

Bayesian Neural Net.

- 1) W : no longer deterministic. $\{m_w, \Sigma_w\}$. ss
 2) Given $D = \{x^i, t^i\}_{i=1}^N$

- 3) Likelihood

$$p(D|w, \beta) \propto \exp(-\beta E_D(w)) \text{, where } E_D(w) = \frac{1}{2} \sum (f(x^i, w) - t^i)^2$$

- 4) Prior

$$p(w|A) \propto \exp(-E_w(w)) \text{, where } E_w(w) = \frac{1}{2} w^T A w$$

- 5) Posterior

$$p(w|D, A, \beta) \propto \exp(-\beta E_D(w) - E_w(w))$$

- 6) Evaluation at x_*

$$\langle f(x_*) \rangle = \int_w f(x_*, w) p(w|D, A, \beta) dw$$

- 7)

- Laplace's Apprx.

- MCMC

- Ensemble Learning \Leftarrow Variation Inference.

Introduce $Q(w)$ that approximates $p(w|.)$

$$\ln p(D|A, \beta) = \ln \int p(D|w, A, \beta) p(w|A) dw$$

$$= \ln \int \left(\frac{p(D|w, A, \beta) p(w|A)}{Q(w)} \right) Q(w) dw$$

$$Q(\mathbb{E}[x]) \leq \mathbb{E}[Q(x)]$$

$$\geq \int \ln \left(\frac{p(D|w, A, \beta) p(w|A)}{Q(w)} \right) Q(w) dw \quad (*)$$

$$= \mathbb{E}[Q]$$

maximize $Q(w)$

$\partial \ln p(w|.)$ $\nwarrow \searrow$

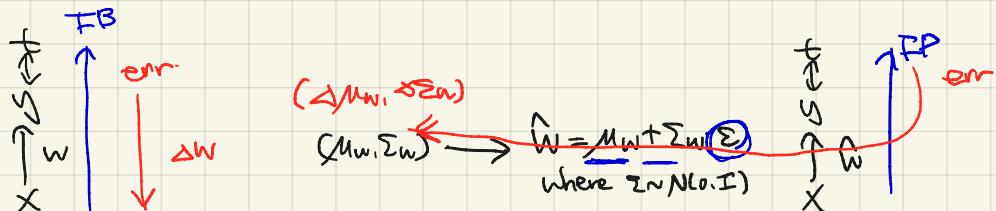
$$\frac{\text{KL}(Q||P)}{\mathbb{E}[Q]}$$

$$\ln p(D|A, \beta)$$

$$\begin{aligned}
8) \ln p(D) &= \int \ln p(D) q(w|\theta) dw \\
&= \int q(w|\theta) \ln \frac{p(D)p(w|D)}{p(w|D)} dw \\
&= \int q(w|\theta) \ln \frac{p(w)p(D|w)}{p(w|D)} dw \\
&= \int q(w|\theta) \ln \frac{q(w|\theta)p(w)p(D|w)}{q(w|\theta)p(w|D)} dw \\
&= \int q(w|\theta) \ln \frac{q(w|\theta)}{p(w|D)} dw + \int q(w|\theta) \ln \frac{p(D|w)p(w)}{q(w|\theta)} dw \\
&= \boxed{\text{KL}(q(w|\theta) || p(w|D))} + \boxed{\mathbb{E}_{q(\cdot)}[\ln p(D)]} \uparrow \\
&\quad \text{KLD b/w } q(\cdot) \& \text{ posterior.} \downarrow \quad \text{ELBO}
\end{aligned}$$

$$\begin{aligned}
9) \boxed{\mathbb{E}_{q(\cdot)}[\ln p(D)]} &= \int q(w|\theta) \ln \frac{p(D|w)p(w)}{q(w|\theta)} dw \\
&= \int \left(-\ln p(D|w) + \ln \frac{p(w)}{q(w|\theta)} \right) q(w|\theta) dw \\
&= \boxed{\mathbb{E}_{q(w|\theta)}[-\ln p(D|w)]} - \boxed{\text{KL}(q(w|\theta) || p(w))} \downarrow \\
&\quad \text{Expectation of LogLik. under } q. \quad \text{prior fitting term,}
\end{aligned}$$

$$10) \ln p(D) = \boxed{\text{KL}(q(w|\theta) || p(w|D))} + \boxed{\mathbb{E}_{q(\cdot)}[\ln p(D)]} \uparrow$$



11) Weight Uncertainty in Neural Networks. (Being Bayesian by Backpropagation) (BBB)

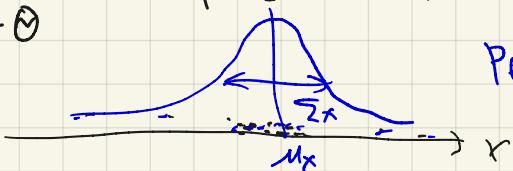
$$\begin{aligned}
\mathbb{E}_{q(\cdot)}[\ln p(D)] &= \int q(w|\theta) \ln \frac{p(D|w)p(w)}{q(w|\theta)} dw \\
&= \mathbb{E}_q \left[\ln p(D|w) + \ln p(w) - \ln q(w|\theta) \right] \\
&\approx \left(\underbrace{\ln p(D|\tilde{w})}_{\text{log lik}} + \underbrace{\ln p(\tilde{w})}_{\text{log prior}} - \underbrace{\ln q(\tilde{w}|\theta)}_{\text{tractable.}} \right), \quad \tilde{w} \sim q(w|\theta)
\end{aligned}$$

Summary: $X \xrightarrow{w} Y$, $p(w|D) \sim Q(w) = N(w; \underline{\mu}_w, \Sigma_w)$
 optimizel.

(2) Deep Latent Variable Model (DLVM), e.g., VAE.

- $D = \{x_i\}_{i=1}^N \rightarrow p_\theta(x)$ to sample $\tilde{x} \sim p_\theta(x)$

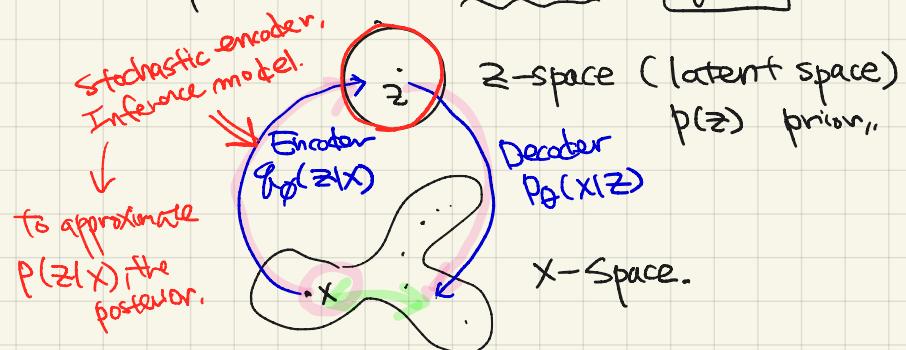
- How? $\theta \leftarrow \underset{\theta \in \Theta}{\operatorname{argmax}} \ln p(D|\theta)$
 example)



$$p_\theta(x) = N(x; \mu_x, \Sigma_x)$$

$$\theta = \{\mu_x, \Sigma_x\}$$

- $P(z) : z \rightarrow x : \boxed{p_\theta(x|z)}$, $p_\theta(x) = \int p_\theta(x|z)p(z)dz$



Goal is to find $q_\phi(z|x)$ that approximates $p(z|x)$.

$$\text{Similarly, } \ln p_\theta(x) = \mathbb{E}_{q_\phi(z|x)} [\ln p_\theta(x)]$$

$$= \mathbb{E}_{q_\phi(z|x)} \left[\ln \frac{p_\theta(x,z)}{p_\theta(z|x)} \right]$$

$$= \mathbb{E}_{q_\phi(z|x)} \left[\ln \frac{\frac{p_\theta(x,z)}{q_\phi(z|x)}}{\frac{p_\theta(z|x)}{p_\theta(z|x)}} \right]$$

$$= \mathbb{E}_{q_\phi(z|x)} \left[\ln \frac{p_\theta(x|z)}{q_\phi(z|x)} \right] + \mathbb{E}_q \left[\ln \frac{q_\phi(z|x)}{p_\theta(z|x)} \right]$$

$$\text{D}\text{KL}(q_\phi(z|x) \parallel p_\theta(z|x))$$

Intracode.

$$\text{ELBO} = \mathbb{E}_{q_\phi(z|x)} \left[\ln \frac{p_\theta(x,z)}{q_\phi(z|x)} \right]$$

$$= \int \ln \frac{\frac{p_\theta(x|z)p(z)}{q_\phi(z|x)}}{q_\phi(z|x)} q_\phi(z|x) dz$$

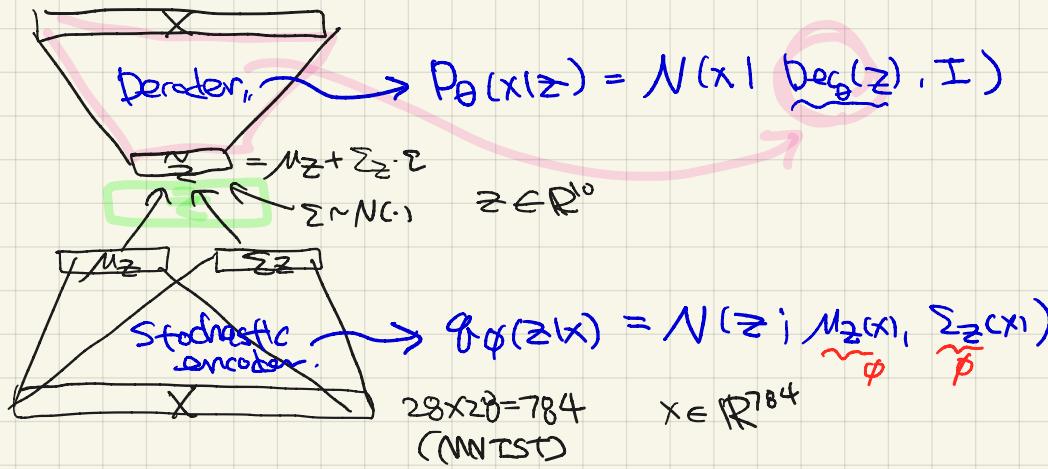
$$= \underbrace{\mathbb{E}_{q_\phi(z|x)} [p_\theta(x|z)]}_{\text{reconstruction term}} - \underbrace{\text{D}\text{KL}(q_\phi(z|x) \parallel p(z))}_{\text{prior fitly}}$$

$$X \xrightarrow{\text{encoder}} \tilde{z} \xrightarrow{\text{decoder}} \tilde{x}$$

VAE.

prior fitly,

(3)



$$KL(q_{\phi}(z|x) || N(0, I)) = \frac{1}{2} \sum_{i=1}^D \left(\frac{\sigma_{z_i}^2}{\sigma_{\phi(z_i)}^2} + M_{\phi(z_i)}^2 - \ln(\frac{\sigma_{z_i}^2}{\sigma_{\phi(z_i)}^2}) - 1 \right)$$

network output of stochastic encoder.

Big Limitation!! p(z)

Given

1) BNN: $D = \{(x_i, y_i)\}_{i=1}^N$

prior

$p(w)$

variational distribution (goal)

$Q(w)$

2) VAE: $D = \{x_i\}_{i=1}^N$

$p(z)$

$q_{\phi}(z|x) \& P_{\theta}(x|z)$

Objectives of VAE:

- 1) Reconstruction
- 2) KLD (prior fitting)

$$KL(Q_\phi(z|x) \parallel N(0, I))$$

$$= \frac{1}{2} \sum_{i=1}^D (\underline{\sigma_{z_i}^2} + \underline{m_{z_i}^2} - \ln(\underline{\sigma_{z_i}^2}) - 1) \quad - (*)$$

* outputs of neural networks.

WAE₁₁ → (WAE-GAN)

* GAN: Generative Adversarial Net.

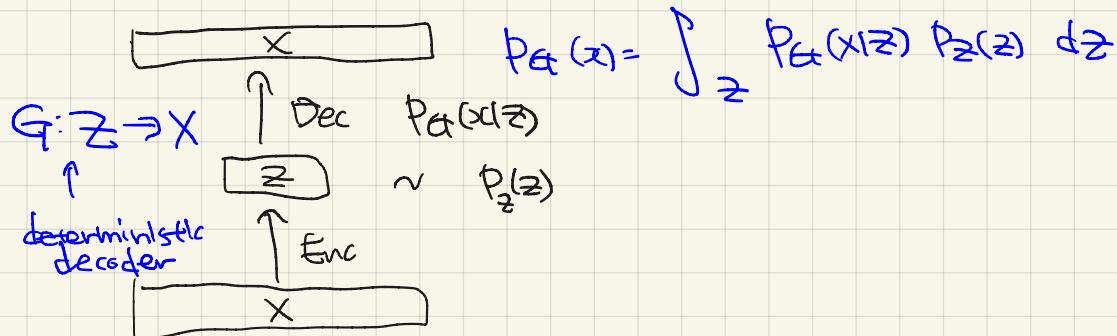
Optimal transport and its dual formulation,

$$W_C(P_X, P_G) = \inf_{\substack{\text{cost} \uparrow \\ \text{data dist.}}} \sup_{\substack{\text{Generativ dist.}}} \mathbb{E}_{(x,y) \sim P} [c(x,y)]$$

If $c(x,y) = d(x,y)$, then the following Kantorovich-Rubinstein duality hold.

$$W_1(P_X, P_G) = \sup_{f \in F_L} \left(\mathbb{E}_{x \sim P_X} [f(x)] - \mathbb{E}_{y \sim P_G} [f(y)] \right) \quad - (**)$$

L-Lipsh.



Third. For P_G as defined above with deterministic $P_G(x|z)$ and any function $G: Z \rightarrow X$

$$\inf_{P \in P(X \times P_X, Z \times P_G)} \mathbb{E}_{(x,y) \sim P} [c(x,y)] \Leftarrow W.D.$$

$$= \inf_{Q: Q_z = P_Z} \mathbb{E}_P \mathbb{E}_Q [c(x, G(z))] \quad (\text{fixed})$$

data encoder decoder
Reconstruction,

* marginal prior fitting.

where Q_z is the marginal distribution of Z when $x \sim P_X$ and $z \sim Q(z|x)$.

From $(\star \star)$,

Recon.

$$D_{WAE}(P_x, P_{\mathcal{G}}) = \inf_{Q(z|x) \in Q} \mathbb{E}_{P_x} \mathbb{E}_{Q(z|x)} [c(x, \mathcal{G}(z))] + \lambda D_z(Q_z, P_z)$$

↑ Prior fitting term,
total divergence.

\Downarrow GAN-based.

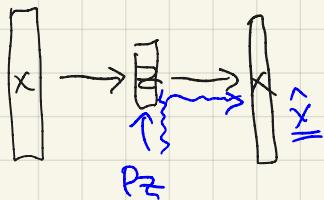
$$D_z(Q_z, P_z) = D_{JS}(Q_z, P_z) \leftarrow \text{adversarial training.}$$

*GAN

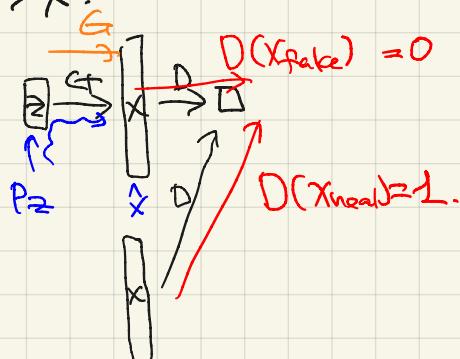
$$\min_{\mathcal{G}} \max_D \left\{ \log(D(x)) + \log(1 - D(\mathcal{G}(z))) \right\} \quad (1)$$

where $D: x \rightarrow \mathbb{R}$, $\mathcal{G}: z \rightarrow x$.

VAE:



GAN:



i) Fix \mathcal{G} , Optimize D ,

$$\begin{aligned} & \min_{\mathcal{G}} \max_D \left\{ \mathbb{E}_{x \sim P_x} [\log(D(x))] + \mathbb{E}_{z \sim P_z} [\log(1 - D(\mathcal{G}(z)))] \right\} \\ &= \min_{\mathcal{G}} \max_D V(\mathcal{G}, D) \end{aligned}$$

Since \mathcal{G} is fixed, let $y = \mathcal{G}(z)$

$$\Rightarrow \mathbb{E}_{z \sim P_z} [\log(1 - D(\mathcal{G}(z)))] = \mathbb{E}_{y \sim P_y} [\log(1 - D(y))]$$

\downarrow

$$\begin{aligned} V(\mathcal{G}, D) &= \mathbb{E}_{x \sim P_x} [\log(D(x))] + \mathbb{E}_{y \sim P_y} [\log(1 - D(y))] \\ &= \int_X P_{\text{data}}(x) \log D(x) dx + \int_Y P_g(y) \log(1 - D(y)) dy \\ &= \int_X P_{\text{data}}(x) \log D(x) dx + P_g(x) \log(1 - D(x)) dx \end{aligned}$$

$$D_{\mathcal{G}}^*(x) = \arg \max_D V(\mathcal{G}, D) = \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_{\mathcal{G}}(x)}$$

($\because y \rightarrow a \log(y) + b \log(1-y)$ achieves its maximum in $[0,1]$ at $\frac{a}{a+b}$)

ii) Then, fix $D = D^*_G$,

$$\begin{aligned}
 V(G, D^*) &= \mathbb{E}_{x \sim P_{\text{data}}} [\log(D^*(x))] + \mathbb{E}_{x \sim P_g(x)} [\log(1 - D^*(x))] \\
 C(G) &= \mathbb{E}_{x \sim P_{\text{data}}} \left[\log \left(\frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \right) \right] + \\
 &\quad \mathbb{E}_{x \sim P_g(x)} \left[\log \left(\frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)} \right) \right] \\
 &= -\log(q) + \text{KL}(P_{\text{data}} \parallel \frac{P_{\text{data}} + P_g}{2}) + \text{KL}(P_g \parallel \frac{P_{\text{data}} + P_g}{2}) \\
 &= -\log(q) + 2 \cdot \text{JSD}(P_{\text{data}} \parallel P_g)
 \end{aligned}$$

\therefore GAN objective minimizes JS'D between P_{data} and P_g .