

Variational Inference & Variational Autoencoders

Casey Meehan & Mary Anne Smart

CSE 254

May 2019

1 / 45

Overview

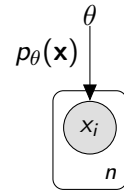
- 1 Variational Inference: A Review for Statisticians (Blei, Kucukelbir, & McAuliffe)
 - When and why use VI?
 - The Evidence Lower Bound (ELBO)
 - Mean Field Variational Inference
 - Comparison with older methods
- 2 Auto-Encoding Variational Bayes (Kingma & Welling)

2 / 45

Bayesian Modeling

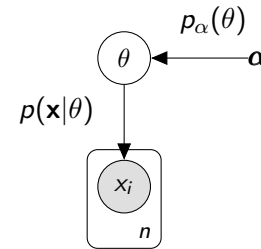
- Define a probabilistic model for your data in which the parameters are random variables
- Consider a simple coin-toss model, $x_i \in \{0, 1\}$:

Non-Bayesian



$$X_i \sim \text{Bern}(\theta)$$

Bayesian



$$X_i \sim \text{Bern}(\theta)$$
$$\theta \sim \text{Beta}(\alpha_1, \alpha_2)$$

3 / 45

MLE, MAP, Full Bayesian

Max Likelihood Est. (MLE)

- Latent θ is a numerical parameter, not random variable.
- Choose the θ^* that best explains the data
- Maximize log likelihood:

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \log p_{\theta}(x_i) = \frac{\sum_{i=1}^n x_i}{n}$$

- likely to overfit without enough data:
- seeing 1 sample, $x_1 = 1$ implies $\theta^* = 1$

4 / 45

MLE, MAP, Full Bayesian

Maximum A Posteriori Estimate (MAP)

- Latent $\theta \sim \text{Beta}(\alpha_1, \alpha_2)$ is a random variable
- Choose the θ^* with maximum probability
- Maximize the log posterior:

$$\begin{aligned}\theta^* &= \arg \max_{\theta} \left[\log p(\theta|\mathbf{x}) \right] \\ &= \arg \max_{\theta} \left[\log \frac{p(\mathbf{x}|\theta)}{p(\mathbf{x})} \right] \\ &= \frac{\alpha_1 - 1 + \sum_{i=1}^n x_i}{n + \alpha_1 + \alpha_2 - 2}\end{aligned}$$

- less likely to overfit:
- Say, $\alpha_1 = \alpha_2 = 2$. Then seeing 1 sample, $x_1 = 1$ implies $\theta^* = \frac{2}{3}$

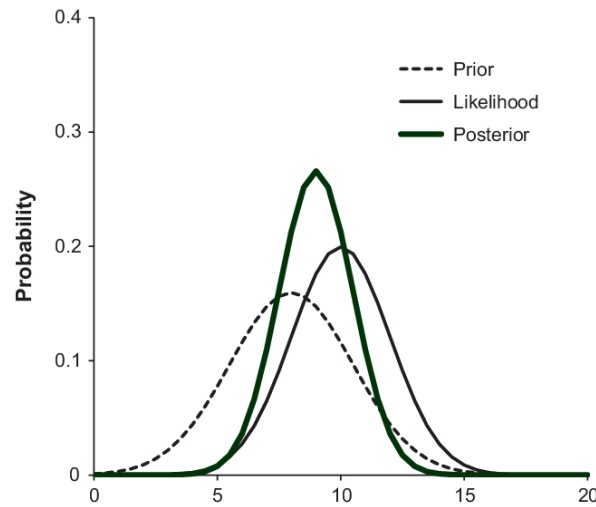
5 / 45

MLE, MAP, Full Bayesian

Full Bayesian

(goal of variational inference)

- Latent $\theta \sim \text{Beta}(\alpha_1, \alpha_2)$ is a random variable
- Find the full distribution of $\theta|\mathbf{x}, \alpha$. Don't pick a single value



6 / 45

A VI Story...

- 1 we have a sample of n data points, $\mathbf{x} = \{x_i : i \in [1, n]\}$
- 2 we have latent vars, $\mathbf{z} = \{z_i : i \in [1, m]\}$,
- 3 we have a probabilistic model with a likelihood and a prior distribution

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

- 4 we want to be full Bayesian, but the posterior is intractable!

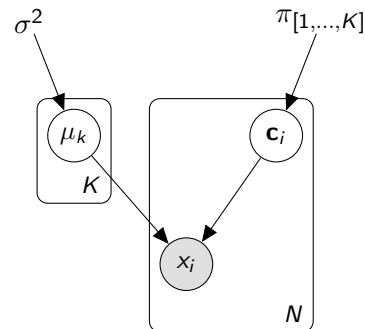
$$p(\mathbf{z}|\mathbf{x})$$

- 5 so, we pick a simpler class of distributions, \mathcal{Q} , to approximate $p(\mathbf{z}|\mathbf{x})$
- 6 we choose the optimal $q_\phi(\mathbf{z}|\mathbf{x}) \in \mathcal{Q}$ by minimizing the reverse KL divergence from $p(\mathbf{z}|\mathbf{x})$
- 7 ...which is equivalent to maximizing the ELBO

7 / 45

Latent Variables can be Difficult to Infer

- Consider this Gaussian mixture model of N samples and K clusters:



$$\begin{aligned}\mu_k &\sim \mathcal{N}(0, \sigma^2) \\ \mathbf{c}_i &\sim \text{Cat}(\pi_1, \dots, \pi_K) \\ x_i | \mathbf{c}_i, \boldsymbol{\mu} &\sim \mathcal{N}(\mathbf{c}_i^\top \boldsymbol{\mu}, \mathbf{I})\end{aligned}$$

8 / 45

Latent Variables can be Difficult to Infer

- The posterior probability requires computing the evidence...

$$p(\mathbf{x}) = p(x_1, x_2, \dots, x_n)$$

$$p(\boldsymbol{\mu}, \mathbf{c} | \mathbf{x}) = \frac{p(\boldsymbol{\mu}, \mathbf{c}, \mathbf{x})}{p(\mathbf{x})} = \frac{\prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(\mathbf{c}_i) p(x_i | \mathbf{c}_i, \boldsymbol{\mu})}{p(\mathbf{x})}$$

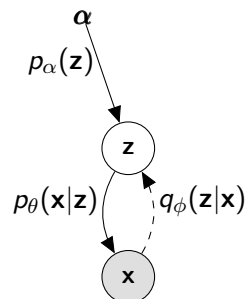
- ...and it's a pain to compute

$$\begin{aligned} p(\mathbf{x}) &= \int p(\boldsymbol{\mu}) \prod_{i=1}^n p(x_i | \boldsymbol{\mu}) d\boldsymbol{\mu} \\ &= \int p(\boldsymbol{\mu}) \prod_{i=1}^n \sum_{\mathbf{c}_i} p(\mathbf{c}_i) p(x_i | \mathbf{c}_i, \boldsymbol{\mu}) d\boldsymbol{\mu} \implies \mathcal{O}(K^n) \end{aligned}$$

9 / 45

A general framework for variational inference

- Generally, let \mathbf{z} be our set of latent variables, \mathbf{x} our observed data
- Variational inference approximates $p(\mathbf{z} | \mathbf{x})$ with some $q(\mathbf{z} | \mathbf{x})$



- In the case of the GMM above...

$$\mathbf{z} = \{\mu_1, \dots, \mu_K, \mathbf{c}_1, \dots, \mathbf{c}_n\}$$

$$\mathbf{x} = \{x_1, \dots, x_n\}$$

$$\boldsymbol{\alpha} = \{\sigma^2, \pi_1, \dots, \pi_K\}$$

$$\boldsymbol{\theta} = \{\emptyset\}$$

$$\phi \text{ not defined yet}$$

10 / 45

The ELBO: choosing $q_\phi(\mathbf{z}|\mathbf{x})$ with KL divergence

- $q_\phi(\mathbf{z}|\mathbf{x})$ belongs to set of distributions \mathcal{Q} parameterized by $\phi \in \Phi$.
- VI chooses $q_\phi(\mathbf{z}|\mathbf{x})$ by maximizing the Evidence Lower Bound
- This is akin to finding parameters $\phi \in \Phi$ that minimize the KL distance between $p(\mathbf{z}|\mathbf{x})$ and $q_\phi(\mathbf{z}|\mathbf{x})$.

$$\min_{\phi} \left[D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) \right]$$

11 / 45

Deriving the ELBO

Factor the KL divergence; let $\mathbb{E}_q = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}$

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) \tag{1}$$

$$= \mathbb{E}_q [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_q [\log p(\mathbf{z}|\mathbf{x})] \tag{2}$$

$$= \mathbb{E}_q [\log q_\phi(\mathbf{z}|\mathbf{x})] - \left(\mathbb{E}_q [\log p_\theta(\mathbf{x}|\mathbf{z})] + \mathbb{E}_q [\log p_\alpha(\mathbf{z})] - \log p(\mathbf{x}) \right) \tag{3}$$

$$= D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\alpha(\mathbf{z})) - \mathbb{E}_q [\log p_\theta(\mathbf{x}|\mathbf{z})] + \log p(\mathbf{x}) \tag{4}$$

So the evidence is bounded by the **ELBO**:

$$\begin{aligned} \log p(\mathbf{x}) &= \mathbb{E}_q [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\alpha(\mathbf{z})) + D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}|\mathbf{x})) \\ &\geq \mathbb{E}_q [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\alpha(\mathbf{z})) \end{aligned}$$

12 / 45

Anatomy of the ELBO?

$$\begin{aligned}\text{ELBO} &= \mathbb{E}_q[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\alpha(\mathbf{z})) \\ &= \text{likelihood} \quad + \quad \text{regularizer}\end{aligned}$$

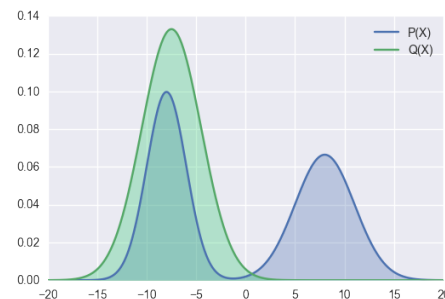
- The first term encourages $q_\phi(\mathbf{z}|\mathbf{x})$ to place mass on \mathbf{z} values that explain \mathbf{x} according to the likelihood, $p_\theta(\mathbf{x}|\mathbf{z})$
- The second term encourages $q_\phi(\mathbf{z}|\mathbf{x})$ to be close to the prior distribution, $p_\alpha(\mathbf{z})$, effectively regularizing the approximate posterior

13 / 45

Reverse KL (used) vs. Forward KL (not used)

Reverse KL zero 'forces':

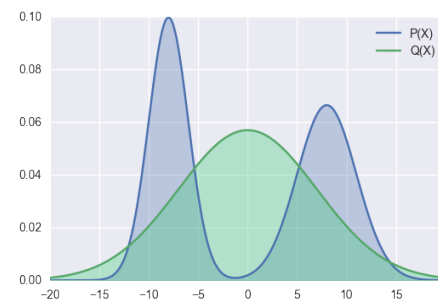
$$\begin{aligned}D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})) \\ = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z}|\mathbf{x})} \right]\end{aligned}$$



img credit: Eric Jang

Forward KL zero 'avoids':

$$\begin{aligned}D_{KL}(p_\theta(\mathbf{z}|\mathbf{x})||q_\phi(\mathbf{z}|\mathbf{x})) \\ = \mathbb{E}_{p_\theta(\mathbf{z}|\mathbf{x})} \left[\log \frac{p_\theta(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]\end{aligned}$$



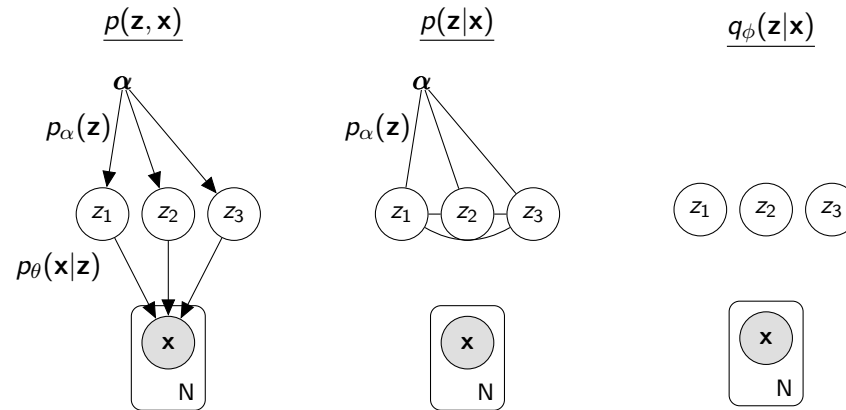
Consequence: support of posterior may be under-represented

14 / 45

Mean Field Variational Inference

- Mean-field variational inference assumes $q_\phi(\mathbf{z}|\mathbf{x})$ completely factors:

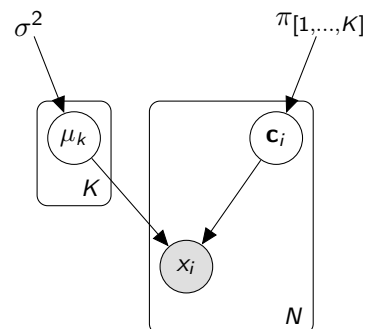
$$q_\phi(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^M q_i(z_i|\mathbf{x})$$



15 / 45

Mean Field on GMM example

- A MF approximation for our GMM applies gaussians for the μ_k 's and categoricals for the \mathbf{c}_i 's:



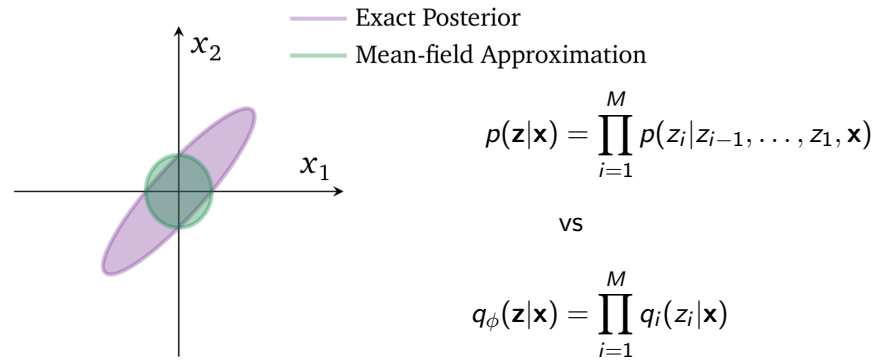
$$q_\phi(\mu_k|\mathbf{x}) = \mathcal{N}(\mu_k; m_k, s_k^2)$$

$$q_\phi(\mathbf{c}_i|\mathbf{x}) = \text{Cat}(\mathbf{c}_i; \psi_1, \dots, \psi_k)$$

- these are optimal forms of mean-field for GMM

16 / 45

Mean Field captures marginal, not covariate characteristics



17 / 45

Coordinate Ascent Variational Inference (CAVI)

- One method to optimize our variational bound is by updating each $q_i(z_i|\mathbf{x})$ while keeping all other Z_{-i} constant.
- The update will look a lot like EM, except we're updating an approximation:

$$q_i(z_i|\mathbf{x}) \propto \exp \left\{ \mathbb{E}_{q_{-i}} \left[\log p(z_i, \mathbf{Z}_{-i}, \mathbf{x}) \right] \right\}$$

- Note that the expectation over q_{-i} only needs to include the markov blanket of i , not all z_i 's

18 / 45

Coordinate Ascent Variational Inference (CAVI)

- CAVI iteratively updates each coordinate: $q_i(z_i|\mathbf{x})$

$$\text{ELBO} = \mathbb{E}_q[\log p(\mathbf{z}, \mathbf{x})] - \mathbb{E}_q[\log q_\phi(\mathbf{z}|\mathbf{x})] \quad (1)$$

$$= \log p(\mathbf{x}) + \mathbb{E}_q[\log p(\mathbf{z}|\mathbf{x})] - \sum_{i=1}^M \mathbb{E}_{q_i}[\log q_i(z_i|\mathbf{x})] \quad (2)$$

$$= \text{const} + \sum_{i=1}^M \left(\mathbb{E}_q[\log p(z_i|z_{1:(i-1)}, \mathbf{x})] - \mathbb{E}_{q_i}[\log q_i(z_i|\mathbf{x})] \right) \quad (3)$$

$$\text{ELBO}_i = \mathbb{E}_{q_i}[\mathbb{E}_{q_{-i}}[\log p(z_i|z_{-i}, \mathbf{x})]] - \mathbb{E}_{q_i}[\log q_i(z_i|\mathbf{x})] \quad (i)$$

$$= \int q_i(z_i) \left(\mathbb{E}_{q_{-i}}[\log p(z_i|z_{-i}, \mathbf{x})] - \log q_i(z_i|\mathbf{x}) \right) dz_i \quad (ii)$$

19 / 45

Coordinate Ascent Variational Inference (CAVI)

$$\text{ELBO}_i = \mathbb{E}_{q_i}[\mathbb{E}_{q_{-i}}[\log p(z_i|z_{-i}, \mathbf{x})]] - \mathbb{E}_{q_i}[\log q_i(z_i|\mathbf{x})] \quad (i)$$

$$= \int q_i(z_i) \left(\mathbb{E}_{q_{-i}}[\log p(z_i|z_{-i}, \mathbf{x})] - \log q_i(z_i|\mathbf{x}) \right) dz_i \quad (ii)$$

$$\frac{d}{dq_i(z_i)} \text{ELBO}_i = \mathbb{E}_{q_{-i}}[\log p(z_i|z_{-i}, \mathbf{x})] - \log q_i(z_i|\mathbf{x}) - 1 = 0 \quad (iii)$$

(5)

solve with Lagrange multipliers:

$$\begin{aligned} q_i(z_i|\mathbf{x}) &\propto \exp \left\{ \mathbb{E}_{q_{-i}}[\log p(z_i|z_{-i}, \mathbf{x})] \right\} \\ &\propto \exp \left\{ \mathbb{E}_{q_{-i}}[\log p(z_i, \mathbf{z}_{-i}, \mathbf{x})] \right\} \end{aligned}$$

20 / 45

How Does VI Compare to Older Methods like EM?

- Doesn't Expectation Maximization (EM) compute latent variables?

$$\text{ELBO} = \mathbb{E}_q[\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\alpha(\mathbf{z}))$$

- EM sets $q_\phi(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$ and iterates to find fixed latent params, \mathbf{z}
- VI, does not assume posterior is tractable – needed for bayesian settings, where latent params are not fixed, but distributed

21 / 45

How Does VI Compare to Older Methods like MCMC?

- MCMC: produce samples that asymptotically approach $p(\mathbf{z}|\mathbf{x})$
 - More accurate, with better guarantees, but computationally very expensive
- VI: run optimization to analytically approximate $p(\mathbf{z}|\mathbf{x})$
 - VI is generally faster, and scales nicely for large datasets
 - Take advantage of stochastic and distributed optimization methods

22 / 45

Open Questions in VI

- VI is not as well understood as MCMC, open areas include
 - Theory to bound accuracy of VI for certain classes of models
 - Improved optimization methods for massive data
 - Developing VI approximations that work for a wide class of models
 - Improving posterior representation without overly complex family of q_ϕ 's

23 / 45

Questions on General Variational Inference?

「_(ツ)_/」

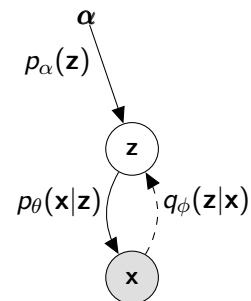
24 / 45

- SGVB (Stochastic Gradient Variational Bayes) estimator
- AEVB (Auto-Encoding Variational Bayes) algorithm
- Variational Auto-encoder

25 / 45

VAE Terminology

- $q_\phi(\mathbf{z}|\mathbf{x})$ is our **encoder**. Given some \mathbf{x} in our observed space, how is \mathbf{z} distributed in latent space?
- $p_\theta(\mathbf{x}|\mathbf{z})$ is our **decoder**. Given some \mathbf{z} in latent space, how is \mathbf{x} distributed in the original space?



26 / 45

Experiments

- $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{x}|\mathbf{z})$ are learned jointly using a neural network.
- Once we are done training the network, we can generate samples from the learned distribution $p_\theta(\mathbf{x}|\mathbf{z})$.
- Datasets
 - MNIST: training set of 60,000 28x28 pixel images of handwritten digits
 - Frey Face: almost 2000 20x28 images of Brendan Frey's face, taken from sequential video frames

27 / 45

MNIST Generated Samples

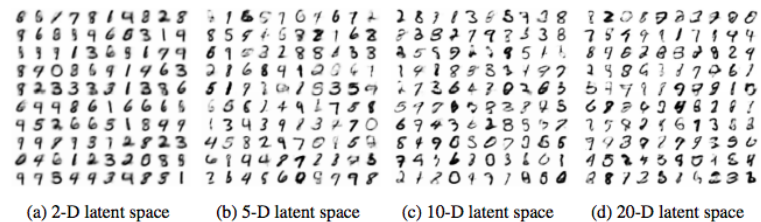


Figure 5: Random samples from learned generative models of MNIST for different dimensionalities of latent space.

Figure: (Kingma & Welling, 2013)

28 / 45

Visualizing Learned 2-D Manifolds

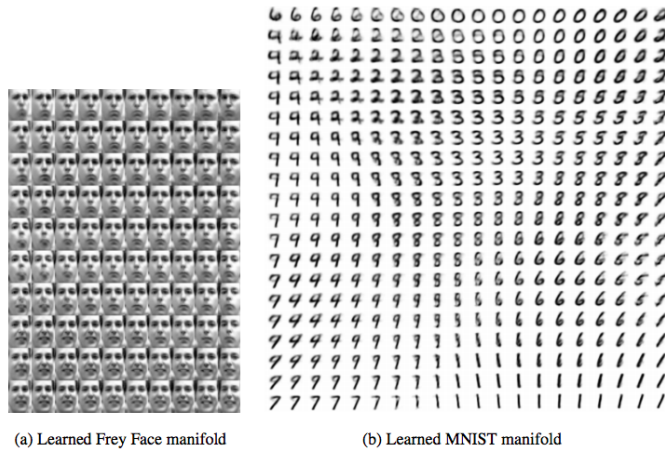


Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables \mathbf{z} . For each of these values \mathbf{z} , we plotted the corresponding generative $p_{\theta}(\mathbf{x}|\mathbf{z})$ with the learned parameters θ .

29 / 45

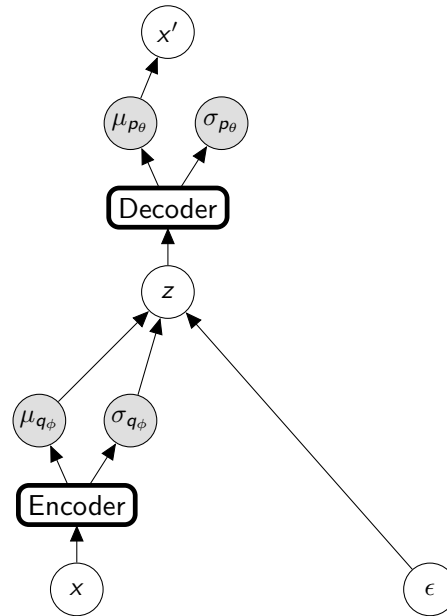
VAE Example

- Let the prior over the latent variables be the centered isotropic multivariate Gaussian:
$$p_{\alpha}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}).$$
- Let $q_{\phi}(\mathbf{z}|\mathbf{x})$ be a multivariate Gaussian with a diagonal covariance structure:
$$q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)}\mathbf{I})$$

where μ and σ are outputs of our encoder network, a neural network with one hidden layer.
- Let $p_{\theta}(\mathbf{x}|\mathbf{z})$ also be a multivariate Gaussian parameterized by a neural network (our decoder network) with one hidden layer.

30 / 45

VAE Diagram



31 / 45

Reparameterization Trick

- We can often reparameterize $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$ as $\mathbf{z} = g(\mathbf{x}, \epsilon)$, where g is some deterministic, differentiable function, and $\epsilon \sim p(\epsilon)$ is a noise variable.
- Gaussian example:
 - We want to sample $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})$.
 - Let $\epsilon \sim \mathcal{N}(0, \mathbf{I})$.
 - Let $g(\mathbf{x}, \epsilon) = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon$.
 - We can write: $\mathbf{z} = g(\mathbf{x}, \epsilon)$
- Instead of sampling \mathbf{z} from q_ϕ directly, we sample ϵ from $\mathcal{N}(0, \mathbf{I})$ and apply the differentiable function g .

32 / 45

Applying the Reparameterization Trick

$$\begin{aligned} E_{q_\phi(z|x)}[f(z)] &= \int q_\phi(z|x) f(z) dz \\ &= \int p(\epsilon) f(z) d\epsilon \\ &= \int p(\epsilon) f(g(x, \epsilon)) d\epsilon \\ &= E_{p(\epsilon)}[f(g(x, \epsilon))] \end{aligned}$$

33 / 45

When can we apply the reparameterization trick?

- *When q belongs to any location-scale family*

e.g. Gaussian, Laplace

- *When q has a tractable inverse CDF*

Let g be the inverse CDF and let $\epsilon \sim \mathcal{U}(0, 1)$.

e.g. Exponential, Cauchy, Logistic

- *When we can express $z \sim q$ as some compositional transformation of auxiliary variables*

e.g. Log-Normal, Dirichlet

34 / 45

Recall that in VI, we select $q_\phi(\mathbf{z}|\mathbf{x})$ to maximize the ELBO (Evidence Lower Bound):

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\alpha(\mathbf{z})) + E_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})}[\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})]$$

35 / 45

Why differentiate the ELBO?

- If we want to train the neural network from the earlier example through backprop, we need a differentiable objective function.
- Even if we don't intend to parameterize p_θ and q_ϕ with a neural network, the ability to use stochastic gradient ascent techniques to optimize the ELBO can help VI scale to large datasets.

36 / 45

Auto-Encoding Variational Bayes Algorithm

Algorithm 1 Minibatch version of the AEVB algorithm.

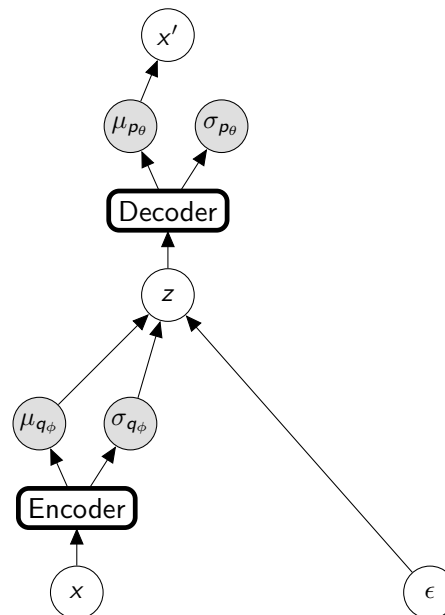
- 1: $\theta, \phi \leftarrow$ Initialize parameters
- 2: **repeat**
- 3: $\mathbf{X}^M \leftarrow$ Random minibatch of M datapoints
- 4: $\epsilon \leftarrow$ Random samples from $p(\epsilon)$
- 5: $\mathbf{g} \leftarrow \nabla_{\theta, \phi} \tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon)$
- 6: $\theta, \phi \leftarrow$ Update parameters using gradients (e.g. SGD or Adagrad)
- 7: **until** convergence of parameters (θ, ϕ)
- 8: **return** θ, ϕ

$$\tilde{\mathcal{L}}^M(\theta, \phi; \mathbf{X}^M, \epsilon) = \frac{N}{M} \sum_{i=1}^M \tilde{\mathcal{L}}(\theta, \phi; x^{(i)})$$

where $\tilde{\mathcal{L}}$ is the differentiable estimator of the ELBO that we will derive next. ($\tilde{\mathcal{L}}^M$ is the minibatch version.)

37 / 45

VAE Diagram



38 / 45

- ELBO (Evidence Lower Bound)
$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) = -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\alpha(\mathbf{z})) + E_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})}[\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})]$$
- How to differentiate with respect to ϕ and θ ?
- In practice, we can often obtain closed form expression for the first term (KL-divergence).
- For the second term we need to do some kind of estimation by sampling.

39 / 45

Deriving the SGVB Estimator

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

where $\mathbf{z}^{(i,l)} = g_\phi(\mathbf{x}^{(i)}, \epsilon^{(i,l)})$ and $\epsilon^{(l)} \sim p(\epsilon)$

This gives us our SGVB estimator:

$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)}) \approx -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\alpha(\mathbf{z})) + \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})$$

40 / 45

Back to Our Original Example

- Let the prior over the latent variables be the centered isotropic multivariate Gaussian:

$$p_\alpha(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}).$$

- Let $q_\phi(\mathbf{z}|\mathbf{x})$ be a multivariate Gaussian with a diagonal covariance structure:

$$q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)}\mathbf{I})$$

where μ and σ are outputs of our encoder network, a neural network with one hidden layer.

- Claim:

$$\begin{aligned} & -D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\alpha(\mathbf{z})) \\ &= \frac{1}{2} \sum_{j=1}^J (1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2) \end{aligned}$$

41 / 45

Deriving an Expression for the KL-Divergence

$$\begin{aligned} -D_{KL}(Q||P) &= -\int q(z) \log \frac{q(z)}{p(z)} dz \\ &= \int (q(z) \log p(z) dz - \int q(z) \log q(z) dz) \\ &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) - \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \\ &= [-\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\mu_j^2 + \sigma_j^2)] - [-\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2)] \\ &= \frac{1}{2} \sum_{j=1}^J (1 + \log((\sigma_j)^2) - (\mu_j)^2 - (\sigma_j)^2) \end{aligned}$$

Note: J denotes the dimensionality of z.

42 / 45

Final Result

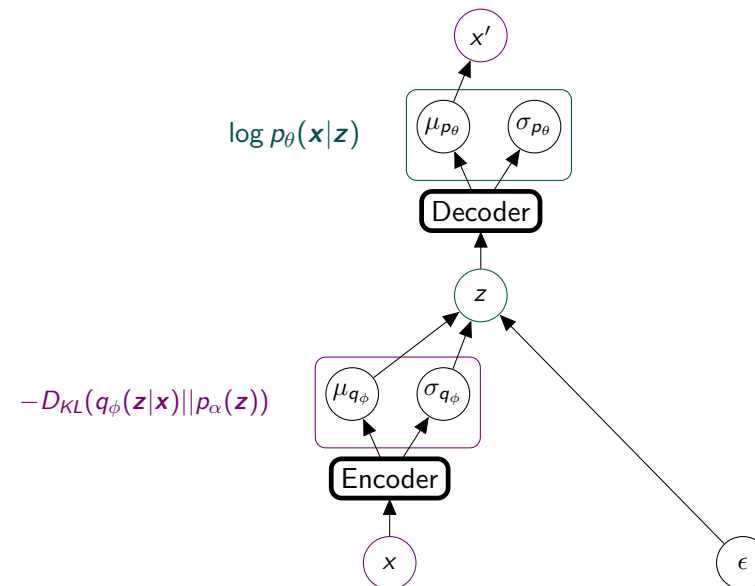
$$\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$$

$$\approx \frac{1}{2} \sum_{j=1}^J (1 + \log((\sigma_j)^2) - (\mu_j)^2 - (\sigma_j)^2) + \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)})$$

where $\mathbf{z}^{(i,l)} = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \boldsymbol{\epsilon}^{(l)}$ and $\boldsymbol{\epsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

43 / 45

VAE Diagram Revisited



44 / 45

See Also

- These are the specific papers we talked about:

Kingma, Diederik P and Welling, Max. Auto-Encoding Variational Bayes. In *The 2nd International Conference on Learning Representations (ICLR)*, 2013.

David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: a review for statisticians. *arXiv:1601.00670*, 2016.

- This is a helpful tutorial that we also consulted:

C. Doersch. Tutorial on variational autoencoders. *arXiv:1606.05908*, 2016.

- Many variations of the standard VAE exist. For example:

- Conditional VAEs
- VAE-GANs