

VI & BNN (1)

Stochastic Variational Inference (SVI) and Variational Autoencoder (VAE)

Keywords : Variational Inference, Scalable Variational Inference, ELBO,
Probabilistic Deep Learning, Variational Autoencoder

Seunghan Lee

February 01, 2021

Contents

1. Variational Inference

1. MCMC vs VI
2. Evidence Lower Bound (ELBO)
3. Mean Field Variational Inference (MFVI)

2. Stochastic Variational Inference

3. Variational Auto Encoder (VAE)

1. Structure of VAE
2. Update Decoder
3. Update Encoder (Log-derivative Trick & Reparameterization Trick)
4. Implementation using Pytorch

Paper List (+ references)

- Practical variational inference for neural networks (2011)
- Stochastic variational inference (2013)
- Auto-Encoding Variational Bayes (2013)
- Variational Inference : A review for statisticians (2017)
- Advances in variational inference (2018)
- Deep Bayes (<https://deepbayes.ru/>)

1. Variational Inference

1. Variational Inference

1-1. MCMC vs Variational Inference

Two main approaches to find the **(intractable) posterior** in Bayesian Inference!

(1) MCMC : sampling from the unnormalized posterior

- (Pros) Unbiased
- (Cons) High computational cost

(2) Variational Inference : Approximating target distn with a simpler distn

- (Pros) Faster & Scalable
- (Cons) Biased

1. Variational Inference

1-1. MCMC vs Variational Inference

Will be going to focus on **Variational Inference**

Before getting on...

- **KL divergence**

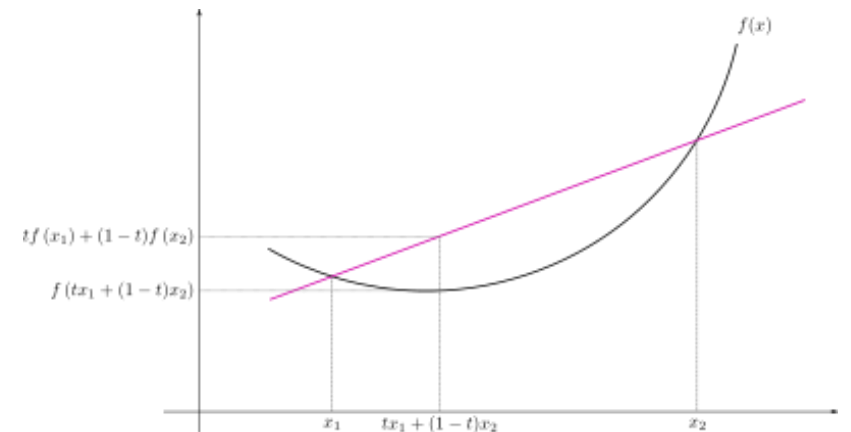
$$KL(q(\theta) \| p(\theta | x)) = \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta$$

- $KL(q \| p) \geq 0$
- $KL(q \| p) = 0 \Leftrightarrow q = p$
- $KL(q \| p) \neq KL(p \| q)$

- **Jensen's Inequality**

If $g(x)$ is a convex function on R_X and $E[g(X)]$ and $g(E[X])$ are finite,

$$E[g(X)] \geq g(E[X])$$



1. Variational Inference

1-2. Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x) &= \int q(\theta) \log p(x) d\theta = \int q(\theta) \log \frac{p(x, \theta)}{p(\theta | x)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta) q(\theta)}{p(\theta | x) q(\theta)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta = \\ &= \mathcal{L}(q(\theta)) + KL(q(\theta) \| p(\theta | x)) \text{ (Non-negative)}\end{aligned}$$

1. Variational Inference

1-2. Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x) &= \int q(\theta) \log p(x) d\theta = \int q(\theta) \log \frac{p(x, \theta)}{p(\theta | x)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta) q(\theta)}{p(\theta | x) q(\theta)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta = \\ &= \boxed{\mathcal{L}(q(\theta))} + \underbrace{KL(q(\theta) \| p(\theta | x))}_{\text{(Non-negative) }}\end{aligned}$$

$$\begin{aligned}\log p(x) &= \log \int_z p(x, z) \\ &= \log \int_z p(x, z) \frac{q(z)}{q(z)} \\ &= \log \left(\mathbb{E}_q \left[\frac{p(x, z)}{q(z)} \right] \right) \\ &\geq \mathbb{E}_q [\log p(x, z)] - \mathbb{E}_q [\log q(z)]\end{aligned}$$

Jensen's Inequality 

ELBO (Evidence Lower Bound)

Why called like that? Think of "Evidence" in Bayes rule!

$$\log p(x) \geq \mathcal{L}(q(\theta))$$

1. Variational Inference

1-2. Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x) &= \int q(\theta) \log p(x) d\theta = \int q(\theta) \log \frac{p(x, \theta)}{p(\theta | x)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta) q(\theta)}{p(\theta | x) q(\theta)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta = \\ &= \boxed{\mathcal{L}(q(\theta))} + \underbrace{KL(q(\theta) \| p(\theta | x))}_{\text{Non-negative}}\end{aligned}$$

Minimizing KL-divergence

ELBO (Evidence Lower Bound)

Why called like that? Think of "Evidence" in Bayes rule!

$$\log p(x) \geq \mathcal{L}(q(\theta))$$

1. Variational Inference

1-2. Evidence Lower Bound (ELBO)

$$\begin{aligned}\log p(x) &= \int q(\theta) \log p(x) d\theta = \int q(\theta) \log \frac{p(x, \theta)}{p(\theta | x)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta) q(\theta)}{p(\theta | x) q(\theta)} d\theta = \\ &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta + \int q(\theta) \log \frac{q(\theta)}{p(\theta | x)} d\theta = \\ &= \boxed{\mathcal{L}(q(\theta))} + \underbrace{KL(q(\theta) \| p(\theta | x))}_{\text{(Non-negative) }}\end{aligned}$$

Maximizing ELBO

ELBO (Evidence Lower Bound)

Why called like that? Think of "Evidence" in Bayes rule!

$$\log p(x) \geq \mathcal{L}(q(\theta))$$

1. Variational Inference

1-2. Evidence Lower Bound (ELBO)

Interpretation of ELBO

Rewrite ELBO as below

$$\begin{aligned}\mathcal{L}(q(\theta)) &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta = \int q(\theta) \log \frac{p(x | \theta)p(\theta)}{q(\theta)} d\theta = \\ &= \int q(\theta) \log p(x | \theta) d\theta + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta = \\ &= \mathbb{E}_{q(\theta)} \log p(x | \theta) - KL(q(\theta) \| p(\theta))\end{aligned}$$

1. Variational Inference

1-2. Evidence Lower Bound (ELBO)

Interpretation of ELBO

Rewrite ELBO as below

$$\begin{aligned}\mathcal{L}(q(\theta)) &= \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta = \int q(\theta) \log \frac{p(x | \theta)p(\theta)}{q(\theta)} d\theta = \\ &= \int q(\theta) \log p(x | \theta) d\theta + \int q(\theta) \log \frac{p(\theta)}{q(\theta)} d\theta = \\ &= \underbrace{\mathbb{E}_{q(\theta)} \log p(x | \theta)}_{\text{Term 1)} \text{ Encourage good fit!}} - \underbrace{KL(q(\theta) \| p(\theta))}_{\text{Term 2)} \text{ Regularize! Encourage posterior to be close to prior}}\end{aligned}$$

Term 1) Encourage good fit!

Term 2) Regularize! Encourage posterior to be close to prior

1. Variational Inference

1-3. Mean Field Variational Inference (MFVI)

We have to optimize w.r.t ELBO, but how?

$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$

Mean Field Assumption


(For simplicity, variational parameters are factorized as below, with an Independent Assumption)

$$q(\theta) = \prod_{j=1}^m q_j(\theta_j), \quad \theta = [\theta_1, \dots, \theta_m]$$

Due to the assumption, its flexibility is limited!


1. Variational Inference

1-3. Mean Field Variational Inference (MFVI)

$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$
$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) = q_1(\theta_1) \cdots q_m(\theta_m)}$$


1. Variational Inference

1-3. Mean Field Variational Inference (MFVI)

$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) \in \mathcal{Q}}$$
$$\mathcal{L}(q(\theta)) = \int q(\theta) \log \frac{p(x, \theta)}{q(\theta)} d\theta \rightarrow \max_{q(\theta) = q_1(\theta_1) \cdots q_m(\theta_m)}$$


Solution (proof on the next page)

$$q_j(\theta_j) = r_j(\theta_j) = \frac{1}{Z_j} \exp \left(\mathbb{E}_{q_{i \neq j}} \log p(x, \theta) \right)$$

Coordinate-ascent method

CAVI (Coordinate Ascent Variational Inference)

1. Variational Inference

1-3. Mean Field Variational Inference (MFVI)

$$q_j(\theta_j) = r_j(\theta_j) = \frac{1}{Z_j} \exp(\mathbb{E}_{q_{i \neq j}} \log p(x, \theta))$$

Proof)

$$\begin{aligned} \mathcal{L}(q(\theta)) &= \mathbb{E}_{q(\theta)} \log p(x, \theta) - \mathbb{E}_{q(\theta)} \log q(\theta) = \\ &= \mathbb{E}_{q(\theta)} \log p(x, \theta) - \sum_{k=1}^m \mathbb{E}_{q_k(\theta_k)} \log q_k(\theta_k) = \\ &= \mathbb{E}_{q_j(\theta_j)} [\mathbb{E}_{q_{i \neq j}} \log p(x, \theta)] - \mathbb{E}_{q_j(\theta_j)} \log q_j(\theta_j) + Const = \\ &= \left\{ r_j(\theta_j) = \frac{1}{Z_j} \exp(\mathbb{E}_{q_{i \neq j}} \log p(x, \theta)) \right\} = \\ &= \mathbb{E}_{q_j(\theta_j)} \log \frac{r_j(\theta_j)}{q_j(\theta_j)} + Const = -KL(q_j(\theta_j) \| r_j(\theta_j)) + Const \end{aligned}$$

1. Variational Inference

1-3. Mean Field Variational Inference (MFVI)

Algorithm

Initialize $q(\theta) = \prod_{j=1}^m q_j(\theta_j)$

Iterations:

- Update each factor q_1, \dots, q_m :

$$q_j(\theta_j) = \frac{1}{Z_j} \exp(\mathbb{E}_{q_{i \neq j}} \log p(x, \theta))$$

- Compute ELBO $\mathcal{L}(q(\theta))$

Repeat until convergence of ELBO

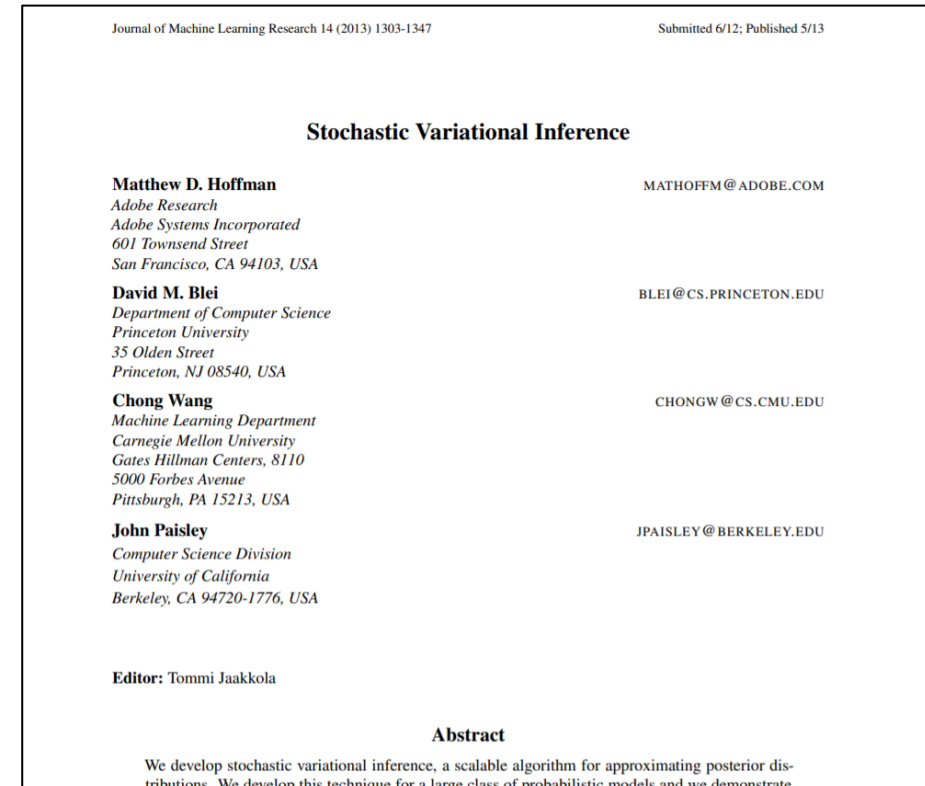
2. Stochastic Variational Inference

2. Stochastic Variational Inference

Stochastic Variational Inference (SVI)?

a stochastic optimization algorithm for mean-field variational inference,
that can **handle massive dataset (scalability)**

- Mean Field Variational Inference
- Stochastic Optimization



2. Stochastic Variational Inference

Mean Field Variational Inference

(This time, we will set hidden variables into 2 parts)

Factorize joint distribution

$$p(x, z, \beta | \alpha) = p(\beta | \alpha) \prod_{n=1}^N p(x_n, z_n | \beta)$$

Approximate using MFVI

$$p(z, \beta | x) = \frac{p(x, z, \beta)}{\int p(x, z, \beta) dz d\beta} \longleftarrow$$

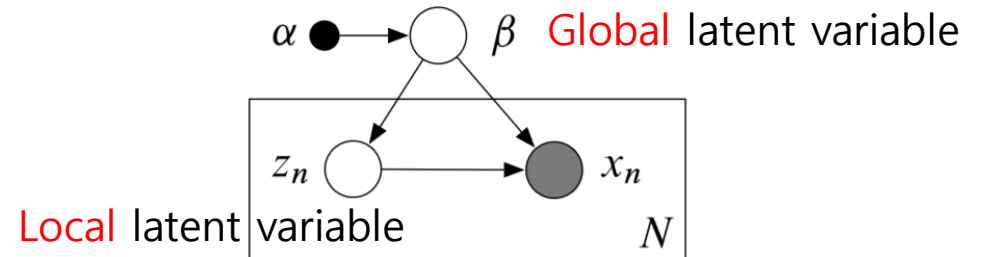


Figure 2: A graphical model with observations $x_{1:N}$, local hidden variables $z_{1:N}$ and global hidden variables β . The distribution of each observation x_n only depends on its corresponding local variable z_n and the global variables β . (Though not pictured, each hidden variable z_n , observation x_n , and global variable β may be a collection of multiple random variables.)

Approximating distn :

$$q(z, \beta) = q(\beta | \lambda) \prod_{n=1}^N \prod_{j=1}^J q(z_{nj} | \phi_{nj})$$

2. Stochastic Variational Inference

Mean Field Variational Inference

$$q(z, \beta) = q(\beta | \lambda) \prod_{n=1}^N \prod_{j=1}^J q(z_{nj} | \phi_{nj})$$

Set $q(\beta | \lambda)$ and $q(z_{nj} | \phi_{nj})$ to be in the same exponential family

- $q(\beta | \lambda) = h(\beta) \exp\{\lambda^\top t(\beta) - a_g(\lambda)\}$
- $q(z_{nj} | \phi_{nj}) = h(z_{nj}) \exp\{\phi_{nj}^\top t(z_{nj}) - a_\ell(\phi_{nj})\}$

$$\nabla_\lambda \mathcal{L} = \nabla_\lambda^2 a_g(\lambda) (\mathbb{E}_q[\eta_g(x, z, \alpha)] - \lambda) \quad \longrightarrow \quad \lambda = \mathbb{E}_q[\eta_g(x, z, \alpha)]$$

$$\nabla_{\phi_{nj}} \mathcal{L} = \nabla_{\phi_{nj}}^2 a_\ell(\phi_{nj}) (\mathbb{E}_q[\eta_\ell(x_n, z_{n,-j}, \beta)] - \phi_{nj}) \quad \longrightarrow \quad \phi_{nj} = \mathbb{E}_q[\eta_\ell(x_n, z_{n,-j}, \beta)]$$

2. Stochastic Variational Inference

```
1: Initialize  $\lambda^{(0)}$  randomly.
2: repeat
3:   for each local variational parameter  $\phi_{nj}$  do
4:     Update  $\phi_{nj}$ ,  $\phi_{nj}^{(t)} = \mathbb{E}_{q^{(t-1)}}[\eta_{\ell,j}(x_n, z_{n,-j}, \beta)]$ .
5:   end for
6:   Update the global variational parameters,  $\lambda^{(t)} = \mathbb{E}_{q^{(t)}}[\eta_g(z_{1:N}, x_{1:N})]$ .
7: until the ELBO converges
```

1. Variational Inference

1-3. Mean Field Variational Inference (MFVI)

Algorithm

Initialize $q(\theta) = \prod_{j=1}^m q_j(\theta_j)$

Iterations:

- Update each factor q_1, \dots, q_m :

$$q_j(\theta_j) = \frac{1}{Z_j} \exp(\mathbb{E}_{q_{i \neq j}} \log p(x, \theta))$$

- Compute ELBO $\mathcal{L}(q(\theta))$

Repeat until convergence of ELBO

Figure 3: Coordinate ascent mean-field variational inference.

2. Stochastic Variational Inference

```
1: Initialize  $\lambda^{(0)}$  randomly.
2: repeat
3:   for each local variational parameter  $\phi_{nj}$  do
4:     Update  $\phi_{nj}$ ,  $\phi_{nj}^{(t)} = \mathbb{E}_{q^{(t-1)}}[\eta_{\ell,j}(x_n, z_{n,-j}, \beta)]$ .
5:   end for
6:   Update the global variational parameters,  $\lambda^{(t)} = \mathbb{E}_{q^{(t)}}[\eta_g(z_{1:N}, x_{1:N})]$ .
7: until the ELBO converges
```

Each data has its OWN local variational parameter

1. Variational Inference

1-3. Mean Field Variational Inference (MFVI)

Algorithm

Initialize $q(\theta) = \prod_{j=1}^m q_j(\theta_j)$

Iterations:

- Update each factor q_1, \dots, q_m :

$$q_j(\theta_j) = \frac{1}{Z_j} \exp(\mathbb{E}_{q_{i \neq j}} \log p(x, \theta))$$

- Compute ELBO $\mathcal{L}(q(\theta))$

Repeat until convergence of ELBO

Figure 3: Coordinate ascent mean-field variational inference.

INEFFICIENT for large data sets!

(should optimize the local variational params for each data, before re-estimating the global variational params)

SVI uses “**stochastic optimization**” to fit global variational parameters.

2. Stochastic Variational Inference

Algorithm

- 1: Initialize $\lambda^{(0)}$ randomly.
- 2: Set the step-size schedule ρ_t appropriately.
- 3: **repeat**
- 4: Sample a data point x_i uniformly from the data set.
- 5: Compute its local variational parameter,

$$\phi = \mathbb{E}_{\lambda^{(t-1)}} [\eta_g(x_i^{(N)}, z_i^{(N)})].$$

- 6: Compute intermediate global parameters as though x_i is replicated N times,

$$\hat{\lambda} = \mathbb{E}_{\phi} [\eta_g(x_i^{(N)}, z_i^{(N)})].$$

- 7: Update the current estimate of the global variational parameters,

$$\lambda^{(t)} = (1 - \rho_t) \lambda^{(t-1)} + \rho_t \hat{\lambda}.$$

- 8: **until** forever

Simple! Just think it as

(1) SGD + (2) VI

$$\begin{aligned} \lambda^{(t)} &= \lambda^{(t-1)} + \rho_t (\hat{\lambda}_t - \lambda^{(t-1)}) \\ &= (1 - \rho_t) \lambda^{(t-1)} + \rho_t \hat{\lambda}_t. \end{aligned}$$

Figure 4: Stochastic variational inference.

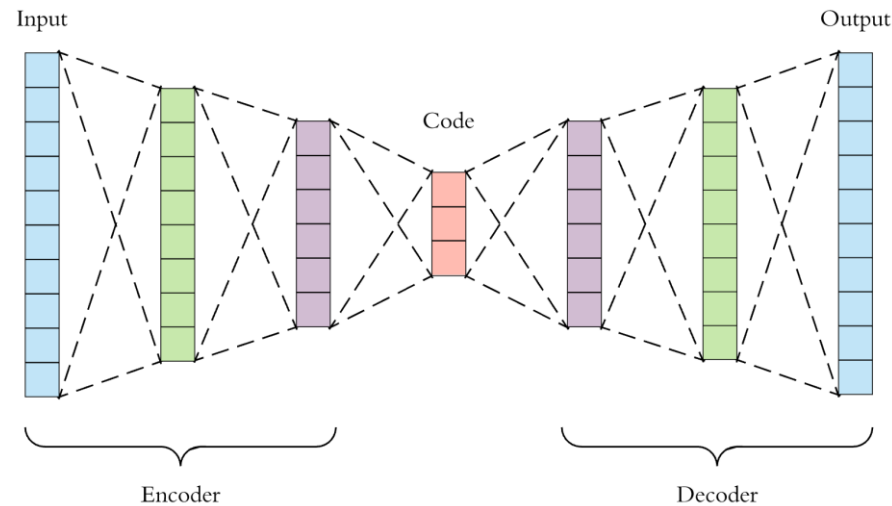
3. Variational Auto Encoder (VAE)

3. Variational Auto Encoder

3-1. Structure of VAE

Auto Encoder : "The aim of an **autoencoder** is to learn a **representation (encoding)** for a set of data "

What's the difference between AE & VAE?



Auto Encoder (AE)

3. Variational Auto Encoder

3-1. Structure of VAE

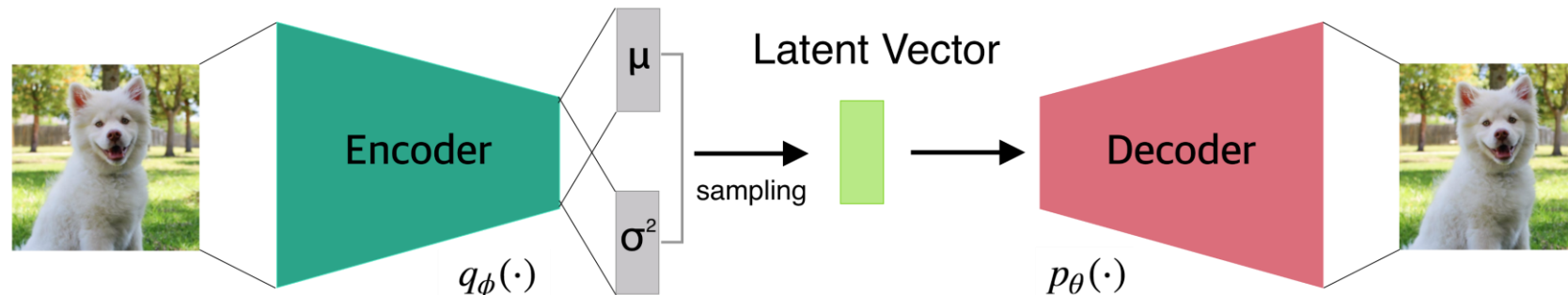
(1) Encoder

- inference network
- input : x , output : z
- $q(z \mid x, \phi)$

(2) Decoder

- generative network
- input : z , output : as closely as x
- $p(x \mid z, \theta)$

Variational Autoencoder



3. Variational Auto Encoder

3-1. Structure of VAE

Variational Inference : approximate with variational distribution!

$$p(z_i | x_i, \theta) \approx q(z_i | x_i, \phi) = \prod_{j=1}^d N(z_{ij} | \underline{\mu_j(x_i)}, \underline{\sigma_j^2(x_i)})$$

(use Neural Network! Flexible!)

Objective Function (ELBO)

- Minimize KL-divergence $q(Z | X, \phi) = \arg \min_{\phi} KL(q(Z | X, \phi) || p(Z | X, \theta))$

= Maximize ELBO $\mathcal{L}(\phi, \theta) = \int q(Z | X, \phi) \log \frac{p(X|Z, \theta)p(Z)}{q(Z|X, \phi)} dZ \rightarrow \max_{\phi, \theta}$

3. Variational Auto Encoder

3-1. Structure of VAE

Update the parameters of **Encoder & Decoder**.

Due to flexible & complex model (Neural Network), it seems hard to solve...

But using some techniques (+ tricks), we can solve it!

Stochastic Optimization

- 1) Mini-batch
- 2) Monte Carlo Estimation

Tricks

- 1) **Log-derivative trick**
- 2) **Reparameterization trick**

3. Variational Auto Encoder

3-2. Update Decoder (parameter : θ)

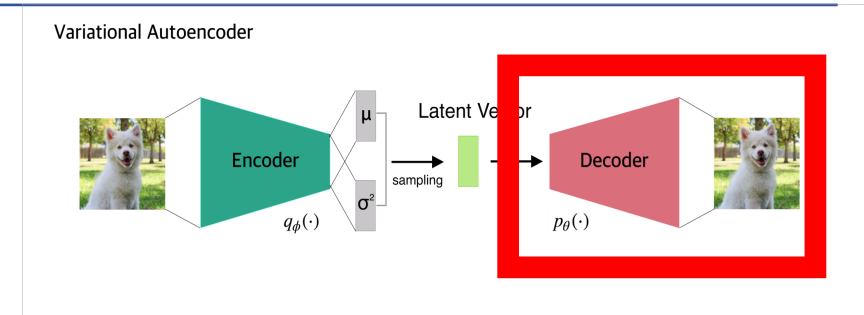
Not that hard to solve!

(1) Mini-batch

$$\begin{aligned}\nabla_{\theta} \mathcal{L}(\phi, \theta) &= \nabla_{\theta} \sum_{i=1}^n \int q(z_i | x_i, \phi) \log \frac{p(x_i | z_i, \theta) p(z_i)}{q(z_i | x_i, \phi)} dz_i \\ &= \sum_{i=1}^n \int q(z_i | x_i, \phi) \nabla_{\theta} \log p(x_i | z_i, \theta) dz_i \\ &\approx n \int q(z_i | x_i, \phi) \nabla_{\theta} \log p(x_i | z_i, \theta) dz_i, \quad i \sim \mathcal{U}\{1, \dots, n\}\end{aligned}$$

(2) Monte-Carlo estimation

$$n \int q(z_i | x_i, \phi) \nabla_{\theta} \log p(x_i | z_i, \theta) dz_i \approx n \nabla_{\theta} \log p(x_i | z_i^*, \theta), \quad z_i^* \sim q(z_i | x_i, \phi)$$



3. Variational Auto Encoder

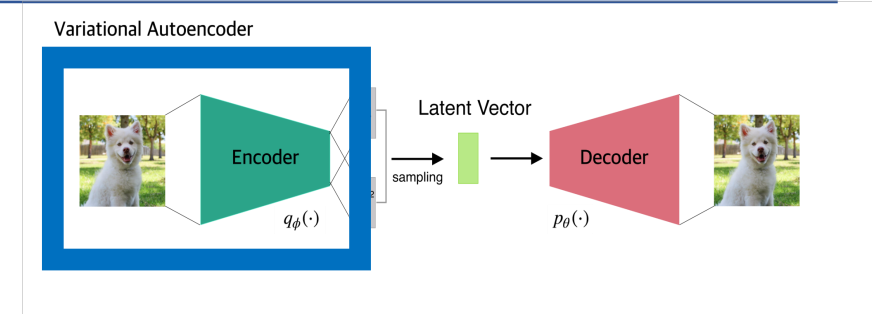
3-3. Update Encoder (parameter : ϕ)

A bit trickier than updating Encoder... ☹

$$\nabla_{\phi} \mathcal{L}(\phi, \theta) = \nabla_{\phi} \sum_{i=1}^n \int q(z_i | x_i, \phi) \log \frac{p(x_i | z_i, \theta) p(z_i)}{q(z_i | x_i, \phi)} dz_i$$

$$\neq \sum_{i=1}^n \int q(z_i | x_i, \phi) \nabla_{\phi} \log p(x_i | z_i, \theta) dz_i$$

We need some tricks to solve this!



3. Variational Auto Encoder

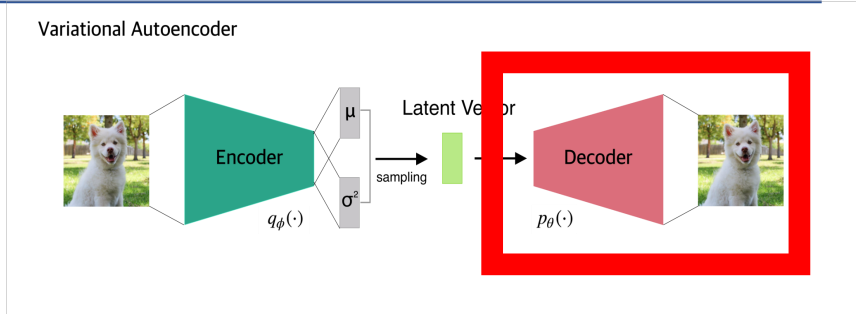
3-3. Update Encoder (parameter : ϕ)

a. Log Derivate Trick

$$\frac{\partial}{\partial x} p(y | x) = p(y | x) \frac{\partial}{\partial x} \log p(y | x)$$

For question like solving $E_{y|x} h(x, y)$

$$\begin{aligned} \frac{\partial}{\partial x} \int p(y | x) h(x, y) dy &= \int \frac{\partial}{\partial x} (p(y | x) h(x, y)) dy \\ &= \int \left(h(x, y) \frac{\partial}{\partial x} p(y | x) + p(y | x) \frac{\partial}{\partial x} h(x, y) \right) dy \\ &= \int p(y | x) \frac{\partial}{\partial x} h(x, y) dy + \int h(x, y) \frac{\partial}{\partial x} p(y | x) dy \end{aligned}$$



3. Variational Auto Encoder

3-3. Update Encoder (parameter : ϕ)

a. Log Derivate Trick

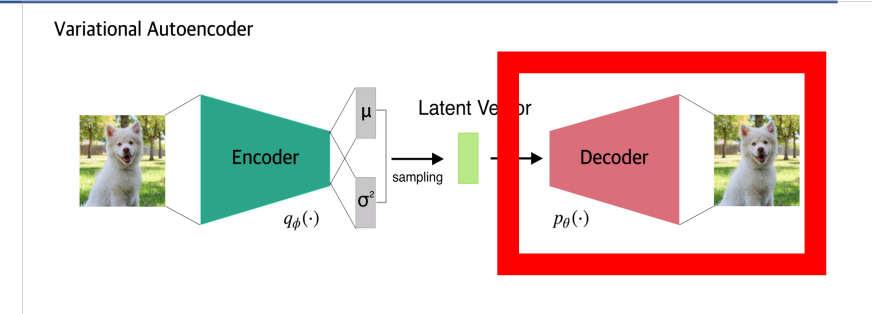
$$\begin{aligned}\frac{\partial}{\partial x} \int p(y | x) h(x, y) dy &= \int \frac{\partial}{\partial x} (p(y | x) h(x, y)) dy \\ &= \int \left(h(x, y) \frac{\partial}{\partial x} p(y | x) + p(y | x) \frac{\partial}{\partial x} h(x, y) \right) dy \\ &= \int \underbrace{p(y | x) \frac{\partial}{\partial x} h(x, y) dy}_{\text{First term}} + \int \underbrace{h(x, y) \frac{\partial}{\partial x} p(y | x) dy}_{\text{Second term}}\end{aligned}$$

First term

- $\int p(y | x) \frac{\partial}{\partial x} h(x, y) dy$ - easy to solve (using Monte Carlo estimation)

Second term

- $\int h(x, y) \frac{\partial}{\partial x} p(y | x) dy$ - hard to solve - but can be solved using "Log-derivative trick"

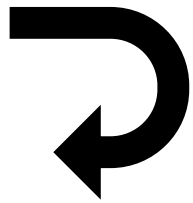


3. Variational Auto Encoder

3-3. Update Encoder (parameter : ϕ)

a. Log Derivate Trick

Apply it to our derivative of ELBO!

$$\begin{aligned}\frac{\partial}{\partial \phi} \mathcal{L}(\phi, \theta) &= \frac{\partial}{\partial \phi} \int q(Z | X, \phi) \log p(X | Z, \theta) dZ \\ &\approx n \log p(x_i | z_i^*, \theta) \frac{\partial}{\partial \phi} \log q(z_i^* | x_i, \phi)\end{aligned}$$


1. Mini batching
2. Log-derivative Trick

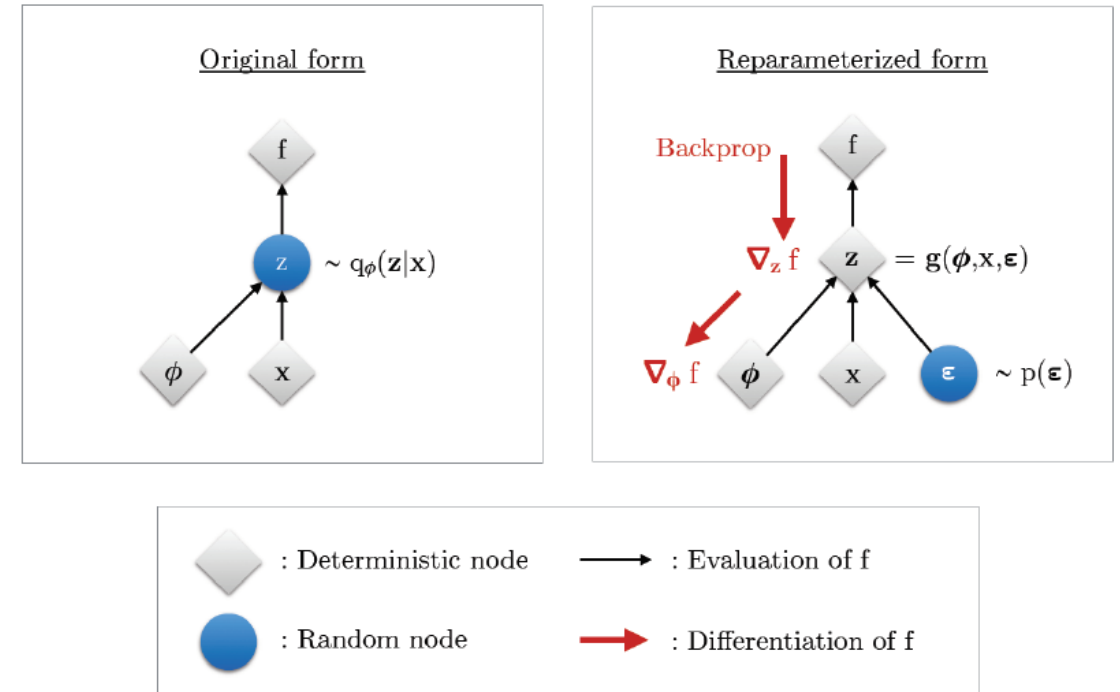
3. Variational Auto Encoder

3-3. Update Encoder (parameter : ϕ)

b. Reparameterization Trick

We can not backpropagate through "random" variable!
It should be deterministic!

$$\begin{aligned}\frac{\partial}{\partial x} \int p(y | x) h(x, y) dy &= \frac{\partial}{\partial x} \int r(\epsilon) h(x, g(\epsilon, x)) d\epsilon \\ &\approx \frac{d}{dx} h(x, g(x, \hat{\epsilon})) \\ &= \frac{\partial}{\partial x} h(x, g(x, \hat{\epsilon})) + \frac{\partial}{\partial g} h(x, g(x, \hat{\epsilon})) \frac{\partial}{\partial x} g(x, \hat{\epsilon}) \quad \text{where } \hat{\epsilon} \sim r(\epsilon)\end{aligned}$$



3. Variational Auto Encoder

3-3. Update Encoder (parameter : ϕ)

b. Reparameterization Trick

Apply it to our derivative of ELBO!

$$\begin{aligned} n \frac{\partial}{\partial \phi} \int q(z_i | x_i, \phi) \log p(x_i | z_i, \theta) &= n \frac{\partial}{\partial \phi} \int r(\epsilon) \log p(x_i | g(\epsilon, x_i, \phi), \theta) d\epsilon \\ &\approx n \frac{\partial}{\partial \phi} \log p(x_i | g(\hat{\epsilon}, x_i, \phi), \theta), \quad \text{where } \hat{\epsilon} \sim r(\epsilon) \end{aligned}$$

(simple example)

$$q_{\phi}(z_i | x_i) = N(\mu_i, \sigma_i^2 \mathbf{I})$$

$$\Rightarrow z_i = \mu_i + \sigma_i \odot \epsilon_i \quad \text{where } \epsilon_i \sim N(\mathbf{0}, \mathbf{I})$$

3. Variational Auto Encoder

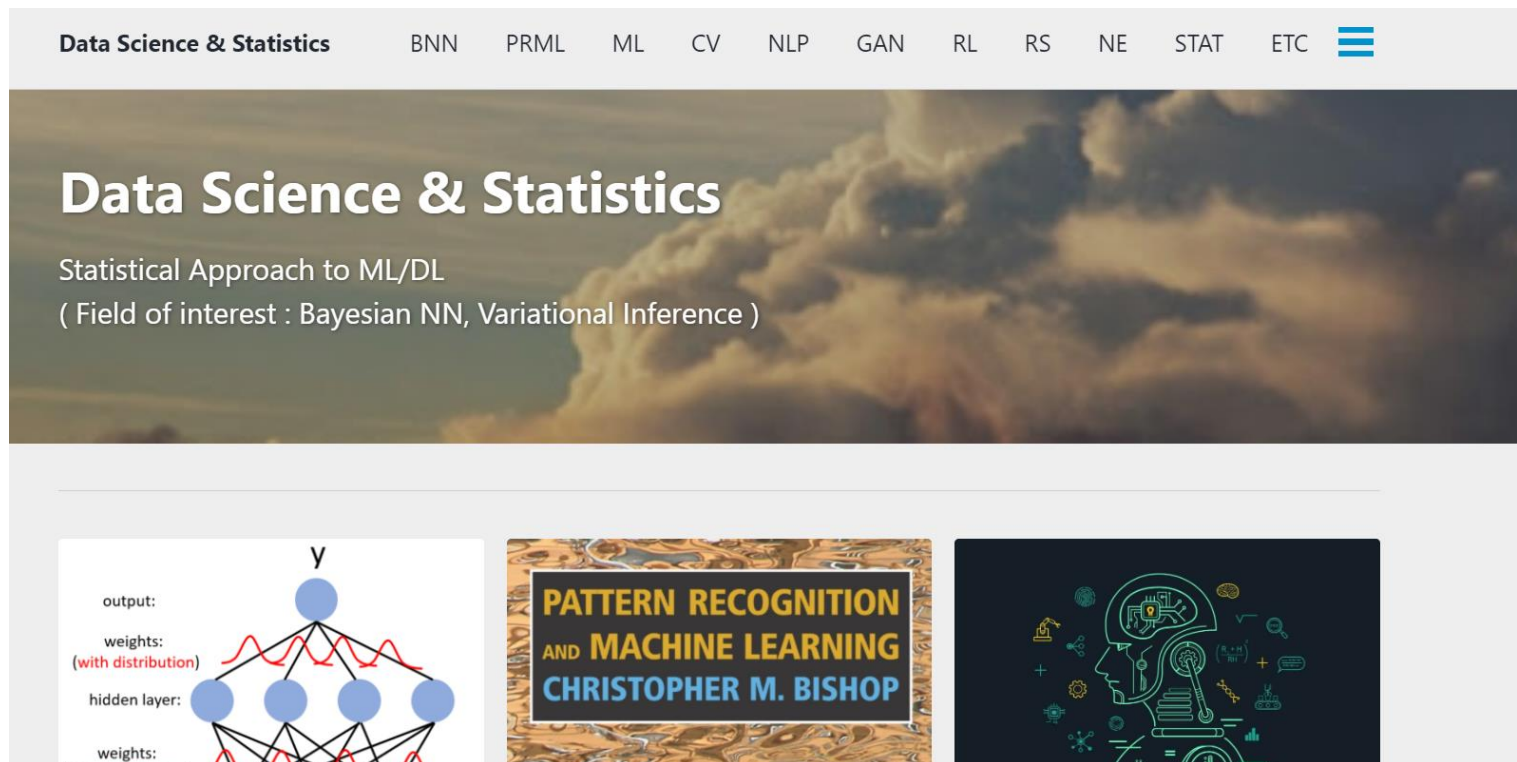
Summary of VAE

- step 1) sample $i \sim \mathcal{U}\{1, \dots, n\}$
- step 2) compute stochastic gradient of ELBO (w.r.t θ and ϕ)
 - Update θ (decoder parameter)
 - stoch. grad $_{\theta} \mathcal{L}(\phi, \theta) = n \frac{\partial}{\partial \theta} \log p(x_i | z_i^*, \theta)$
where $z_i^* \sim q(z_i | x_i, \phi)$
 - Update ϕ (encoder parameter)
 - stoch.grad $_{\phi} \mathcal{L}(\phi, \theta) = n \frac{\partial}{\partial \phi} \log p(x_i | g(\hat{\epsilon}, x_i, \phi), \theta) - \frac{\partial}{\partial \phi} KL(q(z_i | x_i, \phi) || p(z_i))$
where $\hat{\epsilon} \sim r(\epsilon)$
- Update until stopping criterion reaches

3. Variational Auto Encoder

3-4. Implementation using Pytorch

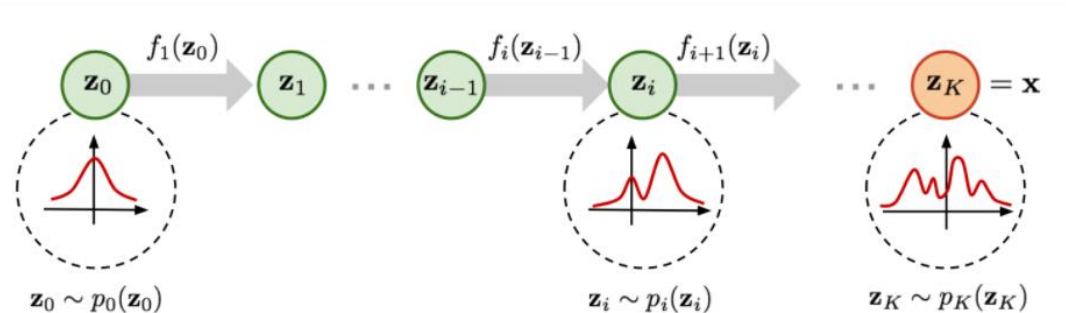
<https://seunghan96.github.io/stat/gan/bnn/code-6.Variational-Auto-Encoder/>



Summary

Have dealt with **basic concepts** to know before reading papers about various **variational inference methods**, **Bayesian NN**.

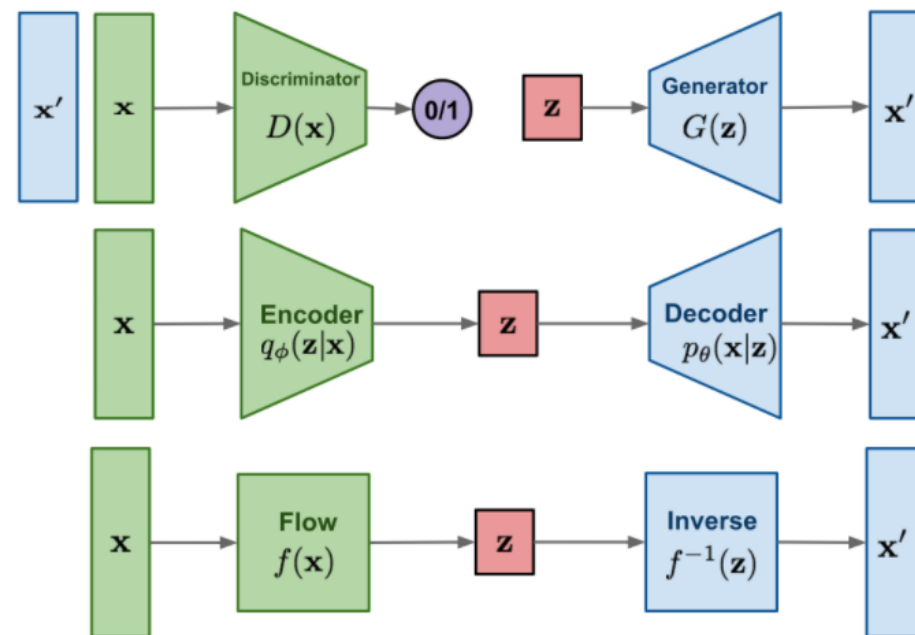
Next Presentation : **Normalizing Flow**



GAN: minimax the classification error loss.

VAE: maximize ELBO.

Flow-based generative models: minimize the negative log-likelihood



Thank you !

Review of Papers regarding VI/BNN

(+ some Statistical Models / Machine Learning / Deep Learning)

can be found in my github blog below :)

<https://seunghan96.github.io/>

(for more about VI/BNN <https://seunghan96.github.io/categories/bnn/>)