

Variational Inference

1

References

These notes are based on the following papers:

- Blei, D. M., Kucukelbir, A. and McAuliffe, J. D. (2017). Variational Inference: A Review for Statisticians. Journal of the American Statistical Association, 112:518,859-877.
- Jordan, M. I., Ghahramani, Z., Jaakkola, T. S. and Saul, L. K. (1999). An Introduction to Variational Methods for Graphical Models. Machine Learning, 37, 183-233.

2

Introduction

MCMC

- Theoretical guarantees (asymptotically) of sampling from the target density
- Computationally intensive but conceptually simple
- Handles multi-modal posterior distributions

Variational Inference

- No theoretical guarantees
- Good for big data and complex models
- Faster than MCMC but requires derivation of variational updates
- Can have problems with multi-modal posteriors

3

Introduction

- Variational methods based on **calculus of variations**
- Complex problem turned to a simpler one by decoupling degrees of freedom in the original problem
- Decoupling done by extending the original problem with additional variational parameters

4

Intuition

We first develop some intuition about variational methods using a simple example.

Write the log function as:

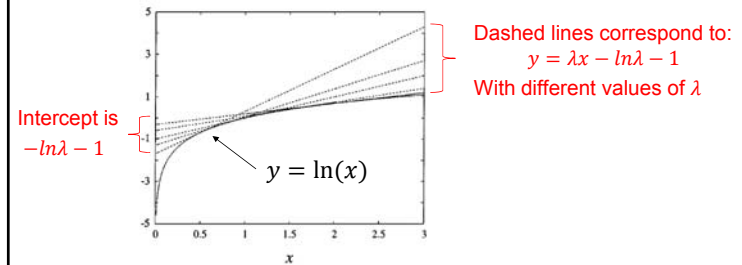
$$\ln(x) = \min_{\lambda} \{\lambda x - \ln \lambda - 1\}$$

Note:

- λ is the variational parameter
- For each value of x , we need to compute the minimization of λ .

5

Intuition



- Varying λ produces a series of upper bounds:

$$\ln(x) \leq \lambda x - \ln \lambda - 1$$
- Minimizing λ produces the exact value for $\ln(x)$
- Note: $\ln(x)$ is a concave function

6

Intuition

Why did we write $\ln(x) = \min_{\lambda} \{\lambda x - \ln \lambda - 1\}$?

- Comes from **convex duality**: a concave function $f(x)$ can be represented by a dual function as

$$f(x) = \min_{\lambda} \{\lambda^T x - f^*(\lambda)\}$$

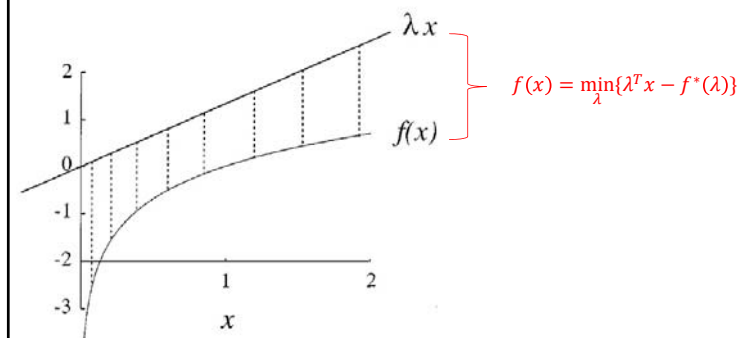
Where

$$f^*(\lambda) = \min_x \{\lambda^T x - f(x)\}$$

- Applies to convex functions as well but you get a lower bound

7

Intuition



8

Variational Inference

- Let X be a set of observed variables (e.g. evidence variables)
- Let Z be a set of latent variables
- Given inference query $P(Z|X) = \frac{P(X,Z)}{P(X)}$
- Need to compute $P(X) = \sum_Z P(X, Z)$ if Z is discrete or $\int P(X, Z) dZ$ if continuous
- The denominator is typically very expensive to compute

9

The ELBO

- Goal: choose a density $q(z) \in Q$ which is the closest approximation to $p(z|x)$
- Here Q is a family of densities over the latent variables
- Need to solve the following optimization problem:

$$q^*(z) = \operatorname{argmin}_{q(z) \in Q} KL(q(z) || p(z|x))$$

10

The ELBO

$$\begin{aligned} KL(q(z)||p(z|x)) &= \int q(z) \log \frac{q(z)}{p(z|x)} dz = \int q(z) \log \frac{q(z)p(x)}{p(x,z)} dz \\ &= \int q(z) [\log q(z) + \log p(x) - \log p(x,z)] dz \\ &= \int [q(z) \log q(z) + q(z) \log p(x) - q(z) \log p(x,z)] dz \\ &= \int q(z) \log q(z) dz + \int q(z) \log p(x) dz - \int q(z) \log p(x,z) dz \\ &= E_{q(z)}[\log q(z)] + E_{q(z)}[\log p(x)] - E_{q(z)}[\log p(x,z)] \end{aligned}$$

Doesn't involve $q(z)$ so it can be taken out of the expectation

11

The ELBO

$$\begin{aligned} KL(q(z)||p(z|x)) &= E_{q(z)}[\log q(z)] - E_{q(z)}[\log p(x,z)] + \log p(x) \end{aligned}$$

Remember that this is very hard to compute because $p(x) = \int p(x,z) dz$

Rewrite as:

$$\begin{aligned} \log p(x) &= KL(q(z)||p(z|x)) + E_{q(z)}[\log p(x,z)] - E_{q(z)}[\log q(z)] \\ &= KL(q(z)||p(z|x)) + ELBO(q) \end{aligned}$$

Where

$$ELBO(q) = E_{q(z)}[\log p(x,z)] - E_{q(z)}[\log q(z)]$$

12

The ELBO

- Because KL divergence is ≥ 0

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

$$\Rightarrow \log p(\mathbf{x}) \geq ELBO(q)$$

- $p(\mathbf{x})$ is the probability of the evidence, hence this is an evidence lower bound (ELBO)
- Instead of minimizing the KL divergence, we maximize the ELBO

$$ELBO(q) = E_{q(\mathbf{z})}[\log p(\mathbf{x}, \mathbf{z})] - E_{q(\mathbf{z})}[\log q(\mathbf{z})]$$

13

The ELBO

Another point about $ELBO(q)$

$$= E_{q(\mathbf{z})}[\log p(\mathbf{x}, \mathbf{z})] - E_{q(\mathbf{z})}[\log q(\mathbf{z})]$$

$$= E_{q(\mathbf{z})}[\log p(\mathbf{x}|\mathbf{z})] + E_{q(\mathbf{z})}[\log p(\mathbf{z})] - E_{q(\mathbf{z})}[\log q(\mathbf{z})]$$

$$= E_{q(\mathbf{z})}[\log p(\mathbf{x}|\mathbf{z})] + \int q(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} - \int q(\mathbf{z}) \log q(\mathbf{z}) d\mathbf{z}$$

$$= E_{q(\mathbf{z})}[\log p(\mathbf{x}|\mathbf{z})] + \int q(\mathbf{z}) \log \frac{p(\mathbf{z})}{q(\mathbf{z})} d\mathbf{z}$$

$$= E_{q(\mathbf{z})}[\log p(\mathbf{x}|\mathbf{z})] + KL(q(\mathbf{z})||p(\mathbf{z}))$$

This is an expected likelihood.

Places mass of $q(\mathbf{z})$ on configurations of the latent variables \mathbf{z} that explain the observed data \mathbf{x} .

This makes $q(\mathbf{z})$ resemble the prior $p(\mathbf{z})$

14

Mean Field

In the mean-field variational family, the latent variables are:

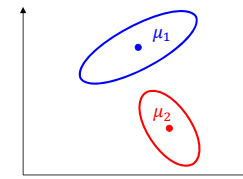
- Mutually independent
- Each has its own factor (and parameters) in the variational family

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

15

Bayesian Mixture of Gaussians

- A Gaussian mixture model assumes there are K Gaussians that generate the data, each with its own mean μ_k and variance σ_k



- We can make this a Bayesian model by putting a prior on the means of the K Gaussians

16

Bayesian Mixture of Gaussians

The Generative model:

- $\mu_k \sim N(\mathbf{0}, \sigma^2)$ for $k = 1, \dots, K$
- $c_i \sim \text{Categorical}(\frac{1}{K}, \dots, \frac{1}{K})$ for $i = 1, \dots, n$
- $x_i | c_i, \mu \sim N(c_i^T \mu, 1)$ for $i = 1, \dots, n$

The joint density is:

$$p(\mu, c, x) = p(\mu) \prod_{i=1}^n p(c_i) p(x_i | c_i, \mu)$$

17

Bayesian Mixture of Gaussians

Prior for the mean of each mixture component

σ is a given hyperparameter

- $\mu_k \sim N(\mathbf{0}, \sigma^2)$ for $k = 1, \dots, K$
 - $c_i \sim \text{Categorical}(\frac{1}{K}, \dots, \frac{1}{K})$
 - $x_i | c_i, \mu \sim N(c_i^T \mu, 1)$
- Cluster assignment that takes values $1, \dots, K$. Encoded as an indicator K-vector, with 0s everywhere except for a 1 in the position corresponding to the cluster c_i belongs to.

The joint density is:

$$\mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_k \end{bmatrix}$$

$$p(\mu, c, x) = p(\mu) \prod_{i=1}^n p(c_i) p(x_i | c_i, \mu)$$

18

Bayesian Mixture of Gaussians

- Computing the evidence requires:

$$p(x) = \int p(\mu) \prod_{i=1}^n \sum_{c_i} p(c_i) p(x_i | c_i, \mu) d\mu$$

- This K-dimensional integral takes $O(K^n)$ time to compute.
- Variational Inference to the rescue!

19

Bayesian Mixture of Gaussians

Joint density:

$$p(\mu, c, x) = p(\mu) \prod_{i=1}^n p(c_i) p(x_i | c_i, \mu)$$

Mean-field variational family:

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

ELBO:

$$E_{q(\mu, c)}[\log p(\mu, c, x)] - E_{q(\mu, c)}[\log q(\mu, c)]$$

20

Bayesian Mixture of Gaussians

$$\begin{aligned}
 ELBO(\mathbf{m}, \mathbf{s}^2, \boldsymbol{\varphi}) &= E_{q(\boldsymbol{\mu}, \mathbf{c})}[\log p(\boldsymbol{\mu}, \mathbf{c}, \mathbf{x})] - E_{q(\boldsymbol{\mu}, \mathbf{c})}[\log q(\boldsymbol{\mu}, \mathbf{c})] \\
 &= E_{q(\boldsymbol{\mu}, \mathbf{c})} \left[\log \left(p(\boldsymbol{\mu}) \prod_{i=1}^n p(c_i) p(x_i | c_i, \boldsymbol{\mu}) \right) \right] \\
 &\quad - E_{q(\boldsymbol{\mu}, \mathbf{c})} [\log (\prod_{k=1}^K q(\mu_k; m_k, s_k^2) \prod_{i=1}^n q(c_i; \varphi_i))] \\
 &= \sum_{k=1}^K E_{q(\boldsymbol{\mu}, \mathbf{c})} [\log p(\mu_k); m_k, s_k^2] \\
 &\quad + \sum_{i=1}^n (E_{q(\boldsymbol{\mu}, \mathbf{c})} [\log p(c_i); \varphi_i] + E_{q(\boldsymbol{\mu}, \mathbf{c})} [\log p(x_i; c_i, \boldsymbol{\mu}); \varphi_i, \mathbf{m}, \mathbf{s}^2]) \\
 &\quad - \sum_{i=1}^n E_{q(\boldsymbol{\mu}, \mathbf{c})} [\log q(c_i; \varphi_i)] - \sum_{k=1}^K E_{q(\boldsymbol{\mu}, \mathbf{c})} [\log q(\mu_k; m_k, s_k^2)]
 \end{aligned}$$

21

Bayesian Mixture of Gaussians

- With the ELBO, we now need to optimize the variational parameters
- One way to do this is coordinate ascent variational inference (CAVI) (Bishop 2006)
- Works by optimizing each parameter while keeping the others fixed
- Need to come up with updates for φ_{ik}, m_k, s_k
- Done iteratively until ELBO converges

22

Bayesian Mixture of Gaussians

For CAVI:

- Uses the **complete conditional** of z_j i.e. $p(z_j | \mathbf{z}_{-j}, \mathbf{x})$
- Optimization uses the following:

$$q_j^*(z_j) \propto \exp\{E_{-j}[\log p(z_j, \mathbf{z}_{-j}, \mathbf{x})]\}$$

This expectation is over all the other variational factors being fixed i.e. $\prod_{l \neq j} q_l(z_l)$

- We won't go through the derivation. See (Bishop 2006) for details

23

Bayesian Mixture of Gaussians

For Bayesian Mixture of Gaussians:

1. Compute update for mixture assignments.
2. Compute update for mixture component means and variances.

24

Bayesian Mixture of Gaussians

1. Computing update for φ_{ik}

$$q^*(c_i; \varphi_i) \propto \exp\{\underbrace{\log p(c_i)} + \underbrace{E[\log P(x_i|c_i, \boldsymbol{\mu}); \mathbf{m}, \mathbf{s}^2]}\}$$

$$\log p(c_i) = -1/K$$

$$\varphi_{ik} \propto \exp\{E[\mu_k; m_k, s_k^2]x_i - E[\mu_k^2; m_k, s_k^2]/2\}$$

(derivation left as an exercise)

Note: Expectation will be over $\prod_{k=1}^K q(\mu_k; m_k, s_k^2) \prod_{j \neq i} q(c_j; \varphi_j)$

25

Bayesian Mixture of Gaussians

2. Computing update for m_k, s_k

$$q(\mu_k; m_k, s_k^2) \propto \exp\left\{\log p(\mu_k) + \sum_{i=1}^n E[\log p(x_i|c_i, \boldsymbol{\mu}); \varphi_i, \mathbf{m}_{-k}, \mathbf{s}_{-k}^2]\right\}$$

Note: Expectation will be over $\prod_{i \neq k} q(\mu_i; m_i, s_i^2) \prod_{i=1}^n q(c_i; \varphi_i)$

This leads to update equations: (derivation left as an exercise)

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

$$s_k^2 = \frac{1}{\frac{1}{\sigma^2} + \sum_i \varphi_{ik}}$$

26

Concluding remarks

- It takes some work to derive variational inference equations
- Generic variational updates have been derived for special cases e.g. when complete conditional is in the exponential family

27