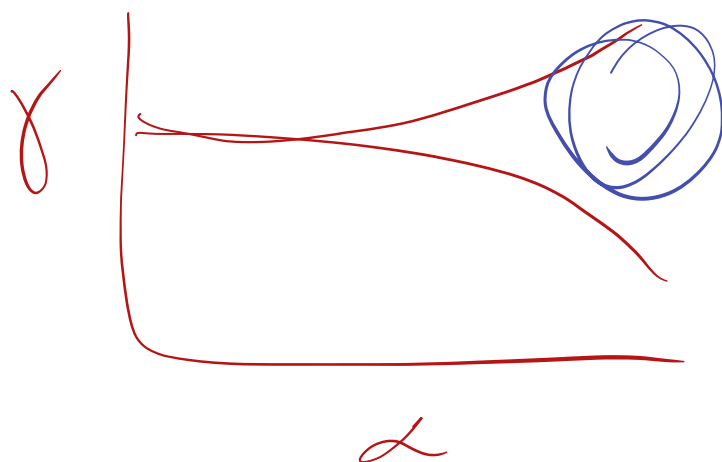# Approximate Inference

- HW 5 (last one), due 2 weeks

- Reading B.7 (Variational Inference)

- Reminder: submit Project "Proposals"

# Approximate Inference

- MCMC can be very slow in high dimensional problems

- Idea: find a distribution that is easy to sample from which closely approximate $p(\theta \mid y)$

- A couple of examples
    - Laplace Approximation — *Normal Approximation*
    - INLA — *Iterated Nested Laplace Approx.*
    - Variational Bayes *(13.7)*

# Mean-Field Variational Bayes

- Solve an optimization problem:

$$\underset{\lambda}{\operatorname{argmin}} \ \operatorname{dist}(p(\theta \mid y), g_\lambda(\theta))$$

*Want: 1. $g$ close to $P(\theta|y)$*

*2. $g$ easy to sample from*

- In mean-field inference we restrict $g(\theta_1, \ldots, \theta_d) = \prod_i g(\theta_i)$

- For VB, we usually use KL-divergence to measure the "distance" between probability distributions

  - KL is not a true distance since $KL(p|q) \neq KL(q|p)$.

*Kullback-Leibler Divergence*

# Mean-Field Variational Bayes

<span style="color:red">*Reverse- KL*</span>

- Optimize $KL(q(\theta)||p(\theta \mid y)) = -E_q \log(\frac{p(\theta|y)}{q(\theta)})$

- Problem: expectation under $q$ means that the variational distribution usually underestimates posterior uncertainty

*Approximating Distn*

$KL(q||p) \geq 0$

$KL(q||p) = 0$ iff $q = p$

$q(\theta) \ll p(\theta|y)$

$\text{supp}(q) \subseteq \text{supp}(p(\theta|y))$

$p(\theta|y) = 0 \implies q(\theta) = 0$

$$\text{Fwd:} \quad KL(P \| q) = -E_P \log \frac{q}{P}$$

$$KL(q(\theta) \| P(\theta|y)) =$$

$$-E_q\left(\log \frac{P(\theta|y)}{q(\theta)}\right) = E_q \log q(\theta) - E_q\left[\log \overbrace{P(\theta|y)}^{\frac{P(\theta,y)}{P(y)}}\right]$$

$$= E_q \log q(\theta) - E_q \log P(\theta,y) + \underbrace{\log P(y)}_{\text{model evidence}}$$

$$\geq 0$$

$$\log P(y) \geq \underbrace{E_q \log P(\theta,y) - E_q \log q}_{\text{Evidence Lower Bound (ELBO)}}$$

**ELBO:** $E_q \log P(\theta, y) - E_q \log q =$

$$E_q \log(P(y|\theta)P(\theta)) - E_q \log q =$$

$$E_q \log P(y|\theta) + \underbrace{E_q \log P(\theta) - E_q \log q}_{KL(q(\theta) \| P(\theta))}$$

$$= \underbrace{E_q \log P(y|\theta)}_{\substack{\text{Favors } q(\theta) \\ \text{which explain} \\ \text{the data.}}} + \underbrace{KL(q(\theta) \| P(\theta))}_{\substack{\text{Favor } q \text{ close} \\ \text{to the prior}}}$$

## Mean-Field VI.

$$q(\theta_1, \dots \theta_d) = q(\theta_1) \dots q(\theta_d)$$

Coordinate Updates:

Update one $q(\theta_i)$ at a time

conditional on the rest.

$$E_q \log \frac{P(\theta, y)}{q(\theta)} = \int_{\theta_d} \cdots \int_{\theta_j} \log\left(\frac{P(\theta, y)}{q(\theta_1)\cdots q(\theta_d)}\right) q(\theta_1)\cdots q(\theta_d) \, d\theta_1, \ldots d\theta_d$$

$$= \int_{\theta_{-j}} \left( \int_{\theta_j} \log\left(\frac{P(\theta, y)}{q(\theta_j)}\right) q(\theta_j) \, d\theta_j \right) d\theta_{-j} -$$

$$\int_{\theta_{-j}} q(\theta_{-j}) \log \theta_{-j} \, d\theta_{-j}$$

Ex: $J = 3$;

$$q(\theta_{-j}) = q_1(\theta_1) q_2(\theta_2) q_4(\theta_4) \cdots q_d(\theta_d)$$

To update $q_j(\theta_j)$ (fixing the rest)

$$\max_{q_j} \int_{\theta_{-j}} \int_{\theta_j} \log\left(\frac{P(\theta, y)}{q(\theta_j)}\right) q(\theta_j) \, d\theta_j -$$

$$\max_{q_j} \int_{\Theta_j} E_{q_{-j}(\Theta_{-j})} \log\left(\frac{P(\Theta, y)}{q_j(\Theta_j)}\right) q_j(\Theta_j) \, d\Theta_j$$

$$\max_{q_j} E_{q_j} E_{q_{-j}} \log P(\Theta, y) - E_{q_j} \log q_j$$

$$\text{Call: } \tilde{P}_j(\Theta_j, y) = E_{q_{-j}} \log P(\Theta, y)$$

$$\max_{q_j} E_{q_j} \frac{\tilde{P}_j(\Theta_j, y)}{q_j} = -KL\left(q_j, \tilde{P}_j\right)$$

$$\log q_j(\Theta) = E_{q_{-j}} \log P(\Theta, y) + \text{const}$$

Mean Field Algo:

until ELBO converges:

  for j in 1:D {

   $q_j(\Theta) \leftarrow$ distribution $\propto E_{q_{-j}} \log P(\Theta, y)$

  }

$$Y_i \sim N(\theta_i, \sigma^2)$$

$$\theta_i \sim N(\mu, \tau^2)$$

$$P(\mu) \propto \text{const}$$

$$\log P(\theta_1, \ldots \theta_8, \mu, y_1, \ldots y_8) =$$

$$\log \prod_{i=1}^{8} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\left(\frac{y_i - \theta_i}{2\sigma^2}\right)^2} \frac{1}{\sqrt{2\pi\tau^2}} e^{-\left(\frac{\theta_i - \mu}{2\tau^2}\right)^2}$$

until convergence:

    update $q_1(\theta_1)$ given $q_{-1}$

      $\vdots$

      $q_8(\theta_8)$

      $q_0(\mu)$      given others

$$q_j(\vartheta_j) \propto \underset{q_{-j}}{E}\left\{-\frac{(y_j-\vartheta_j)^2}{2\sigma^2} - \frac{(\vartheta_j-\mu)^2}{2\tau^2}\right\}$$

$$\propto e^{-\frac{(y_j-\vartheta_j)^2}{2\sigma^2} - \underset{q_{-j}}{E}\left[\frac{\vartheta_j^2-2\vartheta_j\mu+\mu^2}{2\tau^2}\right]}$$

$$\propto e^{-\frac{(y_j-\vartheta_j)^2}{2\sigma^2} - \frac{\vartheta_j^2}{2\tau^2} - \frac{2\vartheta_j \underset{q(\mu)}{E}\mu + \underset{q(\mu)}{E}[\mu^2]}{2\tau^2}}$$

$$q_\mu(\mu) \propto e^{\underset{q_{-\mu}}{E}\sum_{i=1}^{b}\frac{(\vartheta_i-\mu)^2}{2\tau^2}}$$

$$\propto e^{\underset{q_{-\mu}}{E}\sum\frac{\vartheta_i^2-2\mu\vartheta_i+\mu^2}{2\tau^2}}$$

$$\propto e^{\sum_{i>1}^{b}\frac{\underset{q(\vartheta_i)}{E}[\vartheta_i^2]-2\underset{q(\vartheta_i)}{E}[\mu]+\mu^2}{2\tau^2}}$$

# Automatic Differentiation Variational Inference

See: https://arxiv.org/pdf/1603.00788.pdf (Blei et al)

ADVI implemented in STAN.

Reminder: $q \ll p(\theta | y)$

$$supp(q(\theta)) \subseteq supp(p(\theta|y))$$

Define 1-1 transformation $T$,

from $supp(\theta) \longrightarrow \mathbb{R}^u$

$$\phi = T(\theta), \quad \phi \in \mathbb{R}^d$$

$$KL\big(q(\phi) \,\|\, P(\phi \,|\, y)\big)$$

$$P(\phi \,|\, y) = P_\theta\big(T^{-1}(\phi) \,|\, y\big)\,\big|\, J_{T^{-1}}(\phi)\big|$$

$$q(\phi) \sim \mathcal{N}\big(\mu_d, \Sigma_{d \times d}\big)$$

$$\min_{\mu, \Sigma} \quad KL\big(q(\phi) \,\|\, P(\phi \,|\, y)\big)$$

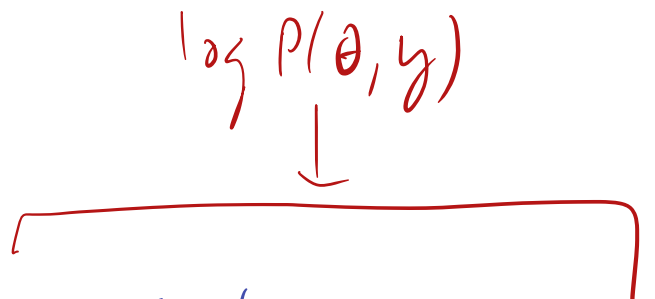$$\underline{ELBO} = \mathbb{E}_{q_{(\mu, \Sigma)}} \log P_\theta\big(T^{-1}(\phi), y\big) + \log\big|\bar{J}_{T^{-1}}(\phi)\big|$$

$$- \mathbb{E}_q \log q(\phi)$$

$$\mu^*, \Sigma^* = \underset{\mu, \Sigma}{\text{argmax}} \; ELBO(\mu, \Sigma)$$

$$n = S(\phi) = \Sigma^{-1/2}(\phi - \mu)$$

$$\Rightarrow q(n) \sim N(0, I)$$

ELBO:

$$E_{N(0,I)} \log \overbrace{P\left(T^{-1}\left(\underset{\mu, \Sigma}{S^{-1}}(n)\right), y\right)}^{\log P(\theta, y)} +$$

$$\log | J_{T^{-1}}(S^{-1}(n))| - E_{q(\phi)} \log q(\phi)$$