

Variational Auto-Encoders

Mauro Camara Escudero

School of Mathematics, University of Bristol

Table of contents

1. Background
2. AEVB Objective Function
3. ELBO Optimization
4. Marginal Likelihood Estimation
5. Variational Auto-Encoders
6. Relationship between EM and VAE

Background

Latent Variable Models (LVMs)

Notation: \mathbf{x} observed, \mathbf{z} latent, θ parameter of interest.

Goals:

- Generative Modelling

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

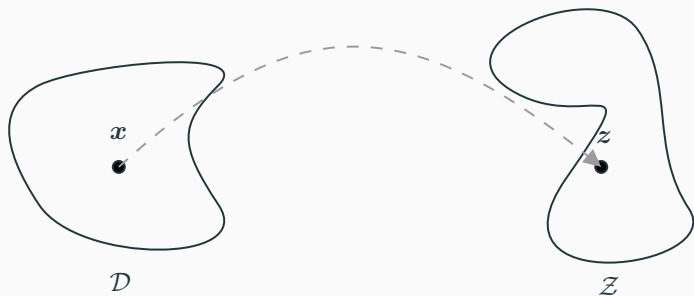
- Posterior Inference

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

- Parameter Estimation

$$\theta^* = \arg \max_{\theta} \prod_{\mathbf{x} \in \mathcal{D}} p_{\theta}(\mathbf{x})$$

Space Diagram



EM Algorithm for MLE

Initialize $\theta^{(0)}$ and $t = 0$.

- Compute conditional distribution of latent given observations

$$\{p_{\theta^{(t)}}(\mathbf{z} \mid \mathbf{x}) : \mathbf{x} \in \mathcal{D}\}$$

- Choose new parameter value $\theta^{(t+1)}$ so that it maximises

$$\sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{p_{\theta^{(t)}}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})]$$

Problem: Breaks if $p_{\theta^{(t)}}(\mathbf{z} \mid \mathbf{x})$ are **intractable**.

Mean-Field VI for Parameter Posterior

Define **factorized variational** distribution

$$\prod_{i=1}^{|\mathcal{Z}|} q_{\phi_i}(\mathbf{z}_i) \approx \prod_{\mathbf{z} \in \mathcal{Z}} p_{\theta}(\mathbf{z} \mid \mathbf{x})$$

Minimize KL-divergence for each data point / variational parameter

$$\min_{\phi} \text{KL}(q_{\phi_i}(\mathbf{z}) \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x}))$$

Problem: Breaks when $\mathbb{E}_{q_{\phi_i}}[\cdot]$ are **intractable** and doesn't **scale** to big data.

- **What it is used for:** Inference and Generative Modelling in LVMs.
- **How it works:** Optimization of an unbiased estimator of the ELBO (objective function) using SGD.
- **What's a VAE:** AEVB where probability distributions in LVM are parametrized by Neural Networks.

AEVB Objective Function

Objective Function: ELBO Derivation

Introduce **recognition model** $q_\phi(\mathbf{z} \mid \mathbf{x})$ with one parameter vector ϕ shared across data points.

$$\begin{aligned}\text{KL}(q_\phi(\mathbf{z} \mid \mathbf{x}) \parallel p_\theta(\mathbf{z} \mid \mathbf{x})) &= \mathbb{E}_{q_\phi} [\log q_\phi(\mathbf{z} \mid \mathbf{x}) - \log p_\theta(\mathbf{z} \mid \mathbf{x})] \\ &= -\mathbb{E}_{q_\phi} \left[\log \left(\frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z} \mid \mathbf{x})} \right) \right] + \log p_\theta(\mathbf{x})\end{aligned}$$

Objective Function: ELBO Derivation

Introduce **recognition model** $q_\phi(\mathbf{z} | \mathbf{x})$ with one parameter vector ϕ shared across data points.

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

First term on RHS is **Evidence Lower Bound** (ELBO) denoted $\mathcal{L}_{\theta, \phi}(\mathbf{x})$

$$\sum_{\mathbf{x} \in \mathcal{D}} \log p_\theta(\mathbf{x}) = \sum_{\mathbf{x} \in \mathcal{D}} (\text{KL}(q_\phi \parallel p_\theta(\mathbf{z} | \mathbf{x})) + \mathcal{L}_{\theta, \phi}(\mathbf{x})) \geq \sum_{\mathbf{x} \in \mathcal{D}} \mathcal{L}_{\theta, \phi}(\mathbf{x})$$

Since $\text{KL}(\cdot \parallel \cdot) \geq 0$.

Objective Function: Two Birds, One Stone

$$\max_{\theta, \phi} \mathcal{L}_{\theta, \phi}(\mathbf{x}) \implies \begin{cases} \max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}} \log p_{\theta}(\mathbf{x}) & \text{as } \log p_{\theta}(\mathbf{x}) \geq \mathcal{L}_{\theta, \phi}(\mathbf{x}) \\ \min_{\phi} \sum_{\mathbf{x} \in \mathcal{D}} \text{KL} & \text{as } \log p_{\theta}(\mathbf{x}) - \text{KL} = \mathcal{L}_{\theta, \phi}(\mathbf{x}) \end{cases}$$

Therefore maximizing ELBO leads to:

- $p_{\theta}(\mathbf{x})$ (i.e. **generative model**) improving.
- q_{ϕ} becoming a better **approximation**.

Objective Function: Alternative ELBO

Rewrite ELBO as **expected reconstruction error** regularized by penalizing approximate posteriors $q_\phi(\mathbf{z} | \mathbf{x})$ that are far from the prior $p_\theta(\mathbf{z})$.

$$\begin{aligned}\mathcal{L}_{\theta,\phi}(\mathbf{x}) &= \log p_\theta(\mathbf{x}) - \text{KL}(q_\phi || p_\theta(\mathbf{z} | \mathbf{x})) \\ &= \mathbb{E}_{q_\phi} [\log (p_\theta(\mathbf{x} | \mathbf{z})p_\theta(\mathbf{z})) - \log q_\phi(\mathbf{z} | \mathbf{x})] \\ &= \underbrace{\mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{x} | \mathbf{z})]}_{\text{Expected Log-Likelihood}} - \underbrace{\text{KL}(q_\phi(\mathbf{z} | \mathbf{x}) || p_\theta(\mathbf{z}))}_{\text{Regularization Term}}\end{aligned}$$

Notice: Optimize ELBO using stochastic gradient optimization requires $\nabla_{\theta,\phi}\mathcal{L}_{\theta,\phi}(\mathbf{x})$

ELBO Optimization

Intractability of Objective Function Gradients

ELBO gradient $\nabla_{\phi} \mathbb{E}_{q_{\phi}}[\cdot]$ difficult to approximate with Monte Carlo as we cannot exchange gradient and expectation $\nabla_{\phi} \mathbb{E}_{q_{\phi}}[\cdot] \neq \mathbb{E}_{q_{\phi}}[\nabla_{\phi}]$

$$\mathcal{L}_{\theta, \phi}(\mathbf{x}) = \begin{cases} \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z} | \mathbf{x})] \\ \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x} | \mathbf{z})] - \text{KL}(q_{\phi} || p_{\theta}(\mathbf{z})) \end{cases}$$

Can we write expectation with respect to a distribution independent of ϕ ?

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}}[\cdot] \stackrel{?}{=} \mathbb{E}_{p(\epsilon)}[\nabla_{\phi}]$$

Reparametrization Trick

We know how to write a general MVN in terms of a Standard MVN

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \implies \mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^\top)$$

Thus expectations can be written as

$$\mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^\top)} [f(\mathbf{z})] = \mathbb{E}_{\mathcal{N}(\mathbf{0}, \mathbf{I})} [f(\boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon})]$$

In general write $\mathbf{z} = g_\phi(\boldsymbol{\epsilon}, \mathbf{x})$ as a **deterministic** and **differentiable** function of \mathbf{x} and $\boldsymbol{\epsilon}$, where $p(\boldsymbol{\epsilon})$ is independent of ϕ .

$$\nabla_\phi \mathbb{E}_{q_\phi} [f(\mathbf{z})] = \mathbb{E}_{p(\boldsymbol{\epsilon})} [\nabla_\phi f(g_\phi(\boldsymbol{\epsilon}, \mathbf{x}))]$$

SGC on Objective Function Unbiased Estimators

Obtain two **unbiased** estimators for ELBO based on $\epsilon^{(i)} \stackrel{\text{i.i.d.}}{\sim} p(\epsilon)$.

$$\tilde{\mathcal{L}}_{\theta, \phi}(\mathbf{x}) = \begin{cases} \frac{1}{L} \sum_{i=1}^L \left[\log p_{\theta}(\mathbf{x}, g_{\phi}(\epsilon^{(i)}, \mathbf{x})) - \log q_{\phi}(g_{\phi}(\epsilon^{(i)}, \mathbf{x}) \mid \mathbf{x}) \right] \\ \frac{1}{L} \sum_{i=1}^L \left[\log p_{\theta}(\mathbf{x} \mid g_{\phi}(\epsilon^{(i)}, \mathbf{x})) \right] - \underbrace{\text{KL}(q_{\phi} \parallel p_{\theta}(\mathbf{z}))}_{\text{Often available in closed form}} \end{cases}$$

SGD randomly samples a **minibatch** of data $M \subseteq \mathcal{D}$ and uses mini-batch gradients

$$\frac{1}{|\mathcal{M}|} \sum_{\mathbf{x} \in \mathcal{M}} \nabla_{\theta, \phi} \tilde{\mathcal{L}}_{\theta, \phi}(\mathbf{x})$$

Marginal Likelihood Estimation

Estimating $\log p_{\theta}(\mathcal{D})$

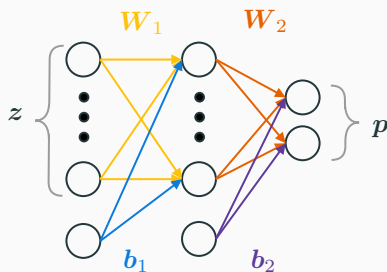
After training we can estimate the log marginal likelihood using **importance sampling**.

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} \\ &= \log \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\frac{p_{\theta}(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z} | \mathbf{x})} \right] \\ &\approx \log \frac{1}{L} \sum_{i=1}^L \frac{p_{\theta}(\mathbf{x}, \mathbf{z}^{(i)})}{q_{\phi}(\mathbf{z}^{(i)} | \mathbf{x})} \quad \mathbf{z}^{(i)} \stackrel{\text{i.i.d.}}{\sim} q_{\phi}(\mathbf{z} | \mathbf{x})\end{aligned}$$

Variational Auto-Encoders

Parametrizing Distributions via Neural Networks

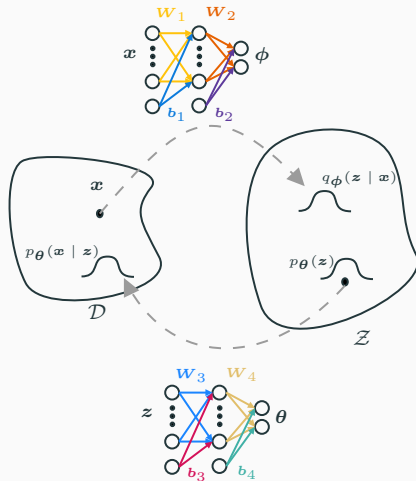
Suppose \mathbf{x} is binary vector of Bernoulli trials. Then $p_{\theta}(\mathbf{x} \mid \mathbf{z})$ is parametrized by a vector of probabilities \mathbf{p} which can be constructed via an MLP.



The log-likelihood then becomes

$$\log p_{\theta}(\mathbf{x} \mid \mathbf{z}) = \sum_j x_j \log p_j + (1 - x_j) \log(1 - p_j)$$

VAE = AEVB + NN



Relationship between EM and VAE

Variational EM

Recall the EM algorithm:

- Compute approximate posteriors $\{p_{\theta^{(t)}}(\mathbf{z} \mid \mathbf{x}) : \mathbf{x} \in \mathcal{D}\}$
- Find optimal parameter $\theta^{(t+1)} = \arg \max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}} \mathbb{E}_{p_{\theta^{(t)}}(\mathbf{z} \mid \mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})]$

Consider ELBO as a functional of q_{ϕ} and a function of θ .

$$\mathcal{L}_{\mathbf{x}}(\theta, q_{\phi}) = \begin{cases} \log p_{\theta}(\mathbf{x}) - \text{KL}(q_{\phi} \parallel p_{\theta}(\mathbf{z} \mid \mathbf{x})) & (1) \\ \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] - \mathbb{E}_{q_{\phi}} [\log q_{\phi}] & (2) \end{cases}$$

E-step: Maximize (1) wrt q_{ϕ} (KL is zero and the bound is **tight**).

$$\left\{ p_{\theta^{(t)}}(\mathbf{z} \mid \mathbf{x}) = \arg \max_{q_{\phi}} \mathcal{L}_{\mathbf{x}}(\theta^{(t)}, q_{\phi}) : \mathbf{x} \in \mathcal{D} \right\}$$

M-step: Maximize (2) wrt θ

$$\theta^{(t+1)} = \arg \max_{\theta} \sum_{\mathbf{x} \in \mathcal{D}} \mathcal{L}_{\mathbf{x}}(\theta, p_{\theta^{(t)}}(\mathbf{z} \mid \mathbf{x}))$$

In summary, EM algorithm and VAE optimize the **same objective**.

- When expectations are in closed-form, one should use EM, which uses **coordinate ascent**.
- When expectations are intractable, VAE uses **stochastic gradient ascent** on an unbiased estimator of the objective function.

Thank you

