

MAP estimate on GLMs

DATA: y_n output, $x_n \in \mathbb{R}^D$ input vector, $n=1, 2, \dots, N$

LINEAR MODEL: $y_n = x_n^T z + \epsilon_n$ \leftarrow Noise $\Rightarrow p(y_n/x_n^T z) = \mathcal{N}(x_n^T z, \sigma^2)$
 \downarrow weights (let $\epsilon_n \sim \mathcal{N}(0, \sigma^2)$)

GENERALIZED: (GLM) Replace p by any Exp-Family distribution.

MAXIMUM LIKELIHOOD: \rightarrow Estimate $z!$
(MLE)

$$z_{MLE} = \arg \max_z \sum_{n=1}^N \log p(y_n/x_n^T z) - \frac{1}{2} \|z\|_2^2$$

MAP:

$$z_{MAP} = \arg \max_z \sum_{n=1}^N \log p(y_n/x_n^T z) + \log \mathcal{N}(z/0, 1/\lambda I)$$

$$z_{MAP} = \arg \max_z \sum_{n=1}^N \log p(y_n/z^T x_n) + \log p(z)$$

Data term Regularizer

Stochastic Gradient Descent (SGD)

$$\underline{z}_{\text{MAP}} = \arg \max_{\underline{z}} \overbrace{\sum_{n=1}^N \log p(y_n / \underline{z}^T \underline{x}_n) + \log p(\underline{z})} := L(\underline{z})$$

Randomly sample an n and compute an "unbiased" estimate

$$\hat{L}(\underline{z}) = N \log p(y_n / \underline{z}^T \underline{x}_n) + \log p(\underline{z}), \quad \mathbb{E}[\hat{L}(\underline{z})] = L(\underline{z})$$

$$\frac{\partial \hat{L}}{\partial \underline{z}} = N \frac{\partial}{\partial \underline{z}} \text{"} + \frac{\partial}{\partial \underline{z}} \text{"}$$

scalar step-size > 0

SGD UPDATE:
$$\underline{z}_{t+1} = \underline{z}_t + \alpha_t \left. \frac{\partial \hat{L}}{\partial \underline{z}} \right|_{\underline{z}=\underline{z}_t}$$

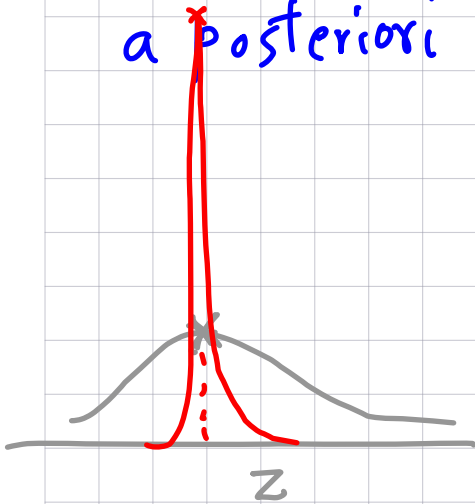
Q: converge?

MAP to Bayesian Inference

MAP:

$$\underline{z}_{\text{MAP}} = \arg \max_{\underline{z}} \sum_{n=1}^N \log p(y_n / \underline{z}^T \underline{x}_n) + \log p(\underline{z})$$

↓
Maximum
a posteriori



$$\log \left[\prod_{n=1}^N p(y_n / \underline{z}^T \underline{x}_n) \times p(\underline{z}) \right]$$

$$= \arg \max_{\underline{z}} \log P(\underline{y}, \underline{z})$$

$$= \text{Mode of } P(\underline{y}, \underline{z}) = \text{Mode of } p(\underline{z} / \underline{y})$$

$$= \text{Max. of a Posterior}$$

Bayes rule:

$$p(\underline{z} | \underline{y})$$

Posterior

$$= \frac{p(\underline{y}, \underline{z})}{\int p(\underline{y}, \underline{z}) d\underline{z}}$$

$$= \frac{p(\underline{y}, \underline{z})}{p(\underline{y})}$$

← Normalizing constant

INTRACTABLE INTEGRAL

Computational issues

GLM:

$$p(\underline{z}|\underline{y}) \propto \prod_{n=1}^N p(y_n | \underline{z}^T \underline{x}_n) p(\underline{z})$$

$\xrightarrow{p(y,z)}$
 \downarrow Likelihoods \downarrow Prior

CONJUGATE PRIORS:

$$\int \prod_{n=1}^N \text{Gaussian}(y_n | \underline{z}^T \underline{x}_n, \sigma^2) \times \text{Gaussian}(\underline{z} | 0, \mathbf{I}) d\underline{z}$$

Complete the Square!

$$= c \times \int \text{Gaussian}(\underline{z} | \cdot, \cdot) d\underline{z} = \log \det(\Sigma)$$

A Covariance matrix!

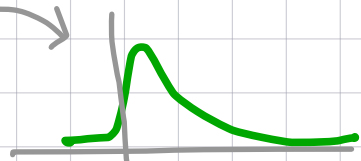
Q: Comp. cost for Gaussian?

NON-CONJUGATE :

$$\int \prod_{n=1}^N p(y_n | \dots) \times \text{Gaussian}(\underline{z} | 0, \mathbf{I}) d\underline{z}$$

other than Gaussian, usually no closed-form expression

$p(\underline{z}|\underline{y})$

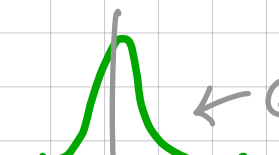


Logistic

\propto



\times



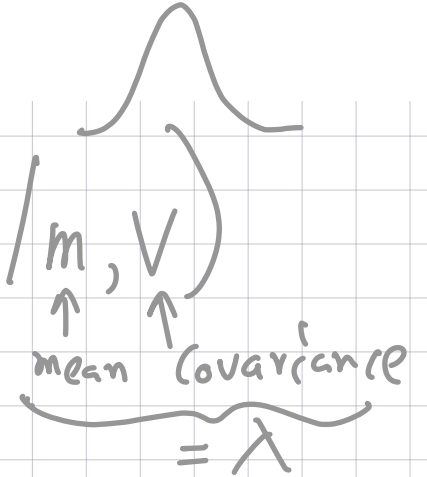
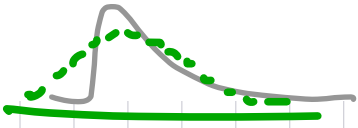
Gaussian

WE WILL CONVERT INTEGRATION \rightarrow OPTIMIZATION

Question

Why is Bayesian inference computationally challenging?

Variational Inference



MAIN IDEA: $P(z|y) \approx Q(z; \lambda)$ e.g. $= N(z | m, V)$

LOWER BOUND:

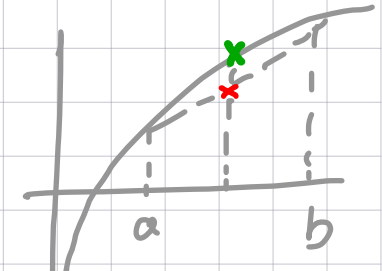
$$\log P(y) = \log \int P(y, z) dz = \log \int Q(z; \lambda) \frac{P(y, z)}{Q(z; \lambda)} dz$$

$$\geq \int Q(z) \log \frac{\prod P(y_n | \cdot) P(z)}{Q(z; \lambda)} dz \quad (\text{Jensen's inequality})$$

$$= \sum_{n=1}^N \int Q(z) \log P(y_n | \cdot) dz - \int Q(z) \log \frac{Q(z; \lambda)}{P(z)} dz$$

$$\log [ta + (1-t)b] \geq t \log a + (1-t) \log b$$

(green)
(red)



$$VI: \max_{\lambda} \sum_{n=1}^N \mathbb{E}_{Q(z; \lambda)} [\log p(y_n | x_n^T z)] - D_{KL} [Q(z; \lambda) || P(z)]$$

Data-term

Regularizer

Q: Equality?

Question

How does variational inference solve the computational problems with Bayesian Inference?

VI by using SGD

(SEE "VI in five lines in python")

$$\text{MAP: } \max_z \sum_{n=1}^N \log p(y_n / z^T x_n) + \log p(z)$$

$$\text{VI: } \max_{\lambda} \sum_{n=1}^N \mathbb{E}_{q(z; \lambda)} [\log p(y_n | x_n^T z)] - D_{KL} [q(z; \lambda) \parallel P(z)] = \mathcal{L}(\lambda)$$

How to compute an unbiased stochastic gradient?

Doubly-stochastic estimate: Sample an "n" at random
 Sample a $z_t^* \sim q(z; \lambda_t)$

$$L(\lambda) \approx \hat{L}(\lambda) = N \log p(y_n / x_n^T z_t^*) - D_{KL} [q(z_t^* | \lambda) \parallel P(z_t^*)]$$

$$\text{SGD: } \lambda_{t+1} = \lambda_t + \alpha_t \left. \frac{\partial \hat{L}}{\partial \lambda} \right|_{\lambda = \lambda_t}$$

gradient w.r.t. λ

Reparameterization trick: say $\lambda = \{m, \sigma^2\}$ then, $z^* = m + \sigma \epsilon$ where $\epsilon \sim \mathcal{N}(0, 1)$

Q: Issues?

Question

Describe one method to perform variational inference.

Discuss computational challenges associated with it.

Mean-Field Inference

ISSUE : size of λ could be $>$ size of z !

EXAMPLE :

SOLUTION :

$$\text{Mean-field : } q(z; \lambda) = \prod_{i=1}^D q_i(z_i; \lambda_i)$$

Q: Issues?

Variational Auto-Encoder

Kingma & Welling (ICLR 2014)

Consider unsupervised learning with data $y_n \in \mathbb{R}^D, n=1, \dots, N$

Variational Auto-Encoder

Variational Lower bound: $\mathbb{E}_{q_{\phi}}[\log p_{\theta}(y/z)] - D_{KL}[q_{\phi} \parallel p(z)]$

"SGD + variance reduction" works quite well!

SHOW RESULTS

Question

How does variational auto-encoder solve the computational problem with the standard mean-field inference?

Assignment

Question 1: Why is Bayesian inference computationally challenging?

Question 2: How does variational inference solve the computational problems with Bayesian Inference?

Question 3: Describe one method to perform variational inference. Discuss computational challenges associated with it?

~~Question 4: How does variational auto-encoder solve the computational problem with the standard mean-field inference?~~