

Variational Autoencoders (VAEs)

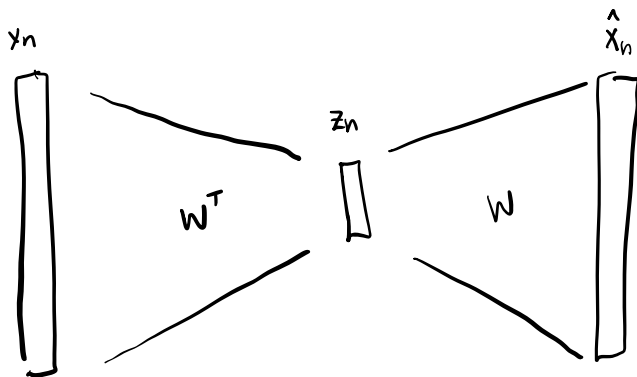
STATS 305C: Applied Statistics

Scott Linderman

May 2, 2022

PCA as a linear autoencoder

Recall from Lecture 5 that PCA could be motivated as a linear autoencoder trained to minimize reconstruction error subject to having orthogonal weights.



$$W = \operatorname{argmin} \|x_n - \hat{x}_n\|_2^2$$
$$\text{s.t. } W^T W = I$$

Deep autoencoders

Why restrict ourselves to **linear** autoencoders? The neural network community has used **deep autoencoders** (a.k.a. autoassociative networks) for nonlinear dimensionality reduction [LeCun, 1987, Bourlard and Kamp, 1988, Hinton and Zemel, 1993, Hinton and Salakhutdinov, 2006, Vincent et al., 2010]. See also, Goodfellow et al. [2016, Ch. 14].

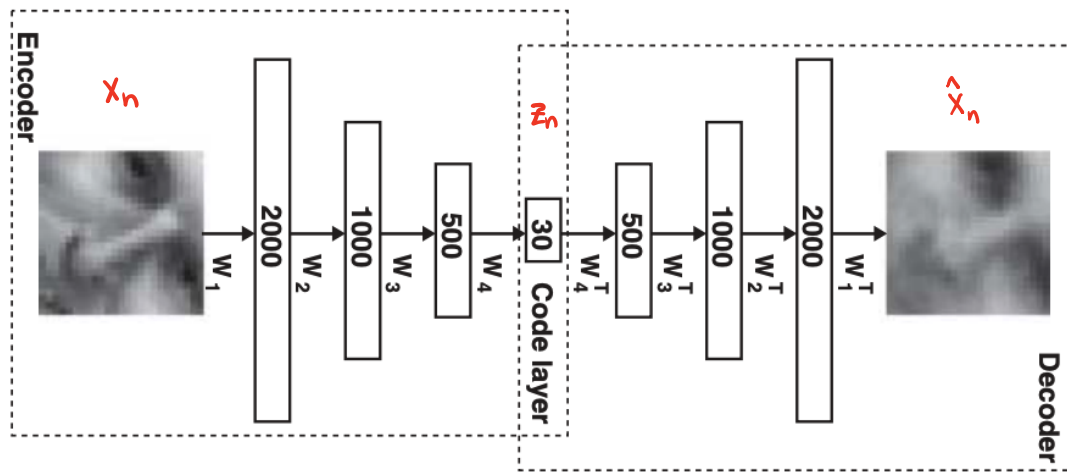
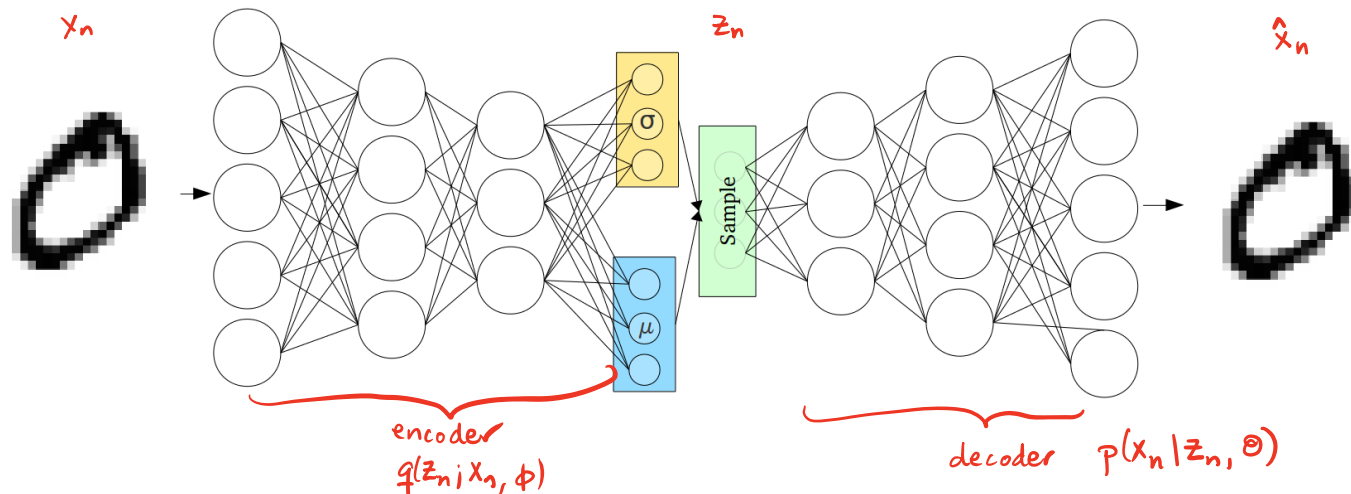


Figure: Figure from Hinton and Salakhutdinov [2006]

Variational autoencoders as deep, stochastic, regularized autoencoders

Kingma and Welling [2014] and Rezende et al. [2014] concurrently developed what we now call **variational autoencoders**. The idea is to treat the hidden codes as random variables. As we will see, VAEs can be viewed as **deep generative models** combined with **amortized variational inference**.



$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{q(z_n; x_n, \phi)} [\log p(x_n | z_n; \theta)] - D_{\text{KL}}(q(z_n; x_n, \phi) \| p(z_n)) \quad (1)$$

Outline

- ▶ The generative model
- ▶ Learning via variational expectation maximization
- ▶ Stochastic gradient ascent
- ▶ Unbiased gradient estimators
- ▶ Amortized inference

The generative model

VAEs start with a “deep” but conceptually simple generative model,

$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{0}, I) \quad \mathbb{R}^H \quad (2)$$

$$\mathbf{x}_n \sim \mathcal{N}(g(\mathbf{z}_n; \boldsymbol{\theta}), I) \quad (3)$$

where $g : \mathbb{R}^H \rightarrow \mathbb{R}^D$ is a **nonlinear mapping** from $\mathbf{z}_n \in \mathbb{R}^H$ to $\mathbb{E}[\mathbf{x}_n] \in \mathbb{R}^D$, parameterized by $\boldsymbol{\theta}$.

We will assume g is a simple **feedforward neural network** (a.k.a. multilayer perceptron) of the form,

$$g(\mathbf{z}; \boldsymbol{\theta}) = g_L(g_{L-1}(\cdots g_1(\mathbf{z}) \cdots)) \quad (4)$$

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

$$g_\ell(\mathbf{u}_\ell) = \text{relu}(\mathbf{W}_\ell \mathbf{u}_\ell + \mathbf{b}_\ell); \quad \text{relu}(a) = \max(0, a). \quad (5)$$

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

Two goals

The **learning goal** is to find the parameters that **maximize the marginal probability of the data**,

$$\theta^* = \arg \max_{\theta} p(\mathbf{X}; \theta) \quad (6)$$

$$= \arg \max_{\theta} \prod_{n=1}^N \int p(\mathbf{x}_n | \mathbf{z}_n; \theta) p(\mathbf{z}_n; \theta) d\mathbf{z}_n \quad (7)$$

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

$$p(\mathbf{z}_n | \mathbf{x}_n; \theta) = \frac{p(\mathbf{x}_n | \mathbf{z}_n; \theta) p(\mathbf{z}_n; \theta)}{\int p(\mathbf{x}_n | \mathbf{z}'_n; \theta) p(\mathbf{z}'_n; \theta) d\mathbf{z}'_n} \quad (8)$$

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.

"ELBO Surgery" (Hoffman + Johnson, —)

$$\log p(\mathbf{X}; \boldsymbol{\theta}) = \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) \quad (9)$$

$$\begin{aligned} \textcircled{1} \quad & \mathbb{E}_{q(\mathbf{z}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})] \\ & - \beta \text{KL}(q(\mathbf{z}_n) \parallel p(\mathbf{z}_n)) \end{aligned} \geq \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) - D_{\text{KL}}(q_n(\mathbf{z}_n) \parallel p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta})) \quad [\text{for any } q_n] \quad (10)$$

$$\begin{aligned} \textcircled{2} \quad & \mathbb{E}_{q(\mathbf{z}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta})] \\ & + \mathbb{H}[q(\mathbf{z}_n)] \end{aligned} \quad \star = \sum_{n=1}^N \underbrace{\mathbb{E}_{q_n(\mathbf{z}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q_n(\mathbf{z}_n)]}_{\text{"local ELBO"}} \quad (11)$$

$$\underline{\beta\text{-VAE (Higgins, 2017?)}} \quad \triangleq \sum_{n=1}^N \mathcal{L}_n(q_n, \boldsymbol{\theta}) \quad \leftarrow \quad (12)$$

$$= \mathcal{L}(q, \boldsymbol{\theta}) \quad (13)$$

where $q = \{q_n\}_{n=1}^N$. Here, I've written the ELBO as a sum of "local ELBOs" $\mathcal{L}_n(q_n, \boldsymbol{\theta})$.

Optimal variational posterior

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

$$q_n(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{x}_n, \theta). \quad (14)$$

Question: The deep generative model above 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?

$$\begin{aligned} p(\mathbf{z}_n | \mathbf{x}_n; \theta) &\propto p(\mathbf{z}_n) p(\mathbf{x}_n | \mathbf{z}_n; \theta) \\ &= \mathcal{N}(\mathbf{z}_n; \mathbf{0}, \mathbf{I}) \mathcal{N}(\mathbf{x}_n | g(\mathbf{z}_n, \theta), \mathbf{I}) \\ &= \exp \left\{ -\frac{1}{2} \mathbf{z}_n^\top \mathbf{z}_n - \frac{1}{2} (\mathbf{x}_n - g(\mathbf{z}_n, \theta))^\top (\mathbf{x}_n - g(\mathbf{z}_n, \theta)) \right\} \\ &= \exp \left\{ -\frac{1}{2} \mathbf{z}_n^\top \mathbf{z}_n - \frac{1}{2} g(\mathbf{z}_n, \theta)^\top g(\mathbf{z}_n, \theta) + \dots \right\} \end{aligned}$$

Fixed form variational inference

Nevertheless, we can still constrain q_n to be Gaussian and seek the best Gaussian approximation to the posterior. This is sometimes called **fixed-form variational inference**.

For example, let,

$$\mathcal{Q} = \{q : q(\mathbf{z}) = \mathcal{N}(\mathbf{z} \mid \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) \text{ for some } \boldsymbol{\mu} \in \mathbb{R}^H, \boldsymbol{\sigma}^2 \in \mathbb{R}_+^H\} \quad (15)$$

Then, for fixed parameters $\boldsymbol{\theta}$, the best q_n in this **variational family** is,

$$q_n^* = \arg \max_{q_n \in \mathcal{Q}} \mathcal{L}_n(q_n, \boldsymbol{\theta}) \quad (16)$$

$$= \arg \min_{q_n \in \mathcal{Q}} D_{\text{KL}}(q_n(\mathbf{z}_n) \parallel p(\mathbf{z}_n \mid \mathbf{x}_n; \boldsymbol{\theta})). \quad (17)$$

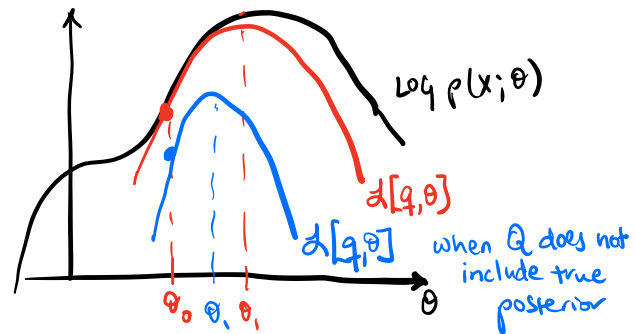
Variational expectation-maximization (vEM)

Now we can introduce a new algorithm: **variational expectation maximization**.

Repeat until either the ELBO or the parameters converges:

1. **M-step:** Set $\theta \leftarrow \arg \max_{\theta} \mathcal{L}(q, \theta)$
2. **E-step:** For $n = 1, \dots, N$
 - Set $q_n \leftarrow \arg \max_{q_n \in \mathcal{Q}} \mathcal{L}_n(q_n, \theta)$
3. Compute the ELBO $\mathcal{L}(q, \theta) \leq \log p(x; \theta)$ for any q

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



Generic M-step with gradient ascent

For exponential family mixture models and simple factor analysis, the M-steps had closed form.

For deep generative models, however, we need a more general approach.

If the parameters are unconstrained and the ELBO is differentiable wrt θ , we can use **gradient ascent**.

Repeat:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \mathcal{L}(q, \theta) = \theta + \alpha \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n)} \left[\underbrace{\nabla_{\theta} \left[-\frac{1}{2} (\mathbf{x}_n - g(\mathbf{z}_n; \theta))^T (\mathbf{x}_n - g(\mathbf{z}_n; \theta)) \right]}_{\nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n; \theta)} \right] + \dots \quad (18)$$

with **step size** α . Typically, you decrease the step size over iterations so that $\alpha_1 \geq \alpha_2 \geq \dots$

Generic M-step with *stochastic* gradient ascent (SGD)

- For more complex models, the expected gradient may be intractable as well. In those cases, we can use Monte Carlo to obtain an **unbiased estimate** of the gradient,

$$\nabla_{\theta} \mathcal{L}(q, \theta) \approx \widehat{\nabla}_{\theta} \mathcal{L}(q, \theta) = \sum_{n=1}^N \left[\frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n^{(m)}; \theta) \right] \quad (19)$$

where $\mathbf{z}_n^{(m)} \stackrel{\text{iid}}{\sim} q(\mathbf{z}_n)$ for $m = 1, \dots, M$.

- While we're at it, we can use Monte Carlo to approximate the sum over data points as well,

$$\widehat{\nabla}_{\theta} \mathcal{L}(q, \theta) = N \mathbb{E}_{n \sim \text{Unif}(1, \dots, N)} \left[\frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n^{(m)}; \theta) \right] \quad (20)$$

$$\approx \frac{N}{M} \sum_{m=1}^M \nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n^{(m)}; \theta) \quad (21)$$

where $n \sim \text{Unif}(1, \dots, N)$ and $\mathbf{z}_n^{(m)} \stackrel{\text{iid}}{\sim} q(\mathbf{z}_n)$ for $m = 1, \dots, M$.

Generic M-step with *stochastic* gradient ascent (SGD)

- For more complex models, the expected gradient may be intractable as well. In those cases, we can use Monte Carlo to obtain an **unbiased estimate** of the gradient,

$$\nabla_{\theta} \mathcal{L}(q, \theta) \approx \widehat{\nabla}_{\theta} \mathcal{L}(q, \theta) = \sum_{n=1}^N \left[\frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n^{(m)}; \theta) \right] \quad (19)$$

where $\mathbf{z}_n^{(m)} \stackrel{\text{iid}}{\sim} q(\mathbf{z}_n)$ for $m = 1, \dots, M$.

- While we're at it, we can use Monte Carlo to approximate the sum over data points as well,

$$\widehat{\nabla}_{\theta} \mathcal{L}(q, \theta) = N \mathbb{E}_{n \sim \text{Unif}(1, \dots, N)} \left[\frac{1}{M} \sum_{m=1}^M \nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n^{(m)}; \theta) \right] \quad (20)$$

$$\approx \frac{N}{M} \sum_{m=1}^M \nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n^{(m)}; \theta) \quad \text{[automatic diff.]} \quad (21)$$

where $n \sim \text{Unif}(1, \dots, N)$ and $\mathbf{z}_n^{(m)} \stackrel{\text{iid}}{\sim} q(\mathbf{z}_n)$ for $m = 1, \dots, M$.

SGD convergence and extensions

When does SGD work? This is a well studied problem in stochastic optimization [Bottou et al., 1998, Robbins and Siegmund, 1971].

Under relatively mild conditions, SGD converges to a **local minimum** if the step sizes obey the **Robbins-Monro conditions**,

$$\sum_{i=0}^{\infty} \alpha_i = \infty \quad \text{and} \quad \sum_{i=0}^{\infty} \alpha_i^2 < \infty \quad (22)$$

There have been dozens of extensions to basic SGD including,

- ▶ SGD with momentum
- ▶ AdaGrad [Duchi et al., 2011]
- ▶ RMSProp
- ▶ Adam [Kingma and Ba, 2014]

Note: we still need to compute the gradient $\nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n^{(m)}; \theta)$. We'll come back to this

Variational expectation-maximization (vEM)

Now we can introduce a new algorithm: **variational expectation maximization**.

Repeat until either the ELBO or the parameters converge:

1. **M-step:** Set $\theta \leftarrow \arg \max_{\theta} \mathcal{L}(q, \theta)$ SGD, unbiased $\hat{\nabla}_{\theta} \mathcal{L}_n(q_n, \theta)$
2. **E-step:** For $n = 1, \dots, N$
 - Set $q_n \leftarrow \arg \max_{q_n \in \mathcal{Q}} \mathcal{L}_n(q_n, \theta)$
3. Compute the ELBO $\mathcal{L}(q, \theta)$.

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

The variational E-step

As above, assume \mathcal{Q} is the family of Gaussian distributions with diagonal covariance.

Then $q_n(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n \mid \boldsymbol{\mu}_n, \text{diag}(\boldsymbol{\sigma}_n^2))$, with **variational parameters** $\boldsymbol{\mu}_n \in \mathbb{R}^H$ and $\boldsymbol{\sigma}_n^2 \in \mathbb{R}_+^H$.

We know the optimal q_n is,

$$q_n(\mathbf{z}_n) = \arg \max_{q \in \mathcal{Q}} \mathcal{L}_n(q_n, \boldsymbol{\theta}) \quad (23)$$

$$= \arg \max_{q \in \mathcal{Q}} \mathbb{E}_{q_n(\mathbf{z}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q_n(\mathbf{z}_n)] \quad (24)$$

but unfortunately there is no closed form solution.

Stochastic gradient ascent on the local ELBO

- **Idea:** Write the objective in terms of the unconstrained variational parameters

$$\lambda_n \triangleq (\mu_n, \log \sigma_n^2) \in \mathbb{R}^{2H} \quad (25)$$

and then perform stochastic gradient ascent.

- **Note:** Now we will write $q_n(\mathbf{z}_n; \lambda_n)$ and $\mathcal{L}_n(\lambda_n, \theta_n)$ to emphasize that the variational posterior and hence the local ELBO are both parameterized by λ_n .
- To perform SGD, we need an unbiased estimate of the gradient of the local ELBO,

$$\nabla_{\lambda_n} \mathcal{L}_n(\lambda_n, \theta) = \nabla_{\lambda_n} \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \theta) - \log q(\mathbf{z}_n)] \quad (26)$$

$$\neq \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} [\nabla_{\lambda_n} (\log p(\mathbf{x}_n, \mathbf{z}_n; \theta) - \log q(\mathbf{z}_n))]. \quad (27)$$

- **Question:** Why can't we simply bring the gradient inside the expectation like we did for the M-step?

Stochastic gradient ascent on the local ELBO

- **Idea:** Write the objective in terms of the unconstrained variational parameters

$$\lambda_n \triangleq (\mu_n, \log \sigma_n^2) \in \mathbb{R}^{2H} \quad (25)$$

and then perform stochastic gradient ascent.

- **Note:** Now we will write $q_n(\mathbf{z}_n; \lambda_n)$ and $\mathcal{L}_n(\lambda_n, \theta)$ to emphasize that the variational posterior and hence the local ELBO are both parameterized by λ_n .
- To perform SGD, we need an unbiased estimate of the gradient of the local ELBO,

$$\nabla_{\lambda_n} \mathcal{L}_n(\lambda_n, \theta) = \nabla_{\lambda_n} \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \theta) - \log q(\mathbf{z}_n)] \quad (26)$$

$$\neq \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} [\nabla_{\lambda_n} (\log p(\mathbf{x}_n, \mathbf{z}_n; \theta) - \log q(\mathbf{z}_n))]. \quad (27)$$

- **Question:** Why can't we simply bring the gradient inside the expectation like we did for the M-step?

Stochastic gradient ascent on the local ELBO

- **Idea:** Write the objective in terms of the unconstrained variational parameters

$$\boldsymbol{\lambda}_n \triangleq (\boldsymbol{\mu}_n, \log \boldsymbol{\sigma}_n^2) \in \mathbb{R}^{2H} \quad (25)$$

and then perform stochastic gradient ascent.

- **Note:** Now we will write $q_n(\mathbf{z}_n; \boldsymbol{\lambda}_n)$ and $\mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}_n)$ to emphasize that the variational posterior and hence the local ELBO are both parameterized by $\boldsymbol{\lambda}_n$.
- To perform SGD, we need an unbiased estimate of the gradient of the local ELBO,

$$\nabla_{\boldsymbol{\lambda}_n} \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\lambda}_n} \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n)] \quad (26)$$

$$\neq \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\nabla_{\boldsymbol{\lambda}_n} (\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n))]. \quad (27)$$

- **Question:** Why can't we simply bring the gradient inside the expectation like we did for the M-step?

Stochastic gradient ascent on the local ELBO

- **Idea:** Write the objective in terms of the unconstrained variational parameters

$$\boldsymbol{\lambda}_n \triangleq (\boldsymbol{\mu}_n, \log \boldsymbol{\sigma}_n^2) \in \mathbb{R}^{2H} \quad (25)$$

and then perform stochastic gradient ascent.

- **Note:** Now we will write $q_n(\mathbf{z}_n; \boldsymbol{\lambda}_n)$ and $\mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}_n)$ to emphasize that the variational posterior and hence the local ELBO are both parameterized by $\boldsymbol{\lambda}_n$.
- To perform SGD, we need an unbiased estimate of the gradient of the local ELBO,

$$\nabla_{\boldsymbol{\lambda}_n} \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\lambda}_n} \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n)] \quad (26)$$

$$\neq \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\nabla_{\boldsymbol{\lambda}_n} (\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n))]. \quad (27)$$

- **Question:** Why can't we simply bring the gradient inside the expectation like we did for the M-step?

The “score function” gradient estimator I

The basic problem is that the variational parameters λ_n determine the distribution we are taking an expectation under. However, there are a few ways to obtain unbiased estimates of the gradient.

One approach is called the **score function gradient estimator** or the **REINFORCE estimator** [Williams, 1992]. It is based on the following identity,

$$\nabla_{\lambda} \log q(\mathbf{z}; \lambda) = \frac{\nabla_{\lambda} q(\mathbf{z}; \lambda)}{q(\mathbf{z}; \lambda)} \quad (28)$$

where the l.h.s. is called the score function of distribution q .

The “score function” gradient estimator II

We can use this identity to obtain an unbiased estimate of the gradient of an expectation,

$$\nabla_{\lambda} \mathbb{E}_{q(\mathbf{z}; \lambda)} [h(\mathbf{z})] = \nabla_{\lambda} \int q(\mathbf{z}; \lambda) h(\mathbf{z}) d\mathbf{z} \quad (29)$$

$$= \int (\nabla_{\lambda} q(\mathbf{z}; \lambda)) h(\mathbf{z}) d\mathbf{z} \quad (30)$$

$$= \int (q(\mathbf{z}; \lambda) \nabla_{\lambda} \log q(\mathbf{z}; \lambda)) h(\mathbf{z}) d\mathbf{z} \quad (31)$$

$$= \mathbb{E}_{q(\mathbf{z}; \lambda)} [(\nabla_{\lambda} \log q(\mathbf{z}; \lambda)) h(\mathbf{z})] \quad (32)$$

From this identity, we can obtain an unbiased Monte Carlo estimate,

$$\widehat{\nabla}_{\lambda} \mathbb{E}_{q(\mathbf{z}; \lambda)} [h(\mathbf{z})] = \frac{1}{M} \sum_{m=1}^M [\nabla_{\lambda} \log q(\mathbf{z}^{(m)}; \lambda) h(\mathbf{z}^{(m)})]; \quad \mathbf{z}^{(m)} \stackrel{\text{iid}}{\sim} q(\mathbf{z}; \lambda) \quad (33)$$

The “score function” gradient estimator III

Notes:

1. The exchange of the gradient and the integral is allowed as long as the dominated convergence theorem holds, and it usually does for ML applications.
2. The score function gradient estimator is broadly applicable; e.g. it works for discrete and continuous latent variables \mathbf{z} . We just need the log density to be continuously differentiable wrt $\boldsymbol{\lambda}$ and to be able to sample from q .
3. If h is a function of both \mathbf{z} and $\boldsymbol{\lambda}$, you need to apply the product rule. This gives another term,

$$\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} [h(\mathbf{z}, \boldsymbol{\lambda})] = \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} [(\nabla_{\boldsymbol{\lambda}} \log q(\mathbf{z}; \boldsymbol{\lambda})) h(\mathbf{z}, \boldsymbol{\lambda})] + \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} [\nabla_{\boldsymbol{\lambda}} h(\mathbf{z}, \boldsymbol{\lambda})] \quad (34)$$

Control variates

Though broadly applicable, the score function estimator is often too high variance to be useful. This problem can often be mitigated with **control variates**. [c.f. Owen textbook]

Recall that the expectation of the score is zero,

$$\mathbb{E}_{q(\mathbf{z}; \lambda)} [\nabla_{\lambda} \log q(\mathbf{z}; \lambda)] = \int q(\mathbf{z}; \lambda) \nabla_{\lambda} \log q(\mathbf{z}; \lambda) d\mathbf{z} \quad (35)$$

$$= \int \nabla_{\lambda} q(\mathbf{z}; \lambda) d\mathbf{z} \quad (36)$$

$$= \nabla_{\lambda} \int q(\mathbf{z}; \lambda) d\mathbf{z} \quad (37)$$

$$= \nabla_{\lambda} 1 = 0. \quad (38)$$

Thus, we can subtract off any **baseline** from the function of interest without changing the expectation,

$$\mathbb{E}_{q(\mathbf{z}; \lambda)} [h(\mathbf{z}) \nabla_{\lambda} \log q(\mathbf{z}; \lambda)] = \mathbb{E}_{q(\mathbf{z}; \lambda)} [(h(\mathbf{z}) - b) \nabla_{\lambda} \log q(\mathbf{z}; \lambda)]. \quad (39)$$

The pathwise gradient estimator

Suppose $q(\mathbf{z}; \boldsymbol{\lambda}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$, as in our variational posteriors (again $\boldsymbol{\lambda} = (\boldsymbol{\mu}, \log \boldsymbol{\sigma}^2)$). Then,

$$\mathbf{z} \sim q(\mathbf{z}; \boldsymbol{\lambda}) \iff \mathbf{z} = r(\boldsymbol{\lambda}, \boldsymbol{\epsilon}) \quad (40)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I) \quad (41)$$

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

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

$$\mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} [h(\mathbf{z}, \boldsymbol{\lambda})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I)} [h(r(\boldsymbol{\lambda}, \boldsymbol{\epsilon}), \boldsymbol{\lambda})] \quad (42)$$

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

$$\nabla_{\boldsymbol{\lambda}} \mathbb{E}_{q(\mathbf{z}; \boldsymbol{\lambda})} [h(\mathbf{z}, \boldsymbol{\lambda})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I)} [\nabla_{\boldsymbol{\lambda}} h(\overset{r(\boldsymbol{\lambda}, \boldsymbol{\epsilon})}{\cancel{\boldsymbol{\mu} + \boldsymbol{\sigma} \boldsymbol{\epsilon}}}, \boldsymbol{\lambda})] \quad (43)$$

and **use Monte Carlo to obtain an unbiased estimate of the final expectation.**

Empirically comparing estimator variances

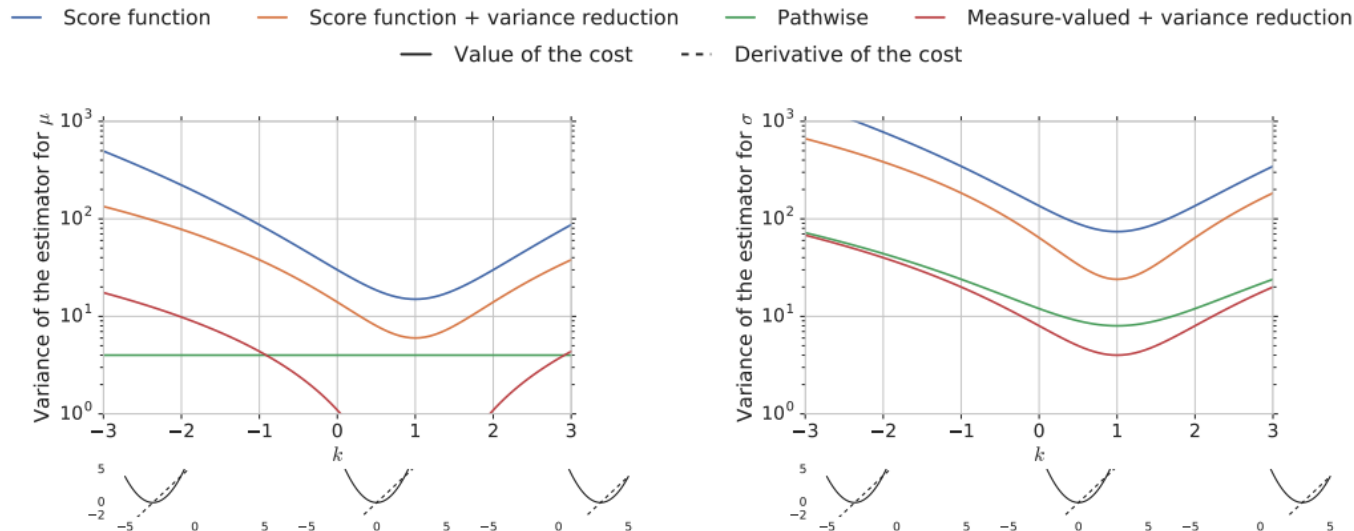


Figure 2: Variance of the stochastic estimates of $\nabla_{\theta} \mathbb{E}_{\mathcal{N}(x|\mu, \sigma^2)} [(x - k)^2]$ for $\mu = \sigma = 1$ as a function of k for three different classes of gradient estimators. Left: $\theta = \mu$; right: $\theta = \sigma$. The graphs in the bottom row show the function (solid) and its gradient (dashed) for $k \in \{-3, 0, 3\}$.

Variational expectation-maximization (vEM)

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

Repeat until either the ELBO or the parameters converges:

1. **M-step:** Set $\theta \leftarrow \arg \max_{\theta} \mathcal{L}(q, \theta)$ [with stochastic gradient ascent on the ELBO]
2. **E-step:** For $n = 1, \dots, N$
 - ▶ Set $q_n \leftarrow \arg \max_{q_n \in \mathcal{Q}} \mathcal{L}_n(q_n, \theta)$
 - ▶ Set $\lambda_n \leftarrow \arg \max_{\lambda_n} \mathcal{L}_n(\lambda_n, \theta)$
[with stochastic gradient ascent on the local ELBO using either the score function estimator or the pathwise gradient estimator]
3. Compute the ELBO $\mathcal{L}(q, \theta)$. [with Monte Carlo]

Amortized inference with recognition networks

- Note that vEM involves a costly E-step to find the variational parameters λ_n for each data point. This could involve many steps of stochastic gradient descent inside just the E-step!
- With a finite computational budget, we might be better off doing more gradient steps on θ and fewer on the local variational parameters.
- Note that the optimal variational parameters are just a function of the data point and the model parameters,

$$\lambda_n^* = \arg \min_{\lambda_n} D_{\text{KL}}(q(\mathbf{z}_n; \lambda_n) \parallel p(\mathbf{z}_n \mid \mathbf{x}_n, \theta)) \triangleq f^*(\mathbf{x}_n, \theta). \quad (44)$$

for some implicit and generally nonlinear function f^* .

$$\hat{\lambda}_n = \hat{f}(\mathbf{x}_n; \phi)$$

Amortized inference with recognition networks II

- ▶ VAEs learn an approximation to $f^*(\mathbf{x}_n, \theta)$ with an **inference network**, a.k.a. **recognition network** or **encoder**.
- ▶ The inference network is (yet another) neural network that takes in a data point \mathbf{x}_n and outputs variational parameters \mathbf{z}_n ,

$$\lambda_n \approx f(\mathbf{x}_n, \phi), \quad (45)$$

where ϕ are the weights of the network.

- ▶ The advantage is that the inference network is very fast; in the E-step, we simply need to pass a data point through the network to obtain the variational parameters.
- ▶ 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 buys us more updates of θ .

Amortization and approximation gaps

Cremer et al. [2018] consider the relative effects of the **amortization gap** and the **approximation gap** on variational EM.

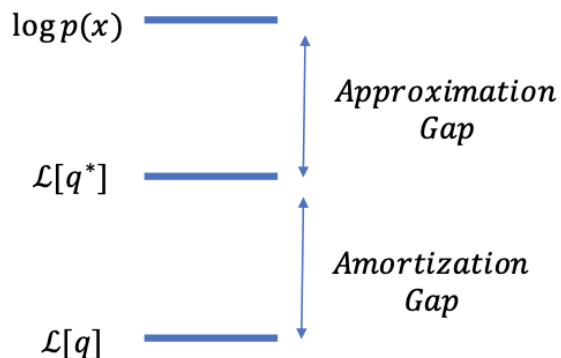


Figure 1. Gaps in Inference

Linear VAEs

Question: What does the optimal encoder network look like for a VAE with a linear generative model,

$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (46)$$

$$\mathbf{x}_n \sim \mathcal{N}(\mathbf{W}\mathbf{z}_n + \mathbf{b}, \mathbf{I}) \quad (47)$$

$$p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta}) = \mathcal{N}(\mathcal{J}^{-1}\mathbf{h}, \mathcal{J}^{-1}) \quad p(\mathbf{z}_n, \mathbf{x}_n) = \exp\left\{-\frac{1}{2}\mathbf{z}_n^T\mathbf{z}_n - \frac{1}{2}(\mathbf{x}_n - \mathbf{W}\mathbf{z}_n - \mathbf{b})^T(\mathbf{x}_n - \mathbf{W}\mathbf{z}_n - \mathbf{b})\right\}$$

$$\mathcal{J} = (\mathbf{W}^T\mathbf{W} + \mathbf{I})$$

$$\mathbf{h} = \mathbf{W}^T(\mathbf{x}_n - \mathbf{b})$$

$$= \exp\left\{-\frac{1}{2}\mathbf{z}_n^T \underbrace{(\mathbf{I} + \mathbf{W}^T\mathbf{W})}_{\mathcal{J}} \mathbf{z}_n + \underbrace{(\mathbf{x}_n - \mathbf{b})^T \mathbf{W}}_{\mathbf{h}} \mathbf{z}_n\right\}$$

$$\mu = \mathcal{J}^{-1}\mathbf{h} = \underline{(\mathbf{W}^T\mathbf{W} + \mathbf{I})^{-1}\mathbf{W}^T(\mathbf{x}_n - \mathbf{b})} = \mathbf{A}\mathbf{x}_n + \bar{\mathbf{b}} \quad ; \quad \mathbf{A} = (\mathbf{W}^T\mathbf{W} + \mathbf{I})^{-1}\mathbf{W}^T$$

$$\Sigma_z = \mathcal{J}^{-1} = (\mathbf{W}^T\mathbf{W} + \mathbf{I})^{-1}$$

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. Repeat until either the ELBO or the parameters converges:

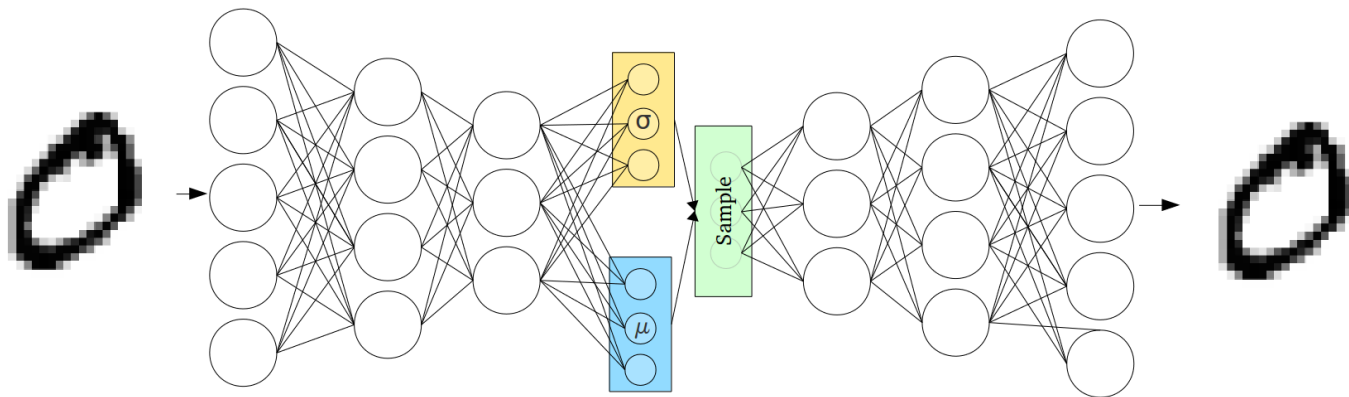
1. Sample data point $n \sim \text{Unif}(1, \dots, N)$. [Or a minibatch of data points.]
2. Estimate the local ELBO $\mathcal{L}_n(\phi, \theta)$ with Monte Carlo. [Note: it is a function of ϕ instead of λ_n .]
3. Compute unbiased Monte Carlo estimates of the gradients $\widehat{\nabla}_{\theta} \mathcal{L}_n(\phi,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

$$\theta \leftarrow \theta + \alpha_i \widehat{\nabla}_{\theta} \mathcal{L}_n(\phi, \theta) \quad (48)$$

$$\phi \leftarrow \phi + \alpha_i \widehat{\nabla}_{\phi} \mathcal{L}_n(\phi, \theta) \quad (49)$$

with step size α_i decreasing over iterations i according to a valid schedule.

VAEs from an autoencoder perspective



From <https://towardsdatascience.com/intuitively-understanding-variational-autoencoders-1bfe67eb5daf>

References I

Yann LeCun. *Modeles connexionnistes de l'apprentissage*. PhD thesis, These de Doctorat, Universite Paris, 1987.

Hervé Bourlard and Yves Kamp. Auto-association by multilayer perceptrons and singular value decomposition. *Biological cybernetics*, 59(4):291–294, 1988.

Geoffrey E Hinton and Richard Zemel. Autoencoders, minimum description length and helmholtz free energy. *Advances in neural information processing systems*, 6, 1993.

G E Hinton and R R Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, July 2006.

Pascal Vincent, Hugo Larochelle, Isabelle Lajoie, Yoshua Bengio, Pierre-Antoine Manzagol, and Léon Bottou. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of machine learning research*, 11(12), 2010.

References II

Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016.

<http://www.deeplearningbook.org>.

D P Kingma and M Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2014.

Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. January 2014.

Léon Bottou et al. Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9):142, 1998.

Herbert Robbins and David Siegmund. A convergence theorem for non negative almost supermartingales and some applications. In *Optimizing methods in statistics*, pages 233–257. Elsevier, 1971.

John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.

References III

- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte Carlo gradient estimation in machine learning. *Journal of Machine Learning Research*, 21(132):1–62, 2020.
- Chris Cremer, Xuechen Li, and David Duvenaud. Inference suboptimality in variational autoencoders. In *International Conference on Machine Learning*, pages 1078–1086. PMLR, 2018.