

Approximate Inference Variational Inference

Task: eval. the post. distribution

$$P(Z|X) \quad \text{could be} \quad P(Z|X; \theta)$$

\uparrow Z observed
 \uparrow latent variables / params

Often the case that

- space in which Z lives is very, evaluating all possible (enumerating all) Z 's
- posterior doesn't have a nice analytic form
 - * continuous vars: integrations might not be closed form
 - * discrete

alternative MCMC.

Variational Inference Applied to the Bayesian inference prob.

(section 10.1 of PRML)

Z latent vars & params (set)

X set of observed vars

} might be N iid obs

$$X = \{x_1, \dots, x_N\}$$

$$Z = \{z_1, \dots, z_N\}$$

Probabilistic model

$$P(X, Z)$$

Goal:

Find Posterior Dist $P(Z|X)$ and evidence $P(X)$

$$\ln p(X) = \mathcal{L}(q) + \text{KL}(q \| p)$$

where

$$\mathcal{L}(q) = \int q(z) \cdot \ln \left\{ \frac{P(X, z)}{q(z)} \right\} dz$$

$$\text{KL}(q \| p) = - \int q(z) \ln \left\{ \frac{P(z|X)}{q(z)} \right\} dz$$

multidimensional int's dimensionality

dim is # Z vars and dimensionality of each

$$\begin{aligned}
\ln p(x) &= \mathcal{L}(q) + KL(q||p) \\
&= \int q(z) \ln \left\{ \frac{P(x, z)}{q(z)} \right\} dz - \int q(z) \ln \left\{ \frac{P(z|x)}{q(z)} \right\} dz \\
&= \int q(z) \ln \left\{ \frac{P(z|x) P(x)}{q(z)} \right\} dz - \int q(z) \ln P(z|x) dz + \int q(z) \ln q(z) dz \\
&= \int q(z) \ln P(z|x) dz + \int q(z) \ln P(x) dz - \int q(z) \ln q(z) dz \\
&\quad - \int q(z) \ln P(z|x) dz + \int q(z) \ln q(z) dz \\
&= \int q(z) \ln p(x) dz = \ln p(x) \int q(z) dz = \ln p(x) \cdot 1
\end{aligned}$$

- Differs EM story only in that θ no longer appears

- Game: maximize $\mathcal{L}(q)$ (lower bound on the evidence) - equivalent to minimizing the KL div.

If all q 's possible then $q(z) = p(z|x)$ is minimum but $p(z|x)$ is complicated.

Approach

Restrict the family of distribution $q(z)$ to "simple" distributions, and then to seek the member of this family that most closely approx. $p(x)$

Note:

- choice of $q(z)$ all about tractability
- more complex $q(z)$'s are limited by computation & no overfitting

Choices for $q(z)$

- Parameterized $q(z|\omega)$ is governed by params ω
- Factorized $q(z) = \prod_{i=1}^M q_i(z_i)$ } kind-of "independence" assumption

$$q(z) = \prod_{i=1}^M q_i(z_i)$$

family of approximately dist's

no restrictions on the form of indiv. q_i dist's

Amongst all dist's in this family, which makes $\mathcal{L}(q)$ the largest?

$$\mathcal{L}(q) = \int q(z) \ln \left\{ \frac{P(x, z)}{q(z)} \right\} dz$$

use family $q(z) = \prod_{i=1}^M q_i(z_i)$ where z_i is a subset of the latent vars

call $q_i(z_i) \equiv q_i$

$$\mathcal{L}(q) = \int \left(\prod_i q_i \right) \ln \left\{ \frac{P(x, z)}{\prod_i q_i} \right\} dz$$

$q \in$ factorized family

EM-like objective, find conditions at optimal $\mathcal{L}(q)$ for each q_i .

Split out a single term q_i . Max $\mathcal{L}(z)$ w.r.t. a single term?

$$\mathcal{L}(q) = \int \left(\prod_i q_i \right) \left(\ln P(x, z) - \sum_i \ln q_i \right) dz$$

extract single factor q_j

$$= \int q_j \left(\prod_{i \neq j} q_i \right) \ln P(x, z) dz - \int \left(\prod_i q_i \right) \left(\sum_i \ln q_i \right) dz$$

remember $q_j = q_j(z_j)$, $z_j \subseteq z$

$$= \int q_j \left\{ \int \dots \int \prod_{i \neq j} q_i \ln P(x, z) dz_1 \dots dz_{j-1} dz_{j+1} \dots dz_m \right\} dz_j$$

$$\left[- \int q_j \left(\prod_{i \neq j} q_i \right) \left(\ln q_j + \sum_{i \neq j} \ln q_i \right) dz \right]$$

operating on boxed quantity above

$$= - \int q_j \left(\prod_{i \neq j} q_i \right) \ln q_j dz + \int q_j \left(\prod_{i \neq j} q_i \right) \cdot \left(\sum_{i \neq j} \ln q_i \right) dz$$

$$\int q_j \ln q_j dz_j$$

integrate to 1

const w.r.t. to q_j because
 $\int q_j dz_j \left\{ \prod_{i \neq j} q_i \right\} dz_j$
 const w.r.t. q_j

Mess - clean - it up

$$\mathcal{L}(q) = \int q_j \left\{ \int \dots \int \prod_{i \neq j} q_i \ln P(x, z) dz_1 \dots dz_{j-1} dz_{j+1} \dots dz_m \right\} dz_j - \int q_j \ln q_j dz_j$$

If we define

$$\begin{aligned} \ln \tilde{p}(x, z_j) &= \mathbb{E}_{i \neq j} [\ln p(x, z)] \\ &= \int \dots \int \ln P(x, z) dz_1 \dots dz_{j-1} dz_{j+1} \dots dz_m \end{aligned}$$

$$\mathcal{L}(q) = \int q_j \ln \tilde{p}(x, z_j) dz_j - \int q_j \ln q_j dz_j + \text{const}$$

$$L(q) = \int q_j \ln \tilde{p}(x, z_j) dz_j - \int q_j \ln q_j dz_j + \text{const}$$

Goal: component-wise maximization of $L(q)$,
 in particular, right now max $L(q)$ w.r.t. q_j

Recognize the $L(q)$ is a negative KL div.
 between $\tilde{p}(x, z_j)$ and $q_j(z_j)$

Minimize KL divergence between $\tilde{p}(x, z_j)$ and $q_j(z_j)$.

The optimal $q_j^*(z_j)$ is given by

forms of
 basis
 all, approaches to
 variational
 inference

$$\ln q_j^*(z_j) = \mathbb{E}_{i \sim j} [\ln p(x, z)] + \text{const}$$

- Condition of optimality

- Think about this for a second, the "coordinate/component" max is a function of z_j (subset of params or lat. vars) - The right hand side is the joint dist of all obs & latent variables, but with all latent vars besides z_j integrated out - this leaves a function of same vars z_j on right hand side as well.

- The specific dist. form of $q_j^*(z_j)$ will often emerge from this rule.

- $p(x, z)$ probably has interesting cond. independencies to exploit.

- in expectation - (actually log) many (most of the terms will be absorbed into the constant but not all)

- coupling between approximating factors, say

z_k and z_m $k \neq m$.

- No closed form sol'n in general.

$$\ln q_j^*(z_j) = \mathbb{E}_{i \neq j} [\ln p(x, z)] + \text{const}$$

lacks normalization (const)

- only yields q_j^* up to a multiplicative factor
- one can normalize this distribution by either
 - (usually)
 - inspection (will be become clear)
 - or by explicit normalization

$$q_j^*(z_j) = \frac{\exp(\mathbb{E}_{i \neq j} [\ln p(x, z)])}{\int \exp(\mathbb{E}_{i \neq j} [\ln p(x, z)]) dz_j}$$

- Set of the eqns for all q_i , $i=1..M$ is a set of "consistency" conditions for the max. They are not an explicit soln, not in general closed form, and have to be cycled through until numerical convergence.

State without proof convergence of these interdependent updates is guaranteed because the objective is convex.

Teaching Example

Variational approximation to a full covariance Gaussian (z^D)

Remember, in general $\vec{z} \in \mathbb{R}^2 \Rightarrow \vec{z} \sim \mathcal{N}(\mu, \Sigma)$

$$p(\vec{z}) \neq p(z_1)p(z_2) \text{ unless } \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \text{ and } z_1 \perp z_2$$

Goal: find the independent Gaussian dist (diagonal Gaussian) that best approximates $p(\vec{z})$

Factorization $q(z_1)q(z_2)$