

Variational Inference and Deep Learning

Can Yang

Department of Mathematics
The Hong Kong University of Science and Technology

October 29, 2018

Learning from data

- ▶ Assume the observed variable \mathbf{x} is a random sample from an **unknown underlying process**, whose true distribution $p^*(\mathbf{x})$ is unknown. We attempt to approximate this underlying process with a chosen model $p_{\theta}(\mathbf{x})$, with parameter θ :

$$\mathbf{x} \sim p_{\theta}(\mathbf{x})$$

- ▶ Learning is, most commonly, the process of searching for a value of the parameter θ such that the probability distribution function given by the model, $p_{\theta}(\mathbf{x})$, approximates the true distribution of the data, denoted by $p^*(\mathbf{x})$, such that for any observed \mathbf{x} :

$$p_{\theta}(\mathbf{x}) \approx p^*(\mathbf{x})$$

Conditional models

- In the case of classification or regression problems, we are not interested in learning an unconditional model $p_{\theta}(\mathbf{x})$, but a conditional model $p_{\theta}(\mathbf{y}|\mathbf{x})$ that approximates the underlying conditional distribution $p^*(\mathbf{y}|\mathbf{x})$:

$$p_{\theta}(\mathbf{y}|\mathbf{x}) \approx p^*(\mathbf{y}|\mathbf{x}).$$

- \mathbf{x} is an image, and \mathbf{y} is the image's class (labeled by a human). In this case, $p_{\theta}(\mathbf{y}|\mathbf{x})$ is typically chosen to be a categorical distribution.
- Examples include (multi-class) logistic regression and neural networks:

$$\mathbf{p} = \text{NeuralNet}(\mathbf{x})$$

$$p_{\theta}(y|\mathbf{x}) = \text{Categorical}(y; \mathbf{p})$$

Learning in fully observed models

- ▶ We often collect a dataset \mathcal{D} consisting of $N \geq 1$ datapoints

$$\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\} = \{\mathbf{x}^{(i)}\}_{i=1}^N.$$

- ▶ The observations $\{\mathbf{x}^{(i)}\}_{i=1}^N$ are said to be *i.i.d.*, for independently and identically distributed.
- ▶ With i.i.d. dataset \mathcal{D} of size $N_{\mathcal{D}}$, the maximum likelihood criterion is

$$\log p_{\boldsymbol{\theta}}(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\boldsymbol{\theta}}(\mathbf{x}).$$

- ▶ Gradient methods can be used to find a local optimum of the ML objective: $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t+1)} + s_t \nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\mathbf{x})$.

Stochastic gradient descent (SGD)

- ▶ Consider a minibatches of data $\mathcal{M} \subset \mathcal{D}$ of size $N_{\mathcal{M}}$. With such minibatches, we can form an unbiased estimator of the ML criterion:

$$\frac{1}{N_{\mathcal{D}}} \log p_{\theta}(\mathcal{D}) \simeq \frac{1}{N_{\mathcal{M}}} \log p_{\theta}(\mathcal{M}) = \frac{1}{N_{\mathcal{M}}} \sum_{\mathbf{x} \in \mathcal{M}} \log p_{\theta}(\mathbf{x})$$

- ▶ The unbiased estimator $\log p_{\theta}(\mathcal{M})$ is differentiable, yielding the unbiased stochastic gradients:

$$\frac{1}{N_{\mathcal{D}}} \nabla_{\theta} \log p_{\theta}(\mathcal{D}) \simeq \frac{1}{N_{\mathcal{M}}} \nabla_{\theta} \log p_{\theta}(\mathcal{M}) = \frac{1}{N_{\mathcal{M}}} \sum_{\mathbf{x} \in \mathcal{M}} \nabla_{\theta} \log p_{\theta}(\mathbf{x})$$

Learning and inference in deep latent variable models

- ▶ Latent variables are variables that are part of the model, but which we don't observe, and are therefore not part of the data. We typically use \mathbf{z} to denote such latent variables.
- ▶ The marginal distribution over the observed variables is given by

$$p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\theta}(\mathbf{z}) d\mathbf{z}$$

- ▶ Such an implicit distribution over \mathbf{x} can be quite flexible. If \mathbf{z} is discrete and $p_{\theta}(\mathbf{x}|\mathbf{z})$ is a Gaussian distribution, then $p_{\theta}(\mathbf{x})$ is a mixture-of-Gaussians distribution. For continuous \mathbf{z} , $p_{\theta}(\mathbf{x})$ can be seen as an infinite mixture, which are potentially more powerful than discrete mixtures.
- ▶ As we shall see in variational auto encode (VAE), the continuous version is generally more efficient to work with due to the re-parameterization trick.

Deep Latent Variable Model (DLVM)

- ▶ Joint distribution

$$p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{z})p_{\theta}(\mathbf{x}|\mathbf{z}),$$

where $p_{\theta}(\mathbf{z})$ is the prior distribution over \mathbf{z} and $p_{\theta}(\mathbf{x}|\mathbf{z})$ can be specified via deep neural networks.

- ▶ Generative model with spherical Gaussian latent space and Bernoulli observation:

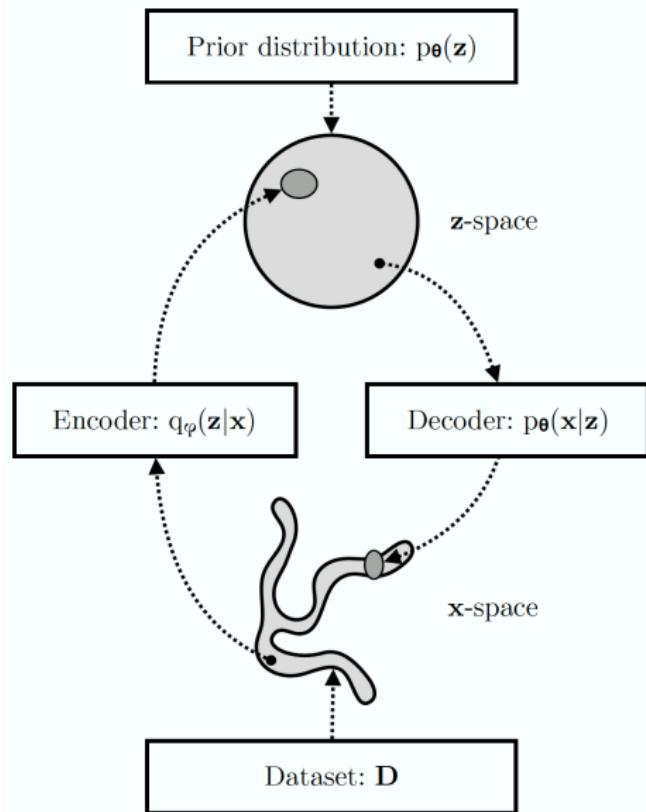
$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$$

$$\mathbf{p} = \text{DecoderNeuralNet}_{\theta}(\mathbf{z})$$

$$\begin{aligned}\log p(\mathbf{x}|\mathbf{z}) &= \sum_{j=1}^D \log p(x_j|\mathbf{z}) = \sum_{j=1}^D \log \text{Bernoulli}(x_j; p_j) \\ &= \sum_{j=1}^D x_j \log p_j + (1 - x_j) \log(1 - p_j)\end{aligned}\tag{1}$$

where $\forall p_j \in \mathbf{p} : 0 \leq p_j \leq 1$ and $\sum_j p_j = 1$.

- ▶ Intractability: $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z})d\mathbf{z}$ is often intractable.



Variational Autoencoders (VAE)

- ▶ To turn the DLVM intractable posterior inference and learning problems into tractable problems, an inference model $q_{\phi}(\mathbf{z}|\mathbf{x})$ is introduced so that

$$q_{\phi}(\mathbf{z}|\mathbf{x}) \approx p_{\theta}(\mathbf{z}|\mathbf{x}),$$

where ϕ is variational parameter.

- ▶ Like the generative model, $q_{\phi}(\mathbf{z}|\mathbf{x})$ can be parameterized using deep neural network. For example:

$$\begin{aligned} (\boldsymbol{\mu}, \log \boldsymbol{\sigma}) &= \text{EncoderNeuralNet}_{\phi}(\mathbf{x}) \\ q_{\phi}(\mathbf{z}|\mathbf{x}) &= \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma})) \end{aligned} \tag{2}$$

in this case, the variational parameter ϕ include the weights coefficients of the neural network.

Objective: Evidence Lower Bound (ELBO)

- ▶ For any choice of $q_\phi(\mathbf{z}|\mathbf{x})$,

$$\begin{aligned}\log p_{\boldsymbol{\theta}}(\mathbf{x}) &= \mathbb{E}_{q_\phi} [\log p_{\boldsymbol{\theta}}(\mathbf{x})] \\ &= \mathbb{E}_{q_\phi} \left[\log \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})} \right] \right] \\ &= \mathbb{E}_{q_\phi} \left[\log \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z}) q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x}) p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})} \right] \right] \\ &= \mathbb{E}_{q_\phi} \left[\log \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \right] + \mathbb{E}_{q_\phi} \left[\log \left[\frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})} \right] \right] \\ &= \mathcal{L}_{\boldsymbol{\theta}, \phi}(\mathbf{x}) + \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}))\end{aligned}\tag{3}$$

Connection with EM

- ▶ For standard EM algorithms, the posterior is often known, $q(\mathbf{z}|\mathbf{x}) = p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$.
- ▶ Then the KL term becomes zero.

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}) = \mathcal{L}_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbb{E}_q[\log p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})] - \mathbb{E}_q[\log q(\mathbf{z}|\mathbf{x})].$$

- ▶ The above step is indeed the E-step in the standard EM algorithm.
- ▶ The M-step is

$$\boldsymbol{\theta}_{\text{new}} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}_{\boldsymbol{\theta}}(\mathbf{x}).$$

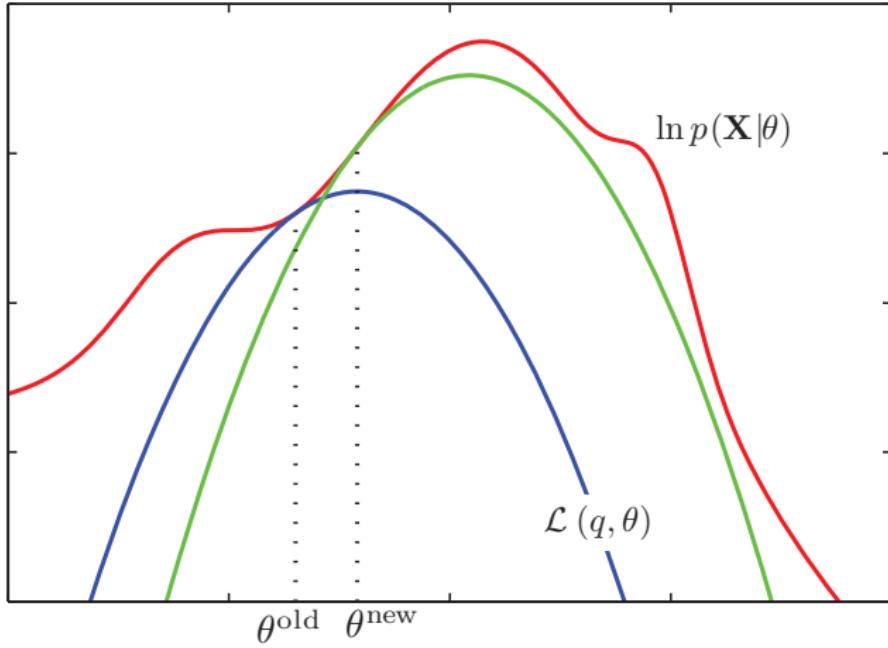


Figure: The EM algorithm involves alternately computing a lower bound on the log likelihood for the current parameter values and then maximizing this bound to obtain the new parameter values.

Stochastic gradient-based optimization

- ▶ Given i.i.d. data $\mathbf{x} \in \mathcal{D}$

$$\mathcal{L}_{\theta, \phi}(\mathcal{D}) = \sum_{\mathbf{x} \in \mathcal{D}} \mathcal{L}_{\theta, \phi}(\mathbf{x}) \quad (4)$$

- ▶ Consider the single-data point ELBO. Its gradient $\nabla_{\theta, \phi} \mathcal{L}_{\theta, \phi}(\mathbf{x})$ is in general intractable. But the unbiased estimator exists.

$$\begin{aligned}\nabla_{\theta} \mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \nabla_{\theta} \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{q_{\phi}} [\nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))] \\ &\simeq \nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})) \\ &= \nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z})),\end{aligned} \quad (5)$$

where \simeq means unbiased estimator. The \mathbf{z} in the last two lines are sampled from $q_{\phi}(\mathbf{z}|\mathbf{x})$

- ▶ The same equation does not apply to ϕ in general

$$\begin{aligned}\nabla_{\phi} \mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \nabla_{\phi} \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &\neq \mathbb{E}_{q_{\phi}} [\nabla_{\phi} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))]\end{aligned} \quad (6)$$

Reparameterization trick

- We express random variable \mathbf{z} , as differentiable and invertible transformation of another random variable $\boldsymbol{\epsilon}$, given \mathbf{x} and $\boldsymbol{\phi}$,

$$\mathbf{z} = g(\boldsymbol{\epsilon}, \boldsymbol{\phi}, \mathbf{x}) \quad (7)$$

such that the distribution of $\boldsymbol{\epsilon}$ is independent of \mathbf{x} or $\boldsymbol{\phi}$.

- Using change of variable, the expectation becomes

$$\mathbb{E}_{q_{\boldsymbol{\phi}}} [f(\mathbf{z})] = \mathbb{E}_{p(\boldsymbol{\epsilon})} [f(\mathbf{z})] \quad (8)$$

- The expectation is independent of $\boldsymbol{\phi}$ and the expectation and the gradient becomes communicative

$$\begin{aligned} \nabla_{\boldsymbol{\phi}} \mathbb{E}_{q_{\boldsymbol{\phi}}} [f(\mathbf{z})] &= \nabla_{\boldsymbol{\phi}} \mathbb{E}_{p(\boldsymbol{\epsilon})} [f(g(\boldsymbol{\epsilon}, \boldsymbol{\phi}, \mathbf{x}))] \\ &= \mathbb{E}_{p(\boldsymbol{\epsilon})} [\nabla_{\boldsymbol{\phi}} f(g(\boldsymbol{\epsilon}, \boldsymbol{\phi}, \mathbf{x}))] \\ &\simeq \nabla_{\boldsymbol{\phi}} f(g(\boldsymbol{\epsilon}, \boldsymbol{\phi}, \mathbf{x})), \end{aligned} \quad (9)$$

where in the last line a random noise sample $\boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$ in $g(\boldsymbol{\epsilon}, \boldsymbol{\phi}, \mathbf{x})$.

Estimate of gradient

- With the reparameterization trick, the ELBO can be written as

$$\begin{aligned}\mathcal{L}_{\theta, \phi}(\mathbf{x}) &= \mathbb{E}_{q_{\phi}}[\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \\ &= \mathbb{E}_{p(\epsilon)}[\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})]\end{aligned}\tag{10}$$

- Using a single noise sample ϵ from $p(\epsilon)$, Monte Carlo estimate of $\mathcal{L}_{\theta, \phi}$ can be obtained:

$$\begin{aligned}\epsilon &\sim p(\epsilon) \\ \mathbf{z} &= g(\epsilon, \phi, \mathbf{x})\end{aligned}\tag{11}$$

$$\tilde{\mathcal{L}}_{\theta, \phi}(\mathbf{x}) = \log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})$$

- Specifically, $\log q_{\phi}(\mathbf{z}|\mathbf{x})$ is computed by

$$\log q_{\phi}(\mathbf{z}|\mathbf{x}) = \log p(\epsilon) - \log |\det(\frac{\partial \mathbf{z}}{\partial \epsilon})|\tag{12}$$

Choosing function $g(\epsilon, \phi, \mathbf{x})$

- ▶ A common choice is simple factorized Gaussian encoder
 $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}))$:

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2) = \text{EncoderNeuralNet}_\phi(\mathbf{x})$$
$$q_\phi(\mathbf{z}|\mathbf{x}) = \prod_i q_\phi(z_i|\mathbf{x}) = \prod_i \mathcal{N}(z_i; \mu_i, \sigma_i^2) \quad (13)$$

- ▶ After re-parameterization,

$$\epsilon \sim p(\epsilon)$$
$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2) = \text{EncoderNeuralNet}_\phi(\mathbf{x}) \quad (14)$$
$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon$$

where \odot denotes element-wise product.

- ▶ Since $\frac{\partial \mathbf{z}}{\partial \epsilon} = \text{diag}(\boldsymbol{\sigma})$, we have $\log |\det(\frac{\partial \mathbf{z}}{\partial \epsilon})| = \sum_i \log \sigma_i$ and the log of posterior is

$$\log q_\phi(\mathbf{z}|\mathbf{x}) = \sum_i \log \mathcal{N}(\epsilon_i; 0, 1) - \log \sigma_i \quad (15)$$

Algorithm 1: Stochastic optimization of the ELBO. Since noise originates from both the minibatch sampling and sampling of $p(\epsilon)$, this is a doubly stochastic optimization procedure. We also refer to this procedure as the *Auto-Encoding Variational Bayes* (AEVB) algorithm.

Data:

\mathcal{D} : Dataset

$q_\phi(\mathbf{z}|\mathbf{x})$: Inference model

$p_\theta(\mathbf{x}, \mathbf{z})$: Generative model

Result:

θ, ϕ : Learned parameters

$(\theta, \phi) \leftarrow$ Initialize parameters

while SGD not converged **do**

$\mathcal{M} \sim \mathcal{D}$ (Random minibatch of data)

$\epsilon \sim p(\epsilon)$ (Random noise for every datapoint in \mathcal{M})

Compute $\tilde{\mathcal{L}}_{\theta, \phi}(\mathcal{M}, \epsilon)$ and its gradients $\nabla_{\theta, \phi} \tilde{\mathcal{L}}_{\theta, \phi}(\mathcal{M}, \epsilon)$

Update θ and ϕ using SGD optimizer

end

Full covariance Gaussian posterior

- ▶ Consider the Gaussian with full covariance:

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (16)$$

with re-parameterization:

$$\begin{aligned}\boldsymbol{\epsilon} &\sim \mathcal{N}(0, \mathbf{I}) \\ \mathbf{z} &= \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon}\end{aligned} \quad (17)$$

where \mathbf{L} is a lower triangular matrix with non-zero entries in the diagonal. The off-diagonal elements define the correlations of the elements in \mathbf{z} .

- ▶ The Jacobian is $\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}} = \mathbf{L}$, hence

$$\log |\det\left(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}}\right)| = \sum_i \log |L_{ii}| \quad (18)$$

Full covariance Gaussian posterior

- ▶ This corresponds to the Cholesky decomposition $\Sigma = \mathbf{L}\mathbf{L}^T$:

$$\begin{aligned}\Sigma &= \mathbb{E}[(\mathbf{z} - \mathbb{E}[\mathbf{z}])(\mathbf{z} - \mathbb{E}[\mathbf{z}])^T] \\ &= \mathbb{E}(\mathbf{L}\boldsymbol{\epsilon}(\mathbf{L}\boldsymbol{\epsilon})^T) = \mathbf{L}\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^T]\mathbf{L}^T \\ &= \mathbf{L}\mathbf{L}^T\end{aligned}\tag{19}$$

- ▶ To build the desired matrix \mathbf{L} , we can

$$\begin{aligned}(\boldsymbol{\mu}, \log \boldsymbol{\sigma}, \mathbf{L}') &\leftarrow \text{EncoderNueralNet}_{\phi}(\mathbf{x}) \\ \mathbf{L} &\leftarrow \mathbf{L}_{mask} \odot \mathbf{L}' + \text{diag}(\boldsymbol{\sigma})\end{aligned}\tag{20}$$

where \mathbf{L}_{mask} is a masking matrix with zeros on and above the diagonal, and ones below the diagonal. Then $\log |\det(\frac{\partial \mathbf{z}}{\partial \boldsymbol{\epsilon}})| = \sum_i \log \sigma_i$.

Algorithm 2: Computation of unbiased estimate of single-datapoint ELBO for example VAE with a full-covariance Gaussian inference model and a factorized Bernoulli generative model. \mathbf{L}_{mask} is a masking matrix with zeros on and above the diagonal, and ones below the diagonal. Note that \mathbf{L} must be a triangular matrix with positive entries on the diagonal.

Data:

\mathbf{x} : a datapoint, and optionally other conditioning information

ϵ : a random sample from $p(\epsilon) = \mathcal{N}(0, \mathbf{I})$

θ : Generative model parameters

ϕ : Inference model parameters

$q_\phi(\mathbf{z}|\mathbf{x})$: Inference model

$p_\theta(\mathbf{x}, \mathbf{z})$: Generative model

Result:

$\tilde{\mathcal{L}}$: unbiased estimate of the single-datapoint ELBO $\mathcal{L}_{\theta, \phi}(\mathbf{x})$

$(\mu, \log \sigma, \mathbf{L}') \leftarrow \text{EncoderNeuralNet}_\phi(\mathbf{x})$

$\mathbf{L} \leftarrow \mathbf{L}_{mask} \odot \mathbf{L}' + \text{diag}(\sigma)$

$\epsilon \sim \mathcal{N}(0, \mathbf{I})$

$\mathbf{z} \leftarrow \mathbf{L}\epsilon + \mu$

$\log q_{\phi}(\mathbf{z}|\mathbf{x}) \leftarrow -\text{sum}(\frac{1}{2}(\epsilon^2 + \log(2\pi) + \log \sigma))$

$\triangleright = q_\phi(\mathbf{z}|\mathbf{x})$

$\log p_{\theta}(\mathbf{z}) \leftarrow -\text{sum}(\frac{1}{2}(\mathbf{z}^2 + \log(2\pi)))$

$\triangleright = p_\theta(\mathbf{z})$

$\mathbf{p} \leftarrow \text{DecoderNeuralNet}_\theta(\mathbf{z})$

$\log p_{\theta}(\mathbf{x}|\mathbf{z}) \leftarrow \text{sum}(\mathbf{x} \odot \log \mathbf{p} + (1 - \mathbf{x}) \odot \log(1 - \mathbf{p}))$

$\triangleright = p_\theta(\mathbf{x}|\mathbf{z})$

$\tilde{\mathcal{L}} = \log p_{\theta}(\mathbf{x}|\mathbf{z}) + \log q_{\phi}(\mathbf{z}|\mathbf{x}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})$

Estimation of marginal likelihood

- ▶ Importance sampling can be adopted in the estimation of marginal likelihood.
- ▶ From

$$\log p_{\theta}(\mathbf{x}) = \log \mathbb{E}_{q_{\phi}} [p_{\theta}(\mathbf{x}, \mathbf{z}) / q_{\phi}(\mathbf{z}|\mathbf{x})],$$

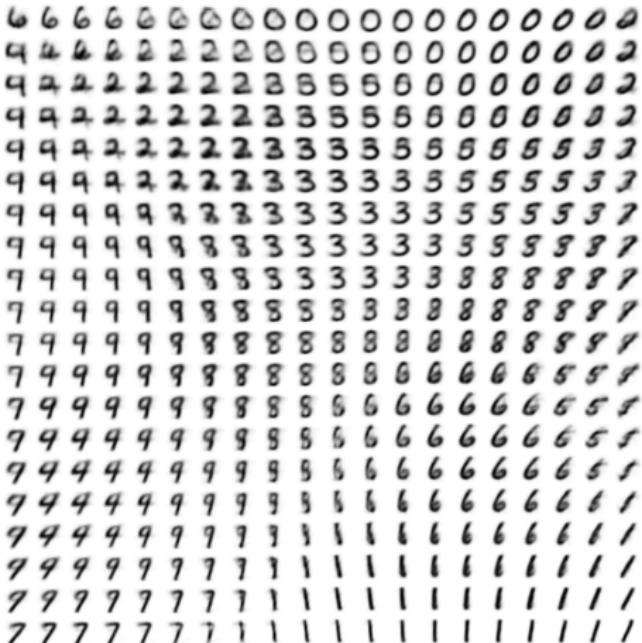
we draw random samples from q_{ϕ} to obtain a Monte Carlo estimator:

$$\log p_{\theta}(\mathbf{x}) \approx \log \frac{1}{L} \sum_{l=1}^L p_{\theta}(\mathbf{x}, \mathbf{z}^l) / q_{\phi}(\mathbf{z}^l | \mathbf{x}^l)$$

- ▶ When $L = 1$, this equals the ELBO estimator of the VAE.



(a) Learned Frey Face manifold



(b) Learned MNIST manifold



(a) 2-D latent space (b) 5-D latent space (c) 10-D latent space (d) 20-D latent space

Figure: Random samples from learned generative models of MNIST for different dimensionalities of latent space.

Plain VAE



(a) Real



(b) Reconstruction

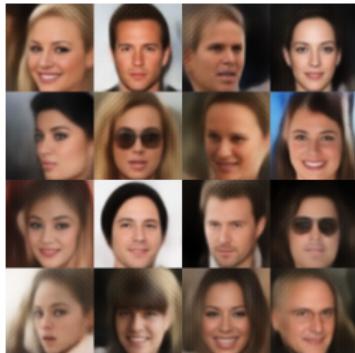


(c) Generation

VAE VGG



(d) Real



(e) Reconstruction



(f) Generation

More on VAE

- ▶ Objective function (population version):

$$\min_{\phi, \theta} \mathbb{E}_{p(\mathbf{x})} [\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})]] + \mathbb{E}_{p(\mathbf{x})} [\text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))] \quad (21)$$

- ▶ Loss(ϕ, θ) = $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [-\log p_\theta(\mathbf{x}|\mathbf{z})]$.
- ▶ $\text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}))$ can be viewed as regularization and it can be evaluated in the closed form (KL of two Gaussians).
- ▶ Loss + Regularization
- ▶ The derivation can be found in paper Tutorial of VAE by C. Doersch.

Adversarial AutoEncoders

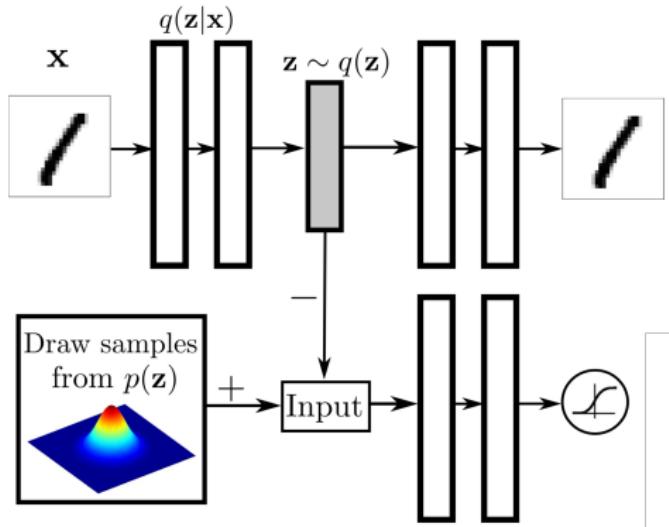


Figure: An AAE example for MNIST

Adversarial AutoEncoders

- ▶ AAE: Plain AE in the data space + GAN in the latent space
- ▶ Goal: learn regularized PAE
- ▶ Latent space: learn to match the prior with adversarial loss
- ▶ Treat the encoder as generator in the latent space
- ▶ Objective function:

$$\min_{\phi, \theta} \left[\text{Loss}(\phi, \theta) + \max_f \left[\mathbb{E}_{p(\mathbf{z})}[\log f(\mathbf{z})] + \mathbb{E}_{q(\mathbf{z})}[\log(1 - f(\mathbf{z}))] \right] \right]$$
$$q(\mathbf{z}) = \int q_\phi(\mathbf{z}|\mathbf{x})p(\mathbf{x})d\mathbf{x} \quad (22)$$

Wasserstein AutoEncoders

- ▶ WAE: the optimal transport view of AAE
- ▶ Goal: learn regularized PAE with optimal transport in data space
- ▶ Latent space: learned by GAN or maximum mean discrepancy (MMD)
- ▶ Objective function:

$$\min_{\phi, \theta} \mathbb{E}_{p(x)} \mathbb{E}_{q_\phi(z|x)} c(x, D_\theta(z)) + \lambda \cdot \mathcal{D}(p(z), q(z))$$
$$q(z) = \int_x q_\phi(z|x)p(x)dx \tag{23}$$

where $c()$ could be the least square loss and $\mathcal{D}(\cdot, \cdot)$ is a distance measure of two density function.

Some New Results

2 0 3 6 7 1 8 4
3 0 8 4 1 4 1 6
3 2 8 2 1 0 7 4
1 1 6 5 6 5 4 3
8 0 1 9 3 0 7 4
8 2 7 6 3 9 9 6
4 1 9 8 5 3 5 7
6 2 0 4 4 8 4 4

(a) Real

2 0 3 6 7 1 8 4
3 0 8 4 0 4 1 6
3 2 8 2 1 0 7 4
1 1 6 5 6 5 4 3
8 0 1 9 3 0 7 4
8 2 7 6 3 9 9 6
4 1 9 1 5 3 5 7
6 2 0 4 4 8 4 4

(b) Reconstruction

4 1 6 2 3 6 8 1
9 4 1 6 4 2 0 6
5 7 5 6 7 8 9 1
3 0 0 1 1 6 9 7
2 0 4 6 9 3 2 7
4 9 6 3 0 7 0 4
0 0 4 0 6 1 1 0
2 7 7 6 6 3 2 3

(c) Generation

Some New Results



(d) Real

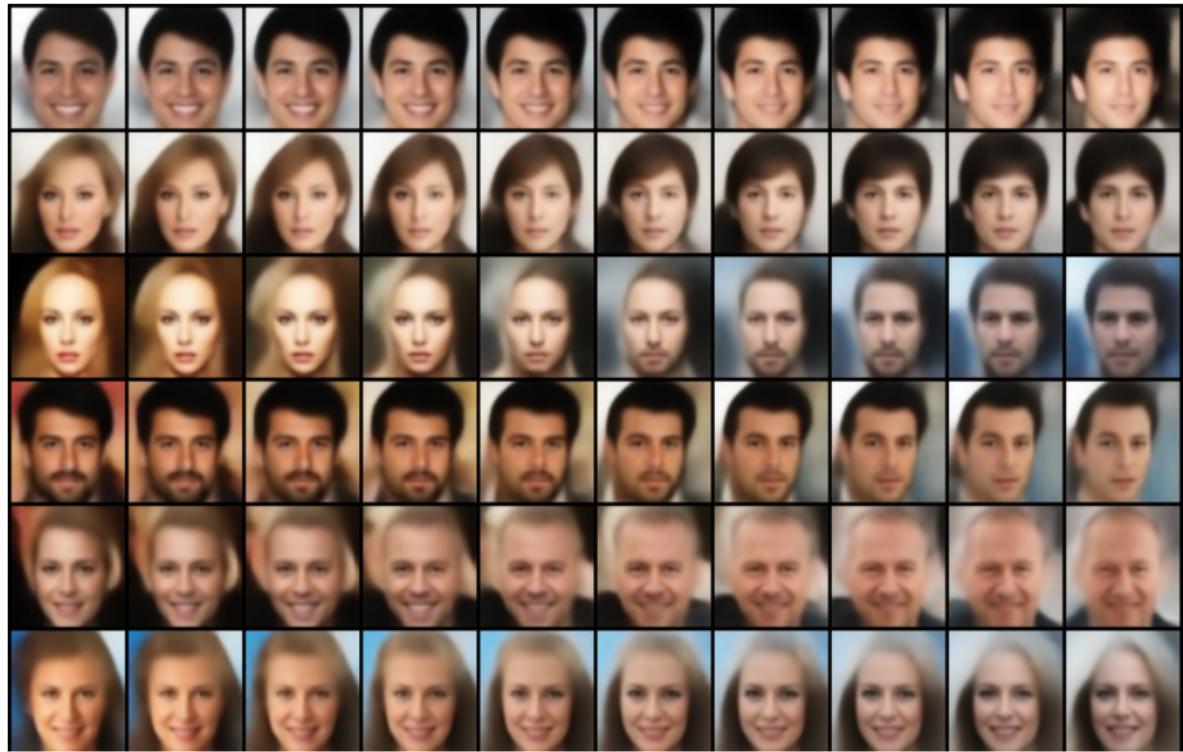


(e) Reconstruction

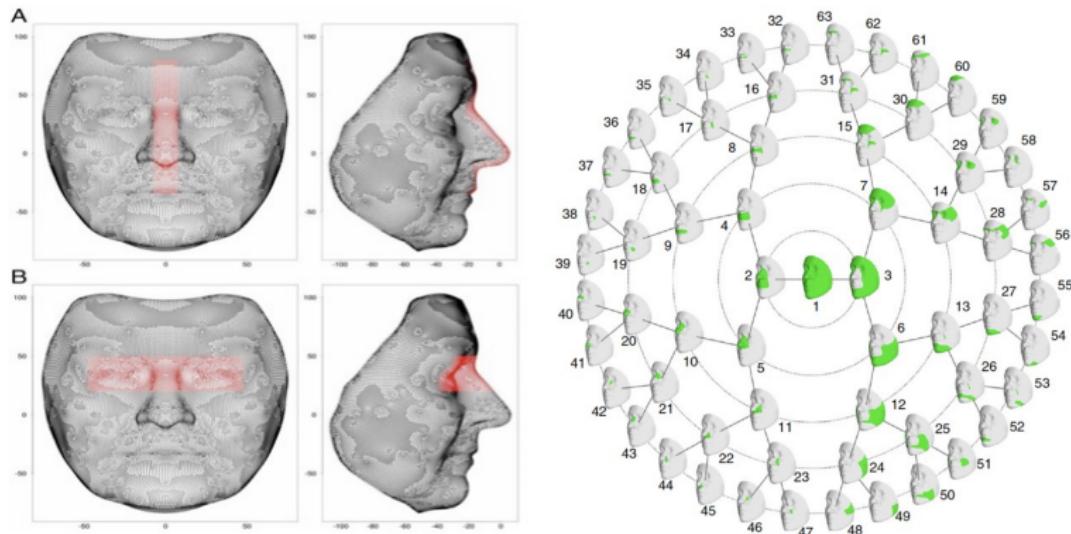


(f) Generation

Some New Results



On-going project: Connecting genome with human face



Thank You!