

Variational Inference

Sargur Srihari

srihari@cedar.buffalo.edu

Plan of Discussion

- Functionals
- Calculus of Variations
 - Maximizing a Functional
- Finding Approximation to a Posterior
 - Minimizing K-L divergence
 - Factorized Distribution
- Example of Bivariate Gaussian

Function versus Functional

- Function takes a value of a variable as input and returns the function value as output
- Functional takes a function as input and returns the functional value as output
- An example is entropy

$$H[p(x)] = \int p(x) \ln p(x) dx$$

Calculus of Variations

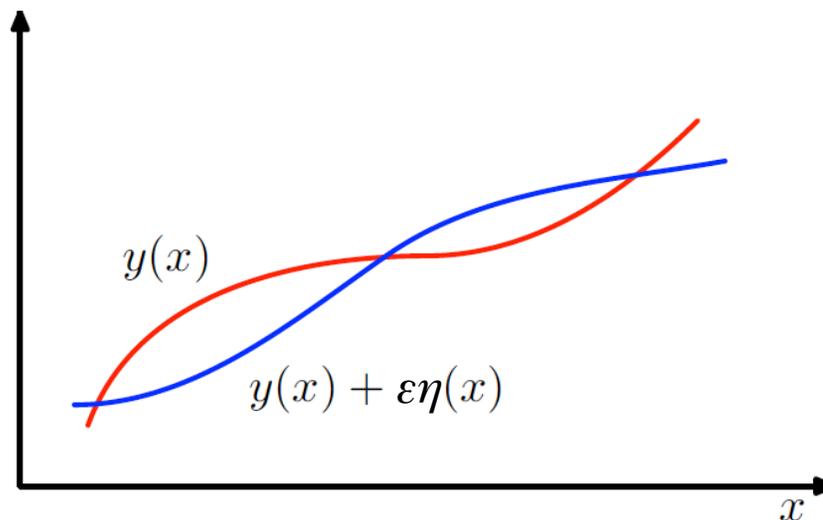
- *Function*: input x is a value, returns a value $y(x)$
- *Functional*: input is a function $y(x)$, returns a value $F[y(x)]$
- *Standard Calculus* concerns *Derivatives* of functions
 - How output changes in response to small changes in input
- *Variational Calculus*: concerns *Functional derivative*
 - How output changes in response to infinitesimal changes in input function

Invented by
Leonhard Euler
Swiss 1707-1783



Functional Derivative

- Defined by considering:
 - How the value of a functional $F[y]$ changes when the function $y(x)$ is changed to $y(x) + \varepsilon\eta(x)$
 - where $\eta(x)$ is an arbitrary function of x



Maximizing a Functional

- Common problem in conventional calculus:
 - Find a value x that maximizes a function $y(x)$
- In calculus of variations find a function $y(x)$ that maximizes a functional $F[y]$
 - E.g., Find a function that maximizes entropy
- Can be used to show that:
 1. Shortest path between two points is a straight line
 2. Distribution that maximizes entropy

$$H[p(x)] = \int p(x) \ln p(x) dx \quad \text{is Gaussian}$$

Variational methods

- Nothing intrinsically approximate
- But naturally lend themselves to approximation
 - By restricting range of functions, e.g.,
 - Quadratic
 - Linear combination of fixed basis functions
 - In which only the coefficients of the linear combination vary
 - Factorization assumptions
 - Method used for probabilistic inference

Application to Probabilistic Inference

- Consider a fully Bayesian model
- All parameters are given prior distributions
- Model may have latent variables as well as parameters
 - Denote all parameters and latent variables by Z
 - Denote set of observed variables by X
 - Given N iid samples for which $X = \{x_1, \dots, x_N\}$ and $Z = \{z_1, \dots, z_N\}$, probabilistic model specifies posterior distribution $p(Z|X)$ as well as model evidence $p(X)$

Lower Bound on $p(X)$ using KL divergence

$$\ln p(X) = L(q) + KL(q \parallel p) \quad \text{Log Marginal Probability}$$

where

$$L(q) = \int q(Z) \ln \left\{ \frac{p(X, Z)}{q(Z)} \right\} dZ$$

and

$$KL\{q \parallel p\} = - \int q(Z) \ln \left\{ \frac{p(Z | X)}{q(Z)} \right\} dZ$$

KL Divergence between proposed q and p (the desired posterior of interest) to be minimized

Functional we wish to maximize

Also applicable to discrete distributions by replacing integrations with summations

Observations on optimization

- Lower bound on $\ln p(X)$ is $L(q)$
- *Maximizing* lower bound $L(q)$ wrt distribution $q(Z)$ is equivalent to *minimizing* KL Divergence
- KL divergence vanishes when $q(Z) = p(Z|X)$

Plan:

- We seek that distribution $q(Z)$ for which $L(q)$ is largest
- Since true posterior is intractable we consider restricted family for $q(Z)$
- Seek member of this family for which KL divergence is minimized

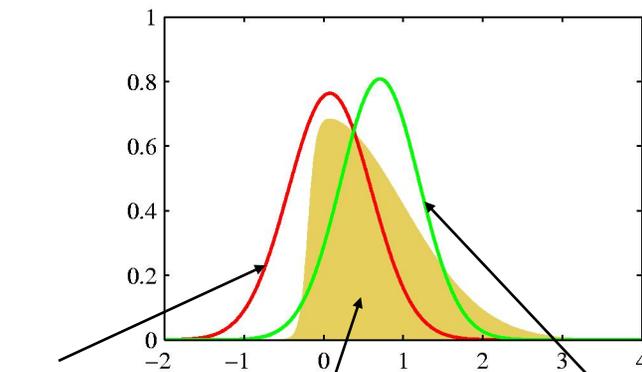
Approximating Family

- Typically true posterior $p(Z|X)$ is intractable
- Pick family of parametric distributions $q(Z|\omega)$ governed by parameters ω
 - Lower bound becomes a function of ω
 - Can exploit standard nonlinear optimization techniques to determine optimal parameters
- Use q with the fitted parameters as a proxy for the posterior
 - e.,g., to predict about future data, or posterior of hidden variables

Variational Approximation Example

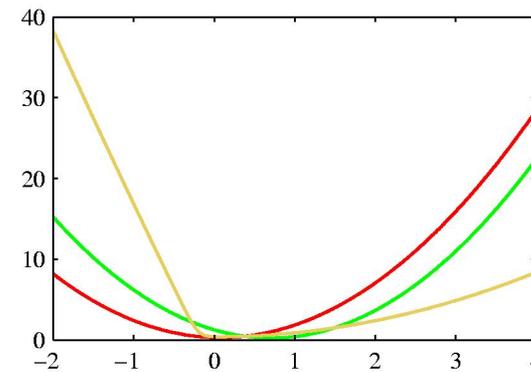
Use parametric distribution $q(Z|\omega)$

Lower bound for the functional $L(q)$ is a function of ω and can use standard nonlinear optimization to determine optimal ω



Laplace

Original
distribution



Negative
Logarithms

Variational distribution is Gaussian
Optimized with respect to mean and variance

Factorized Distribution Approach

- One way to restrict the family of distributions
 - Partition the elements of Z into disjoint groups Z_i , $i=1, \dots, M$, so that

$$q(Z) = \prod_{i=1}^M q_i(Z_i)$$

- No further assumptions made about distribution
 - No restrictions on factors $q_i(Z_i)$
- Method is called *Mean Field Theory* in Physics
 - also *Mean Field Variational Inference*

Factorized Distribution Approach

- Among factorized distributions $q(Z)$
 - we seek one for which lower bound $L(q)$, is largest

$$\begin{aligned}
 L(q) &= \int q(Z) \ln \left\{ \frac{p(X, Z)}{q(Z)} \right\} dZ \\
 &= \int \prod_i q_i \left\{ \ln p(X, Z) - \sum_i \ln q_i \right\} dZ \\
 &= \int q_j \left\{ \int \ln p(X, Z) \prod_{i \neq j} q_i dZ_i \right\} dZ_j - \int q_j \ln q_j dZ_j + const \\
 &= \int q_j \ln \tilde{p}(X, Z_j) dZ_j - \int q_j \ln q_j dZ_j + const
 \end{aligned}$$

Substitute $q(Z) = \prod_{i=1}^M q_i(Z_i)$

Denote $q_j(Z_j)$ by q_j

KL Divergence of $q_j(Z_j)$ and $\tilde{p}(X, Z_j)$

- where we have defined a new distribution $\tilde{p}(X, Z_j)$ as

$$\ln \tilde{p}(X, Z_j) = E_{i \neq j} [\ln p(X, Z)] + const$$

where $E_{i \neq j} [\dots]$ is an expectation wrt the q distribution, so that

$$E_{i \neq j} [\ln p(X, Z)] = \int \ln p(X, Z) \prod_{i \neq j} q_i dZ_i$$

Maximizing $L(q)$ with factorization

- Since we have seen that

$$L(q) = \int q_j \ln \tilde{p}(X, Z_j) dZ_j - \int q_j \ln q_j dZ_j + const$$

- Maximizing $L(q)$ is equivalent to minimizing the Kullback-Leibler divergence between

$$q_j(Z_j) \text{ and } \tilde{p}(X, Z_j)$$

- The minimum occurs when $q_j(Z_j) = \tilde{p}(X, Z_j)$
- Thus we obtain a general expression for the optimal solution $q_j^*(Z_j)$ given by

$$\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + const$$

Optimal Solution Interpretation

- The optimal solution is

$$\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + \text{const}$$

- This provides the basis for application of variational methods
- Log of the optimal solution for factor q_j is obtained by considering
 - log of the joint distribution over all hidden and joint variables and
 - Take expectation wrt all other factors $\{q_i\}$ for $i \neq j$

Normalized Solution

- Additive constant in $\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + const$ is set by normalizing the distribution $q_j^*(Z_j)$
 - Take exponential of both sides and normalize

$$q_j^*(Z_j) = \frac{\exp\left(E_{i \neq j} [\ln p(X, Z)]\right)}{\int \exp\left(E_{i \neq j} [\ln p(X, Z)]\right) dZ_j}$$

- In practice more convenient to work with the form $\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + const$

Obtaining an Explicit Solution

- Set of equations $\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + const$
 - $j=1, \dots, M$ represent a set of consistency conditions for the maximum of the lower bound
 - They do not provide an explicit solution
- Solution obtained by properly initializing $q_i(Z_i)$ and cycling through all other factors
 - Replacing each in turn with revised estimate given by $E_{i \neq j} [\ln p(X, Z)] + const$ evaluated with current estimates for all other factors
 - Convergence guaranteed because bound is convex wrt factors $q_i(Z_i)$

Properties of Factorized Approximations

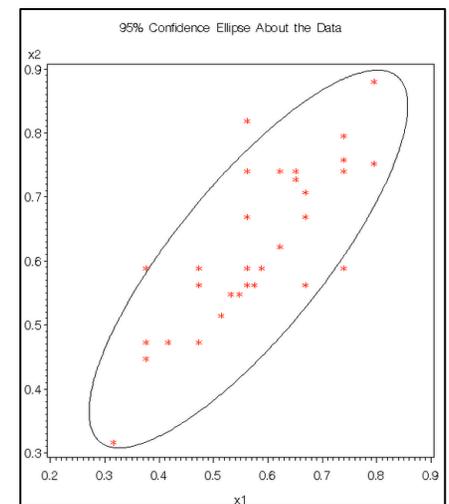
- Our approach to variational inference is based on a factorized approximation to the true posterior distribution
- We will consider approximating a Gaussian distribution using a factorized Gaussian
- It will provide insight into types of inaccuracy introduced in using factorized approximations

Example of Factorized Approximation

- True posterior distribution $p(\mathbf{z})$, $\mathbf{z} = \{z_1, z_2\}$, is Gaussian

$p(\mathbf{z}) = N(\mathbf{z} | \boldsymbol{\mu}, \Lambda^{-1})$ over two correlated variables

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} \quad \Lambda = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix}$$



- We wish to approximate $p(\mathbf{z})$ using a factorized Gaussian of the form $q(\mathbf{z}) = q_1(z_1) q_2(z_2)$

Factor Derivation

- Using the optimal solution for factorized distributions

$$\ln q_j^*(Z_j) = E_{i \neq j} [\ln p(X, Z)] + \text{const}$$

- We find an expression for the optimal factor $q_1^*(z_1)$

Only need terms that have functional dependence on z_1

Other terms are absorbed into normalization constant

$$\begin{aligned} \ln q_1^*(z_1) &= E_{z_2} [\ln p(z)] + \text{const} && \boxed{p(z) = N(z | \mu, \Lambda^{-1})} \\ &= E_{z_2} \left[-\frac{1}{2} (z_1 - \mu_1)^2 \Lambda_{11} - (z_1 - \mu_1) \Lambda_{12} (z_2 - \mu_2) \right] + \text{const} \\ &= -\frac{1}{2} z_1^2 \Lambda_{11} + z_1 \mu_1 \Lambda_{11} - z_1 \Lambda_{12} (E[z_2] - \mu_2) + \text{const} \end{aligned}$$

Since rhs is a quadratic function of z_1 , we can identify $q_1^*(z_1)$ as a Gaussian

Factorized Approximation of Gaussian

- We did not assume $q_1(z_1)$ to be Gaussian
 - but derived it as a variational optimization
 - of KL divergence over all possible distributions
- We can identify mean and precision for the Gaussians as

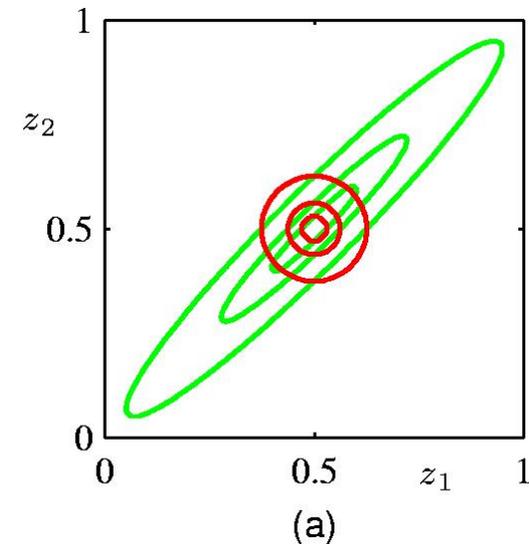
$$\begin{aligned} q_1^*(z_1) &= N(z_1 | m_1, \Lambda_{11}^{-1}) \text{ where } m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (E[z_2] - \mu_2) \\ q_2^*(z_2) &= N(z_2 | m_2, \Lambda_{22}^{-1}) \text{ where } m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (E[z_1] - \mu_1) \end{aligned}$$

– Coupled solution

- $q_1^*(z_1)$ depends on expectation wrt $q^*(z_2)$
- Treated as re-estimation equations and cycling through them in turn

Result of Variational Approximation

- Mean is correctly captured
- Variance is controlled by
 - direction of smallest variance of $p(z)$
 - Variance along orthogonal direction is under-estimated
- In general:
 - factorized variational approximation gives results that are too compact



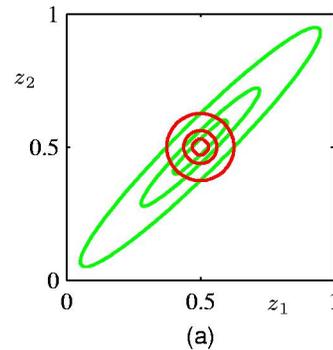
Green: Correlated Gaussian
For 1, 2 and 3 std deviations

Red: $q(z)$ over same variables
given by product
of two independent
univariate Gaussians

Two alternative forms of KL Divergence

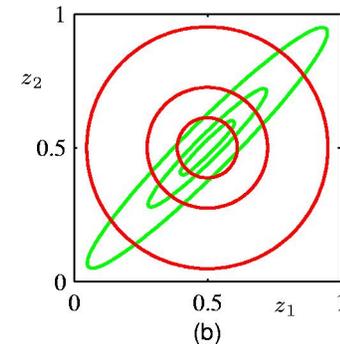
Green: Correlated Gaussian
For 1, 2 and 3 std deviations

Red: $q(z)$ over same variables
given by product
of two independent
univariate Gaussians



Minimization based on
KL Divergence $KL(q||p)$

Mean correctly captured
Variance under-estimated
: too compact

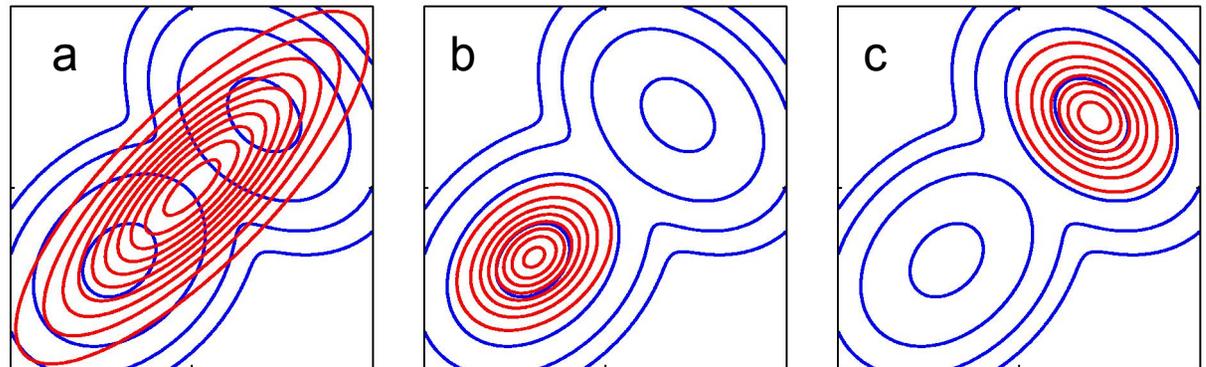


Minimization based on
Reverse KL Divergence
 $KL(p||q)$

Form used in
Expectation Propagation

Alternative KL divergences for Multimodal

- Approximating a multimodal by a unimodal one



Blue Contours: multimodal distribution $p(Z)$

Red Contours in (a): Single Gaussian $q(Z)$ that best approximates $p(Z)$

in (b): single Gaussian $q(Z)$ that best approximates $p(Z)$

in the sense of minimizing $KL(p||q)$

in(c): as in (b) showing different local minimum of KL divergence

Alpha Family of Divergences

- Two forms of divergence are members of the alpha family of divergences

$$D_{\alpha}(p \parallel q) = \frac{4}{1-\alpha^2} \left(1 - \int p(x)^{(1+\alpha)/2} q(x)^{(1-\alpha)/2} dx \right)$$

– where $-\infty < \alpha < \infty$

– $KL(p \parallel Q)$ corresponds to $\alpha \rightarrow 1$

– $KL(q \parallel p)$ corresponds to $\alpha \rightarrow -1$

For all α $D_{\alpha}(p \parallel q) \geq 0$ with equality iff $p(x) = q(x)$

when $\alpha = 0$ we get symmetric divergence

– which is linearly related to the Hellinger Distance (a valid distance measure)

$$D_H(p \parallel q) = \int (p(x)^{1/2} - q(x)^{1/2})^2 dx$$

Example: Univariate Gaussian

- Inferring the parameters of a single Gaussian
 - Mean μ and precision τ given data set $D = \{x_1, \dots, x_N\}$
- Assume Gaussian-Gamma conjugate prior distribution of parameters

- Likelihood function:

$$p(D | \mu, \tau) = \left(\frac{\tau}{2\pi}\right)^{N/2} \exp\left\{-\frac{\tau}{2} \sum_{n=1}^N (x_n - \mu)^2\right\}$$

- Conjugate priors for μ and τ

$$p(\mu | \tau) = N(\mu | \mu_0, (\lambda_0 \tau)^{-1})$$

$$p(\tau) = \text{Gam}(\tau | a_0, b_0)$$

Factorized Variational Approximation

- Posterior distribution $q(\mu, \tau) = q_\mu(\mu)q_\tau(\tau)$
- For $q_\mu(\mu)$ we have a Gaussian
- Optimal solution for factor $q_\tau(\tau)$ is a Gamma

Variational Inference of Univariate Gaussian

- Variational Inference of mean μ and precision τ
- (a) Green: Contours of true posterior $p(\mu, \tau | D)$
Blue: Contours of initial approx $q_\mu(\mu)q_\tau(\tau)$
- (b) after re-estimating the factor $q_\mu(\mu)$.
- (c) after re-estimating the factor $q_\tau(\tau)$.
- (d) Iterative scheme converges to red contours

