

# Variational Autoencoders

## STATS305B: Applied Statistics II

Scott Linderman

February 24, 2025

# Last Time...

## Outline:

- ▶ Principal Components Analysis (PCA)
- ▶ PCA as a linear Gaussian latent variable model
- ▶ Factor analysis
- ▶ Linear Dynamical Systems & the Kalman Filter/Smoothen

# Today...

## Outline:

- ▶ (Probabilistic) PCA as a (stochastic) linear autoencoder
- ▶ Variational Autoencoders (VAEs)
- ▶ Review of Gradient-based Variational Inference
- ▶ Demo

## PCA as a Linear Autoencoder

Last time we introduced PCA as a method for finding dimensions of maximal variance. However, we can arrive at the same model from another perspective: find the linear projection that minimizes the average projection cost.

To formalize this, let

$$W = \begin{bmatrix} | & & | \\ \mathbf{w}_1 & \cdots & \mathbf{w}_M \\ | & & | \end{bmatrix} \in \mathbb{R}^{D \times M} \quad \text{with} \quad W^T W = I$$

be an orthogonal basis for the principal subspace.

We will **encode** each data point by subtracting the mean and projecting onto the principal subspace to obtain  $\mathbf{z}_n = W^T (\mathbf{x}_n - \bar{\mathbf{x}})$ .

## PCA as a Linear Autoencoder

Since  $\mathbf{W}$  is an orthogonal matrix, all we need to do to **decode** the encoded data point is multiply  $\mathbf{W}\mathbf{z}_n$  and add back the mean. That gives us,

$$\hat{\mathbf{x}}_n = \mathbf{W}\mathbf{z}_n + \bar{\mathbf{x}} = \mathbf{W}\mathbf{W}^\top(\mathbf{x}_n - \bar{\mathbf{x}}) + \bar{\mathbf{x}}.$$

**Goal:** Find an orthogonal matrix  $\mathbf{W}$  that *minimizes* the mean squared reconstruction error,

$$\mathcal{L}(\mathbf{W}) = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \hat{\mathbf{x}}_n\|_2^2$$

## PCA as a Linear Autoencoder

We can write this in matrix notation instead. Let  $\mathbf{X} \in \mathbb{R}^{N \times D}$  be the *centered data matrix* with rows  $(\mathbf{x}_n - \bar{\mathbf{x}})^\top$ . Then,

$$\begin{aligned}\mathcal{L}(\mathbf{W}) &= \frac{1}{N} \text{Tr}[(\mathbf{X} - \hat{\mathbf{X}})^\top (\mathbf{X} - \hat{\mathbf{X}})] \\ &= \frac{1}{N} \text{Tr}[(\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^\top)^\top (\mathbf{X} - \mathbf{X}\mathbf{W}\mathbf{W}^\top)] \\ &= \frac{1}{N} \text{Tr}[(\mathbf{X}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top))^\top (\mathbf{X}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top))] \\ &= \frac{1}{N} \text{Tr}[(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)^\top \mathbf{X}^\top \mathbf{X} (\mathbf{I} - \mathbf{W}\mathbf{W}^\top)] \\ &= \text{Tr}[(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)^\top \mathbf{S} (\mathbf{I} - \mathbf{W}\mathbf{W}^\top)]\end{aligned}$$

where  $\mathbf{S} = \frac{1}{N} \mathbf{X}^\top \mathbf{X}$  is the sample covariance matrix.

## PCA as a Linear Autoencoder

Now apply the circular trace property,

$$\mathcal{L}(\mathbf{W}) = \text{Tr}[\mathbf{S}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)^\top]$$

**Question:** What does  $(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)^\top$  equal?

Note that  $\mathbf{I} - \mathbf{W}\mathbf{W}^\top$  is a projection matrix – it projects a vector onto the *nullspace* of  $\mathbf{W}$ . Applying the projection operator twice doesn't change the result. Mathematically,

$$(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)^\top = \mathbf{I} - 2\mathbf{W}\mathbf{W}^\top + \mathbf{W}\mathbf{W}^\top\mathbf{W}\mathbf{W}^\top = \mathbf{I} - \mathbf{W}\mathbf{W}^\top$$

where we used the fact that  $\mathbf{W}^\top\mathbf{W} = \mathbf{I}$  since  $\mathbf{W}$  is an orthogonal matrix.

Thus, the objective simplifies to,

$$\mathcal{L}(\mathbf{W}) = \text{Tr}[\mathbf{S}(\mathbf{I} - \mathbf{W}\mathbf{W}^\top)] = \text{Tr}[\mathbf{S}] - \text{Tr}[\mathbf{S}\mathbf{W}\mathbf{W}^\top] = \text{const} - \text{Tr}[\mathbf{W}^\top\mathbf{S}\mathbf{W}].$$

## PCA as a Linear Autoencoder

Let  $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$  be the eigendecomposition of  $\mathbf{S}$ . (Since it is a covariance matrix, the eigenvectors are orthogonal.) Plugging in,

$$\begin{aligned}\mathcal{L}(\mathbf{W}) &= \text{const} - \text{Tr}[\mathbf{W}^\top \mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top \mathbf{W}] \\ &= \text{const} - \text{Tr}\left[\mathbf{W}^\top \left(\sum_{d=1}^D \lambda_d \mathbf{u}_d \mathbf{u}_d^\top\right) \mathbf{W}\right] \\ &= \text{const} - \sum_{m=1}^M \sum_{d=1}^D \lambda_d \mathbf{w}_m^\top \mathbf{u}_d \mathbf{u}_d^\top \mathbf{w}_m \\ &= \text{const} - \sum_{m=1}^M \sum_{d=1}^D \lambda_d (\mathbf{w}_m^\top \mathbf{u}_d)^2\end{aligned}$$

We want to minimize  $\mathcal{L}(\mathbf{W})$  subject to  $\mathbf{W}$  being orthogonal. What is the solution?  $\mathbf{W} = \mathbf{U}_M$ !



## Probabilistic PCA as a Stochastic Linear Autoencoder

The jump from PCA to probabilistic PCA was to treat the scores,  $\mathbf{z}_n$ , as latent variables. We gave them a Gaussian prior and assumed a generative model,

$$\begin{aligned}\mathbf{z}_n &\stackrel{\text{iid}}{\sim} \mathcal{N}(\mathbf{0}, I) \\ \mathbf{x}_n | \mathbf{z}_n &\sim \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 I),\end{aligned}$$

where  $\mathbf{z}_n \in \mathbb{R}^M$  is a latent variable,  $\mathbf{W} \in \mathbb{R}^{D \times M}$  are the weights,  $\boldsymbol{\mu} \in \mathbb{R}^D$  is the bias parameter, and  $\sigma^2 \in \mathbb{R}_+$  is a variance.

Given the parameters, we showed that the posterior distribution over latent variables is,

$$p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}) = \mathcal{N}(\mathbf{z}_n | (\sigma^2 I + \mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top (\mathbf{x}_n - \boldsymbol{\mu}), \sigma^2 (\sigma^2 I + \mathbf{W}^\top \mathbf{W})^{-1}),$$

and when  $\sigma^2 \rightarrow 0$  and  $\mathbf{W}$  are the PCs, this converges to a delta function on the PCA scores.

## Probabilistic PCA as a Stochastic Linear Autoencoder

To simplify notation, let's rewrite the posterior as,

$$p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}) = \text{N}(\mathbf{z}_n | \mathbf{E}\mathbf{x}_n + \mathbf{d}, \mathbf{G})$$

where

$$\mathbf{E} = (\sigma^2 \mathbf{I} + \mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top$$

$$\mathbf{d} = -(\sigma^2 \mathbf{I} + \mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \boldsymbol{\mu}$$

$$\mathbf{G} = \sigma^2 (\sigma^2 \mathbf{I} + \mathbf{W}^\top \mathbf{W})^{-1}.$$

To “autoencode” the data using probabilistic PCA, we could sample  $\mathbf{z}_n \sim p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta})$  by mapping  $\mathbf{x}_n \mapsto \mathbf{E}\mathbf{x}_n + \mathbf{d}$  and adding Gaussian noise with covariance  $\mathbf{G}$ . Then we can generate a new data point  $\hat{\mathbf{x}}_n \sim \text{N}(\mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ .

Just as PCA could be motivated as finding weights that minimize the reconstruction error, probabilistic PCA can be seen as finding weights that minimize the *expected* reconstruction error (subject to a little regularization).

## Revisiting EM for Probabilistic PCA

Recall that we justified EM by connecting it to variational inference. We showed that EM maximizes the evidence lower bound (ELBO),

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_n \mathbb{E}_{q(\mathbf{z}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \theta) - \log q(\mathbf{z}_n)] \\ &= \sum_n \underbrace{\mathbb{E}_{q(\mathbf{z}_n)} [\log p(\mathbf{x}_n | \mathbf{z}_n; \theta)]}_{\text{expected log likelihood}} - \underbrace{D_{\text{KL}}(q(\mathbf{z}_n) \| p(\mathbf{z}_n))}_{\text{KL to the prior}} \\ &\leq \sum_n \log p(\mathbf{x}_n; \theta)\end{aligned}$$

EM can be seen as coordinate ascent on the ELBO:

1. The E-step sets  $q$  to the posterior over latent variables,  $q(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_n; \theta)$  for the current parameters  $\theta$ .
2. The M-step updates the parameters to maximize the expected log joint probability.

## Revisiting EM for Probabilistic PCA

In probabilistic PCA, maximizing the expected log likelihood is the same as minimizing the expected reconstruction error. Plugging in the definition of the model and the posterior,

$$\begin{aligned}\mathbb{E}_{q(\mathbf{z}_n)} [\log p(\mathbf{x}_n | \mathbf{z}_n; \boldsymbol{\theta})] &= \mathbb{E}_{\mathbf{N}(\mathbf{z}_n | \mathbf{E}\mathbf{x}_n + \mathbf{d}, \mathbf{G})} [\log \mathbf{N}(\mathbf{x}_n | \mathbf{W}\mathbf{z}_n + \boldsymbol{\mu}, \sigma^2 \mathbf{I})] \\ &= -\frac{1}{2\sigma^2} \mathbb{E}_{\mathbf{N}(\mathbf{z}_n | \mathbf{E}\mathbf{x}_n + \mathbf{d}, \mathbf{G})} [\|\mathbf{x}_n - \mathbf{W}\mathbf{z}_n - \boldsymbol{\mu}\|_2^2] + \log \sigma + c\end{aligned}$$

where  $c$  is an additive constant, and  $\mathbf{E}$ ,  $\mathbf{d}$ ,  $\mathbf{G}$  are functions of the parameters, as defined above,

We can't forget about the KL to the prior though! When we maximize the ELBO, we are essentially finding weights that minimize the expected reconstruction error, while also not deviating too far from the standard normal prior on the latent variables.

# Variational Autoencoders

Variational Autoencoders (VAEs) are “deep” but conceptually simple generative models. The generative model is the same as probabilistic PCA, but allowing for **nonlinear** mappings between latent variables and observations.

To sample a data point  $\mathbf{x}_n$ ,

1. First, sample **latent variables**  $\mathbf{z}_n$ ,

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

2. Then sample the data point  $\mathbf{x}_n$  from a conditional distribution with mean,

$$\mathbb{E}[\mathbf{x}_n | \mathbf{z}_n] = g(\mathbf{z}_n; \boldsymbol{\theta}),$$

where  $g : \mathbb{R}^H \rightarrow \mathbb{R}^D$  is a nonlinear mapping parameterized by  $\boldsymbol{\theta}$ .

## Variational Autoencoders

We will assume  $g$  is a simple **feedforward neural network** of the form,

$$g(\mathbf{z}; \boldsymbol{\theta}) = g_L(g_{L-1}(\cdots g_1(\mathbf{z}) \cdots))$$

where each **layer** is a cascade of a linear mapping followed by an element-wise nonlinearity (except for the last layer, perhaps). For example,

$$g_\ell(\mathbf{u}_\ell) = \text{relu}(\mathbf{W}_\ell \mathbf{u}_\ell + \mathbf{b}_\ell); \quad \text{relu}(a) = \max(0, a).$$

The generative parameters consist of the weights and biases,  $\boldsymbol{\theta} = \{\mathbf{W}_\ell, \mathbf{b}_\ell\}_{\ell=1}^L$ .

# Learning and Inference

We have two goals. The **learning goal** is to find the parameters that **maximize the marginal likelihood of the data**,

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(\mathbf{X}; \theta) \\ &= \arg \max_{\theta} \prod_{n=1}^N \int p(\mathbf{x}_n | \mathbf{z}_n; \theta) p(\mathbf{z}_n; \theta) d\mathbf{z}_n\end{aligned}$$

The **inference goal** is to find the **posterior distribution of latent variables**,

$$p(\mathbf{z}_n | \mathbf{x}_n; \theta) = \frac{p(\mathbf{x}_n | \mathbf{z}_n; \theta) p(\mathbf{z}_n; \theta)}{\int p(\mathbf{x}_n | \mathbf{z}'_n; \theta) p(\mathbf{z}'_n; \theta) d\mathbf{z}'_n}$$

Both goals require an integral over  $\mathbf{z}_n$ , but that is intractable for deep generative models.

## The Evidence Lower Bound (ELBO)

**Idea:** Use the ELBO to get a bound on the marginal probability and maximize that instead.

$$\begin{aligned}\log p(\mathbf{X}; \boldsymbol{\theta}) &= \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) \\ &\geq \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{z}_n; \boldsymbol{\lambda}_n) \parallel p(\mathbf{z}_n \mid \mathbf{x}_n; \boldsymbol{\theta})) \\ &= \sum_{n=1}^N \underbrace{\mathbb{E}_{q(\mathbf{z}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n; \boldsymbol{\lambda}_n)]}_{\text{"local ELBO"}} \\ &\triangleq \sum_{n=1}^N \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\lambda}, \boldsymbol{\theta})\end{aligned}$$

where  $\boldsymbol{\lambda} = \{\boldsymbol{\lambda}_n\}_{n=1}^N$ . Here, I've written the ELBO as a sum of *local ELBOs*  $\mathcal{L}_n$ .



## Variational Inference

The ELBO is still maximized (and the bound is tight) when each  $q$  is equal to the true posterior,

$$q(\mathbf{z}_n; \boldsymbol{\lambda}_n) = p(\mathbf{z}_n | \mathbf{x}_n, \boldsymbol{\theta}).$$

Unfortunately, the posterior no longer has a simple, closed form.

**Question:** Suppose  $\mathbf{x}_n \sim N(g(\mathbf{z}_n; \boldsymbol{\theta}), I)$ . This deep generative model has a Gaussian prior on  $\mathbf{z}_n$  and a Gaussian likelihood for  $\mathbf{x}_n$  given  $\mathbf{z}_n$ . Why isn't the posterior Gaussian?

Nevertheless, we can still constrain  $q$  to belong to a simple family. For example, we could constrain it to be Gaussian and seek the best Gaussian approximation to the posterior. This is sometimes called **fixed-form variational inference**. Let,

$$\mathcal{Q} = \{q : q(\mathbf{z}; \boldsymbol{\lambda}) = N(\mathbf{z} | \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \text{ for } \boldsymbol{\lambda} = (\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2) \in \mathbb{R}^{2H}\}$$

## Variational Inference

Then, for fixed parameters  $\theta$ , the best  $q$  in this **variational family** is,

$$\begin{aligned} q^* &= \arg \min_{q \in \mathcal{Q}} D_{\text{KL}}(q(\mathbf{z}_n; \boldsymbol{\lambda}_n) \parallel p(\mathbf{z}_n \mid \mathbf{x}_n; \boldsymbol{\theta})) \\ &= \arg \max_{\boldsymbol{\lambda}_n \in \mathbb{R}^{2H}} \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}). \end{aligned}$$

# Variational Expectation-Maximization (vEM)

Now we can introduce a new algorithm.

**Algorithm** Variational EM (vEM) Repeat until either the ELBO or the parameters converges:

- 1. M-step:** Set  $\theta \leftarrow \arg \max_{\theta} \mathcal{L}(\lambda, \theta)$
- 2. E-step:** Set  $\lambda_n \leftarrow \arg \max_{\lambda_n \in \Lambda} \mathcal{L}_n(\lambda_n, \theta)$  for  $n = 1, \dots, N$
- 3.** Compute (an estimate of) the ELBO  $\mathcal{L}(\lambda, \theta)$ .

In general, none of these steps will have closed form solutions, so we'll have to use approximations.

## Generic M-step with Stochastic Gradient Ascent

For exponential family mixture models, the M-step had a closed form solution. For deep generative models, we need a more general approach.

If the parameters are unconstrained and the ELBO is differentiable wrt  $\theta$ , we can use stochastic gradient ascent.

$$\begin{aligned}\theta &\leftarrow \theta + \alpha \nabla_{\theta} \mathcal{L}(q, \theta) \\ &= \theta + \alpha \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} [\nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n; \theta)]\end{aligned}$$

Note that the expected gradient wrt  $\theta$  can be computed using ordinary Monte Carlo – nothing fancy needed!

## The Variational E-step

Assume  $\mathcal{Q}$  is the family of Gaussian distributions with diagonal covariance:

$$\mathcal{Q} = \{q : q(\mathbf{z}; \boldsymbol{\lambda}) = \text{N}(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \text{ for } \boldsymbol{\lambda} = (\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2) \in \mathbb{R}^{2H}\}$$

This family is indexed by **variational parameters**  $\boldsymbol{\lambda}_n = (\boldsymbol{\mu}_n, \log \boldsymbol{\sigma}_n^2) \in \mathbb{R}^{2H}$ .

To perform SGD, we need an unbiased estimate of the gradient of the local ELBO, but

$$\begin{aligned} \nabla_{\boldsymbol{\lambda}_n} \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) &= \nabla_{\boldsymbol{\lambda}_n} \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n; \boldsymbol{\lambda}_n)] \\ &\neq \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\nabla_{\boldsymbol{\lambda}_n} (\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n; \boldsymbol{\lambda}_n))]. \end{aligned}$$

## Reparameterization Trick

One way around this problem is to use the **reparameterization trick**, aka the **pathwise gradient estimator**. Note that,

$$\mathbf{z}_n \sim q(\mathbf{z}_n; \boldsymbol{\lambda}_n) \iff \mathbf{z}_n = r(\boldsymbol{\lambda}_n, \boldsymbol{\epsilon}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I)$$

where  $r(\boldsymbol{\lambda}_n, \boldsymbol{\epsilon}) = \boldsymbol{\mu}_n + \boldsymbol{\sigma}_n \boldsymbol{\epsilon}$  is a reparameterization of  $\mathbf{z}_n$  in terms of parameters  $\boldsymbol{\lambda}_n$  and noise  $\boldsymbol{\epsilon}$ .

We can use the **law of the unconscious statistician** to rewrite the expectations as,

$$\mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [h(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}, \boldsymbol{\lambda}_n)] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I)} [h(\mathbf{x}_n, r(\boldsymbol{\lambda}_n, \boldsymbol{\epsilon}), \boldsymbol{\theta}, \boldsymbol{\lambda}_n)]$$

where

$$h(\mathbf{x}_n, \mathbf{z}_n, \boldsymbol{\theta}, \boldsymbol{\lambda}_n) = \log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n; \boldsymbol{\lambda}_n).$$

## Reparameterization Trick

The distribution that the expectation is taken under no longer depends on the parameters  $\lambda_n$ , so we can simply take the gradient inside the expectation,

$$\nabla_{\lambda} \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} [h(\mathbf{x}_n, \mathbf{z}_n, \theta, \lambda_n)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I)} [\nabla_{\lambda_n} h(\mathbf{x}_n, r(\lambda_n, \epsilon), \theta, \lambda_n)]$$

Now we can use Monte Carlo to obtain an unbiased estimate of the final expectation!

## Working with mini-batches of data

We can view the ELBO as an expectation over data indices,

$$\begin{aligned}\mathcal{L}(\boldsymbol{\lambda}, \boldsymbol{\theta}) &= \sum_{n=1}^N \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) \\ &= N \mathbb{E}_{n \sim \text{Unif}([N])} [\mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta})].\end{aligned}$$

We can use Monte Carlo to approximate the expectation (and its gradient) by drawing **mini-batches** of data points at random.

In practice, we often cycle through mini-batches of data points deterministically. Each pass over the whole dataset is called an **epoch**.



# Algorithm

Now we can add some detail to our variational expectation maximization algorithm.

For epoch  $i = 1, \dots, \infty$ :

For  $n = 1, \dots, N$ :

1. Sample  $\epsilon_n^{(m)} \stackrel{\text{iid}}{\sim} \mathbf{N}(\mathbf{0}, I)$  for  $m = 1, \dots, M$ .

2. **M-Step:**

2.1 Estimate

$$\hat{\nabla}_{\theta} \mathcal{L}_n(\lambda_n, \theta) = \frac{1}{M} \sum_{m=1}^M \left[ \nabla_{\theta} \log p(\mathbf{x}_n, r(\lambda_n, \epsilon_n^{(m)}); \theta) \right]$$

2.2 Set  $\theta \leftarrow \theta + \alpha_j N \hat{\nabla}_{\theta} \mathcal{L}_n(\lambda_n, \theta)$

3. **E-step:**

# Algorithm

## 3.1 Estimate

$$\hat{\nabla}_{\lambda} \mathcal{L}_n(\lambda_n, \theta) = \frac{1}{M} \sum_{m=1}^M \nabla_{\lambda} \left[ \log p(\mathbf{x}_n, r(\lambda_n, \epsilon_n^{(m)}); \theta) - \log q(r(\lambda_n, \epsilon_n^{(m)}), \lambda_n) \right]$$

3.2 Set  $\lambda_n \leftarrow \lambda_n + \alpha_i \hat{\nabla}_{\lambda} \mathcal{L}_n(\lambda_n, \theta)$ .

## 4. Estimate the ELBO

$$\hat{\mathcal{L}}(\lambda, \theta) = \frac{N}{M} \sum_{m=1}^M \log p(\mathbf{x}_n, r(\lambda_n, \epsilon_n^{(m)}); \theta) - \log q(r(\lambda_n, \epsilon_n^{(m)}), \lambda_n)$$

5. Decay step size  $\alpha_i$  according to schedule.

## Amortized Inference

Note that vEM involves optimizing separate variational parameters  $\lambda_n$  for each data point. For large datasets where we are optimizing using mini-batches of data points, this leads to a strange asymmetry: we update the generative model parameters  $\theta$  every mini-batch, but we only update the variational parameters for the  $n$ -th data point once per epoch. Is there any way to share information across data points?

Note that the optimal variational parameters are just a function of the data point and the model parameters,

$$\lambda_n^* = \arg \min_{\lambda_n} D_{\text{KL}}(q(\mathbf{z}_n; \lambda_n) \parallel p(\mathbf{z}_n \mid \mathbf{x}_n, \theta)) \triangleq f^*(\mathbf{x}_n, \theta).$$

for some implicit and generally nonlinear function  $f^*$ .

VAEs learn an approximation to  $f^*(\mathbf{x}_n, \theta)$  with an **inference network**, a.k.a. **recognition network** or **encoder**.

## Amortized Inference

The inference network is (yet another) neural network that takes in a data point  $\mathbf{x}_n$  and outputs variational parameters  $\lambda_n$ ,

$$\lambda_n \approx f(\mathbf{x}_n, \phi),$$

where  $\phi$  are the weights of the network.

The advantage is that the inference network shares information across data points – it *amortizes* the cost of inference, hence the name. The disadvantage is the output will not minimize the KL divergence. However, in practice we might tolerate a worse variational posterior and a weaker lower bound if it leads to faster optimization of the ELBO overall.

## Putting it all together

Logically, I find it helpful to distinguish between the E and M steps, but with recognition networks and stochastic gradient ascent, the line is blurred.

The final algorithm looks like this.

### Variational EM (with amortized inference)

Repeat until either the ELBO or the parameters converges:

1. Sample data point  $n \sim \text{Unif}(1, \dots, N)$ . [Or a minibatch of data points.]
2. Estimate the local ELBO  $\mathcal{L}_n(\phi, \theta)$  with Monte Carlo. [Note: it is a function of  $\phi$  instead of  $\lambda_n$ .]
3. Compute unbiased Monte Carlo estimates of the gradients  $\widehat{\nabla}_{\theta} \mathcal{L}_n(\phi, \theta)$  and  $\widehat{\nabla}_{\phi} \mathcal{L}_n(\phi, \theta)$ . [The latter requires the reparameterization trick.]

## Putting it all together

### 4. Set

$$\theta \leftarrow \theta + \alpha_i \widehat{\nabla}_{\theta} \mathcal{L}_n(\phi, \theta)$$

$$\phi \leftarrow \phi + \alpha_i \widehat{\nabla}_{\phi} \mathcal{L}_n(\phi, \theta)$$

with step size  $\alpha_i$  decreasing over iterations  $i$  according to a valid schedule.