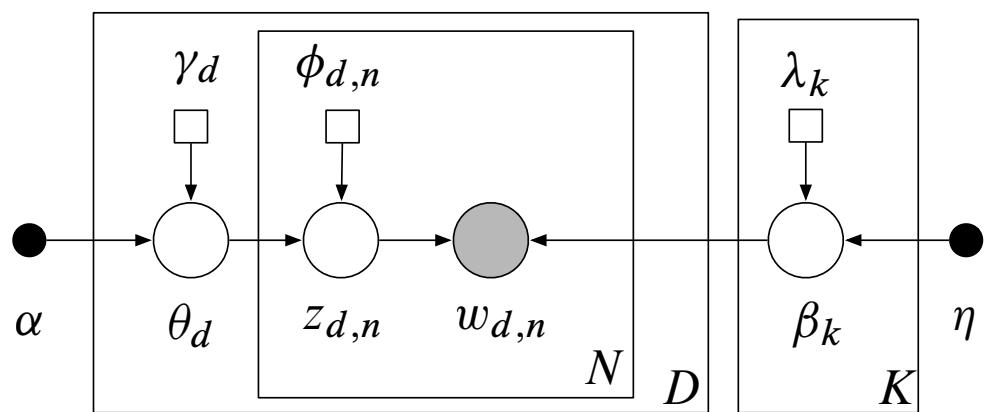
A Generic Class of Models



$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

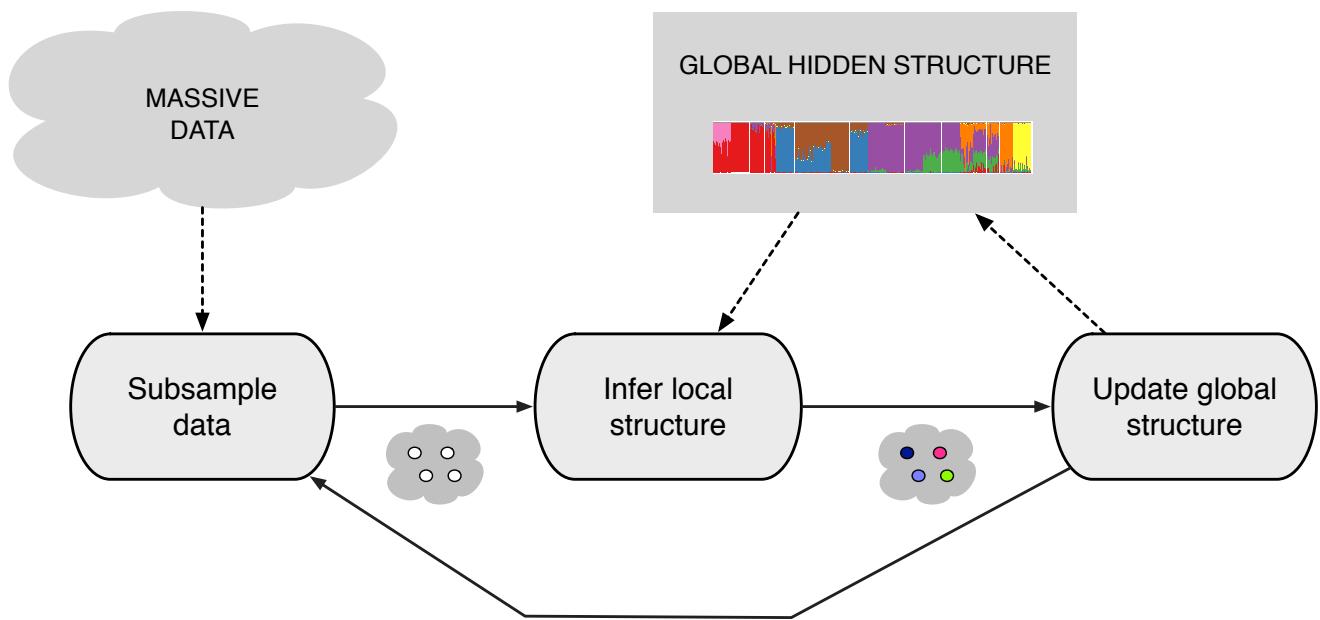
- Bayesian mixture models
- Time series models
(HMMs, linear dynamic systems)
- Factorial models
- Matrix factorization
(factor analysis, PCA, CCA)
- Dirichlet process mixtures, HDPs
- Multilevel regression
(linear, probit, Poisson)
- Stochastic block models
- Mixed-membership models
(LDA and some variants)

Stochastic Variational Inference



- Classical VI is inefficient:
 - Do some local computation *for each data point.*
 - Aggregate these computations to re-estimate global structure.
 - Repeat.
- This cannot handle massive data.
- **Stochastic variational inference** (SVI) scales VI to massive data.

Stochastic Variational Inference



Stochastic Optimization

A STOCHASTIC APPROXIMATION METHOD¹

By HERBERT ROBBINS AND SUTTON MONRO

University of North Carolina

1. Summary. Let $M(x)$ denote the expected value at level x of the response to a certain experiment. $M(x)$ is assumed to be a monotone function of x but is unknown to the experimenter, and it is desired to find the solution $x = \theta$ of the equation $M(x) = \alpha$, where α is a given constant. We give a method for making successive experiments at levels x_1, x_2, \dots in such a way that x_n will tend to θ in probability.



- Replace the gradient with cheaper noisy estimates [Robbins and Monro, 1951]
- Guaranteed to converge to a local optimum [Bottou, 1996]
- Has enabled modern machine learning

Stochastic Optimization

A STOCHASTIC APPROXIMATION METHOD¹

By HERBERT ROBBINS AND SUTTON MONRO

University of North Carolina

1. Summary. Let $M(x)$ denote the expected value at level x of the response to a certain experiment. $M(x)$ is assumed to be a monotone function of x but is unknown to the experimenter, and it is desired to find the solution $x = \theta$ of the equation $M(x) = \alpha$, where α is a given constant. We give a method for making successive experiments at levels x_1, x_2, \dots in such a way that x_n will tend to θ in probability.



- With noisy gradients, update

$$\nu_{t+1} = \nu_t + \rho_t \hat{\nabla}_\nu \mathcal{L}(\nu_t)$$

- Requires unbiased gradients, $\mathbb{E}[\hat{\nabla}_\nu \mathcal{L}(\nu)] = \nabla_\nu \mathcal{L}(\nu)$
- Requires the step size sequence ρ_t follows the Robbins-Monro conditions

Stochastic Variational Inference

- The **natural gradient** of the ELBO [Amari, 1998; Sato, 2001]

$$\nabla_{\lambda}^{\text{nat}} \mathcal{L}(\lambda) = \left(\alpha + \sum_{i=1}^n \mathbb{E}_{\phi_i^*}[t(Z_i, x_i)] \right) - \lambda.$$

- Construct a **noisy natural gradient**,

$$j \sim \text{Uniform}(1, \dots, n)$$

$$\hat{\nabla}_{\lambda}^{\text{nat}} \mathcal{L}(\lambda) = \alpha + n \mathbb{E}_{\phi_j^*}[t(Z_j, x_j)] - \lambda.$$

- This is a good noisy gradient.
 - Its expectation is the exact gradient (*unbiased*).
 - It only depends on optimized parameters of one data point (*cheap*).

Stochastic Variational Inference

Input: data \mathbf{x} , model $p(\beta, \mathbf{z}, \mathbf{x})$.

Initialize λ randomly. Set ρ_t appropriately.

repeat

 Sample $j \sim \text{Unif}(1, \dots, n)$.

 Set local parameter $\phi \leftarrow \mathbb{E}_\lambda [\eta_\ell(\beta, x_j)]$.

 Set intermediate global parameter

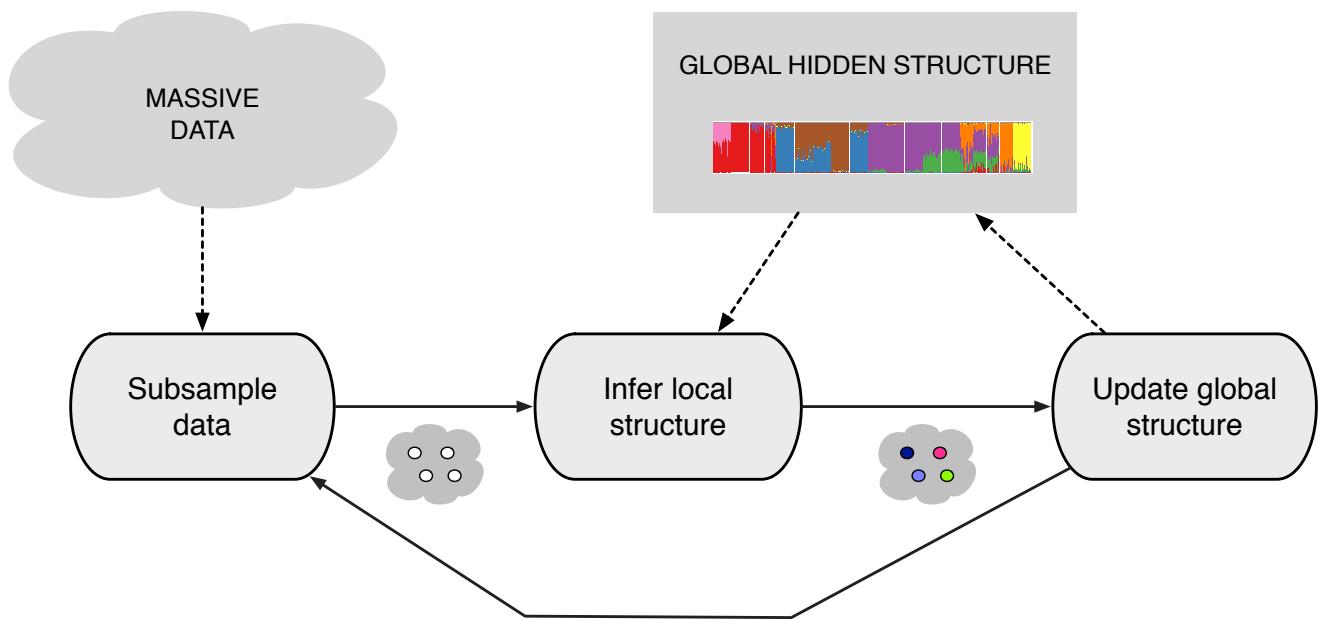
$$\hat{\lambda} = \alpha + n\mathbb{E}_\phi [t(Z_j, x_j)].$$

 Set global parameter

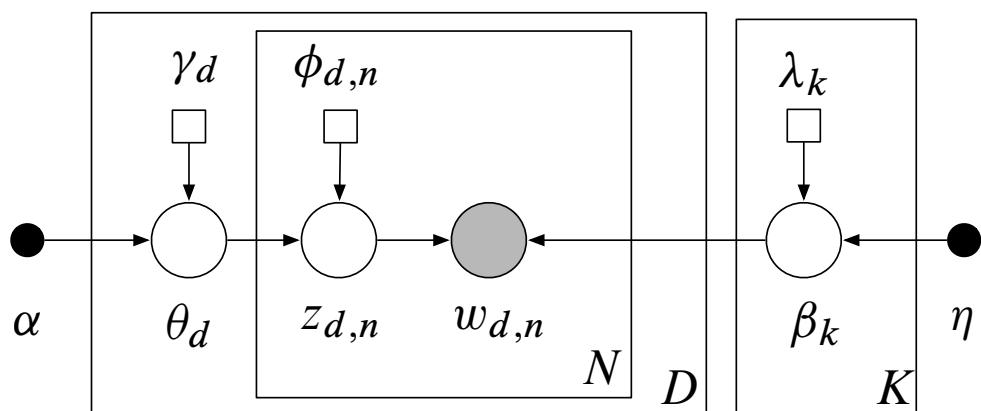
$$\lambda = (1 - \rho_t)\lambda + \rho_t \hat{\lambda}.$$

until *forever*

Stochastic Variational Inference

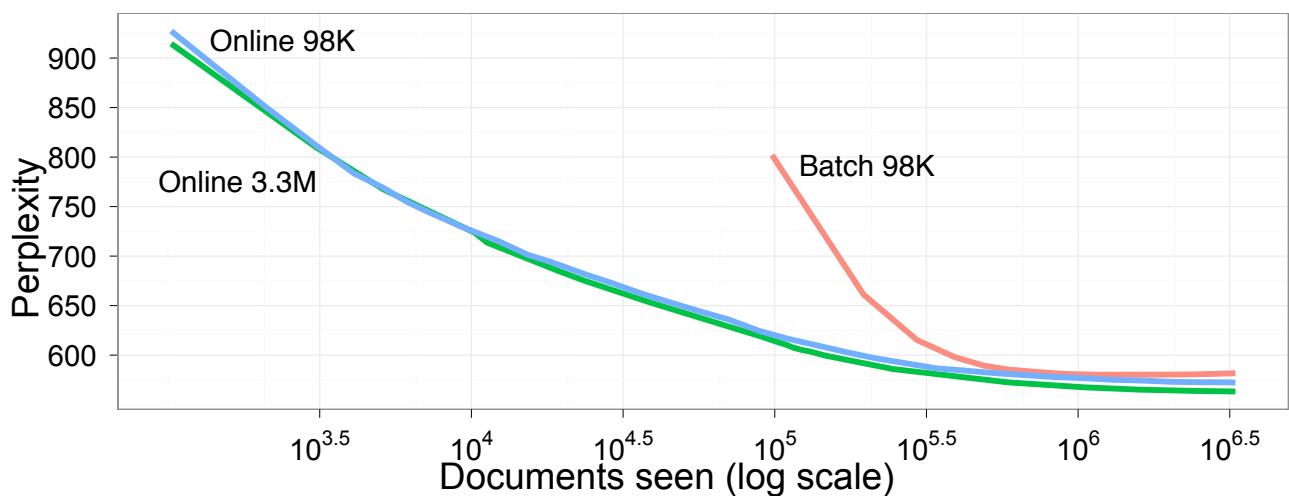


Stochastic Variational Inference in LDA



- Sample a document
- Estimate the local variational parameters using the current topics
- Form intermediate topics from those local parameters
- Update topics as a weighted average of intermediate and current topics

Stochastic Variational Inference in LDA



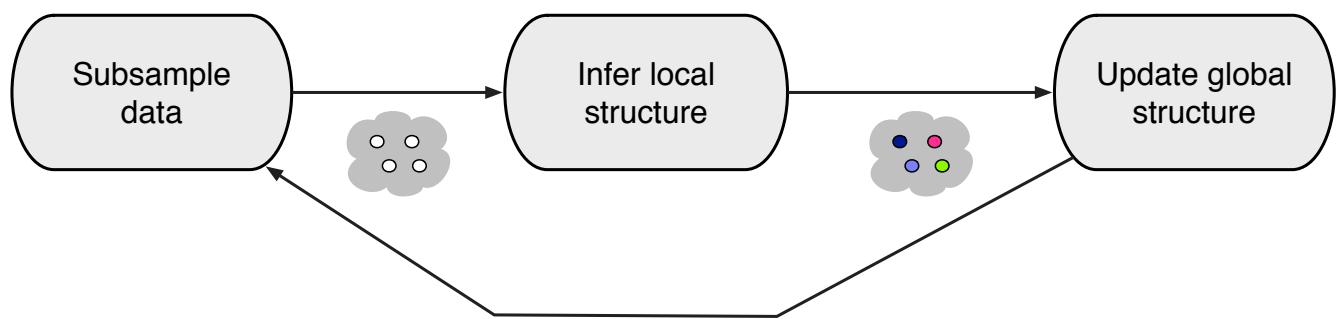
Documents analyzed	2048	4096	8192	12288	16384	32768	49152	65536
Top eight words	systems road made service announced national west language	systems health communication service billion language care road	service systems health companies market communication company billion	service systems companies business company billion health industry	service companies systems business company billion industry market	business service companies industry company management systems services	business service companies industry services company management public	business industry service companies services company management public

[Hoffman et al., 2010]

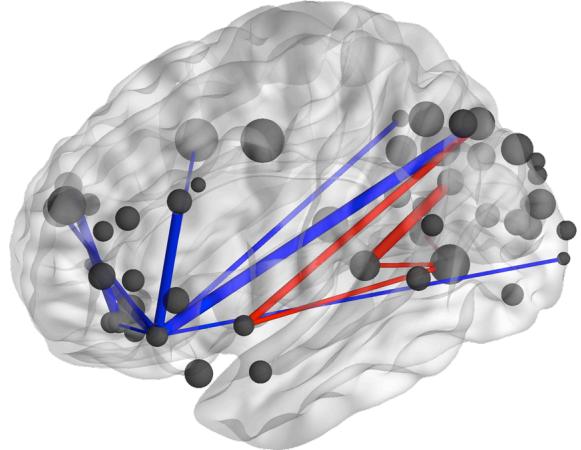
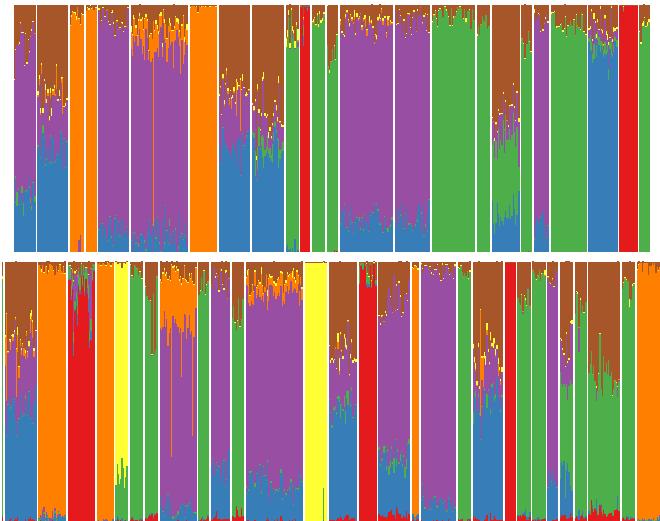
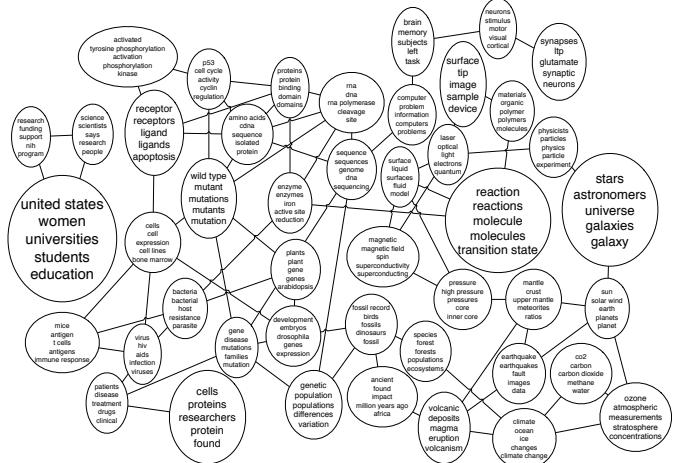
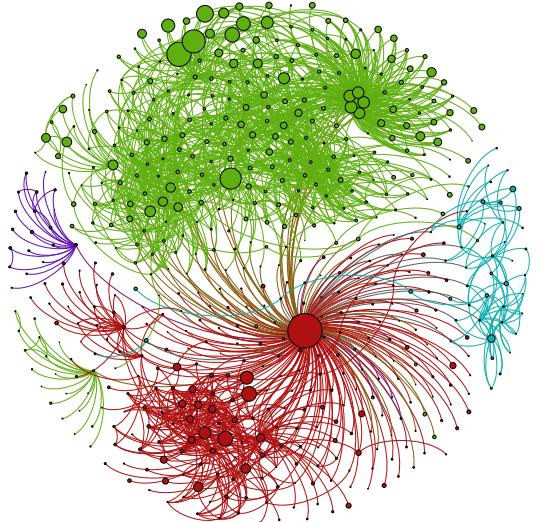


Topics using the HDP, found in 1.8M articles from the New York Times

SVI scales many models



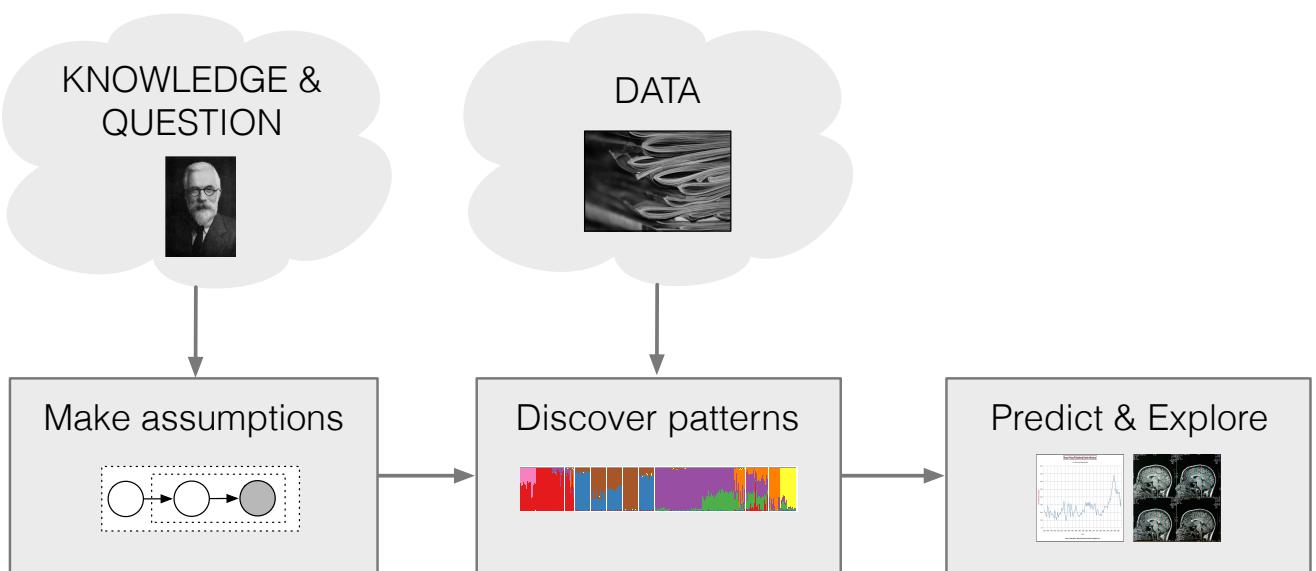
- Bayesian mixture models
- Time series models
(HMMs, linear dynamic systems)
- Factorial models
- Matrix factorization
(factor analysis, PCA, CCA)
- Dirichlet process mixtures, HDPs
- Multilevel regression
(linear, probit, Poisson)
- Stochastic block models
- Mixed-membership models
(LDA and some variants)



PART III

Stochastic Gradients of the ELBO

Review: The Promise



- Realized for conditionally conjugate models
- What about the general case?

The Variational Inference Recipe

Start with a model:

$$p(\mathbf{z}, \mathbf{x})$$



The Variational Inference Recipe

Choose a variational approximation:

$$q(\mathbf{z}; \nu)$$



The Variational Inference Recipe

Write down the ELBO:

$$\mathcal{L}(\nu) = \mathbb{E}_{q(z; \nu)}[\log p(x, z) - \log q(z; \nu)]$$



The Variational Inference Recipe

Compute the expectation(integral):

$$\text{Example: } \mathcal{L}(\nu) = x\nu^2 + \log \nu$$



The Variational Inference Recipe

Take derivatives:

$$\text{Example: } \nabla_{\nu} \mathcal{L}(\nu) = 2x\nu + \frac{1}{\nu}$$



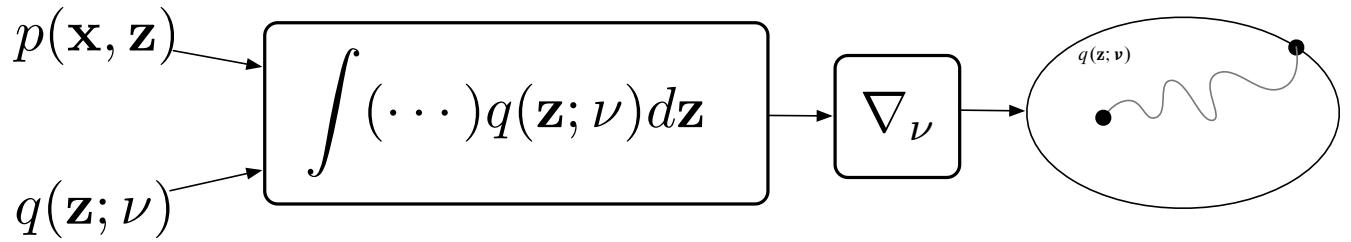
The Variational Inference Recipe

Optimize:

$$\boldsymbol{\nu}_{t+1} = \boldsymbol{\nu}_t + \rho_t \nabla_{\boldsymbol{\nu}} \mathcal{L}$$



The Variational Inference Recipe



Example: Bayesian Logistic Regression

- Data pairs y_i, x_i
- x_i are covariates
- y_i are label
- z is the regression coefficient
- Generative process

$$p(z) \sim N(0, 1)$$
$$p(y_i | x_i, z) \sim \text{Bernoulli}(\sigma(zx_i))$$

VI for Bayesian Logistic Regression

Assume:

- We have one data point (y, x)
- x is a scalar
- The approximating family q is the normal; $\nu = (\mu, \sigma^2)$

The ELBO is

$$\mathcal{L}(\mu, \sigma^2) = \mathbb{E}_q[\log p(z) + \log p(y | x, z) - \log q(z)]$$

VI for Bayesian Logistic Regression

$$\begin{aligned}\mathcal{L}(\mu, \sigma^2) \\ = & \mathbb{E}_q[\log p(z) - \log q(z) + \log p(y|x, z)]\end{aligned}$$

VI for Bayesian Logistic Regression

$$\begin{aligned}\mathcal{L}(\mu, \sigma^2) &= \mathbb{E}_q[\log p(z) - \log q(z) + \log p(y|x, z)] \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q[\log p(y|x, z)] + C\end{aligned}$$

VI for Bayesian Logistic Regression

$$\begin{aligned}\mathcal{L}(\mu, \sigma^2) &= \mathbb{E}_q[\log p(z) - \log q(z) + \log p(y|x, z)] \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q[\log p(y|x, z)] + C \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q[yxz - \log(1 + \exp(xz))]\end{aligned}$$

VI for Bayesian Logistic Regression

$$\begin{aligned}\mathcal{L}(\mu, \sigma^2) &= \mathbb{E}_q[\log p(z) - \log q(z) + \log p(y|x, z)] \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q[\log p(y|x, z)] + C \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q[yxz - \log(1 + \exp(xz))] \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + yx\mu - \mathbb{E}_q[\log(1 + \exp(xz))]\end{aligned}$$

VI for Bayesian Logistic Regression

$$\begin{aligned}\mathcal{L}(\mu, \sigma^2) &= \mathbb{E}_q[\log p(z) - \log q(z) + \log p(y|x, z)] \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q[\log p(y|x, z)] + C \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + \mathbb{E}_q[yxz - \log(1 + \exp(xz))] \\ &= -\frac{1}{2}(\mu^2 + \sigma^2) + \frac{1}{2} \log \sigma^2 + yx\mu - \mathbb{E}_q[\log(1 + \exp(xz))]\end{aligned}$$

We are stuck.

1. We cannot analytically take that expectation.
2. The expectation hides the objectives dependence on the variational parameters. This makes it hard to directly optimize.

Options?

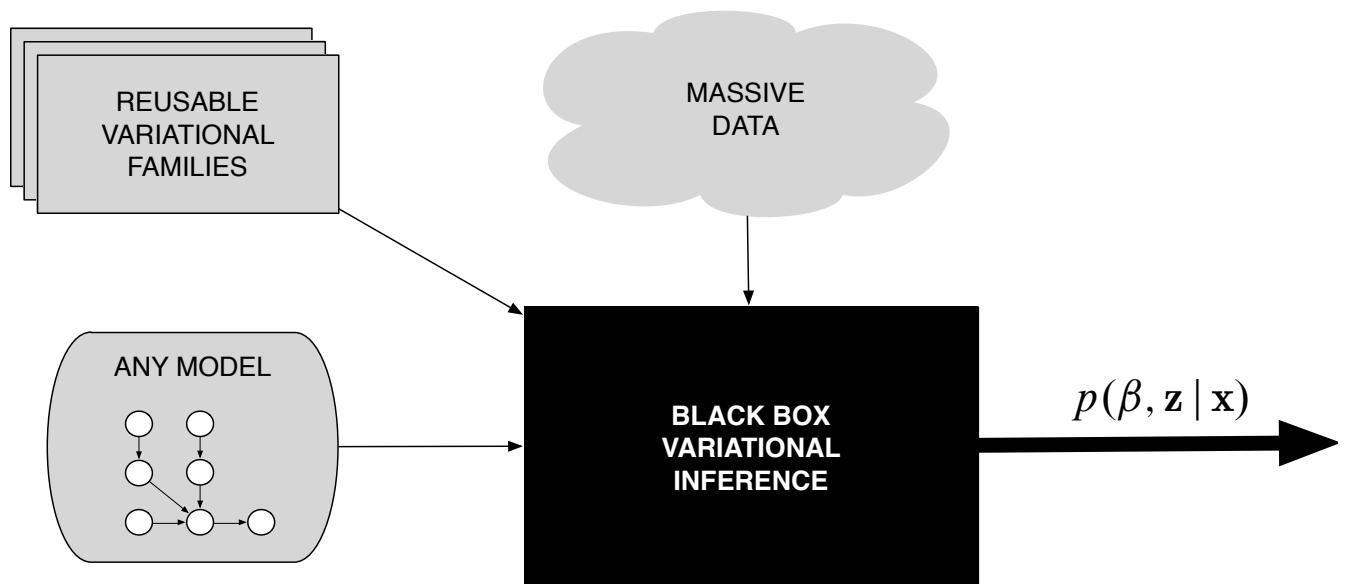
- Derive a model specific bound:
[Jordan and Jaakola; 1996], [Braun and McAuliffe; 2008], others
- More general approximations that require model-specific analysis:
[Wang and Blei; 2013], [Knowles and Minka; 2011]

Nonconjugate Models

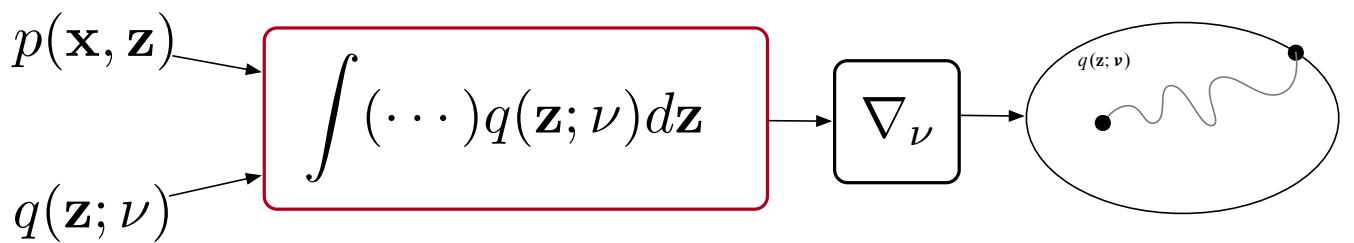
- Nonlinear Time series Models
- Deep Latent Gaussian Models
- Models with Attention
(such as DRAW)
- Generalized Linear Models
(Poisson Regression)
- Stochastic Volatility Models
- Discrete Choice Models
- Bayesian Neural Networks
- Deep Exponential Families
(e.g. Sparse Gamma or Poisson)
- Correlated Topic Model
(including nonparametric variants)
- Sigmoid Belief Network

We need a solution that does not entail model specific work

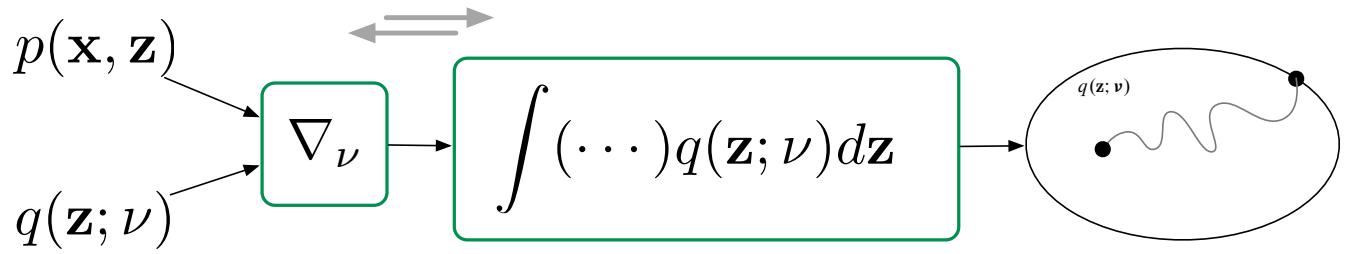
Black Box Variational Inference (BBVI)



The Problem in the Classical VI Recipe



The New VI Recipe



Use stochastic optimization!

Computing Gradients of Expectations

- Define

$$g(\mathbf{z}, \boldsymbol{\nu}) = \log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \boldsymbol{\nu})$$

- What is $\nabla_{\boldsymbol{\nu}} \mathcal{L}$

$$\begin{aligned}\nabla_{\boldsymbol{\nu}} \mathcal{L} &= \nabla_{\boldsymbol{\nu}} \int q(\mathbf{z}; \boldsymbol{\nu}) g(\mathbf{z}, \boldsymbol{\nu}) d\mathbf{z} \\ &= \int \nabla_{\boldsymbol{\nu}} q(\mathbf{z}; \boldsymbol{\nu}) g(\mathbf{z}, \boldsymbol{\nu}) + q(\mathbf{z}; \boldsymbol{\nu}) \nabla_{\boldsymbol{\nu}} g(\mathbf{z}, \boldsymbol{\nu}) d\mathbf{z} \\ &= \int q(\mathbf{z}; \boldsymbol{\nu}) \nabla_{\boldsymbol{\nu}} \log q(\mathbf{z}; \boldsymbol{\nu}) g(\mathbf{z}, \boldsymbol{\nu}) + q(\mathbf{z}; \boldsymbol{\nu}) \nabla_{\boldsymbol{\nu}} g(\mathbf{z}, \boldsymbol{\nu}) d\mathbf{z} \\ &= \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\nu})} [\nabla_{\boldsymbol{\nu}} \log q(\mathbf{z}; \boldsymbol{\nu}) g(\mathbf{z}, \boldsymbol{\nu}) + \nabla_{\boldsymbol{\nu}} g(\mathbf{z}, \boldsymbol{\nu})]\end{aligned}$$

Using $\nabla_{\boldsymbol{\nu}} \log q = \frac{\nabla_{\boldsymbol{\nu}} q}{q}$

Roadmap

- **Score Function Gradients**
- **Pathwise Gradients**
- **Amortized Inference**

Score Function Gradients of the ELBO

Score Function Estimator

Recall

$$\nabla_{\nu} \mathcal{L} = \mathbb{E}_{q(\mathbf{z}; \nu)} [\nabla_{\nu} \log q(\mathbf{z}; \nu) g(\mathbf{z}, \nu) + \nabla_{\nu} g(\mathbf{z}, \nu)]$$

Simplify:

$$\mathbb{E}_q [\nabla_{\nu} g(\mathbf{z}, \nu)] = \mathbb{E}_q [\nabla_{\nu} \log q(\mathbf{z}; \nu)] = 0$$

Gives the gradient:

$$\nabla_{\nu} \mathcal{L} = \mathbb{E}_{q(\mathbf{z}; \nu)} [\nabla_{\nu} \log q(\mathbf{z}; \nu) (\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \nu))]$$

Sometimes called likelihood ratio or REINFORCE gradients

[Glynn 1990; Williams, 1992; Wingate+ 2013; Ranganath+ 2014; Mnih+ 2014]

Noisy Unbiased Gradients

Gradient: $\mathbb{E}_{q(\mathbf{z}; \boldsymbol{\nu})}[\nabla_{\boldsymbol{\nu}} \log q(\mathbf{z}; \boldsymbol{\nu})(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \boldsymbol{\nu}))]$

Noisy unbiased gradients with Monte Carlo!

$$\frac{1}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\nu}} \log q(\mathbf{z}_s; \boldsymbol{\nu})(\log p(\mathbf{x}, \mathbf{z}_s) - \log q(\mathbf{z}_s; \boldsymbol{\nu})),$$

where $\mathbf{z}_s \sim q(\mathbf{z}; \boldsymbol{\nu})$

Basic BBVI

Algorithm 1: Basic Black Box Variational Inference

Input : Model $\log p(\mathbf{x}, \mathbf{z})$,
Variational approximation $q(\mathbf{z}; \boldsymbol{\nu})$

Output : Variational Parameters: $\boldsymbol{\nu}$

while *not converged* **do**

$\mathbf{z}[s] \sim q$ // Draw S samples from q

$\rho = t$ -th value of a Robbins Monro sequence

$\boldsymbol{\nu} = \boldsymbol{\nu} + \rho \frac{1}{S} \sum_{s=1}^S \nabla_{\boldsymbol{\nu}} \log q(\mathbf{z}[s]; \boldsymbol{\nu})(\log p(\mathbf{x}, \mathbf{z}[s]) - \log q(\mathbf{z}[s]; \boldsymbol{\nu}))$

$t = t + 1$

end

The requirements for inference

The noisy gradient:

$$\frac{1}{S} \sum_{s=1}^S \nabla_{\nu} \log q(\mathbf{z}_s; \nu) (\log p(\mathbf{x}, \mathbf{z}_s) - \log q(\mathbf{z}_s; \nu)),$$

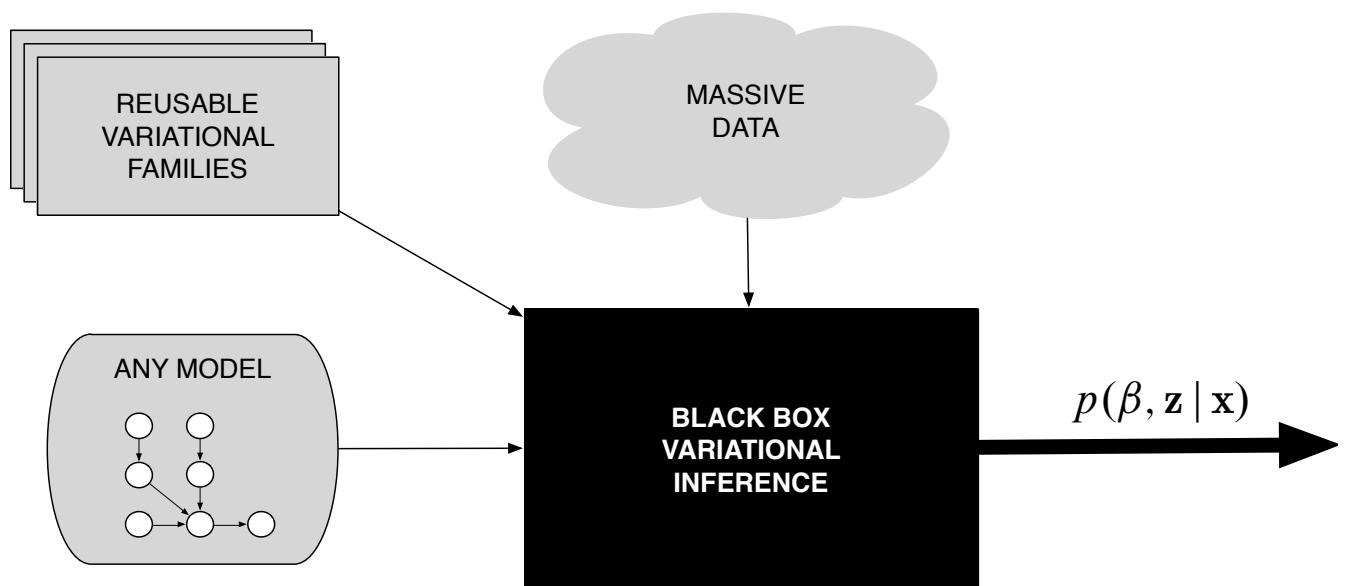
where $\mathbf{z}_s \sim q(\mathbf{z}; \nu)$

To compute the noisy gradient of the ELBO we need

- Sampling from $q(\mathbf{z})$
- Evaluating $\nabla_{\nu} \log q(\mathbf{z}; \nu)$
- Evaluating $\log p(\mathbf{x}, \mathbf{z})$ and $\log q(\mathbf{z})$

There is no model specific work: black box criteria are satisfied

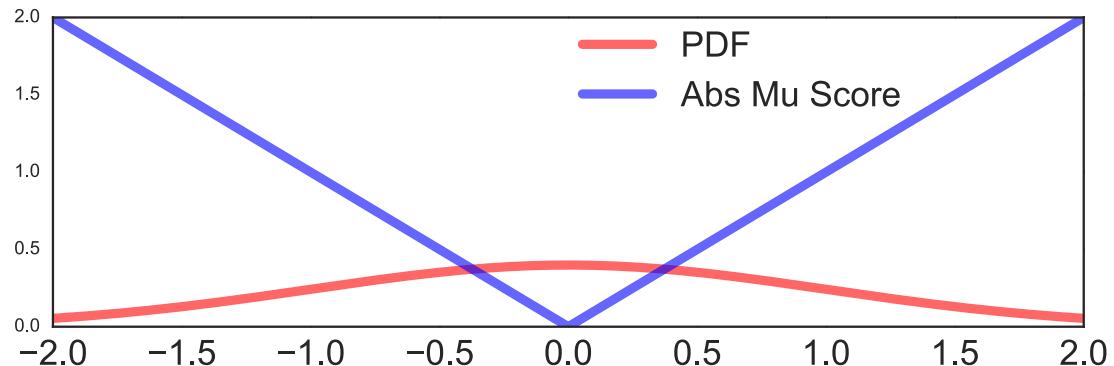
Black Box Variational Inference



Problem: Basic BBVI doesn't work

Variance of the gradient can be a problem

$$\text{Var}_{q(\mathbf{z}; \nu)} = \mathbb{E}_{q(\mathbf{z}; \nu)}[(\nabla_\nu \log q(\mathbf{z}; \nu)(\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \nu)) - \nabla_\nu \mathcal{L})^2].$$



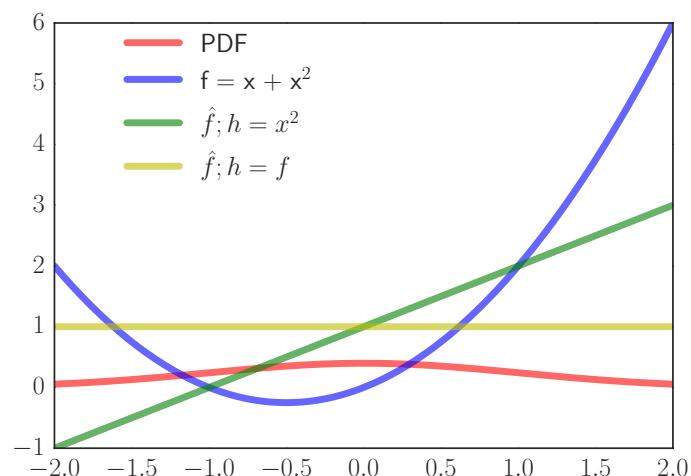
Intuition:

Sampling rare values can lead to large scores and thus high variance

Solution: Control Variates

Replace with f with \hat{f} where $\mathbb{E}[\hat{f}(z)] = \mathbb{E}[f(z)]$. General such class:

$$\hat{f}(z) \triangleq f(z) - a(h(z) - \mathbb{E}[h(z)])$$

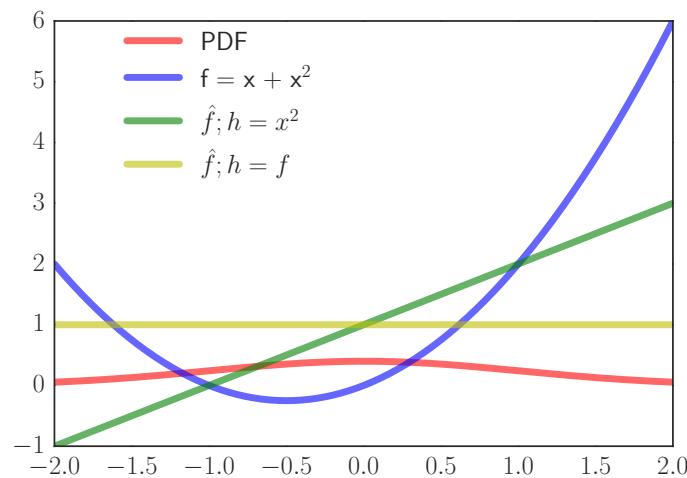


- h is a function of our choice
- a is chosen to minimize the variance
- Good h have high correlation with the original function f

Solution: Control Variates

Replace with f with \hat{f} where $\mathbb{E}[\hat{f}(z)] = \mathbb{E}[f(z)]$. General such class:

$$\hat{f}(z) \triangleq f(z) - a(h(z) - \mathbb{E}[h(z)])$$

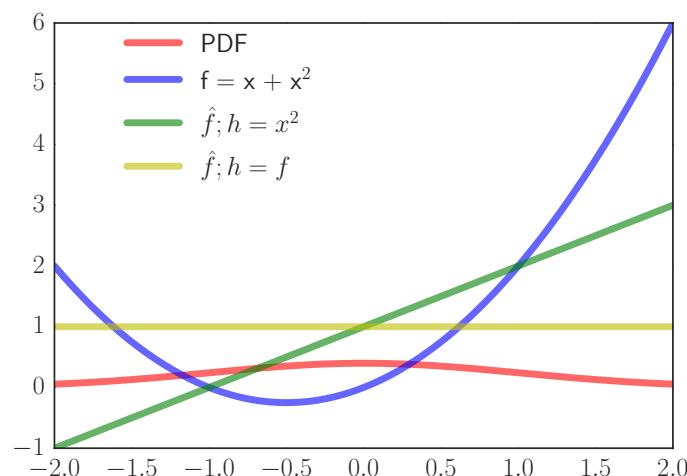


- For variational inference we need functions with known q expectation
- Set h as $\nabla_{\nu} \log q(\mathbf{z}; \nu)$
- Simple as $\mathbb{E}_q[\nabla_{\nu} \log q(\mathbf{z}; \nu)] = 0$ for any q

Solution: Control Variates

Replace with f with \hat{f} where $\mathbb{E}[\hat{f}(z)] = \mathbb{E}[f(z)]$. General such class:

$$\hat{f}(z) \triangleq f(z) - a(h(z) - \mathbb{E}[h(z)])$$



Many of the other techniques from Monte Carlo can help:

- *Importance Sampling, Quasi Monte Carlo, Rao-Blackwellization*

[Ruiz+ 2016; Ranganath+2014; Titsias+2015; Mnih+2016]

Nonconjugate Models

- Nonlinear Time series Models
- Deep Latent Gaussian Models
- Models with Attention
(such as DRAW)
- Generalized Linear Models
(Poisson Regression)
- Stochastic Volatility Models
- Discrete Choice Models
- Bayesian Neural Networks
- Deep Exponential Families
(e.g. Sparse Gamma or Poisson)
- Correlated Topic Model
(including nonparametric variants)
- Sigmoid Belief Network

We can design models based on data rather than inference.

More Assumptions?

The current black box criteria

- Sampling from $q(\mathbf{z})$
- Evaluating $\nabla_{\boldsymbol{\nu}} \log q(\mathbf{z}; \boldsymbol{\nu})$
- Evaluating $\log p(\mathbf{x}, \mathbf{z})$ and $\log q(\mathbf{z})$

Can we make additional assumptions that are not too restrictive?

Pathwise Gradients of the ELBO

Pathwise Estimator

Assume

1. $\mathbf{z} = t(\epsilon, \nu)$ for $\epsilon \sim s(\epsilon)$ implies $\mathbf{z} \sim q(\mathbf{z}; \nu)$

Example:

$$\begin{aligned}\epsilon &\sim \text{Normal}(0, 1) \\ z &= \epsilon\sigma + \mu \\ \rightarrow z &\sim \text{Normal}(\mu, \sigma^2)\end{aligned}$$

2. $\log p(\mathbf{x}, \mathbf{z})$ and $\log q(\mathbf{z})$ are differentiable with respect to \mathbf{z}

Pathwise Estimator

Recall

$$\nabla_{\nu} \mathcal{L} = \mathbb{E}_{q(\mathbf{z}; \nu)} [\nabla_{\nu} \log q(\mathbf{z}; \nu) g(\mathbf{z}, \nu) + \nabla_{\nu} g(\mathbf{z}, \nu)]$$

Rewrite using $\mathbf{z} = t(\epsilon, \nu)$

$$\nabla_{\nu} \mathcal{L} = \mathbb{E}_{s(\epsilon)} [\nabla_{\nu} \log s(\epsilon) g(t(\epsilon, \nu), \nu) + \nabla_{\nu} g(t(\epsilon, \nu), \nu)]$$

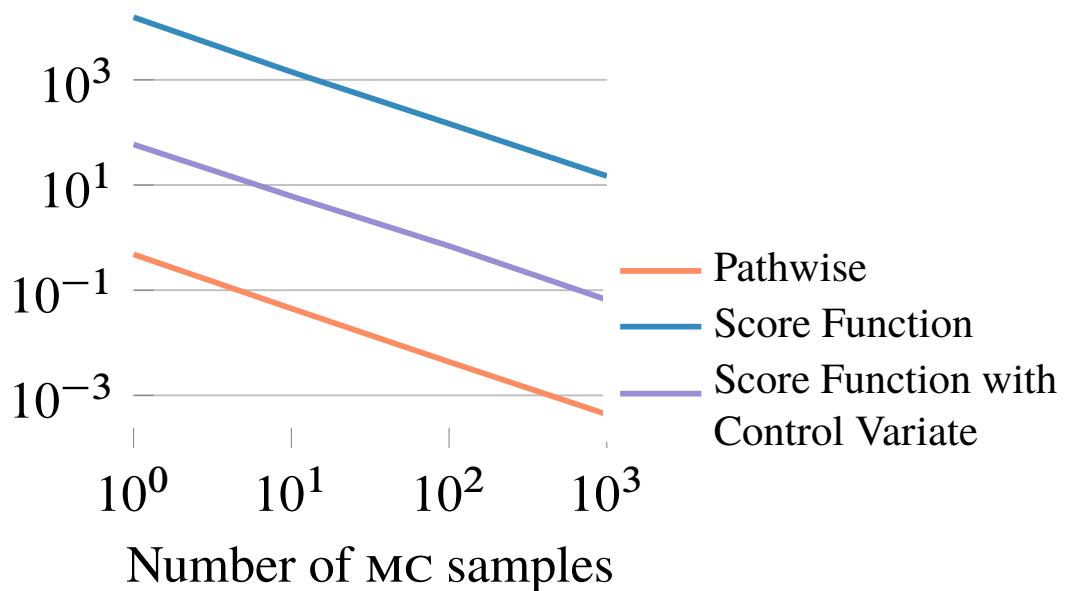
To differentiate:

$$\begin{aligned}\nabla \mathcal{L}(\nu) &= \mathbb{E}_{s(\epsilon)} [\nabla_{\nu} g(t(\epsilon, \nu), \nu)] \\ &= \mathbb{E}_{s(\epsilon)} [\nabla_{\mathbf{z}} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \nu)] \nabla_{\nu} t(\epsilon, \nu) - \nabla_{\nu} \log q(\mathbf{z}; \nu)] \\ &= \mathbb{E}_{s(\epsilon)} [\nabla_{\mathbf{z}} [\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \nu)] \nabla_{\nu} t(\epsilon, \nu)]\end{aligned}$$

This is also known as the reparameterization gradient.

[Glasserman 1991; Fu 2006; Kingma+ 2014; Rezende+ 2014; Titsias+ 2014]

Variance Comparison



[Kucukelbir+ 2016]

Score Function Estimator vs. Pathwise Estimator

Score Function

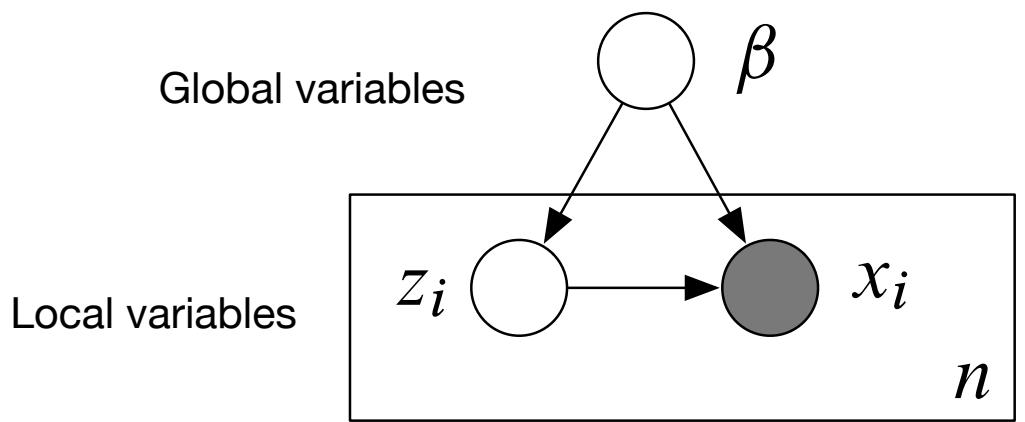
- Differentiates the density $\nabla_{\nu}q(z; \nu)$
- Works for discrete and continuous models
- Works for large class of variational approximations
- Variance can be a big problem

Pathwise

- Differentiates the function $\nabla_z[\log p(\mathbf{x}, \mathbf{z}) - \log q(\mathbf{z}; \nu)]$
- Requires differentiable models
- Requires variational approximation to have form $\mathbf{z} = t(\epsilon, \nu)$
- Generally better behaved variance

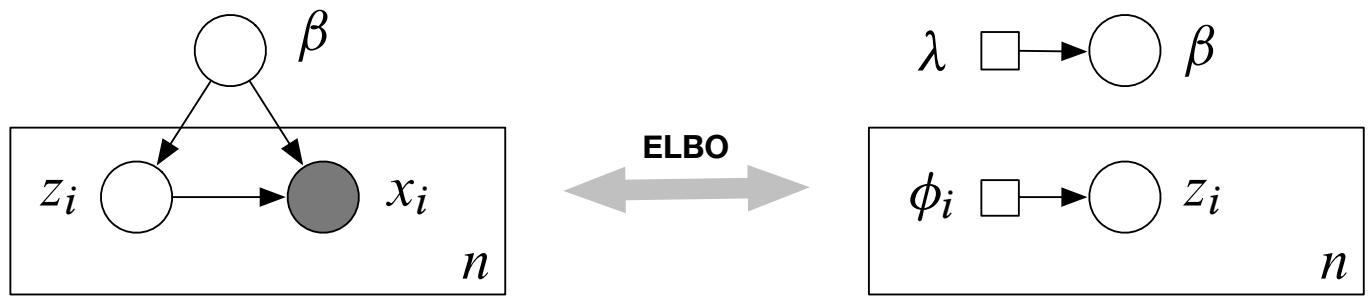
Amortized Inference

Hierarchical Models



$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

Mean Field Variational Approximation



SVI: Revisited

Input: data \mathbf{x} , model $p(\beta, \mathbf{z}, \mathbf{x})$.

Initialize λ randomly. Set ρ_t appropriately.

repeat

 Sample $j \sim \text{Unif}(1, \dots, n)$.

 Set local parameter $\phi \leftarrow \mathbb{E}_\lambda [\eta_\ell(\beta, x_j)]$.

 Set intermediate global parameter

$$\hat{\lambda} = \alpha + n\mathbb{E}_\phi [t(Z_j, x_j)].$$

 Set global parameter

$$\lambda = (1 - \rho_t)\lambda + \rho_t \hat{\lambda}.$$

until *forever*

SVI: The problem

Input: data \mathbf{x} , model $p(\beta, \mathbf{z}, \mathbf{x})$.

Initialize λ randomly. Set ρ_t appropriately.

repeat

 Sample $j \sim \text{Unif}(1, \dots, n)$.

 Set local parameter $\phi \leftarrow \mathbb{E}_\lambda [\eta_\ell(\beta, x_j)]$.

 Set intermediate global parameter

$$\hat{\lambda} = \alpha + n\mathbb{E}_\phi [t(Z_j, x_j)].$$

 Set global parameter

$$\lambda = (1 - \rho_t)\lambda + \rho_t \hat{\lambda}.$$

until *forever*

- These expectations are no longer tractable
- Inner stochastic optimization needed for each data point.

SVI: The problem

Input: data \mathbf{x} , model $p(\beta, \mathbf{z}, \mathbf{x})$.

Initialize λ randomly. Set ρ_t appropriately.

repeat

 Sample $j \sim \text{Unif}(1, \dots, n)$.

 Set local parameter $\phi \leftarrow \mathbb{E}_\lambda [\eta_\ell(\beta, x_j)]$.

 Set intermediate global parameter

$$\hat{\lambda} = \alpha + n\mathbb{E}_\phi [t(Z_j, x_j)].$$

 Set global parameter

$$\lambda = (1 - \rho_t)\lambda + \rho_t \hat{\lambda}.$$

until *forever*

Idea: Learn a mapping f from x_i to ϕ_i

Amortizing Inference

ELBO:

$$\mathcal{L}(\lambda, \phi_{1\dots n}) = \mathbb{E}_q [\log p(\beta, \mathbf{z}, \mathbf{x})] - \mathbb{E}_q \left[\log q(\beta; \lambda) + \sum_{i=1}^n q(z_i; \phi_i) \right]$$

Amortizing the ELBO with *inference network f*:

$$\mathcal{L}(\lambda, \theta) = \mathbb{E}_q [\log p(\beta, \mathbf{z}, \mathbf{x})] - \mathbb{E}_q \left[\log q(\beta; \lambda) + \sum_{i=1}^n q(z_i | \mathbf{x}_i; \phi_i = f_\theta(\mathbf{x}_i)) \right]$$

[Dayan+ 1995; Heess+ 2013; Gershman+ 2014, many others]

Amortized SVI

Input: data \mathbf{x} , model $p(\beta, \mathbf{z}, \mathbf{x})$.

Initialize λ randomly. Set ρ_t appropriately.

repeat

 Sample $\beta \sim q(\beta; \lambda)$.

 Sample $j \sim \text{Unif}(1, \dots, n)$.

 Sample $z_j \sim q(z_j | x_j; \phi_\theta(x_j))$.

 Compute stochastic gradients

$$\hat{\nabla}_\lambda \mathcal{L} = \nabla_\lambda \log q(\beta; \lambda) (\log p(\beta) + n \log p(x_j, z_j | \beta) - \log q(\beta))$$

$$\hat{\nabla}_\theta \mathcal{L} = n \nabla_\theta \log q(z_j | x_j; \theta) (\log p(x_j, z_j | \beta) - \log q(z_j | x_k; \theta))$$

 Update

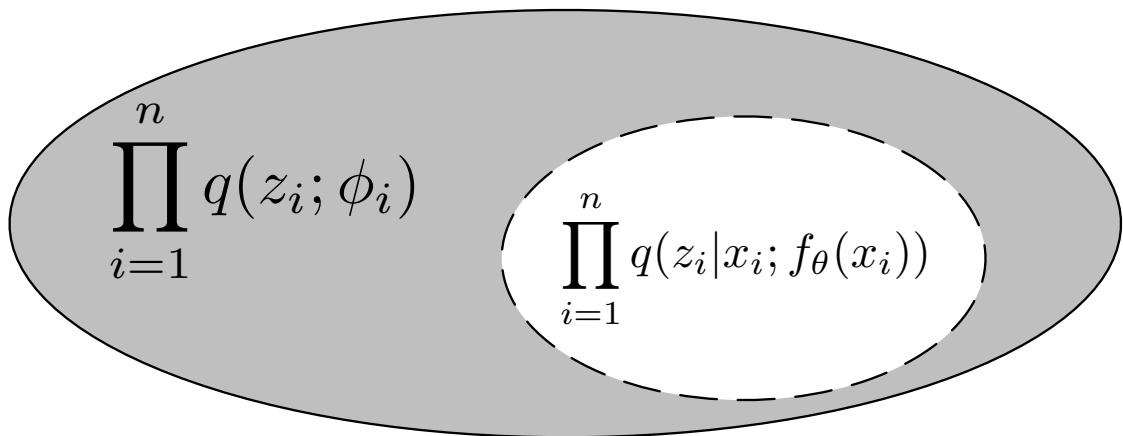
$$\lambda = \lambda + \rho_t \hat{\nabla}_\lambda$$

$$\theta = \theta + \rho_t \hat{\nabla}_\theta.$$

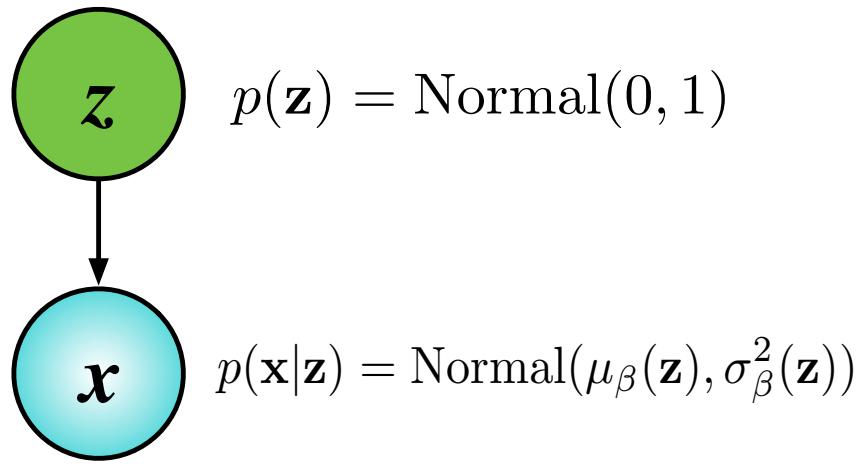
until *forever*

A computational-statistical tradeoff

- Amortized inference is faster, but admits a smaller class of approximations
- The size of the smaller class depends on the flexibility of f



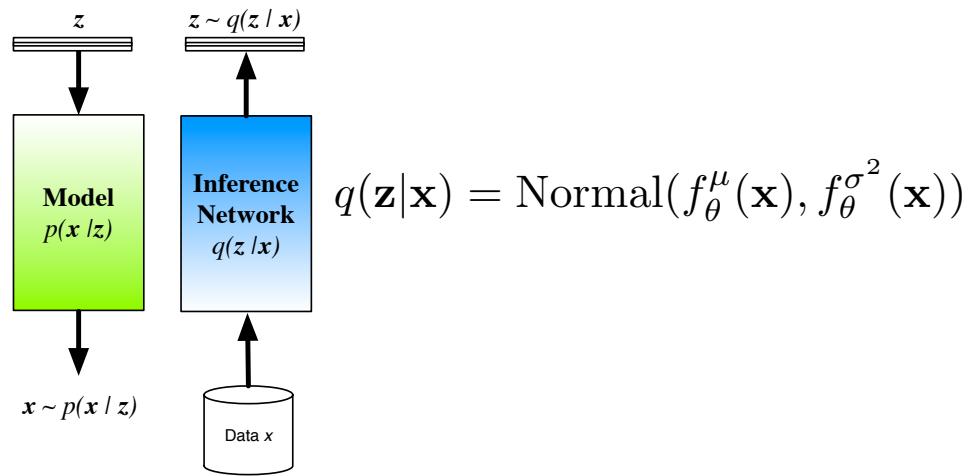
Example: Variational Autoencoder (VAE)



μ and σ^2 are deep networks with parameters β .

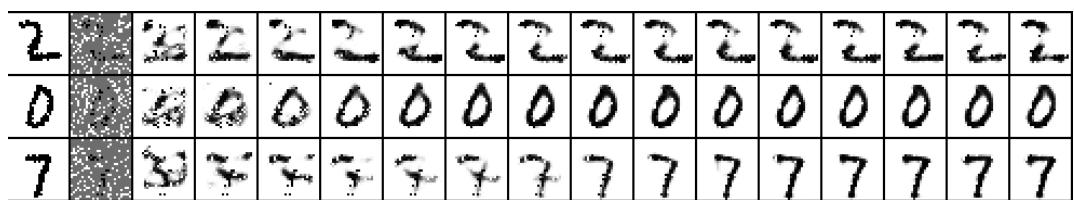
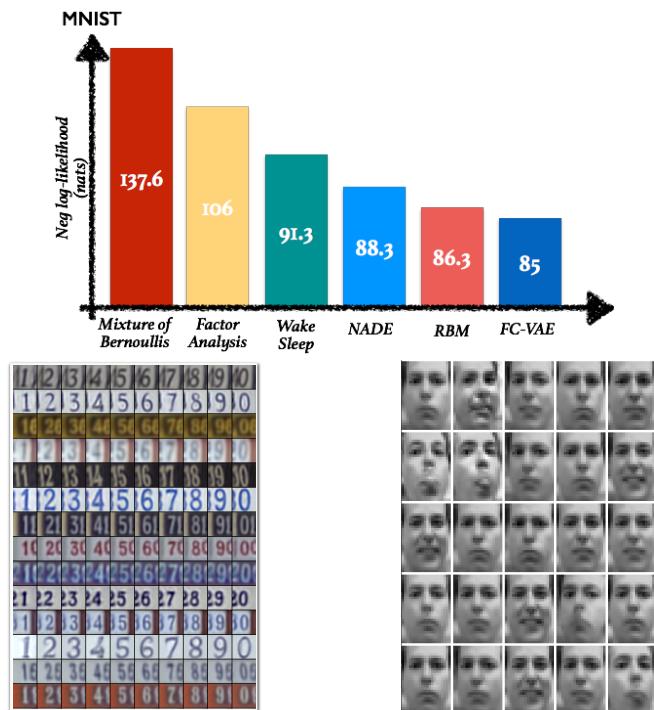
[Kingma+ 2014; Rezende+ 2014]

Example: Variational Autoencoder (VAE)



All functions are deep networks

Example: Variational Autoencoder (VAE)



Rules of Thumb for a New Model

If $\log p(\mathbf{x}, \mathbf{z})$ is \mathbf{z} differentiable

- Try out an approximation q that is reparameterizable

If $\log p(\mathbf{x}, \mathbf{z})$ is not \mathbf{z} differentiable

- Use score function estimator with control variates
- Add further variance reductions based on experimental evidence

Rules of Thumb for a New Model

If $\log p(\mathbf{x}, \mathbf{z})$ is \mathbf{z} differentiable

- Try out an approximation q that is reparameterizable

If $\log p(\mathbf{x}, \mathbf{z})$ is not \mathbf{z} differentiable

- Use score function estimator with control variates
- Add further variance reductions based on experimental evidence

General Advice:

- Use coordinate specific learning rates (e.g. RMSProp, AdaGrad)
- Annealing + Tempering
- Consider parallelizing across samples from q

Software

Systems with Variational Inference:

- Venture, WebPPL, Edward, Stan, PyMC3, Infer.net, Anglican

Good for trying out lots of models

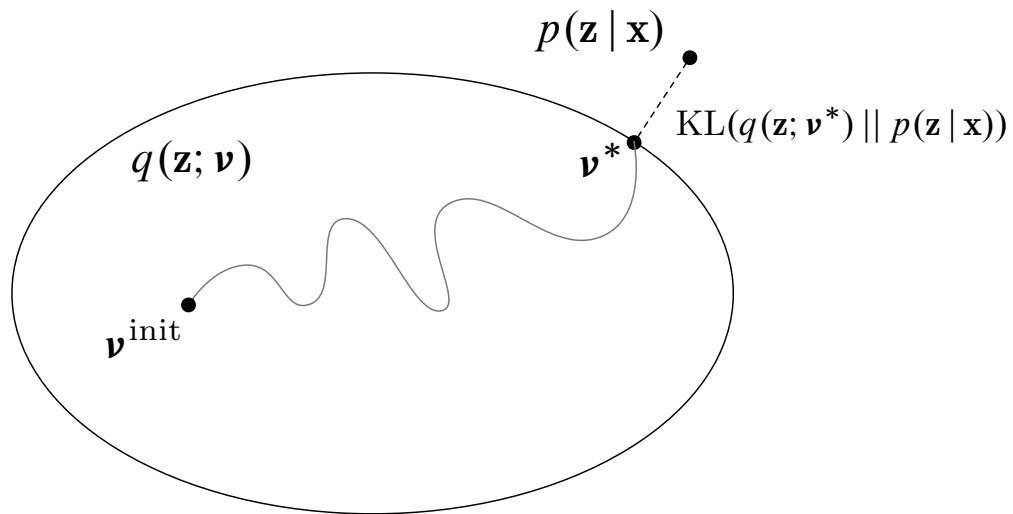
Differentiation Tools:

- Theano, Torch, Tensorflow, Stan Math, Caffe

Can lead to more scalable implementations of individual models

Summary

Variational Inference: Foundations and Modern Methods



VI approximates difficult quantities from complex models.

With **stochastic optimization** we can

- scale up VI to massive data
- enable VI on a wide class of difficult models
- enable VI with elaborate and flexible families of approximations

Bibliography

Introductory Variational Inference

- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S., and Saul, L. K. (1999). An introduction to variational methods for graphical models. *Machine learning*, 37(2), 183-233.
- Beal, Matthew James. Variational algorithms for approximate Bayesian inference. Diss. University of London, 2003.
- Wainwright, Martin J., and Michael I. Jordan. "Graphical models, exponential families, and variational inference." *Foundations and Trends in Machine Learning* 1, no. 1-2 (2008): 1-305.

Bibliography

Applications of Variational Inference

- Frey, Brendan J., and Geoffrey E. Hinton. "Variational learning in nonlinear Gaussian belief networks." *Neural Computation* 11, no. 1 (1999): 193-213.
- Eslami, S. M., Heess, N., Weber, T., Tassa, Y., Kavukcuoglu, K., and Hinton, G. E. Attend, Infer, Repeat: Fast Scene Understanding with Generative Models. NIPS (2016).
- Rezende, Danilo Jimenez, Shakir Mohamed, Ivo Danihelka, Karol Gregor, and Daan Wierstra. "One-Shot Generalization in Deep Generative Models." ICML (2016).
- Kingma, Diederik P., Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. "Semi-supervised learning with deep generative models." In *Advances in Neural Information Processing Systems*, pp. 3581-3589. 2014.

Bibliography

Monte Carlo Gradient Estimation

- Pierre L'Ecuyer, Note: On the interchange of derivative and expectation for likelihood ratio derivative estimators, Management Science, 1995
- Peter W Glynn, Likelihood ratio gradient estimation for stochastic systems, Communications of the ACM, 1990
- Michael C Fu, Gradient estimation, Handbooks in operations research and management science, 2006
- Ronald J Williams, Simple statistical gradient-following algorithms for connectionist reinforcement learning, Machine learning, 1992
- Paul Glasserman, Monte Carlo methods in financial engineering, 2003
- Omiros Papaspiliopoulos, Gareth O Roberts, Martin Skold, A general framework for the parametrization of hierarchical models, Statistical Science, 2007
- Michael C Fu, Gradient estimation, Handbooks in operations research and management science, 2006
- Rajesh Ranganath, Sean Gerrish, and David M. Blei. "Black Box Variational Inference." In AISTATS, pp. 814-822. 2014.
- Andriy Mnih, and Karol Gregor. "Neural variational inference and learning in belief networks." arXiv preprint arXiv:1402.0030 (2014).

Bibliography

Monte Carlo Gradient Estimation (cont.)

- Michalis Titsias and Miguel Lázaro-Gredilla. "Doubly stochastic variational Bayes for non-conjugate inference." (2014).
- David Wingate and Theophane Weber. "Automated variational inference in probabilistic programming." arXiv preprint arXiv:1301.1299 (2013).
- John Paisley, David Blei, and Michael Jordan. "Variational Bayesian inference with stochastic search." arXiv preprint arXiv:1206.6430 (2012).
- Durk Kingma and Max Welling. "Auto-encoding Variational Bayes." ICLR (2014).
- Danilo Jimenez Rezende, Shakir Mohamed, Daan Wierstra. "Stochastic Backpropagation and Approximate Inference in Deep Generative Models." ICML (2014).

Bibliography

Amortized Inference

- Dayan, Peter, Geoffrey E. Hinton, Radford M. Neal, and Richard S. Zemel. "The helmholtz machine." *Neural computation* 7, no. 5 (1995): 889-904.
- Gershman, Samuel J., and Noah D. Goodman. "Amortized inference in probabilistic reasoning." In *Proceedings of the 36th Annual Conference of the Cognitive Science Society*. 2014.
- Heess, Nicolas, Daniel Tarlow, and John Winn. "Learning to pass expectation propagation messages." In *Advances in Neural Information Processing Systems*, pp. 3219-3227. 2013.
- Jitkrittum, Wittawat, Arthur Gretton, Nicolas Heess, S. M. Eslami, Balaji Lakshminarayanan, Dino Sejdinovic, and Zoltán Szabó. "Kernel-based just-in-time learning for passing expectation propagation messages." *arXiv preprint arXiv:1503.02551* (2015).
- Korattikara, Anoop, Vivek Rathod, Kevin Murphy, and Max Welling. "Bayesian dark knowledge." *arXiv preprint arXiv:1506.04416* (2015).