

Variational Autoencoders (VAEs)

STATS 305C: Applied Statistics

Scott Linderman

May 11, 2023

Where are we?

Model	Algorithm	Application
Multivariate Normal Models	Conjugate Inference	Bayesian Linear Regression
Hierarchical Models	MCMC (MH & Gibbs)	Modeling Polling Data
Probabilistic PCA & Factor Analysis	MCMC (HMC)	Images Reconstruction
Mixture Models	EM & Variational Inference	Image Segmentation
Mixed Membership Models	Coordinate Ascent VI	Topic Modeling
Variational Autoencoders	Gradient-based VI	Image Generation
State Space Models	Message Passing	Segmenting Video Data
Bayesian Nonparametrics	Fancy MCMC	Modeling Neural Spike Trains

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.

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].

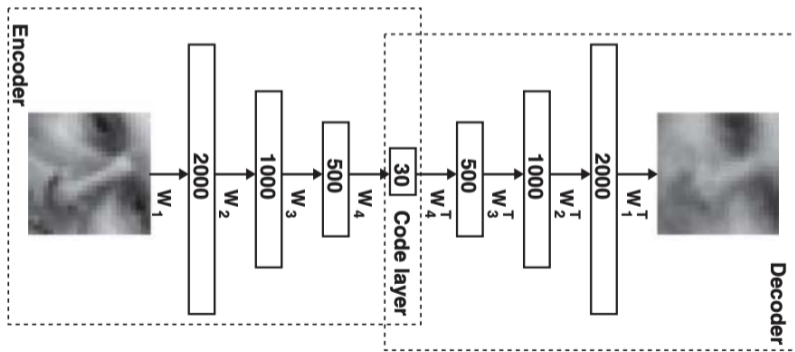
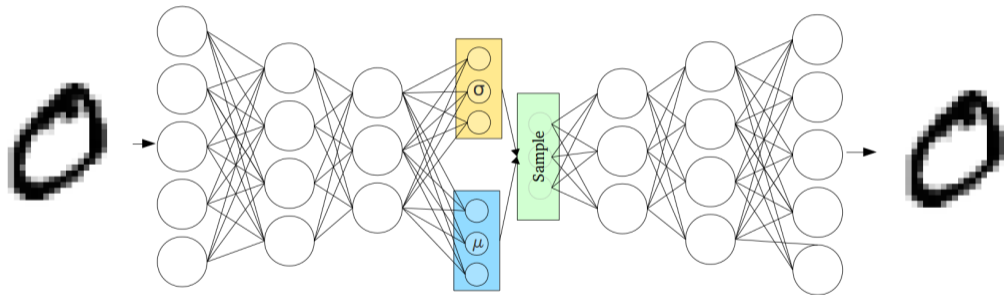


Figure: Figure from Hinton and Salakhutdinov [2006]

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**.



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

Outline

- ▶ The generative model
- ▶ Variational expectation maximization
- ▶ Amortized inference

The generative model

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

$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{0}, I) \tag{2}$$

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

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,

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

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

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

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.

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

$$\geq \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) - D_{\text{KL}}(q_n(\mathbf{z}_n; \boldsymbol{\lambda}_n) \parallel p(\mathbf{z}_n | \mathbf{x}_n; \boldsymbol{\theta})) \quad (10)$$

$$= \sum_{n=1}^N \underbrace{\mathbb{E}_{q_n(\mathbf{z}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q_n(\mathbf{z}_n; \boldsymbol{\lambda}_n)]}_{\text{"local ELBO"}} \quad (11)$$

$$\triangleq \sum_{n=1}^N \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) \quad (12)$$

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

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

Optimal variational posterior

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

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

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?

Review: Gradient-based VI

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, black-box**, or **automatic differentiation VI**.

For example, let,

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

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

$$q_n^* = \arg \min_{q_n \in \mathcal{Q}} D_{\text{KL}}(q_n(\mathbf{z}_n; \boldsymbol{\lambda}_n) \parallel p(\mathbf{z}_n \mid \mathbf{x}_n; \boldsymbol{\theta})) \quad (16)$$

$$= \arg \max_{\boldsymbol{\lambda}_n \in \mathbb{R}^{2H}} \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}). \quad (17)$$

We can maximize the ELBO with **stochastic gradient ascent** using unbiased estimates of the gradient, $\widehat{\nabla}_{\boldsymbol{\lambda}_n} \mathcal{L}(\boldsymbol{\lambda}_n, \boldsymbol{\theta})$, e.g., using the **score-function** or the **pathwise** gradient estimators.

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

Unfortunately, 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 and simple factor analysis, the M-steps had closed form. For deep generative models, we need a more general approach.
- ▶ If the parameters are unconstrained and the ELBO is differentiable wrt θ , we can use **stochastic gradient ascent**.

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \mathcal{L}(q, \theta) = \theta + \alpha \sum_{n=1}^N \mathbb{E}_{q(\mathbf{z}_n; \lambda_n)} [\nabla_{\theta} \log p(\mathbf{x}_n, \mathbf{z}_n; \theta)] \quad (18)$$

- ▶ Note that the expected gradient wrt θ can be computed using ordinary Monte Carlo – no fancy gradient estimators necessary!
- ▶ Likewise, we can use **mini-batches** of data to approximate the sum over data points.

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 $\theta \leftarrow \arg \max_{\theta} \mathcal{L}(q, \theta)$
2. **E-step:** For $n = 1, \dots, N$
 - ▶ Set $\lambda_n \leftarrow \arg \max_{\lambda_n \in \Lambda} \mathcal{L}_n(\lambda_n, \theta)$
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

- ▶ Assume \mathcal{Q} is the family of Gaussian distributions with diagonal covariance:
 $q_n(\mathbf{z}_n) = \mathcal{N}(\mathbf{z}_n \mid \boldsymbol{\mu}_n, \text{diag}(\boldsymbol{\sigma}_n^2))$, with **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

$$\nabla_{\boldsymbol{\lambda}_n} \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\lambda}_n} \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n; \boldsymbol{\lambda}_n)] \quad (19)$$

$$\neq \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\nabla_{\boldsymbol{\lambda}_n} (\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n; \boldsymbol{\lambda}_n))]. \quad (20)$$

- ▶ Last lecture we introduced the score-function and pathwise gradient estimators to tackle this problem. For example,

$$\nabla_{\boldsymbol{\lambda}_n} \mathcal{L}_n(\boldsymbol{\lambda}_n, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\lambda}_n} \mathbb{E}_{q(\mathbf{z}_n; \boldsymbol{\lambda}_n)} [\log p(\mathbf{x}_n, \mathbf{z}_n; \boldsymbol{\theta}) - \log q(\mathbf{z}_n; \boldsymbol{\lambda}_n)] \quad (21)$$

$$= \mathbb{E}_{\boldsymbol{\epsilon}_n \sim \mathcal{N}(\mathbf{0}, I)} [\nabla_{\boldsymbol{\lambda}_n} (\log p(\mathbf{x}_n, r(\boldsymbol{\lambda}_n, \boldsymbol{\epsilon}_n); \boldsymbol{\theta}) - \log q(r(\boldsymbol{\lambda}_n, \boldsymbol{\epsilon}_n); \boldsymbol{\lambda}_n))] \quad (22)$$

where $r(\boldsymbol{\lambda}_n, \boldsymbol{\epsilon}_n) = \boldsymbol{\mu}_n + \boldsymbol{\sigma}_n \boldsymbol{\epsilon}_n$.

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,

$$\lambda_n^* = \arg \min_{\lambda_n} D_{\text{KL}}(q(\mathbf{z}_n; \lambda_n) \parallel p(\mathbf{z}_n \mid \mathbf{x}_n, \theta)) \triangleq f^*(\mathbf{x}_n, \theta). \quad (23)$$

for some implicit and generally nonlinear function f^* .

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 ,

$$\lambda_n \approx f(\mathbf{x}_n, \phi), \quad (24)$$

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.

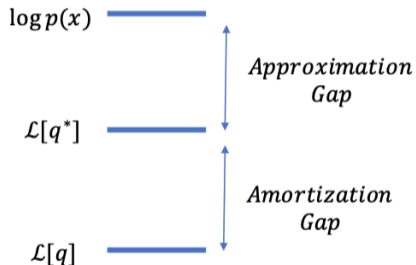


Figure 1. Gaps in Inference

Linear VAEs

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

$$\mathbf{z}_n \sim \mathcal{N}(\mathbf{0}, I) \tag{25}$$

$$\mathbf{x}_n \sim \mathcal{N}(\mathbf{W}\mathbf{z}_n + \mathbf{b}, I) \tag{26}$$

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:

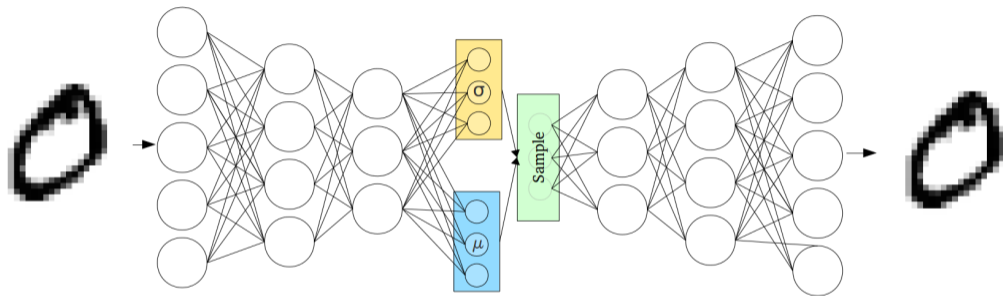
1. Sample data point $n \sim \text{Unif}(1, \dots, N)$. [Or a minibatch of data points.]
2. Estimate the local ELBO $\mathcal{L}_n(\phi, \theta)$ with Monte Carlo. [Note: it is a function of ϕ instead of λ_n .]
3. Compute unbiased Monte Carlo estimates of the gradients $\widehat{\nabla}_{\theta} \mathcal{L}_n(\phi,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

$$\theta \leftarrow \theta + \alpha_i \widehat{\nabla}_{\theta} \mathcal{L}_n(\phi, \theta) \quad (27)$$

$$\phi \leftarrow \phi + \alpha_i \widehat{\nabla}_{\phi} \mathcal{L}_n(\phi, \theta) \quad (28)$$

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.

Chris Cremer, Xuechen Li, and David Duvenaud. Inference suboptimality in variational autoencoders. In *International Conference on Machine Learning*, pages 1078–1086. PMLR, 2018.