

MAP estimate on GLMs

DATA :

LINEAR MODEL:

GENERALIZED :
(GLM)

MAXIMUM LIKELIHOOD :
(MLE)

MAP:

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

Stochastic Gradient Descent (SGD)

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

Randomly sample an n and compute an "unbiased" estimate

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

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

$$\text{Bayes rule: } \underset{\text{Posterior}}{p(\underline{z} | \underline{y})} = \frac{p(\underline{y}, \underline{z})}{\int p(\underline{y}, \underline{z}) d\underline{z}}$$

INTRACTABLE INTEGRAL

Computational issues

GLM:
$$p(\underline{z} | \underline{y}) = \int \prod_{n=1}^N p(y_n / \underline{z}^T \underline{x}_n) p(\underline{z}) d\underline{z}$$

POSTERIOR

The diagram shows the equation for the posterior distribution in a Generalized Linear Model (GLM). The term $p(\underline{z} | \underline{y})$ is labeled as the POSTERIOR. The equation is $p(\underline{z} | \underline{y}) = \int \prod_{n=1}^N p(y_n / \underline{z}^T \underline{x}_n) p(\underline{z}) d\underline{z}$. Red arrows point from the terms $p(y_n / \underline{z}^T \underline{x}_n)$ and $p(\underline{z})$ in the integrand to the word POSTERIOR.

CONJUGATE PRIORS:

Q: Comp. cost for Gaussian?

NON-CONJUGATE :

WE WILL CONVERT INTEGRATION \rightarrow OPTIMIZATION

Q: How to approximate?

Question

Why is Bayesian inference computationally challenging?

Variational Inference

MAIN IDEA: $P(\underline{z} | \underline{y}) \approx$ e.g.

LOWER BOUND:

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

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")

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

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

How to compute an unbiased stochastic gradient?

Doubly-stochastic estimate:

$$\text{SGD: } \lambda_{t+1} = \lambda_t + \alpha_t \widehat{\frac{\partial \mathcal{L}}{\partial \lambda}} \Big|_{\lambda = \lambda_t}$$

other methods

- VMP
- SVI

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?