

VARIATIONAL INFERENCE

$p(x, z)$, with x observable variables
with z latent variables

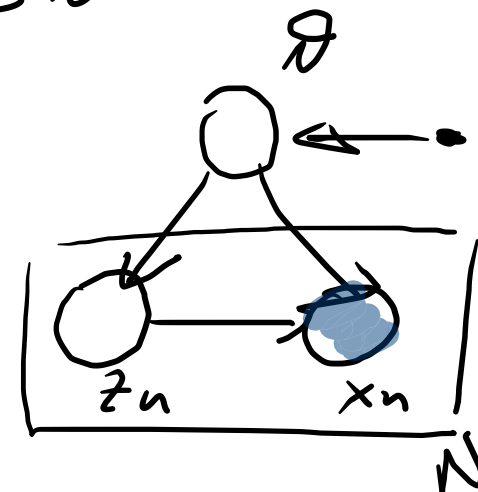
local latent variables z_n
global latent variables θ

$$z = (z_1, \dots, z_N, \theta)$$

We have observation $x_1, \dots, x_N = \underline{x}$ (sometimes x)

$p(z | \underline{x})$ posterior

$p(\underline{x}) = \int p(\underline{x}, z) dz$ evidence



ELBO

$$\underbrace{P(x)}_{\text{fixed}} = \underbrace{f(q)}_{\text{fixed}} + \underbrace{\text{KL}[q \parallel p]}_{\geq 0}$$

$\text{KL}[q(z) \parallel p(z|x)]$

$$f(q) = \int q(z) [\log p(x, z) - \log q(z)] dz = \mathbb{E}_q [\log p(x, z) - \log q(z)]$$

$$\text{KL}[q(z) \parallel p(z|x)] = - \int q(z) [\log p(z|x) - \log q(z)] dz$$

q - VARIATIONAL DISTRIBUTION

FIND THE BEST q THAT APPROXIMATES p(z|x)

minimize $\text{KL}[q \parallel p] \Leftrightarrow$ maximize ELBO $f(q)$

$f(q)$ is maximum when $\text{KL}[q \parallel p] = 0 \Leftrightarrow q(z) = p(z|x)$

Point $p(z|x)$ is INTRACTABLE

Solution: RESTRICT q to a TRactable FAMILY OF DISTRIBUTIONS!

The optimal $q(z)$ likely is s.t. $KL[q||p] \gg 0$

VARIATIONAL INFERENCE \Rightarrow solve an optimization problem
to maximize $f(q)$ in a suitable space of
variational distributions $Q \ni q$

SCENARIO: $Q =$ set of parametric distributions $q(z|\lambda)$ λ parameters M/\mathbb{R}^k

Then we need to find $\underset{\lambda}{\text{argmax}} f(\lambda)$

This is a highly non-linear non-convex optimization

$$q(z|\lambda) = q(\theta|\lambda\theta) \prod_i q(z_i|\lambda_i, \theta) \quad \lambda = (\lambda_1, \lambda_2, \dots, \lambda_n)$$

VARIATIONAL MEAN FIELDS

$P(z|x)$ by $q(z)$. $z = (z_1, \dots, z_M)$ decomposition

$$q(z) = \prod_{i=1}^M q_i(z_i) \rightarrow \text{MEAN FIELD ASSUMPTION}$$

q_i for simple by

$$J(q) = E_q [\log p(x, z) - \log q(z)]$$

$$= \int \prod_i q_i [\log p(x, z) - \sum_i \log q_i] dz$$

factor cent
terms
q_i

$$= \int q_j \left[\int \log p(x, z) \prod_{i \neq j} q_i dz_i \right] dz_j - \int q_j \log q_j dz_j + \text{CONST}$$

$$\rightarrow E_{i \neq j} [\log p(x, z)]$$

defines $\tilde{p}(x, z_j)$

$$= E_{q_j} [\log \tilde{p}(x, z_j) - \log q_j] + \text{CONST}$$

$$\log \tilde{p}(x, z_j) = E_{i \neq j} [\log p(x, z)] + \text{CONST}$$

$$J(q_j) = KL[q_j \parallel \tilde{p}(x, z_j)] + \text{const}$$

$$\tilde{p}(x, z_j) \text{ as}$$

$$\log \tilde{p}(x, z_j) = \mathbb{E}_{i \neq j} [\log p(x, z_i)] + \text{const}$$


This is maximised for q_j , with $q_i, i \neq j$ fixed

$$\text{by } q_j^*(z_j) = \tilde{p}(x, z_j)$$

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

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

$$\int \exp(\mathbb{E}_{i \neq j} [\log p(x, z_i)]) dz_j$$

if we can compute (analytically) $\mathbb{E}_{i \neq j} [\log p(x, z_i)]$ 
 we initialize q_i , then cycle through q_j , optimize them in turn
 (COORDINATE GRADIENT ASCEND), until convergence

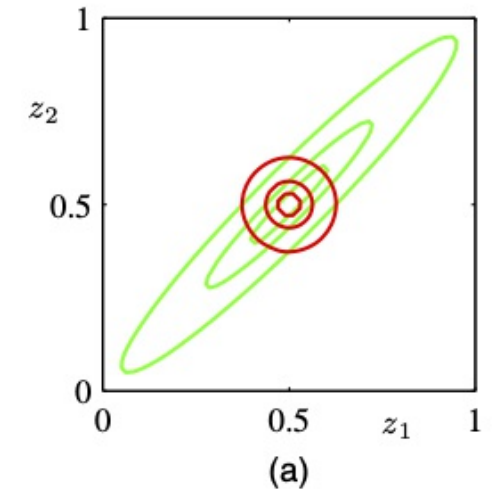
The lower bound is convex on q_j

A SIMPLE EXAMPLE

Gaussian to factorized gaussian $p(z) = \mathcal{N}(z | \mu, \Lambda^{-1})$, $\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}$, Λ precision matrix

$z = (z_1, z_2)$ M.F.V.D $q(z) = q(z_1)q(z_2)$

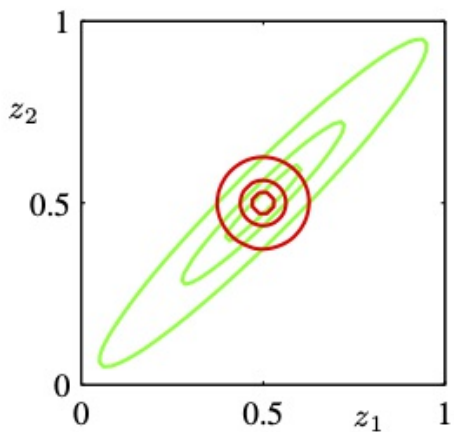
$$\begin{aligned} \log q_1^*(z_1) &= \mathbb{E}_{z_2} [\log p(z)] + \text{const} \\ &= \mathbb{E}_{z_2} \left[-\frac{1}{2} (z_1 - \mu_1)^2 \Lambda_{11} - (z_1 - \mu_1) \Lambda_{12} (z_2 - \mu_2) \right] + \text{const} \\ &= -\frac{1}{2} (z_1 - \mu_1)^2 \Lambda_{11} - (z_1 - \mu_1) \Lambda_{12} (\mathbb{E}[z_2] - \mu_2) + \text{const} \end{aligned}$$



$\Rightarrow q_1^*(z_1)$ is Gaussian $\mathcal{N}(z_1 | m_1, \Lambda_{11}^{-1})$ $m_1 = \mu_1 - \Lambda_{11}^{-1} \Lambda_{12} (\mathbb{E}[z_2] - \mu_2)$

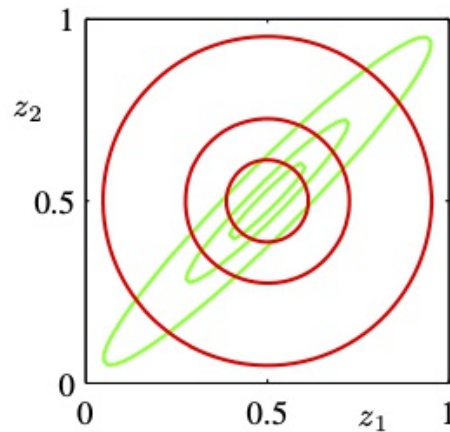
$\Rightarrow q_2^*(z_2) = \mathcal{N}(z_2 | m_2, \Lambda_{22}^{-1})$, $m_2 = \mu_2 - \Lambda_{22}^{-1} \Lambda_{21} (\mathbb{E}[z_1] - \mu_1)$

$m_i = \mathbb{E}[z_i] \Rightarrow m_1 = \mu_1, m_2 = \mu_2$



(a)

$$q(z) = q(z_1)q(z_2)$$



(b)

$$KL[q || p]$$

ZERO FORCING

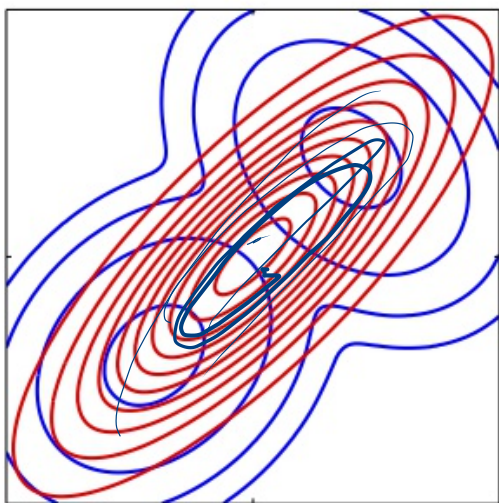
if $p(z) \approx 0$ then $q(z) \approx 0$
one mode

$$KL[p || q]$$

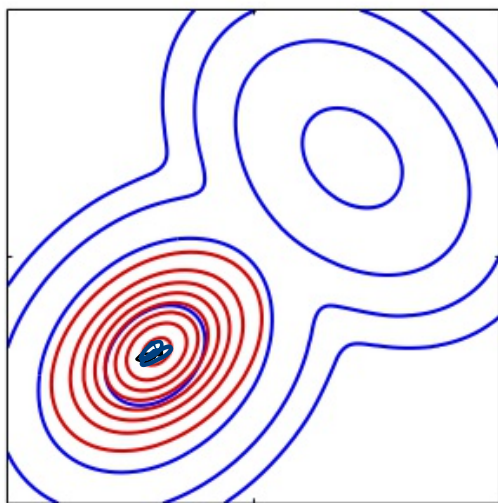
ZERO AVOIDING

$q(z)$ is non-zero where
 $p(z)$ is non-zero
overlaps modes

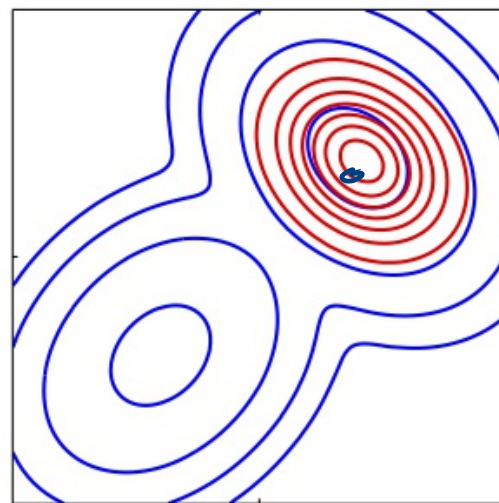
$$KL[p || q]$$



(a)



(b)



(c)

$$KL[q || p]$$

VARIATIONAL LINEAR REGRESSION

$$P(y, \omega, \alpha) = P(y | \omega) P(\omega | \alpha) P(\alpha)$$

$$P(y | \omega) = \prod_{n=1}^N \mathcal{N}(y_n | \omega^T \phi(x_n), \beta^{-1})$$

$$P(\omega | \alpha) = \mathcal{N}(\omega | 0, \alpha^{-1} I)$$

$$P(\alpha) = \text{GAMMA}(\alpha | a_0, b_0) = \frac{1}{\Gamma(a_0)} b_0^{a_0} \alpha^{a_0-1} e^{-b_0 \alpha}$$

$a_0 \propto N$

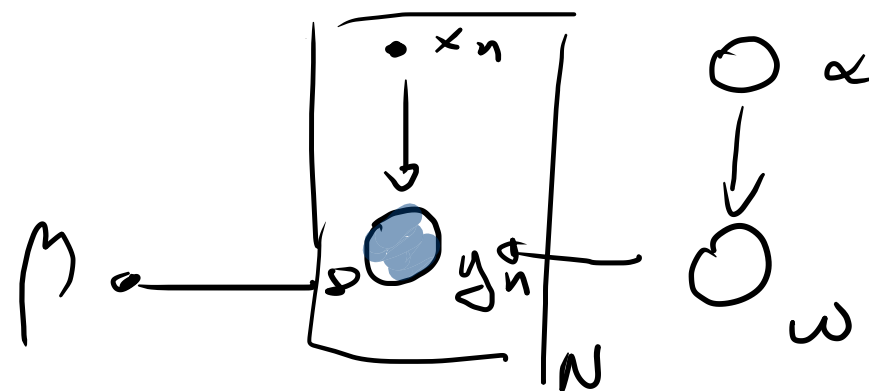
$$P(\omega, \alpha | y)$$

$$q(\omega, \alpha) = q(\omega) q(\alpha)$$

$$q^*(\alpha) = \mathbb{E}_{\omega} [\log P(y, \omega, \alpha)] + \text{CONST}$$

$$\stackrel{!}{=} \log P(\alpha) + \mathbb{E}_{\omega} [\log P(\omega | \alpha)] + \text{CONST}$$

$$\stackrel{!}{=} (a_0 - 1) \log \alpha - b_0 \alpha + \frac{M}{2} \log \alpha - \frac{\alpha}{2} \mathbb{E}[\omega^T \omega] + \text{CONST}$$



$$q^*(\alpha) = \text{GAMMA}(\alpha \mid a_N, b_N)$$

$$a_N = a_0 + \frac{M}{2}$$

$$b_N = b_0 + \frac{1}{2} \mathbb{E}_q[\omega^T \omega]$$

$$\log q^*(\omega) = \log q(y \mid \omega) + \mathbb{E}_\alpha [\log p(\omega \mid \alpha)] + \text{const}$$

$$= -\frac{M}{2} \sum_{n=1}^N [\omega^T \phi(x_n) - y_n]^2 - \frac{1}{2} \mathbb{E}[\alpha] \omega^T \omega + \text{const}$$

$$\Rightarrow q^*(\omega) = \mathcal{N}(\omega \mid m_N, S_N)$$

$$m_N = (M S_N \Phi^T y)$$

$$S_N = (\mathbb{E}[\alpha] \mathbf{I} + M \Phi^T \Phi)^{-1}$$

$$\mathbb{E}[\alpha] = a_N / b_N, \quad \mathbb{E}[\omega^T \omega] = m_N^T m_N + \text{Trace}(S_N)$$

We can compute $\mathcal{L}(q) \approx p(y)$ (approximates model evidence)

BLACK BOX VARIATIONAL INFERENCE

$q(z|\lambda)$ parametric distribution

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

we want to compute $\nabla_{\lambda} f(\lambda)$

SOLUTION 1: Reparametrization trick

if $z = g_{\nu}(\epsilon)$, $\epsilon \sim \hat{q}(\epsilon)$ (indep of λ and fixed), $\bar{\lambda} = (\lambda, \nu)$

$$f(\bar{\lambda}) = \mathbb{E}_{\hat{q}(\epsilon)} [\log p(x, g_{\nu}(\epsilon)) - \log q(g_{\nu}(\epsilon)|\lambda)]$$

$$\nabla_{\bar{\lambda}} f(\bar{\lambda}) = \mathbb{E}_{\hat{q}(\epsilon)} \left[\underbrace{\nabla_{\bar{\lambda}} \log p(x, g_{\nu}(\epsilon)) - \nabla_{\bar{\lambda}} \log q(g_{\nu}(\epsilon)|\lambda)}_{G(\epsilon)} \right]$$

$\epsilon_i \sim \hat{q}(\epsilon)$ $\nabla_{\bar{\lambda}} f(\bar{\lambda}) = \frac{1}{S} \sum_{s=1}^S G(\epsilon_i)$ \leftarrow This is the SGA
stochastic gradient descent

case 2: non-reparametrizable $q(z|\lambda)$

$$\nabla_{\lambda} \mathcal{J}(\lambda) = \nabla_{\lambda} \mathbb{E}_{q(z|\lambda)} [\log p(x, z) - \log q(z|\lambda)]$$

$$= \nabla_{\lambda} \int q(z|\lambda) [\log p(x, z) - \log q(z|\lambda)] dz$$

$$= \int \nabla_{\lambda} [\log p(x, z) - \log q(z|\lambda)] q(z|\lambda) dz + \int \nabla_{\lambda} q(z|\lambda) [\log p(x, z) - \log q(z|\lambda)] dz$$

$$= - \mathbb{E}_q [\nabla_{\lambda} \log q(z|\lambda)] +$$

$$\mathbb{E}_q \left[\frac{\nabla_{\lambda} q(z|\lambda)}{q(z|\lambda)} \right] = \int \nabla_{\lambda} q(z|\lambda) dz = \nabla_{\lambda} \int q(z|\lambda) dz = \nabla_{\lambda} 1 = 0$$

$$\nabla_{\lambda} \log q(z|\lambda) = \frac{\nabla_{\lambda} q(z|\lambda)}{q(z|\lambda)} \Rightarrow \nabla_{\lambda} q(z|\lambda) = \nabla_{\lambda} \log q(z|\lambda) \cdot q(z|\lambda)$$

$$\begin{aligned} \nabla_{\lambda} f(\lambda) &= \int q(z|\lambda) \nabla_{\lambda} \log q(z|\lambda) [\log p(x, z) - \log q(z|\lambda)] dz \\ &= E_q [\nabla_{\lambda} \log q(z|\lambda) [\log p(x, z) - \log q(z|\lambda)]] \end{aligned}$$

$$z_s \sim q(z|\lambda) \Rightarrow \nabla_{\lambda} f_{\lambda} \approx \underbrace{\frac{1}{S} \sum_{s=1}^S \nabla_{\lambda} \log q(z_s|\lambda) [\log p(x, z_s) - \log q(z_s|\lambda)]}_{(*)}$$

VARIANCE REDUCTION

ISSUE: The estimator (*) typically has HIGH VARIANCE

→ RAO-BLACKWELLIZATION

→ CONTROL VARIABLES

RAD-BLACKWELLIZATION

x, y random variables $\mathcal{J}(x, y)$, $E[\mathcal{J}(x, y)]$ is our target

$$a) \hat{\mathcal{J}}(x) = E_y[\mathcal{J}(x, y) | x] \text{ st. } E_x[\hat{\mathcal{J}}(x)] = E_{x,y}[\mathcal{J}(x, y)]$$

$$\text{VAR}[\hat{\mathcal{J}}(x)] = \text{VAR}[\mathcal{J}(x, y)] - E[(\mathcal{J}(x, y) - \hat{\mathcal{J}}(x))^2] < \text{VAR}[\mathcal{J}(x, y)]$$

Mean field factorization of $q(z|\lambda) = \prod_{i=1}^N q(z_i|\lambda_i)$

- $q_{(i)}$: marginal of $q(z|\lambda)$ on the marked block $z_{(i)}$ of z_i ($i \in \mathcal{P}$)

$p_i(x, z_{(i)})$ the product of factors of $p(z|\lambda)$ depending on $z_{(i)}$

$$E_{q_{(i)}} E_{q_{(j)}} \left[\nabla_{\lambda_i} \log q_{(i)} \right] \hat{\nabla}_{\lambda_i} \log q_{(i)} = E_{q_{(i)}} \left[\nabla_{\lambda_i} \log q(z_i|\lambda_i) (\log p_i(x, z_{(i)}) - \log q(z_i|\lambda_i)) \right]$$

$$z_s \sim q_{(i)}(z|\lambda)$$

CONTROL VARIATES

f : $E_q[f]$ our goal. Instead $E_q[\hat{f}]$ s.t. $E_q[\hat{f}] = E[f]$ and $VAR_q[\hat{f}] < VAR_q[f]$

we choose h : $E[h] < \infty$, define

$$\hat{f}_a(z) = f(z) - a(h(z) - E[h(z)])$$

The a^* that minimizes variance is $a^* = \frac{COV(f, h)}{VAR(h)}$

$$f_i(z) = \nabla_{\lambda_i} \log q(z_i | \lambda_i) [\log p_i(x, z_i) - \log q_i(z_i | \lambda_i)]$$

$$h_i(z) = \nabla_{\lambda_i} \log q(z_i | \lambda_i)$$

$E[\cdot] = 0$

$z^{(i)} = n_i$

$$\hat{a}_i^* = \frac{\sum_{d=1}^{n_i} \hat{COV}(f_i^d, h_i^d)}{\sum_{d=1}^{n_i} \hat{VAR}(h_i^d)}$$

$$\hat{\nabla}_{\lambda_i} f = \frac{1}{S} \sum_{s=1}^S \nabla_{\lambda_i} \log q(z_s | \lambda_i) [\log p_i(x, z_s) - \log q_i(z_s | \lambda_i) - \hat{a}_i^*]$$

$z_s \sim q_i(\cdot | \lambda)$