



Computer
Science

CSC696H: Probabilistic Methods in ML

Variational Inference

Prof. Jason Pacheco

Material adapted from: David Blei, NeurIPS 2016 Tutorial

Outline

- Variational Inference
- Stochastic Variational

Outline

- Variational Inference
- Stochastic Variational

Posterior Inference Review

Posterior on latent variable x given data \mathcal{Y} by Bayes' rule:

$$p(x | \mathcal{Y}) = \frac{p(x)p(\mathcal{Y} | x)}{p(\mathcal{Y})}$$

Marginal likelihood given by,

$$p(\mathcal{Y}) = \int p(x)p(\mathcal{Y} | x)dx$$

- Posterior: belief over unknowns, given observed data (knowns)
- Marginal Likelihood: quality of model fit to the observed data

Variational Inference Preview

- Formulate statistical inference as an optimization problem
- Maximize variational lower bound on marginal likelihood

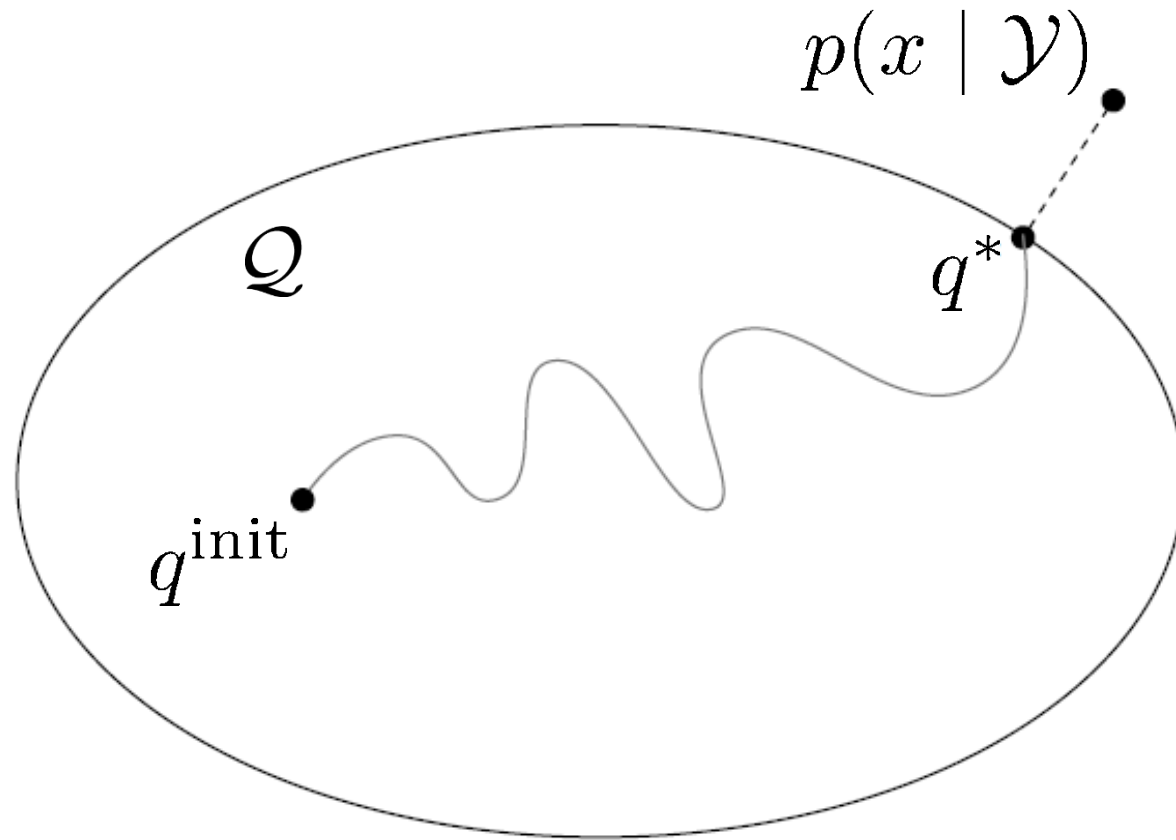
$$\log p(\mathcal{Y}) \geq \max_{q \in \mathcal{Q}} \mathcal{L}(q)$$

- Solution to RHS yields posterior approximation

$$q^* = \arg \max_{q \in \mathcal{Q}} \mathcal{L}(q) \approx p(x | \mathcal{Y})$$

- Constraint set \mathcal{Q} defines tractable family of approximating distributions
- Very often \mathcal{Q} is an *exponential family*

Variational Inference



Expectation Maximization (EM) Lower Bound

Recall EM lower bound of marginal likelihood

$$\log p(\mathcal{Y}) = \log \int p(x)p(\mathcal{Y} | x) dx$$

(Multiply by $q(x)/q(x)=1$)

$$= \log \int p(x)p(\mathcal{Y} | x) \left(\frac{q(x)}{q(x)} \right) dx$$

(Definition of Expected Value)

$$= \log \mathbf{E}_q \left[\frac{p(x)p(\mathcal{Y} | x)}{q(x)} \right]$$

(Jensen's Inequality)

$$\geq \mathbf{E}_q \left[\log \frac{p(x)p(\mathcal{Y} | x)}{q(x)} \right]$$

A Little Information Theory

- The *entropy* is a natural measure of the inherent uncertainty:

$$H(p) = - \int p(x) \log p(x) dx = \mathbb{E}_p[-\log p(x)]$$

- **Interpretation** Difficulty of compression of some random variable

- The *relative entropy* or *Kullback-Leibler (KL) divergence* is a non-negative, but asymmetric, “distance” between a given pair of probability distributions:

$$KL(p||q) = \int p(x) \log \frac{p(x)}{q(x)} dx \quad KL(p||q) \geq 0 \quad KL(p||q) \neq KL(q||p)$$

- The KL divergence equals zero if and only if $p(x) = q(x)$ for all x .

- **Interpretation** The cost of compressing data from distribution $p(x)$ with a code optimized for distribution $q(x)$

EM Lower Bound

$$\mathbf{E}_q \left[\log \frac{p(x)p(\mathcal{Y} | x)}{q(x)} \right] = \mathbf{E}_q \left[\log \frac{p(x)p(\mathcal{Y} | x)}{q(x)} \frac{p(\mathcal{Y})}{p(\mathcal{Y})} \right] \quad (\text{Multiply by 1})$$
$$= \log p(\mathcal{Y}) - \text{KL}(q(x) \| p(x | \mathcal{Y})) \quad (\text{Definition of KL})$$

Bound gap is the Kullback-Leibler divergence $\text{KL}(q \| p)$,

$$\text{KL}(q(x) \| p(x | \mathcal{Y})) = \int q(x) \log \frac{q(x)}{p(x | \mathcal{Y})}$$

Solution to **E-step** is,

$$q^* = \arg \min_q \text{KL}(q(x) \| p(x | \mathcal{Y})) = p(x | \mathcal{Y})$$

This doesn't help us if $p(x | \mathcal{Y})$ is intractable

Variational Lower Bound

Idea Restrict optimization to a set \mathcal{Q} of analytic distributions

$$\log p(\mathcal{Y}) \geq \max_{q \in \mathcal{Q}} \mathcal{L}(q) \equiv \mathbf{E}_q \left[\log \frac{p(x)p(\mathcal{Y} | x)}{q(x)} \right]$$

- If posterior is in set $p(x | \mathcal{Y}) \in \mathcal{Q}$ then exact inference $q(x) = p(x | \mathcal{Y})$
- Otherwise, if $p(x | \mathcal{Y}) \notin \mathcal{Q}$ posterior is closest approximation in KL

$$q^* = \arg \min_{q \in \mathcal{Q}} \text{KL}(q(x) \| p(x | \mathcal{Y}))$$

... and we recover strict lower bound on marginal likelihood with gap

$$\log p(\mathcal{Y}) - \mathcal{L}(q^*) = \text{KL}(q^*(x) \| p(x | \mathcal{Y}))$$

Variational Lower Bound

Two competing terms in variational bound...

$$\begin{aligned}\mathcal{L}(q) &\equiv \mathbb{E}_q \left[\log \frac{p(x)p(\mathcal{Y} | x)}{q(x)} \right] \\ &= \mathbb{E}_q[\log p(x, \mathcal{Y})] - \mathbb{E}_q[\log q(x)] \\ &= \mathbb{E}_q[\log p(x, \mathcal{Y})] + H(q)\end{aligned}$$

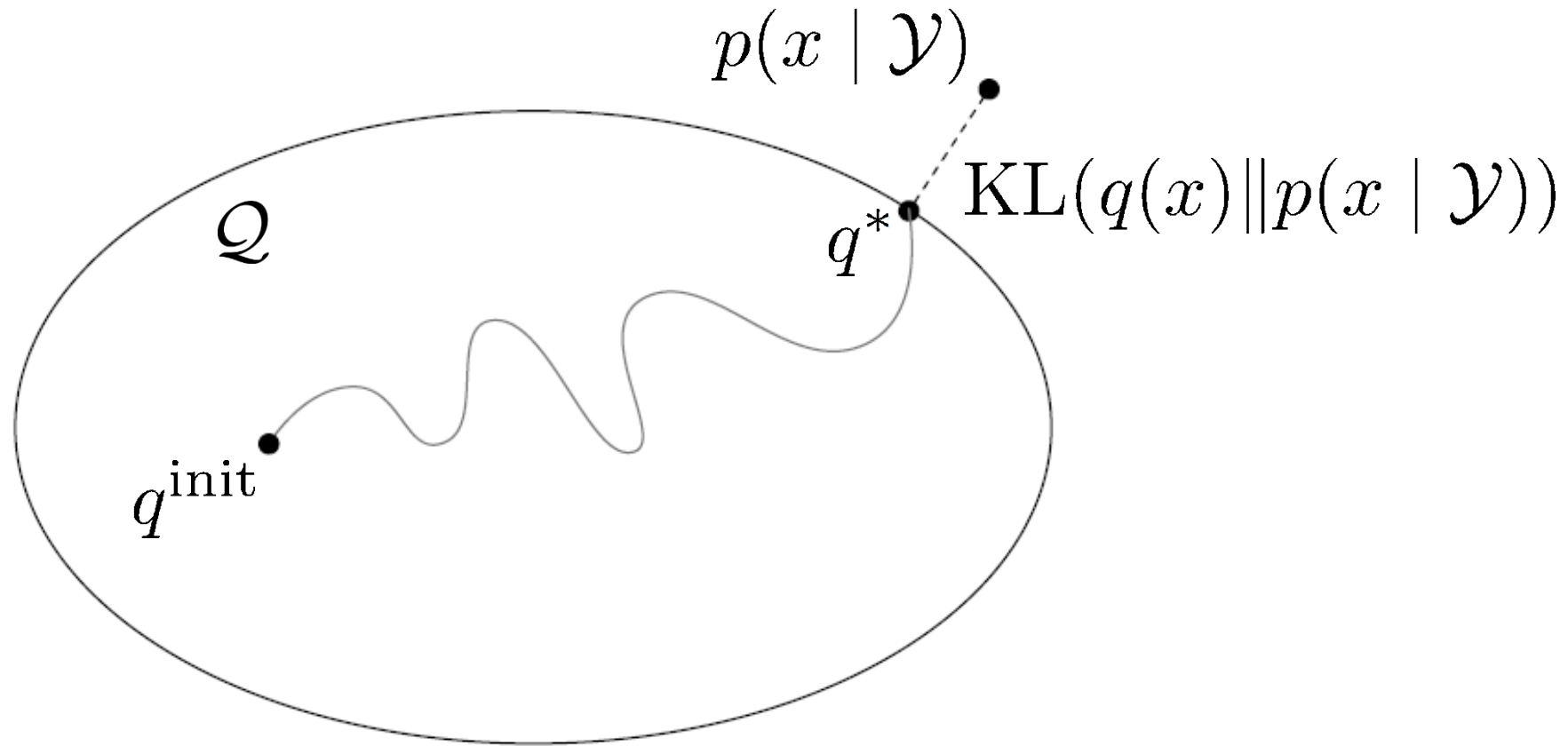
Average (negative) Energy

Encourages $q(x)$ to “agree”
with model $p(x, y)$

Entropy

Encourages $q(x)$ to have
large uncertainty (good for
generalization)

Variational Approximation



Minimize KL between $q(x)$ and posterior $p(x | \mathcal{Y})$.

Relation to EM

- EM is means for approximate *learning*, but we are using it to motivate approximate *inference*
- EM lower bound takes same form as VI lower bound, but with different constraint sets
- Connection with variational inference (VI) is in E-step, which performs inference with fixed parameters

Variational Inference

$$\log p(\mathcal{Y}) \geq \max_{q \in \mathcal{Q}} \mathcal{L}(q) \equiv \mathbb{E}_q[\log p(x, \mathcal{Y})] + H(q)$$

Different sets \mathcal{Q} yield different VI algorithms to optimize bound:

- **Mean Field** Ignore posterior dependencies among variables
- **Loopy BP** *Locally consistent* marginals (exact for tree-structured models)
- **Expectation Propagation (EP)** *Locally consistent moments* (equivalent to Loopy BP for tree-structure exponential families)

Why is it called “variational”?

Differential Calculus

- Typically, we optimize a function $\max_x f(x)$ w.r.t. a **variable** X
- Use standard derivatives/gradients $\nabla_x f(x)$
- Extrema given by zero-gradient conditions $\|\nabla_x f(x)\| = 0$

Calculus of Variations

- Optimize a *functional* (function of a function): $\max_{q(x)} f(q(x))$
- *Functional derivative* characterizes change w.r.t. function $q(x)$
- Extrema given by Euler-Lagrange equation; analogous to zero-gradient condition

In practice, we typically parameterize $q_\mu(x)$ and take standard gradients w.r.t. parameters μ

Summary: Variational Inference

1) Begin with intractable model posterior:

$$p(x | \mathcal{Y}) = \frac{p(x)p(\mathcal{Y} | x)}{p(\mathcal{Y})} \quad \leftarrow \text{Marginal Likelihood}$$

2) Choose a family of approximating distributions \mathcal{Q} that is tractable

3) Maximize variational lower bound on marginal likelihood:

$$\log p(\mathcal{Y}) \geq \max_{q \in \mathcal{Q}} \mathcal{L}(q) \equiv \mathbb{E}_q[\log p(x, \mathcal{Y})] + H(q)$$

4) Maximizer is posterior approximation (in KL divergence)

$$q^* = \arg \max_{q \in \mathcal{Q}} \mathcal{L}(q) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q(x) || p(x | \mathcal{Y}))$$

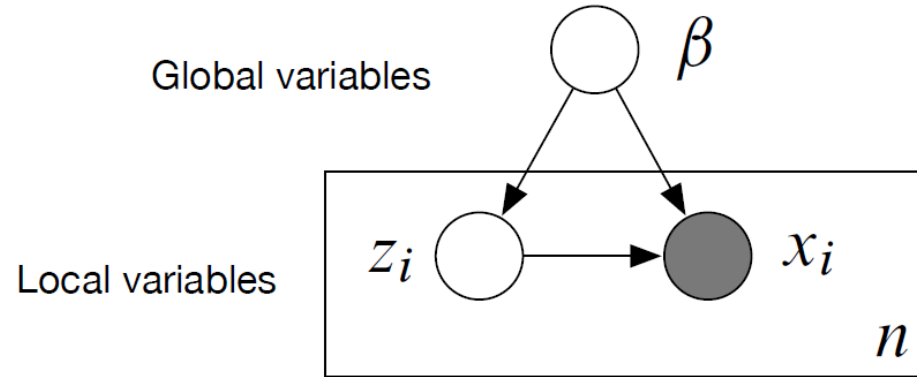
Still need to show...

- a) How to define approximating variational family \mathcal{Q}
- b) How to optimize lower bound

Outline

- Variational Inference
- **Stochastic Variational**

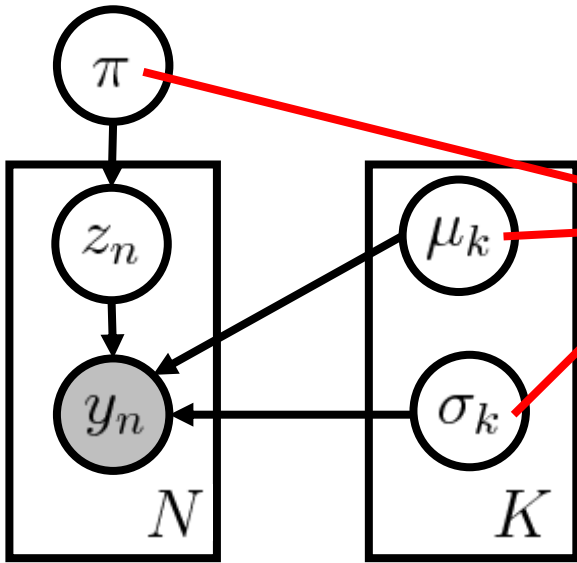
A Generic Class of Directed Models



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

- Bayesian mixture models
- Time series & sequence models (HMMs, Linear dynamical systems)
- Matrix factorization (factor analysis, PCA, CCA)
- Multilevel regression (linear, probit, Poisson)
- Stochastic block models
- Mixed-membership models (Linear discriminant analysis)

Example: Gaussian Mixture Model

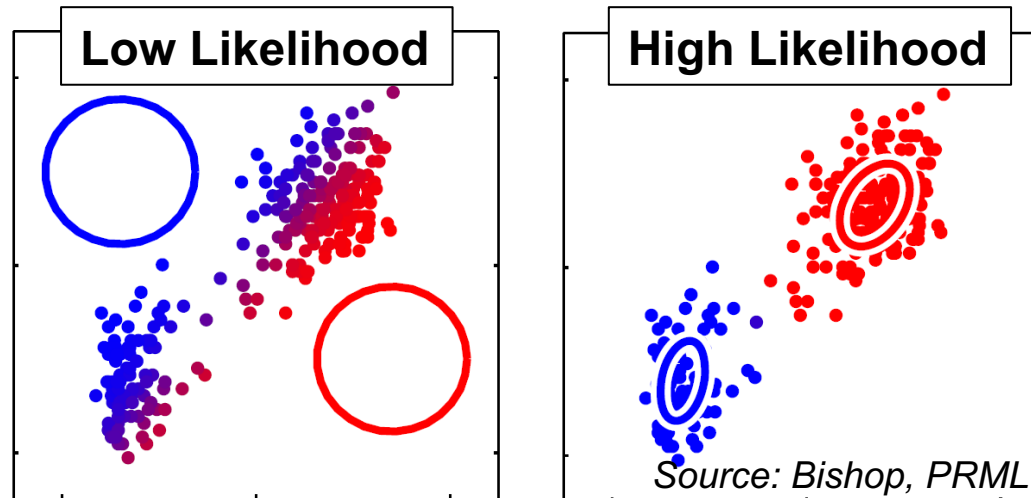


Global variables:

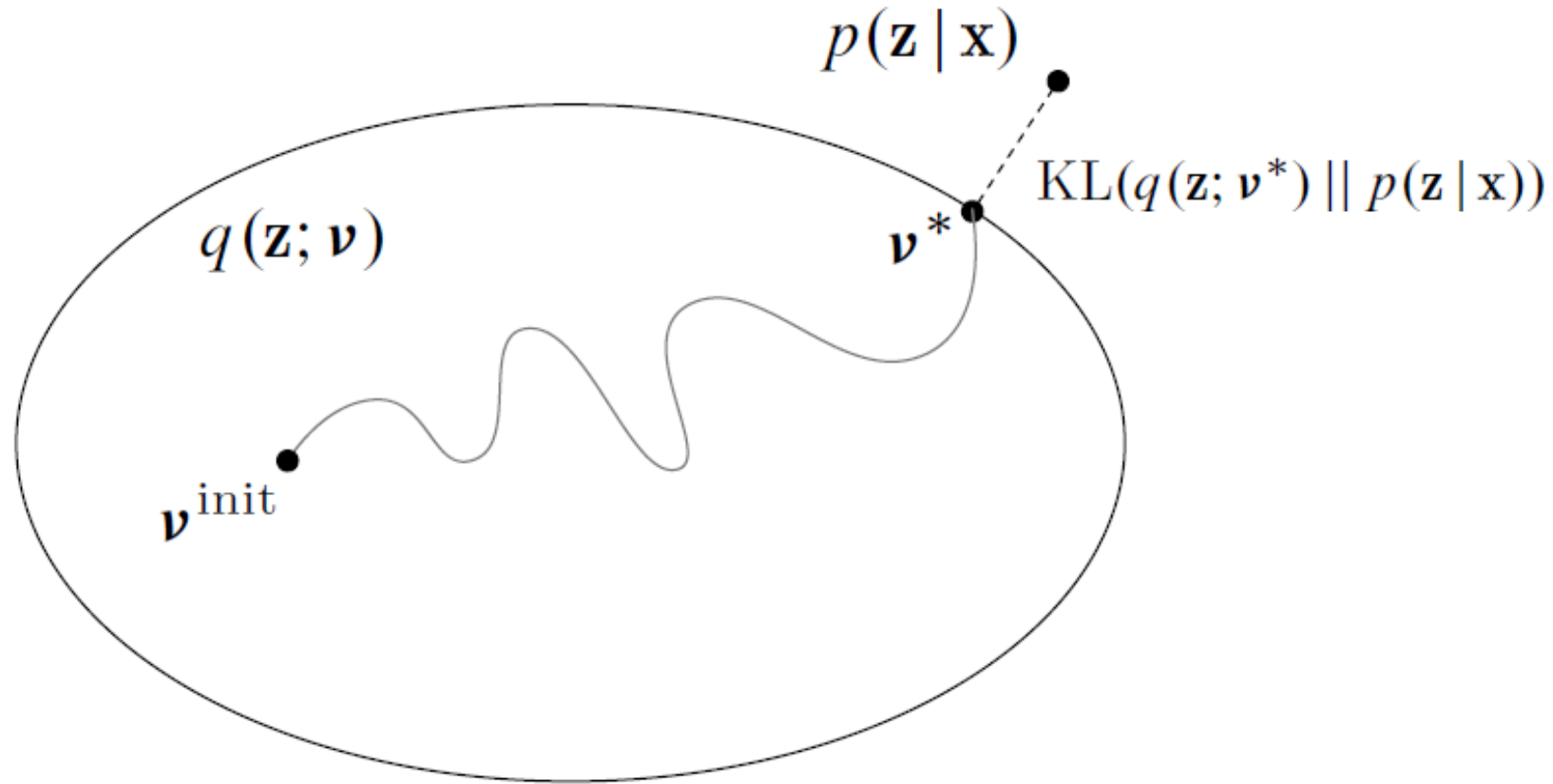
$$\beta = \{\pi, \mu_1, \sigma_1, \dots, \mu_K, \sigma_K\} \quad \mathcal{Y} = \{y_1, \dots, y_N\}$$

Local variables Z control component assignments

GMM



Variational Approximation



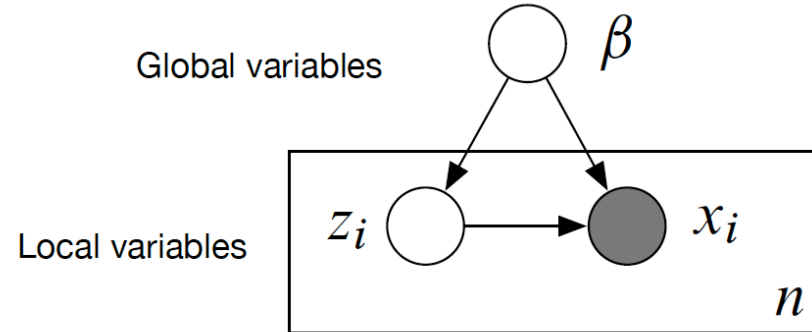
Minimize KL between $q(\beta, \mathbf{z}; \nu)$ and posterior $p(\beta, \mathbf{z} | \mathbf{x})$.

Variational Lower Bound – ELBO

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

- KL is intractable; VI optimizes **evidence lower bound (ELBO)**
 - Lower bounds $\log p(\mathbf{x})$ – marginal likelihood, or *evidence*
 - Maximizing ELBO is equivalent to minimizing KL w.r.t. posterior
- The ELBO trades off two terms
 - The first term prefers $q(\cdot)$ to place mass on the MAP estimate
 - Second term encourages $q(\cdot)$ to be *diffuse* (maximize entropy)
- The ELBO is **non-convex**

Conditionally conjugate models

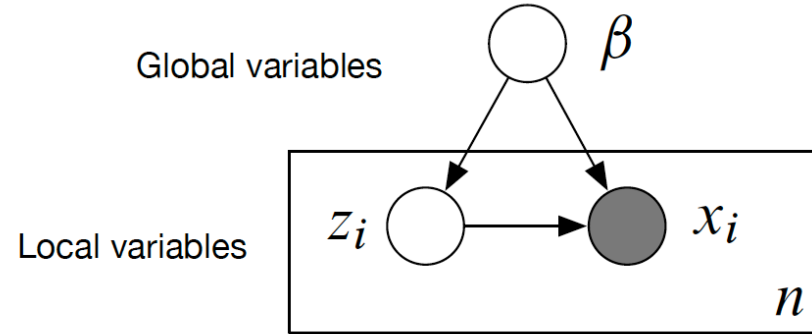


$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

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

Conditionally conjugate 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,

$$p(z_i | \beta, x_i) = \text{expfam}(z_i; \eta_\ell(\beta, x_i))$$
$$p(\beta | \mathbf{z}, \mathbf{x}) = \text{expfam}(\beta; \eta_g(\mathbf{z}, \mathbf{x})),$$

where $\text{expfam}(z; \eta) = h(z) \exp\{\eta^\top z - a(\eta)\}$.

Aside: The exponential family

$$p(x) = h(x) \exp\{\eta^\top t(x) - a(\eta)\}$$

Terminology:

- η the natural parameter
- $t(x)$ the sufficient statistics
- $a(\eta)$ the log normalizer
- $h(x)$ the base density

Aside: The exponential family

$$p(x) = h(x) \exp\{\eta^\top t(x) - a(\eta)\}$$

- The log normalizer is

$$a(\eta) = \log \int \exp\{\eta^\top t(x)\} dx$$

- It ensures the density integrates to one.
- Its gradient calculates the expected sufficient statistics

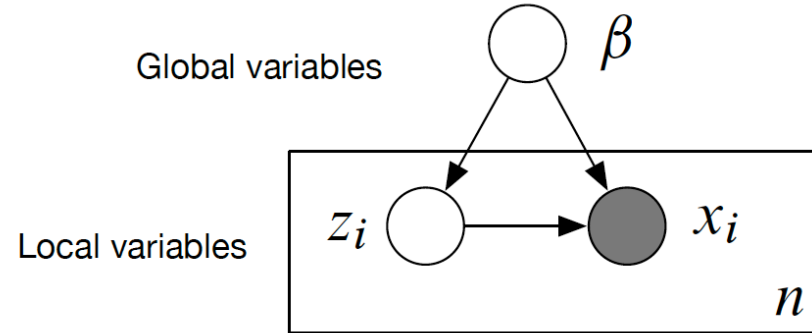
$$\mathbb{E}[t(X)] = \nabla_\eta a(\eta).$$

Aside: The exponential family

$$p(x) = h(x) \exp\{\eta^\top t(x) - a(\eta)\}$$

- Many common distributions are in the exponential family
 - Bernoulli, categorical, Gaussian, Poisson, Beta, Dirichlet, Gamma, etc.
- Outlines the theory around conjugate priors and corresponding posteriors
- Connects closely to variational inference [Wainwright and Jordan, 2008]

Conditionally conjugate models



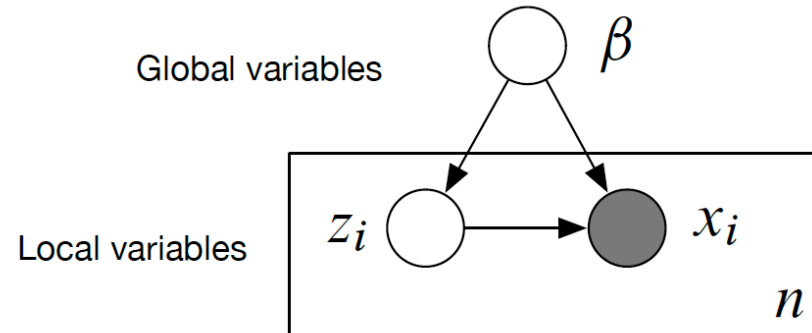
$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

- Each **complete conditional** is in the exponential family.
- 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 $[z_i, x_i]$.

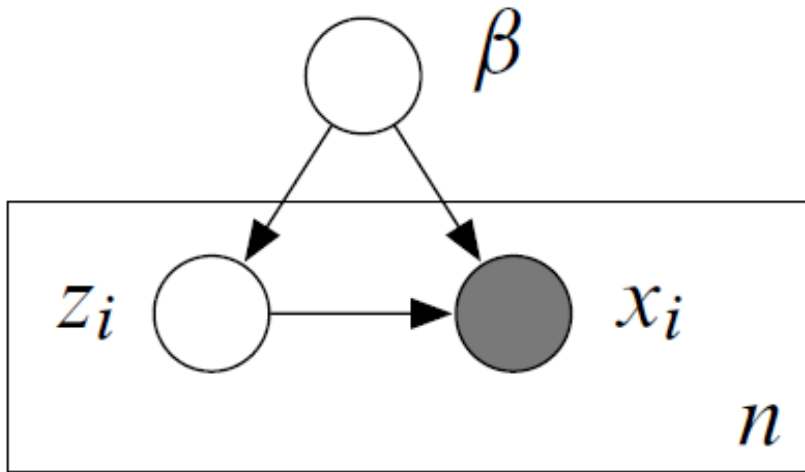
Conditionally conjugate models



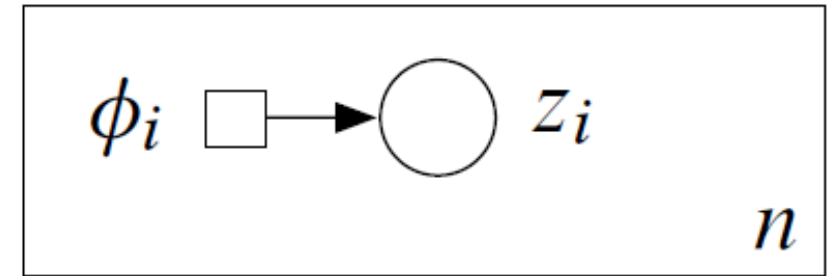
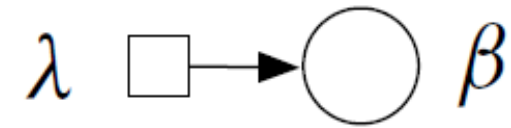
$$p(\beta, \mathbf{z}, \mathbf{x}) = p(\beta) \prod_{i=1}^n p(z_i, x_i | \beta)$$

- 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)

Mean Field for Generic Directed Model



ELBO



PGM of Mean Field Approximation

Recall: mean field family is *fully factorized*

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

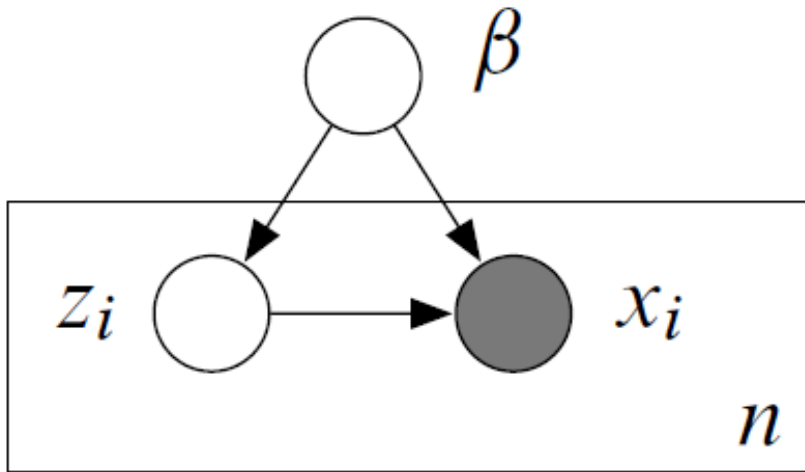
↑ ↑ Variational Parameters

Conditional conjugacy: Each factor is the same expfam as complete conditional

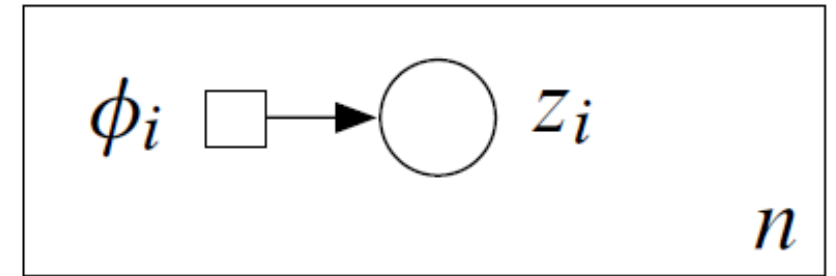
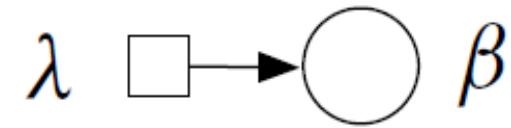
$$p(\beta | \mathbf{z}, \mathbf{x}) = h(\beta) \exp\{\eta_g(\mathbf{z}, \mathbf{x})^\top \beta - a(\eta_g(\mathbf{z}, \mathbf{x}))\}$$

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

Mean Field for Generic Directed Model



ELBO



PGM of Mean Field Approximation

Recall: mean field family is *fully factorized*

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

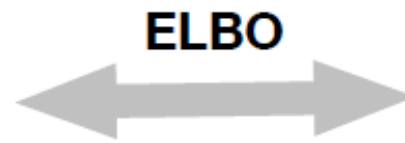
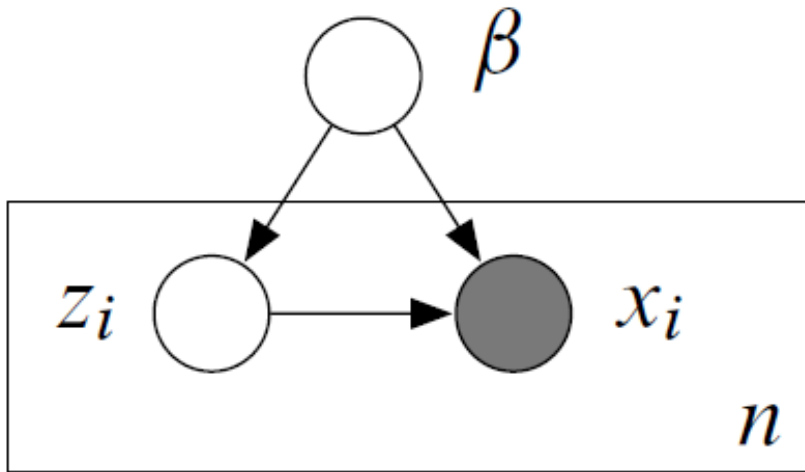
↑ ↑ Variational Parameters

Global parameter ensure conjugacy to (\mathbf{z}, \mathbf{x}) :

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

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

Mean Field for Generic Directed Model



PGM of Mean Field Approximation

Optimize ELBO,

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

Don't forget... entropy decomposes as sum over individual entropies

Traditional VI uses coordinate ascent,

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

Iteratively update each parameter, holding others fixed

- Obvious relationship with Gibbs sampling
- Remember, ELBO is not convex

Coordinate Ascent Mean Field for Generic Model

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*

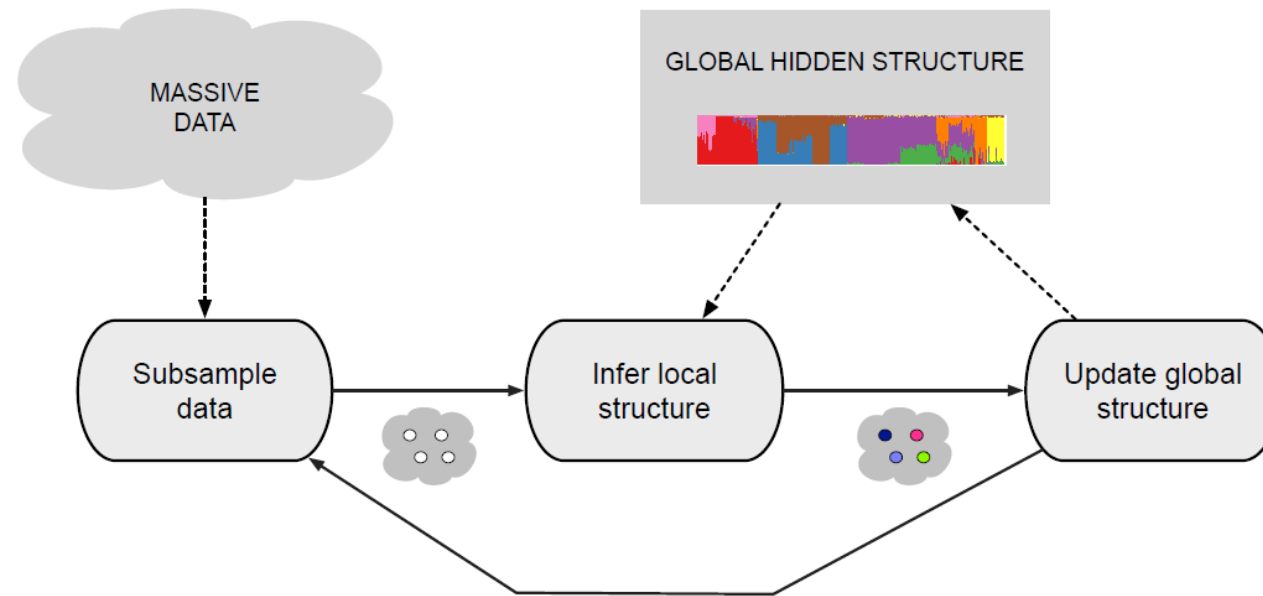
**Need to visit every
data point**



**Need to sum every
data point**



Stochastic (Mean Field) Variational Inference



Classical mean field VI is inefficient for large data

- Do some local computation *for each data point*
- Aggregate computations to re-estimate global structure
- Repeat

Idea visit [random subsets](#) of data to estimate gradient updates on full dataset

Stochastic Gradient Ascent/Descent

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.



- Use cheaper noisy gradient estimates [Robbins and Monro, 1951]
- Guaranteed to converge to local optimum [Bottou, 1996]
- Popular in modern machine learning (e.g. learning deep neural nets)

Stochastic Gradient Ascent/Descent

- Stochastic gradients update:

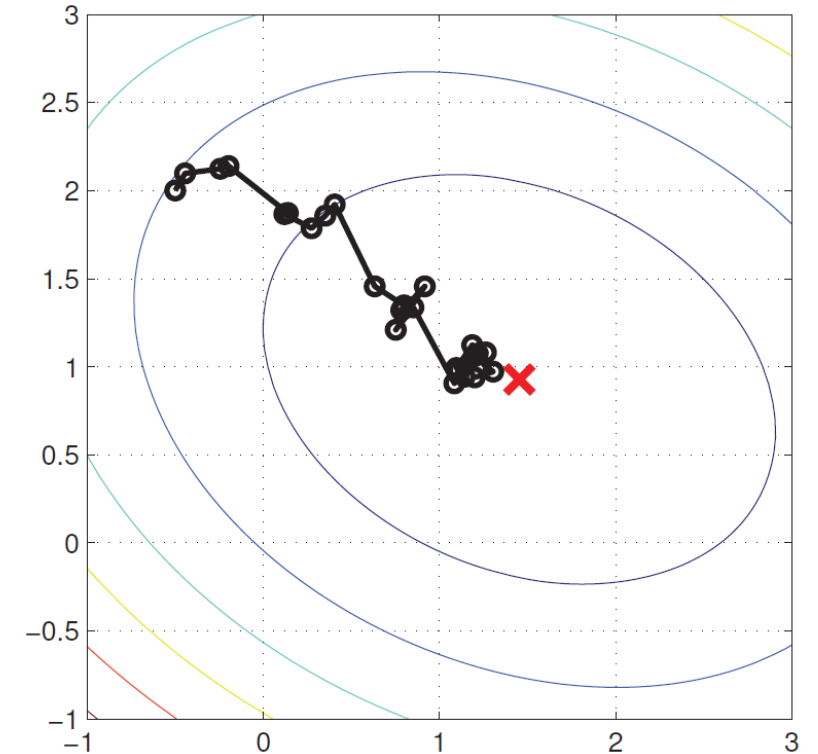
$$\nu_{t+1} = \nu_t + \rho_t \hat{\nabla}_{\nu} \mathcal{L}(\nu_t)$$

- Gradient estimator must be *unbiased*

$$\mathbb{E}[\hat{\nabla}_{\nu} \mathcal{L}(\nu)] = \nabla_{\nu} \mathcal{L}(\nu)$$

- Sequence of step sizes ρ_t must follow **Robbins-Monro conditions**

$$\sum_{t=0}^{\infty} \rho_t = \infty, \quad \sum_{t=0}^{\infty} \rho_t^2 < \infty$$



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

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

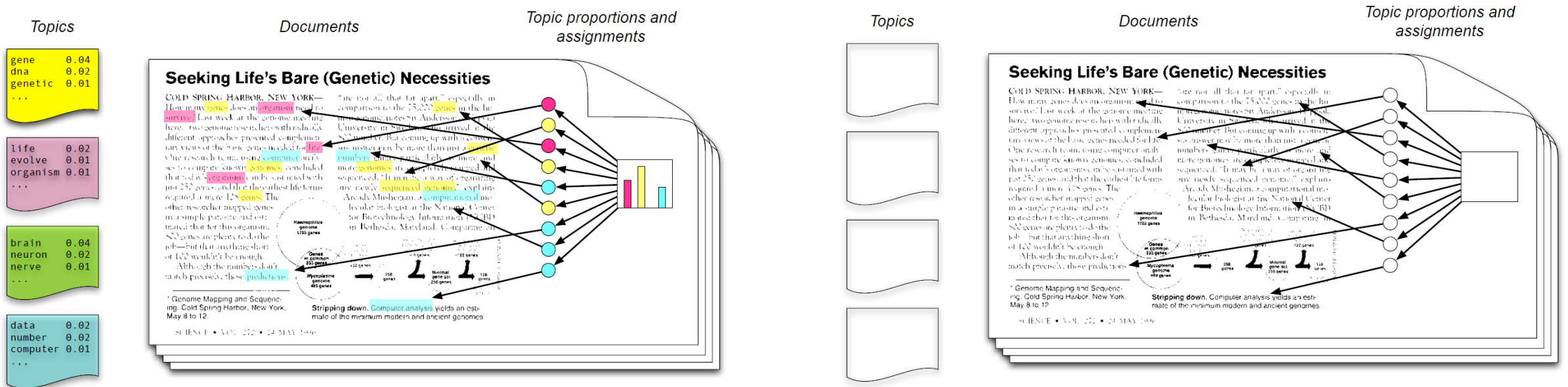
until *forever*

Topic Models



Topic models discover hidden thematic structure in large collections of documents

Topic Models

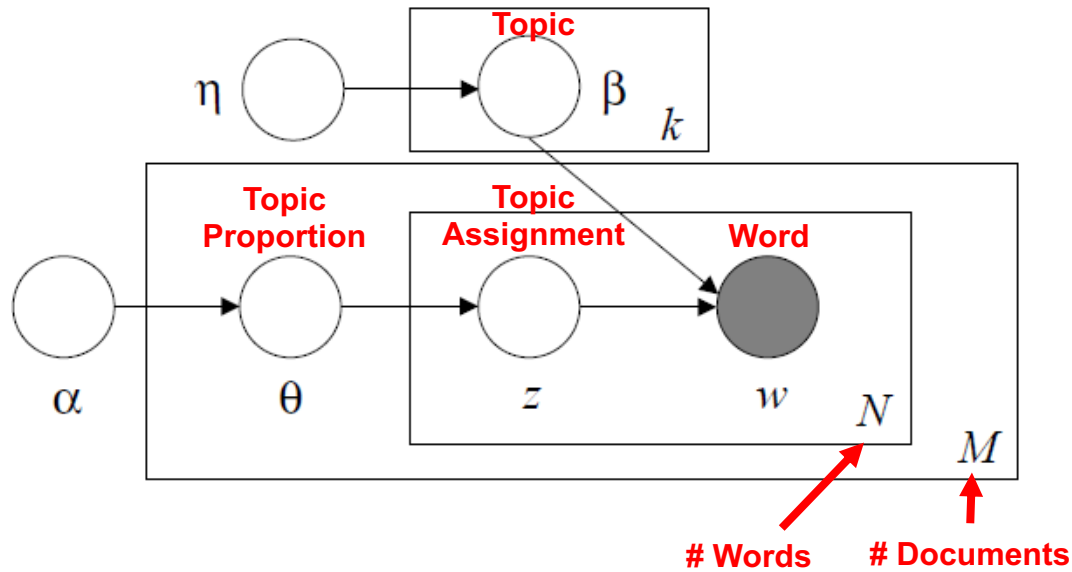


- Each *topic* is a distribution over words (vocabulary)
- Each *document* is a mixture of corpus-wide topics
- Each *word* is drawn from one of the topics (they are distributions)
- But we only observe documents; everything else is hidden (unsupervised learning problem)
- Need to calculate posterior (for millions of documents; billions of latent variables):

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

Topic Models

Latent Dirichlet Allocation (LDA)

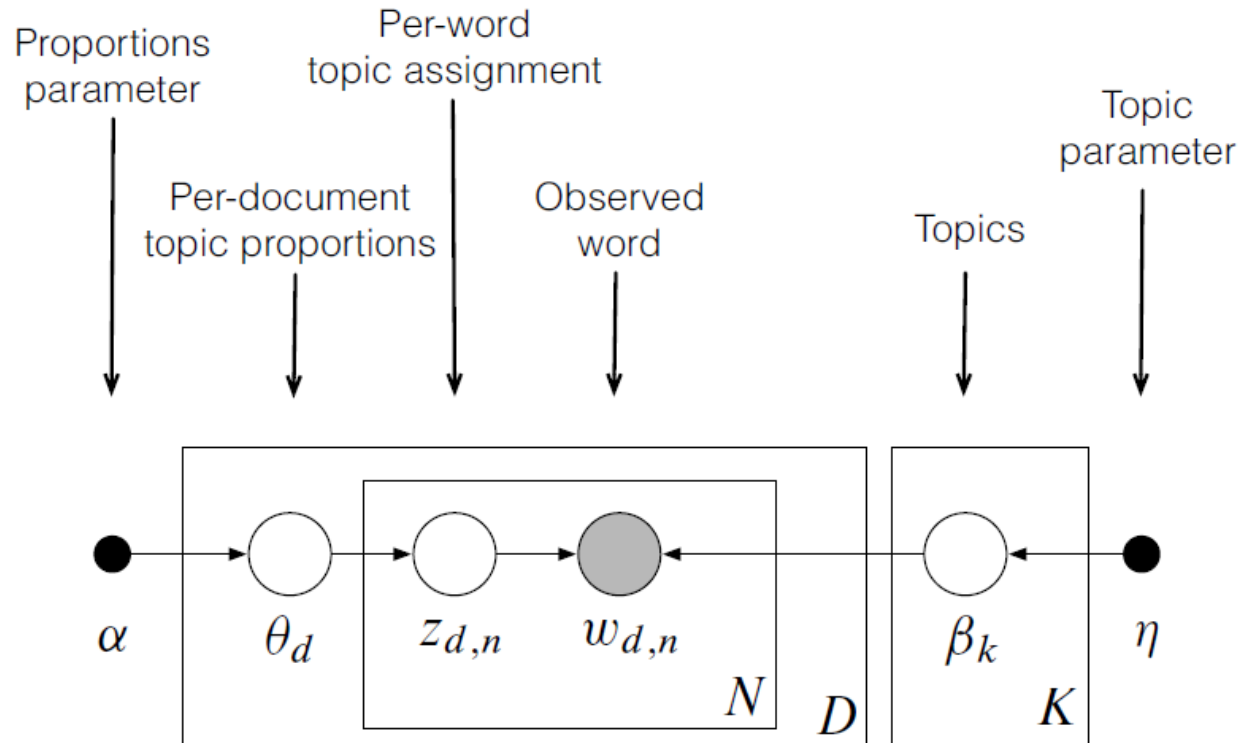


Allows *unsupervised learning* of document corpus via mixture modeling

“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. “Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services,” Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center’s share will be \$200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive \$400,000 each. The Juilliard School, where music and the performing arts are taught, will get \$250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual \$100,000 donation, too.

Example: Latent Dirichlet Allocation



Latent Dirichlet Allocation (LDA):

$$\beta_k \sim \text{Dirichlet}(\eta)$$

$$\theta_d \sim \text{Dirichlet}(\alpha)$$

$$z_{d,n} \mid \theta_d \sim \text{Cat}(\theta_d)$$

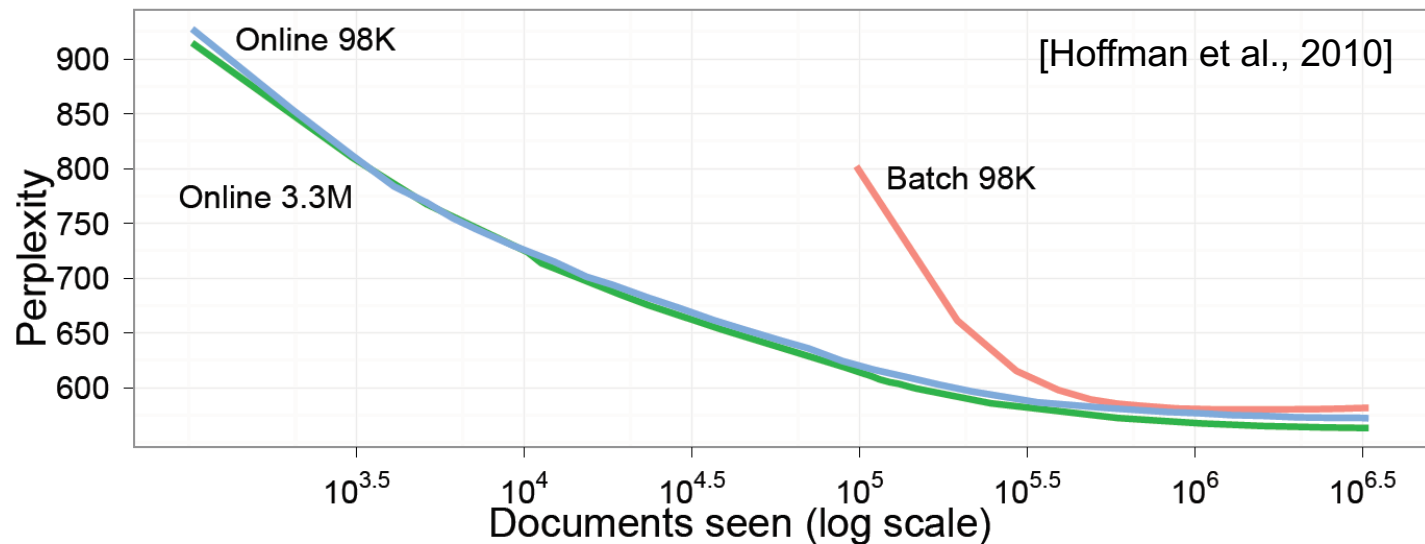
$$w_{d,n} \mid z_{d,n}, \beta \sim \text{Cat}(\beta_{z_{d,n}})$$

- Assumes words are *exchangeable* (“bag-of-words” model)
- Reduces parameters while still yielding useful insights
- Complete conditionals are closed-form (we can do mean field)

Example: Latent Dirichlet Allocation



Topics found in 1.8M articles from the New York Times



- Stochastic VI (online) shows faster learning as compared to standard (batch) updates
- Similar learning rate when dataset increased from 98K to 3.3M documents
- Perplexity measures posterior uncertainty (lower is better)

$$\text{Perplexity} = 2^{H(p)} = 2^{-\sum_x p(x) \log p(x)}$$

Summary: Variational Inference

1) Begin with intractable model posterior:

$$p(x | \mathcal{Y}) = \frac{p(x)p(\mathcal{Y} | x)}{p(\mathcal{Y})} \quad \leftarrow \text{Marginal Likelihood}$$

2) Choose a family of approximating distributions \mathcal{Q} that is tractable

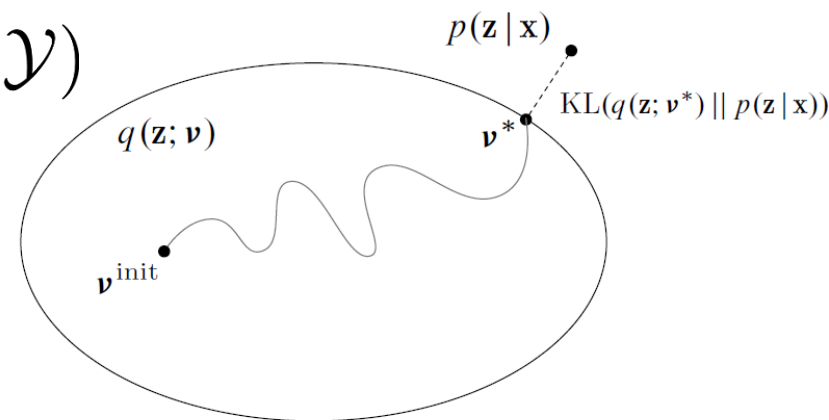
3) Maximize variational lower bound on marginal likelihood:

$$\log p(\mathcal{Y}) \geq \max_{q \in \mathcal{Q}} \mathcal{L}(q) \equiv \mathbb{E}_q[\log p(x, \mathcal{Y})] + H(q)$$

4) Maximizer is posterior approximation (in KL divergence)

$$q^* = \arg \max_{q \in \mathcal{Q}} \mathcal{L}(q) = \arg \min_{q \in \mathcal{Q}} \text{KL}(q(x) || p(x | \mathcal{Y}))$$

Different approximating families \mathcal{Q} lead to different forms of optimizing variational bound



Summary: Mean Field VI

- Mean field family assumes **fully factorized** approximating distribution

$$q(x) = \prod_{s \in \mathcal{V}} q_s(x_s)$$

- Mean field algorithm performs coordinate ascent on lower bound

$$q_s(x_s) \propto \exp \left\{ \mathbb{E}_{q_{\setminus s}} [\log p(x, \mathcal{Y})] \right\}$$

- Coordinate ascent updates require complete conditionals to be conjugate
 - Similar, but stricter, assumption to Gibbs sampling

- MF update takes specific form depending on model $p(\cdot)$, e.g. pairwise MRF:

$$\mu_{sk}^{(i)} \propto \psi_s(k) \exp \left\{ \sum_{t \in \Gamma(s)} \mathbb{E}_{\mu_t^{(i-1)}} [\phi_{st}(k, x_t)] \right\}$$

Summary: Stochastic (Mean Field) VI

- MF coordinate ascent updates require visiting *all data*
 - Doesn't scale to large datasets
- Stochastic VI updates using stochastic gradient ascent
 - Randomly subsample dataset
 - Compute stochastic estimate of full gradient based on subsample
 - Stochastic gradient step on variational parameters (ν here):

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

- Step sizes must decrease over time while satisfying Robbins-Monro conditions

$$\sum_{t=0}^{\infty} \rho_t = \infty, \quad \sum_{t=0}^{\infty} \rho_t^2 < \infty$$

- Often call standard MF “batch” since updates based on full data

