

Probabilistic Graphical Models for Image Analysis - Lecture 4

Stefan Bauer

12 October 2018

Max Planck ETH Center for Learning Systems

Overview

1. Repetition
2. α -Divergence
3. Variational Inference
4. Exponential Families
5. Stochastic Variational Inference
6. Black Box Variational Inference

Repetition

x observations, Z hidden variables, α additional parameters

$$p(z | x, \alpha) = \frac{p(z, x | \alpha)}{\int p(z, x | \alpha)} \quad (1)$$

Idea: Pick family of distributions over latent variables with its own variational parameter

$$q(z | \nu) = \dots?$$

and find variational parameters ν such that q and p are "close".

Mean-field

Assumes that each variable is independent:

$$q(z_1, \dots, z_k) = \prod_{i=1}^k q(z_i) \quad (2)$$

- It does not contain the true posterior since the variables are dependent which can now not be captured by q .
- Offers the possibility to group variables together.
- Use coordinate ascent updates for each $q(z_k)$.
- Here we only specified factorization but **not** the form of the $q(z_k)$.

Using Jensen's inequality to obtain a lower bound

$$\begin{aligned}\log p(x) &= \log \int_Z p(z, x) = \\ &= \log \int_Z p(z, x) \frac{q(z)}{q(z)} \\ &= \log \mathbb{E}_q \left(\frac{p(x, Z)}{q(Z)} \right) \\ &\geq \mathbb{E}_q (\log p(x, Z)) - \mathbb{E}_q (\log q(Z))\end{aligned}$$

Proposal: Choose/design variational family Q such that the expectations are easily computable.

Relation with KL

$$\begin{aligned}\text{KL}[q(z) \parallel p(z | y)] &= \mathbb{E}_q \left[\log \frac{q(Z)}{p(Z | y)} \right] \\ &= \mathbb{E}_q[\log q(Z)] - \mathbb{E}_q[\log p(Z | y)] \\ &= \mathbb{E}_q[\log q(Z)] - \mathbb{E}_q[\log p(Z, y)] + \log p(y) \\ &= -(\mathbb{E}_q[\log p(Z, y)] - \mathbb{E}_q[\log q(Z)]) + \log p(y)\end{aligned}$$

Difference between KL and ELBO is precisely the log normalizer, which does not depend on q and is bounded by the ELBO.

α -Divergence

α -Divergence

$$D_\alpha(p||q) = \frac{\int \alpha p(x) + (1 - \alpha)q(x) - p(x)^\alpha q(x)^{1-\alpha}}{\alpha(1 - \alpha)} dx \quad (3)$$

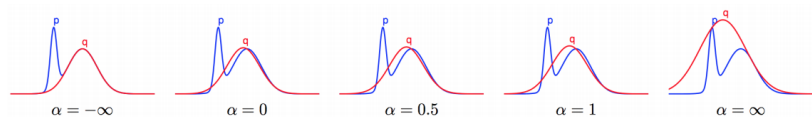
with $\alpha \in (-\infty, \infty)$.

Properties:

- $D_\alpha(p||q)$ is convex with respect to both q and p .
- $D_\alpha(p||q) \geq 0$
- $D_\alpha(p||q) = 0$ when $q = p$
- $\lim_{\alpha \rightarrow 0} D_\alpha(p||q) = \text{KL}(q||p)$
- $\lim_{\alpha \rightarrow 1} D_\alpha(p||q) = \text{KL}(p||q)$

Inclusive v.s Exclusive

Take again Gaussian q approximating two mode Gaussian p from before:



- For $\alpha \geq 1$: **inclusive** since it prefers to stretch across p .
- For $\alpha \leq 0$: **exclusive**, q seeks area of largest total mass, mode seeking.

Expectation Propagation

Instead of minimizing with respect to $\text{KL}(q||p)$ we minimize wrt. $\text{KL}(p||q)$.

Assume q is a member of the exponential family i.e.:

$$q(\mathbf{x}|\eta) = h(\mathbf{x})g(\eta) \exp(\eta^T \mathbf{u}(\mathbf{x}))$$

Then

$$\text{KL}(p||q)(\eta) = -\log g(\eta) - \eta^T \mathbb{E}_{p(z)}[\mathbf{u}(z)] + \text{const.}$$

Minimize KL by setting gradient to zero:

$$-\nabla \log g(\eta) = \mathbb{E}_{p(z)}[\mathbf{u}(z)]$$

from previous slide $-\nabla \log g(\eta) = \mathbb{E}[\mathbf{u}(\mathbf{x})]$, we can then conclude (moment matching)

$$\mathbb{E}_{q(z)}[\mathbf{u}(z)] = \mathbb{E}_{p(z)}[\mathbf{u}(z)]$$

Expectation Propagation - Summary

- The reversed KL is harder to optimized. If the true posterior p factorizes, then we can update single factors iteratively by moment matching.
- Factors are in the exponential family.
- There is no guarantee that the iterations will converge.
- Expectation Propagation aims to preserve the marginals!
- Restriction to exponential family in EP implies: Any product and any division between distributions stays in parametric family and can be done analytically.
- Main applications involve Gaussian Processes and logistic regression (less well suited for GMM).

Variational Inference

Framework

We use variational inference to approximate the posterior distribution

$$\log p(x, \theta) = \text{ELBO}(q, \theta) + \text{KL}(q(z) || p(z|x, \theta)),$$

$$\log p(x, \theta) \geq \mathbb{E}_q[\log p(Z, x)] - \mathbb{E}_q[\log q(Z)]$$

To optimize the lower bound, we can use coordinate ascent!

Problems:

- In each iteration we go over all the data!
- Computing the gradient of the expectations above.

Coordinate Ascent

Given the independence assumption, we can decompose the ELBO as a function of individual $q(z_k)$:

$$\text{ELBO}_j = \int q(z_j) \mathbb{E}_{i \neq j} \log p(z_j | z_{-j}, \mathbf{x}, \theta) dz_j - \int q(z_j) \log q(z_j) dz_j$$

Optimality condition $\frac{d\text{ELBO}}{dq(z_j)} = 0$

Exercise (Lagrange multipliers):

$$q^*(z_j) \propto \exp \mathbb{E}_{-j}[\log p(z_j | Z_{-j}, \mathbf{x})]$$

Coordinate Ascent: Iteratively update each $q(z.)$.

Conditionals

Last slide: $q^*(z_j) \propto \exp \mathbb{E}_{-j}[\log p(z_j|Z_{-j}, \mathbf{x})]$

Assume conditional is in the exponential family i.e.

$$p(z_j|z_{-j}, \mathbf{x}) = h(z_j) \exp(\eta(z_{-j}, \mathbf{x})^\top t(z_j) - a(\eta(z_{-j}, \mathbf{x})))$$

Note: We will see examples in the following Lectures!

Mean-field for exponential family

- Compute log of the conditional

$$\log p(z_j|z_{-j}, \mathbf{x}) = \log(h(z_j)) + \eta(z_{-j}, \mathbf{x})^\top t(z_j) - a(\eta(z_{-j}, \mathbf{x}))$$

- Compute expectation with respect to $q(z_{-j})$

$$\mathbb{E}[\log p(z_j|z_{-j}, \mathbf{x})] = \log(h(z_j)) + \mathbb{E}[\eta(z_{-j}, \mathbf{x})^\top] t(z_j) - \mathbb{E}[a(\eta(z_{-j}, \mathbf{x}))]$$

- Thus $q^*(z_j) \propto h(z_j) \exp(\mathbb{E}[\eta(z_{-j}, \mathbf{x})^\top] t(z_j))$

Conditionals

Note: In the case of an exponential family, the optimal $q(z_j)$ is in the same family as the conditional.

Coordinate Ascent

- Assuming variational parameter ν .

$$q(z_1, \dots, z_k | \nu) = \prod_{i=1}^k q(z_i | \nu_i) \quad (4)$$

- Then each natural variational parameter is set equal to the expectation of the natural conditional parameter given all the other variables and the observations:

$$\nu_j^* = \mathbb{E}[\eta(z_{-j}, \mathbf{x})] \quad (5)$$

Exponential Families

Exponential Families

The exponential family of distributions over x , given parameters η , is defined to be the set of distribution of the form

$$p(x|\eta) = h(x) \exp(\eta^T t(x) - a(\eta)) \quad (6)$$

where:

- η are the so called natural parameters
- $t(x)$ sufficient statistic
- $a(\eta)$ log normalizer
- $\mathbb{E}[t(x)] = \frac{d}{d\eta} a(\eta)$
- Higher moments are next derivatives.

Examples: Bernoulli (last week), Gaussian, Binomial, Multinomial, Poisson, Dirichlet, Beta, Gamma etc.

Conjugacy

Bayesian modeling allows to incorporate priors:

$$\begin{aligned}\eta &\sim F(\cdot|\lambda), \\ x_i &\sim G(\cdot|\eta), \quad \text{for } i \in \{1, \dots, n\}\end{aligned}$$

The posterior distribution of η given the data $x_{1:n}$ is then given by

$$p(\eta|\mathbf{x}, \lambda) \propto F(\eta|\lambda) \prod_{i=1}^n G(x_i|\eta)$$

We say F and G are conjugate if the above posterior belongs to the same functional family as F .

Conjugate prior for the Exponential Family

Exponential family members have a conjugate prior:

$$p(x_i|\eta) = h_l(x) \exp(\eta^T t(x) - a_l(\eta))$$

$$p(\eta|\lambda) = h_c(x) \exp(\lambda_1^T \eta + \lambda_2^T (-a_l(\eta)) - a_c(\lambda))$$

- The natural parameter $\lambda = (\lambda_1, \lambda_2)$ has dimension $\dim(\eta) + 1$
- The sufficient statistics are $(\eta, -a(\eta))$

Exercise: Show the above result.

Conditional conjugacy

Let β be a vector of *global latent variables* (with corresponding global parameters α and z be a vector of local latent variables:

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

Note: For stochastic variational inference, we make an additional assumption:

$$p(z_i, x_i | \beta) = h(z_i, x_i) \exp(\beta^T t(z_i, x_i) - a(\beta))$$

and take an exponential prior on the global variables as the corresponding conjugate prior:

$$p_\alpha(\beta) = h(\beta) \exp(\alpha^T [\beta, -a(\beta)] - a(\alpha)) \text{ with } \alpha = (\alpha_1, \alpha_2).$$

Mean-field for conjugates

Mean-field: $q(z, \beta) = q_\lambda(\beta) \prod_{i=1}^k q_{\varphi_i}(z_i)$

- λ global variational parameter
- φ local variational parameter

Local update $\varphi_i \leftarrow \mathbb{E}_\lambda[\eta_l(\beta, x_i)]$

Global update: $\lambda \leftarrow \mathbb{E}_\varphi[\eta_g(x, z)]$

Note: Coordinate ascent iterates between local and global updates.

Algorithm 1 Mean-field with conjugate family assumption

1: **Input:**

 model p , variational family $q_\varphi(z)$, $q_\lambda(\beta)$

2: **while** ELBO is not converged **do**

4: **for** each data point i **do**

5: Update $\varphi_i \leftarrow \mathbb{E}_\lambda[\eta_l(\beta, \mathbf{x}_i)]$

6: **end for**

7: Update $\lambda \leftarrow \mathbb{E}_\varphi[\eta_g(\mathbf{x}, \mathbf{z})]$

8: **end while**

Stochastic Variational Inference

Gradient Optimization

$$\lambda_{t+1} = \lambda_t + \delta \nabla_{\lambda} f(\lambda_t)$$

or equally

$$\arg \max_{d\lambda} f(\lambda + d\lambda) \text{ st. } \|d\lambda\|^2 \leq \varepsilon \quad (7)$$

Problem: Here it is the euclidean distance, which is not suitable for probability distributions.

Natural Gradient for ELBO:

$$\arg \max_{d\lambda} \text{ELBO}(\lambda + d\lambda) \text{ st. } D_{\text{KL}}^{\text{sym}}(q_{\lambda}, q_{\lambda+d\lambda}) \leq \varepsilon \quad (8)$$

where $D_{\text{KL}}^{\text{sym}}(q, p) = \text{KL}(q||p) + \text{KL}(p||q)$

Gradient Optimization

Riemannian metric $G(\lambda)$ sucht that:

$$d\lambda^\top G(\lambda) d\lambda = D_{\text{KL}}^{\text{sym}}(q_\lambda(\beta), q_{\lambda+d\lambda}(\beta))$$

From information geometry, we know how to calculate the natural gradient:

$$\widehat{\nabla}_\lambda \text{ELBO} = G^{-1}(\lambda) \nabla_\lambda \text{ELBO}$$

where

$$G(\lambda) = \mathbb{E}[(\nabla_\lambda \log q_\lambda(\beta))(\nabla_\lambda \log q_\lambda(\beta))^\top] \quad (9)$$

Gradient Optimization for conjugate models

For our model class, we have:

$$\nabla_{\lambda} \log q_{\lambda}(\beta) = t(\beta) - \mathbb{E}[t(\beta)] \quad (10)$$

Thus

$$G(\lambda) = \nabla_{\lambda}^2 a(\lambda) \quad (11)$$

Recall

$$\nabla_{\lambda} \text{ELBO} = \nabla_{\lambda}^2 a(\lambda) (\mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] - \lambda) \quad (12)$$

then

$$\begin{aligned} \widehat{\nabla}_{\lambda} \text{ELBO} &= G^{-1}(\lambda) \nabla_{\lambda} \text{ELBO} \\ &= G^{-1}(\lambda) \nabla_{\lambda}^2 a(\lambda) (\mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] - \lambda) \\ &= (\mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] - \lambda) \end{aligned}$$

Gradient Step

In each iteration

$$\lambda_t = \lambda_{t-1} + \delta_t \nabla_{\lambda_{t-1}} \text{ELBO}$$

where δ_t is the step size.

Then substituting the above

$$\lambda_t = (\mathbf{1} - \delta_t)\lambda_{t-1} + \delta_t \mathbb{E}[\eta(\mathbf{x}, \mathbf{z})]$$

Algorithm 2 Mean-field with natural gradient

1: **Input:**

 model p , variational family $q_\varphi(z)$, $q_\lambda(\beta)$

2: **while** ELBO is not converged **do**

4: **for** each data point i **do**

5: Update $\varphi_i \leftarrow \mathbb{E}_\lambda[\eta_l(\beta, \mathbf{x}_i)]$

6: **end for**

7: Update $\lambda \leftarrow \lambda + \delta \widehat{\nabla}_\lambda \text{ELBO} = (\mathbf{1} - \delta)\lambda + \delta \mathbb{E}_{q(\varphi)}[\eta_g(\mathbf{x}, z)]$

8: **end while**

Stochastic Optimization

Idea: Maximize a function f using noisy gradients H of that function.

- Noisy gradient H : $\mathbb{E}[H] = \nabla f$
- Step size δ_t
- $x_{t+1} \leftarrow x_t + \delta_t H(x_t)$

Convergence to a local optimum is guaranteed, when:

$$\sum_{t=1}^{\infty} \delta_t = \infty$$
$$\sum_{t=1}^{\infty} \delta_t^2 < \infty$$

Stochastic Gradient Step

Recall

$$\begin{aligned}\widehat{\nabla}_{\lambda} \text{ELBO} &= (\mathbb{E}[\eta(\mathbf{x}, \mathbf{z})] - \lambda) \\ &= \left(\alpha_1 + \sum_{i=1}^n \mathbb{E}_q[t(\mathbf{z}_i, \mathbf{x}_i)], n + \alpha_2 \right) - \lambda\end{aligned}$$

Idea: Construct a noisy natural gradient by sampling

- Sample index $j \sim \text{Uniform}(\mathbf{1}, \dots, n)$
- Rescale

$$\begin{aligned}\widehat{\nabla}_{\lambda} \text{ELBO} &= (\mathbb{E}[\eta(\mathbf{x}_j^{(n)}, \mathbf{z}_j^{(n)})] - \lambda) \\ &= (\alpha_1 + n \mathbb{E}_q[t(\mathbf{z}_j, \mathbf{x}_j)], \mathbf{1} + \alpha_2) - \lambda \\ &=: \widehat{\lambda} - \lambda\end{aligned}$$

Summary gradient step $\lambda_t = (\mathbf{1} - \delta_t) \lambda_{t-1} + \delta_t \widehat{\lambda}$

Algorithm 3 Mean-field with stochastic gradient

1: Input:

model p , variational family $q_\varphi(z)$, $q_\lambda(\beta)$

2: while Stopping criteria is not fulfilled **do**

4: Sample index $j \sim \text{Uniform}(1, \dots, n)$

5: Update $\varphi_i \leftarrow \mathbb{E}_\lambda[\eta_l(\beta, \mathbf{x}_i)]$

6: Compute global parameter estimate $\hat{\lambda} = \mathbb{E}_\varphi[\eta_g(\mathbf{x}_j, \mathbf{z}_j)]$

7: Optimize the global variational parameter $\lambda_{t+1} \leftarrow \lambda_t(1 - \delta_t) + \delta_t \hat{\lambda}$

8: Check step size and update if required!

9: end while

Black Box Variational Inference

Gradient Estimates of the ELBO

$$\text{ELBO} = \mathbb{E}_{q_\nu}[\log p_\theta(z, x)] - \mathbb{E}_q[\log q_\nu(z)]$$

where ν are the parameters of the variational distribution and θ the parameters of the model (as before).

Aim: Maximize the ELBO

Problem: Need unbiased estimates of $\nabla_{\nu, \theta} \text{ELBO}$.

Reparametrization trick^{*}

Simplified notation:

$$\nabla_{\nu} \mathbb{E}_{q_{\nu}} [f_{\nu}(z)]$$

Assume that there exists a fixed reparameterization such that

$$\mathbb{E}_{q_{\nu}} [f_{\nu}(z)] = \mathbb{E}_q [f_{\nu}(g_{\nu}(\varepsilon))]$$

where the expectation on the right does now not depend on ν .
Then

$$\nabla_{\nu} \mathbb{E}_q [f_{\nu}(g_{\nu}(\varepsilon))] = \mathbb{E}_q [\nabla_{\nu} f_{\nu}(g_{\nu}(\varepsilon))]$$

Solution: Obtain unbiased estimates by taking a Monte Carlo estimate of the expectation on the right.

^{*}Will be covered later

Optimizing the ELBO

$$\mathbb{E}_{q(\lambda)}[\log p(z, x) - \log q(z)] =: \mathbb{E}[g(z)]$$

Exercise:

$$\nabla_{\lambda} \text{ELBO} = \nabla_{\lambda} \mathbb{E}[g(z)] = \mathbb{E}[g(z) \nabla \log q(z)] + \mathbb{E}[\nabla g(z)]$$

where $\nabla \log q(z)$ is called the *score function*.

Note: The expectation of the score function is zero for any q
i.e.

$$\mathbb{E}_q[\nabla \log q(z)] = 0$$

Thus, to compute a noisy gradient of the ELBO

- sample from $q(z)$
- evaluate $\nabla \log q(z)$
- evaluate $\log p(x, z)$ and $\log q(z)$

Algorithm 4 Black Box Variational Inference

- 1: **Input:** data x , model $p(x,z)$, variational family $q_\varphi(z)$,
- 2: **while** Stopping criteria is not fulfilled **do**
- 3: Draw L samples $z_l \sim q_\varphi(z)$
- 4: Update variational parameter using the collected samples

$$\varphi \leftarrow \varphi + \delta_t \frac{1}{L} \sum_{l=1}^L \nabla \log q(z_l) (\log p(x, z_l) - \log q(z_l))$$

- 5: Check step size and update if required!
 - 6: **end while**
-

Note: Active research area (problem) is the reduction of the variance of the noisy gradient estimator.

Summary

What we have seen

- Expectation Maximization
- Introduced KL to understand convergence of EM.
- α - Divergence as a more general concept, including Expectation Propagation
- Scalable Stochastic Variational Inference for conjugates in the exponential family
- black box variational inference to overcome exponential family assumption

What is not covered

- Problem for black box is the variance of the noisy gradient
- Extensions black box α -Divergence
- Online inference algorithms

Plan for next week

- Originally planned: Bayesian Non-parametrics
- Most likely: State Space Models

Questions?