

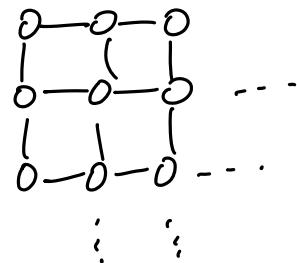
LECTURE 8 - VARIATIONAL INFERENCE

1. Variational lower bounds

Motivation

→ Inference is often intractable!

Ex: • inference in non-tree graphical model
e.g. Ising model



$$P(x) \propto \exp \left\{ \sum_j \theta_j x_j + \sum_{i < j} \theta_{ij} x_i x_j \right\} x_j e^{\{-1\}}$$

$$\mu_j = \bar{\mathcal{E}}_\theta[x_j] \quad \mu_{\hat{j}} = \bar{\mathcal{E}}_\theta[x_{\hat{j}}]$$

- BP does not converge to correct marginals
- exact marginals are hard to compute

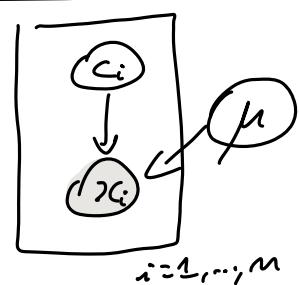
- Complicated posteriors in Bayesian inference

e.g. Bayesian mixture model

$$c_i \sim \text{Unif}\{1, \dots, K\}$$

$$\mu_k \sim \mathcal{N}(0, \sigma^2) \quad (\text{prior})$$

$$x_i \sim \mathcal{N}(\mu_k, 1)$$



Latent variables are $z = (c, \mu)$

posterior : $p(z|x) = \frac{p(c, \mu, x)}{p(x)} \leftarrow \text{OK}$
 $\qquad\qquad\qquad p(x) \leftarrow \text{Not OK!}.$

requires computing intractable integral

$$p(x) = \sum_c p(c) \int \cdots \int p(\mu) \prod_{i=1}^m p(x_i | c_i, \mu) d\mu_1 \cdots d\mu_n$$

$c = (c_1, \dots, c_m)$ takes K^m possible values!

■ Variational lower bound in Exponential families

$$p(x, \theta) = \exp(\langle \theta, \phi(x) \rangle - A(\theta))$$

$$A(\theta) = \log \int \exp(\langle \theta, \phi(x) \rangle) d\nu(x)$$

• recall:

• θ : natural parameter

• $\mu(\theta) = E_\theta[\phi(x)] = \nabla A(\theta)$: mean parameter

• $\theta \mapsto \mu$: inference



• $\mu(\theta)$ achieves maximum in Legendre dual:

$$A(\theta) = \sup_{\mu \in \mathcal{M}} \langle \theta, \mu \rangle - A^*(\mu)$$

$\sup \rightarrow \approx \text{neg-entropy}$

$$\mathcal{M} = \{\mu \in \mathbb{R}^d : \exists q, E_q[\phi(x)] = \mu\}$$

(convex set, polytope if X is discrete)

\rightarrow convex problem in μ



• Lower bound on $A(\theta)$

$$\begin{aligned}
 A(\theta) &= \log \int \exp(\langle \theta, \phi(x) \rangle) d\nu(x) \\
 &= \log \int q(x) \frac{\exp(\langle \theta, \phi(x) \rangle)}{q(x)} d\nu(x) \quad \text{for any } q(x) \\
 &\geq \bar{E}_q \left[\log \frac{\exp(\langle \theta, \phi(x) \rangle)}{q(x)} \right] \quad (\text{Jensen's ineq.}) \\
 &= \bar{E}_q [\langle \theta, \phi(x) \rangle] - \bar{E}_q [\log q] \\
 &= \langle \theta, \mu_q \rangle + H(q) \quad \left(\mu_q := \bar{E}_q [\phi(x)] \right)
 \end{aligned}$$

→ Equality iff $q(x) \propto \exp(\langle \theta, \phi(x) \rangle)$, i.e. $\{q(x) = p(x; \theta)\}$
 $\left\{ \mu_q = \mu(\theta) \right.$
(exact inference)

⇒ "Inference as optimization" / Variational inference (V.I.)

$$A(\theta) = \sup_{q(x)} \langle \theta, \bar{E}_q [\phi(x)] \rangle + H(q)$$

- V.I.:
- consider tractable families of distributions $q \in \mathcal{Q}$ amenable to tractable optimization
 - usually, pick \mathcal{Q} s.t. $\mu_q = \bar{E}_q [\phi(x)]$ is easy to compute for $q \in \mathcal{Q}$
e.g. q factorizes over disconnected graph
 - parameterize by μ_q , express $H(q) = \bar{H}(\mu_q)$
 - not always convex! $dq_1 + (1-d)q_2 \notin \mathcal{Q}$

- Link with KL divergence

Gap in the lower bound:

$$\begin{aligned}
 A(\theta) - \mathbb{E}_q \left[\log \frac{\exp(c\theta, q(x))}{q(x)} \right] &= A(\theta) - \mathbb{E}_q \left[\log \frac{\exp(c\theta, q(x)) - A(\theta)) \exp(A(\theta))}{q(x)} \right] \\
 &= \mathbb{E}_q \left[\log \frac{q(x)}{p(x; \theta)} \right] = \text{KL}(q \parallel p_\theta)
 \end{aligned}$$

Maximize L.B. on $A(\theta)$ \Leftrightarrow Minimize $\text{KL}(q \parallel p_\theta)$

- Variational lower bound for posterior inference

$$\rightarrow p(x) = \int p(x, z) dz$$

↑
 "evidence"

x : observed
 z : latent variables
 parameters in Bayesian inference

- Evidence lower bound (ELBO)

$$\begin{aligned}
 \log p(x) &= \log \int p(x, z) dz \\
 &= \log \mathbb{E}_q \left[\frac{p(x, z)}{q(z)} \right] \\
 &\geq \mathbb{E}_q \left[\log p(x, z) \right] + H(q) =: \text{ELBO}(q)
 \end{aligned}$$

Note: \rightarrow optimal q : $q(z) = p(z|x)$

$\rightarrow \max_q \text{ELBO}(q) \Leftrightarrow \min_q \text{KL}(q \parallel p(\cdot|x))$

\rightarrow similar lower bound used for EM algorithm

[2.] Mean-field Variational Inference

Key idea : use **product distributions** q as tractable families
 i.e. q factorizes on the disconnected graph

■ Ex : Ising model

$$G = (V, E)$$

$$x_i \in \{0, 1\}$$

$$p(x; \theta) \propto \exp \left\{ \sum_{i \in V} \theta_i x_i + \sum_{(i,j) \in E} \theta_{ij} x_i x_j \right\}$$

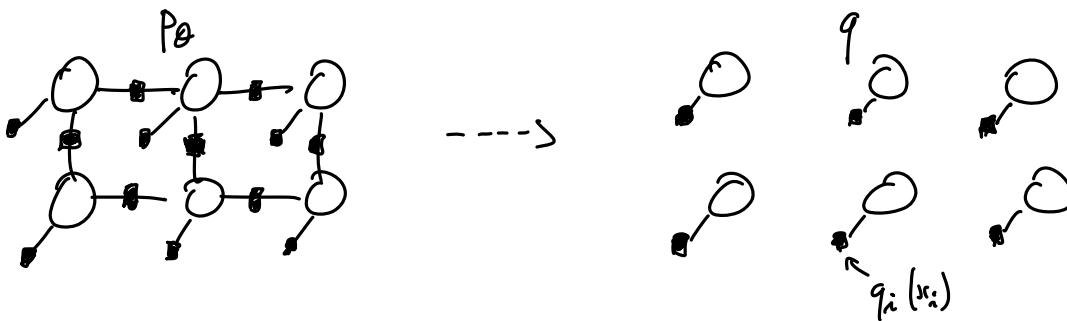
Mean params : $\mu_i = \bar{E}_\theta [x_i]$ $\phi(x) = \begin{pmatrix} (x_i)_i \\ (\delta_{ij} x_j)_{ij} \end{pmatrix}$

$$\mu_{ij} = \bar{E}_\theta [x_i x_j]$$

→ Consider a fully factorized (product) family Q

$$q(x) = \prod_{i \in V} q_i(x_i)$$

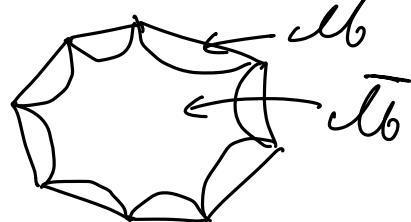
parameterized by $\mu_i := q_i(x_i=1) = \bar{E}_q [x_i]$



Note :- This choice imposes the constraint $\mu_{ij} = \mu_i \mu_j$ on the set of realizable mean params :

$$\bar{\mathcal{M}} = \{ \mu \in \mathcal{M} : \mu_{ij} = \mu_i \mu_j \} \subset \mathcal{M}$$

- \mathcal{B} may be non-convex!



- Variational Lower bound

$$A(\theta) = \sup_q \langle \theta, E_q[\phi(x)] \rangle + H(q)$$

$$\geq \sup_{q \in Q} \langle \theta, E_q[\phi(x)] \rangle + H(q)$$

$$= \max_{(\mu_i)_{i \in V}} \sum_i \theta_i \mu_i + \sum_{(i,j) \in E} \theta_{ij} \mu_i \mu_j + E_q \left[\sum_i \log q_i(x_i) \right]$$

$$= \max_{(\mu_i)_{i \in V}} \sum_i \theta_i \mu_i + \sum_{ij} \theta_{ij} \mu_i \mu_j + \sum_i \bar{H}(\mu_i) =: L(\mu)$$

with $\bar{H}(\mu_i) = \mu_i \log \mu_i + (1-\mu_i) \log(1-\mu_i)$

- Mean Field V.I. algorithm

→ coordinate ascent on $L(\mu)$, i.e., maximize w.r.t. each μ_i separately.

$$\frac{\partial L}{\partial \mu_i} = \theta_i + \sum_{j \in N(i)} \theta_{ij} \mu_j + \log \mu_i + 1 - \log(1-\mu_i) - 1$$

$$\frac{\partial L}{\partial \mu_i} = 0 \Rightarrow \log \frac{\mu_i}{1-\mu_i} = \theta_i + \sum_{j \in N(i)} \theta_{ij} \mu_j$$

$$\frac{\mu_i}{1-\mu_i} = \exp(\theta_i + \sum_j \theta_{ij} \mu_j)$$

$$\mu_i = \sigma(\theta_i + \sum_{j \in N(i)} \theta_{ij} \mu_j)$$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

Algorithm: (MFVI for Ising model)

while $L(\mu)$ has not converged:

for $i \in V$:

$$\mu_i \leftarrow \sigma\left(\theta_i + \sum_{j \in \text{EN}(i)} \theta_{ij} \mu_j\right)$$

Ex: Bayesian Inference

$$p(x) = \int p(x, z) dz \quad \text{where } z \begin{cases} \nearrow \text{hidden variables} \\ \searrow \text{parameters} \end{cases}$$

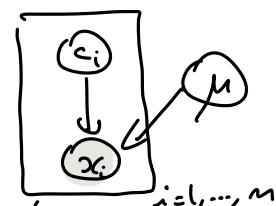
Ex: Bayesian mixture model (BMM): $z = ((c_i)_i, (\mu_k)_k)$

$$c_i \sim \text{Unif}\{1, \dots, K\}$$

\uparrow cluster assignments \uparrow cluster parameters

$$\mu_k \sim \mathcal{N}(0, \sigma^2)$$

$$x_i | c_i = k \sim \mathcal{N}(\mu_k, 1)$$



Factorized variational family:

$$q(z) = \prod_{j=1}^m q_j(z_j)$$

Generic coordinate ascent algorithm:

$$\bar{\text{ELBO}}(q) = \bar{E}_q [\log p(x, z)] + \sum_j \bar{H}(q_j)$$

As a function of q_j :

$$\bar{\text{ELBO}}(q) = E_q [\log p(z_j | z_{-j}, x)] + \bar{H}(q_j) + \text{const}(q_{-j})$$

$$\max \text{ over } q_j \Rightarrow q_j(z_j) \propto \exp\left(\bar{E}_{q_{-j}}[\log p(z_j | z_{-j}, x)]\right)$$

Algorithm: (Coordinate Ascent Variational Inference)

While $\text{ELBO}(q)$ has not converged:

For $j \in \{1, \dots, m\}$:

$$q_j(z_j) \propto \exp\left\{\bar{E}_{z_{-j}}[\log p(z_j | z_{-j}, x)]\right\}$$

- Remark:
- usually, take $q_j(z_j)$ to have same form as $p(z_j | z_{-j}, x)$
 - $p(z_j | z_{-j}, x)$ is often available in closed form when $\rightarrow z_j$ is a discrete latent variable
 $\rightarrow z_j$ is a parameter with conjugate prior to $p(x | z_j)$

Def.: (conjugate prior)

$p(z)$ is a conjugate prior to the likelihood $p(x | z)$ when the posterior $p(z | x) \propto p(x | z)p(z)$ belongs to the same family

Ex:

$$\begin{aligned} p(z) &= \text{Beta}(z, \alpha, \beta) \propto z^{\alpha-1}(1-z)^{\beta-1} \\ p(x | z) &= \text{Ber}(x, z) \propto z^x(1-z)^{1-x} \\ \Rightarrow p(z | x) &= \text{Beta}(z, \alpha+x, \beta+(1-x)) \end{aligned}$$

- other: Dirichlet-Multinomial, Gaussian-Gaussian (fixed covariance)

- **BMM example**: (see [Blei, Kucukelbir, McAuliffe '17])

$$\rightarrow q(z) = q(\mu, c) = \prod_{k=1}^K q(\mu_k; m_k, s_k^2) \prod_{i=1}^m q(c_i; \varphi_i)$$

↗ Gaussian ↗ categorical
 Gaussian categorical

We have:

$$\begin{aligned}
 E[CB](q) &= \sum_{i=1}^m E_q \left\{ \log p(x_i, c_i | \mu) \right\} + \sum_i \bar{H}_i(p_i) \\
 &\quad + \sum_{k=1}^K E_q \left\{ \log p(\mu_k) \right\} + \sum_k \bar{H}_k(m_k, s_k^2)
 \end{aligned}$$

→ max over $\{\varphi_i\}$ is similar to an E-step:

$$q(c_i) \propto \exp \left\{ \bar{E}_\mu [\log p(x_i, c_i | \mu)] \right\}$$

$$\varphi_{ik} \propto \exp \left\{ \bar{E}_\mu [\mu_k] x_i - \frac{1}{2} \bar{E}_\mu [\mu_k^2] \right\}$$

→ max over $\{m_k, s_k\}$ is similar to M-step:

$$m_k = \frac{\sum_i \varphi_{ik} x_i}{\frac{1}{\sigma^2} + \sum_i \varphi_{ik}}, \quad s_k^2 = \frac{1}{\frac{1}{\sigma^2} + \sum_i \varphi_{ik}}$$

Algorithm: ("Variational Bayes" for BMM)

While ELBO has not converged:

(E-step)

For $i = 1 \dots n$

$$\varphi_{ih} \propto \exp \left\{ \bar{\epsilon}_{\mu}[\mu_h] x_i - \frac{1}{2} \bar{\epsilon}_{\mu}[\mu_h^2] \right\}$$

(M-step)

for $h = 1, \dots, K$

$$m_h = \frac{\sum_i \varphi_{ih} x_i}{\frac{1}{\sigma^2} + \sum_i \varphi_{ih}}, \quad s_h^2 = \frac{1}{\frac{1}{\sigma^2} + \sum_i \varphi_{ih}}$$

3. Stochastic VI and Online EM

→ In EM and (Bayesian) VI, the E-step often requires computing expected sufficient statistics across the entire dataset

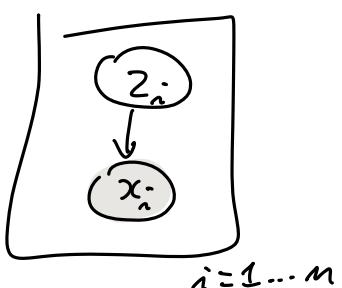
Q: Can we go faster for large datasets?

→ Yes, using **stochastic approximation**!

Ex: **Online EM** (Stochastic VI is similar)

[Neal & Hinton '98], [Cappé & Moulines '08]

[Hoffman et al. '13]



n i.i.d. observations x_i

$$\log p(x; \theta) \geq E_q [\log p(x, z; \theta)] + H(q)$$

$$= \sum_i E_{q_i} [\log p(x_i, z_i; \theta)] + H(q_i)$$

→ In exponential families, the M-step ($\arg \max_{\theta}$) usually depends on "Expected sufficient statistics"

$$S_q = \sum_i E_q [\phi(x_i, z_i)]$$

e.g. for Gaussian Mixtures : $S_q = (S_q^{0,h}, S_q^{1,h})_h$

$$\left\{ \begin{array}{l} S_q^{0,h} = \sum_i E_q [1\{z_i=h\}] \\ S_q^{1,h} = \sum_i E_q [1\{z_i=h\}] x_i \\ S_q^{2,h} = \sum_i E_q [1\{z_i=h\}] x_i x_i^\top \end{array} \right.$$

$$\Rightarrow \text{M-step} : \mu_h = \frac{S_q^{1,h}}{S_q^{0,h}} \quad \Sigma_h = \frac{S_q^{2,h}}{S_q^{0,h}}$$

→ Instead of updating each q_i at every E-step, Online EM updates only one :

(Initialize $S^{(0)} = \sum_i S_i^{(0)}$)

- E-step : pick random $i \in [n]$

$$q_i(z_i) \propto p(z_i | x_i; \theta_t)$$

$$S_i^{(t+1)} = E_{q_i} [\phi(x_i, z_i)]$$

$$S^{(t+1)} = S^{(t)} + (S_i^{(t+1)} - S_i^{(t)})$$

- M-step : update θ_{t+1} using $S^{(t+1)}$

→ Can also use stochastic approximation on infinite data stream $(x_t)_{t \in \mathbb{Z}, \dots}$

$$S^{(t+1)} = (1 - \alpha_t) S^{(t)} + \alpha_t S_t^{(t)}$$