

Latent Variable Models

Muchang Bahng

Spring 2025

Contents

1	Nonlinear Latent Variable Models	2
1.1	Variational Lower Bounds	2
1.2	EM Algorithm	8
1.3	Gaussian Mixture Models	12
1.4	Nonlinear ICA	15
	Bibliography	15

1 Nonlinear Latent Variable Models

Now we will consider ourselves with nonlinear latent variables models, which still defines a simple latent random variable Z with prior $p(z)$, but now a family of nonlinear functions $\{f_\theta(z)\}$ that defines the generative component $f_\theta(x | z)$. In factor models, we have taken linear transformations of random variables and therefore the likelihood had been easy to calculate, differentiate, and therefore optimize.

In the general nonlinear case, we usually deal with f_θ not as a transformation of Z to X , but really $f_\theta(z)$ becomes the *parameters* of $X | Z = z$. This allows to define the *implicitly parameterized* family of distributions $\{p_\theta\}$. Given that the true distribution of the data is $p^*(x)$, we would like to find a distribution $p_\theta(x)$ that is a good approximation.

$$p^*(x) \approx p_\theta(x) \quad (1)$$

To calculate the likelihood $p_\theta(x)$, we must compute the marginal

$$p_\theta(x) = \int p_\theta(x, z) dz = \int p_\theta(x | z) p(z) dz \quad (2)$$

which is known to be computationally intractable due to the integral. At first, it seems like all hope is lost, but statisticians have a few tricks up their sleeves.

1. The first trick is to notice that by Bayes rule, we can compute the likelihood not as an integral, but as

$$p_\theta(x) = \frac{p_\theta(x | z) p(z)}{p_\theta(z | x)} \quad (3)$$

So it suffices to find a good approximation of $p_\theta(z | x)$, which is a probabilistic discriminative model for the latent variable (i.e. we are trying to compute the distribution of z given x as if we were predicting it). We can do MCMC since $p_\theta(z | x) \propto p_\theta(x | z) p(z)$, but often this can be slow to fit.

2. The next trick is called the variational lower bound, which is a lower bound on the log likelihood, and therefore by optimizing it we can hope to optimize the log-likelihood as well. This works well in practice.
3. The next trick is by optimizing the Fisher score, which is the gradient of the log likelihood *with respect to the covariates* (not the parameters!).

1.1 Variational Lower Bounds

We focus on this problem and define a family of distributions $\{q_\phi(z | x)\}_\phi$ and use it to approximate $p_\theta(z | x)$. Therefore, searching for a good ϕ and therefore a good q_ϕ is basically the problem of **variational Bayesian inference**. Essentially we are trying to construct an encoder and a decoder.

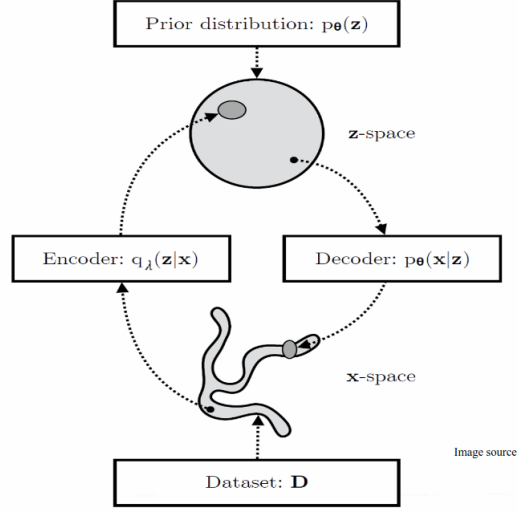


Figure 1: If $q_\phi = p_\theta$, then the diagram commutes, i.e. $p(z)p_\theta(x | z) = p(x)p_\theta(z | x) = p_\theta(x, z)$.

As we have stated before (and in pretty much all density estimation problems), our job is to maximize the log likelihood of the training set:

$$\sum_i \log p(x^{(i)}) \quad (4)$$

In order to do this for this problem, we need a little fact from information theory.

Theorem 1.1 (Log Likelihood vs Conditional Entropy)

The KL divergence can be decomposed to

$$KL(q_\phi(z | x) || p_\theta(z | x)) = \mathbb{E}_{q_\phi(z|x)}[\log q_\phi(z | x)] + \log p_\theta(x) - \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x, z)] \quad (5)$$

and hence

Proof.

Starting with the definition of KL divergence:

$$KL(q_\phi(z | x) || p_\theta(z | x)) = \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{q_\phi(z | x)}{p_\theta(z | x)} \right] \quad (6)$$

$$= \mathbb{E}_{q_\phi(z|x)}[\log q_\phi(z | x)] - \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(z | x)] \quad (7)$$

By Bayes' rule, we know that

$$p_\theta(z | x) = \frac{p_\theta(x, z)}{p_\theta(x)} \quad (8)$$

Substituting this into our equation gives

$$KL(q_\phi(z | x) || p_\theta(z | x)) = \mathbb{E}_{q_\phi(z|x)}[\log q_\phi(z | x)] - \mathbb{E}_{q_\phi(z|x)} \left[\log \frac{p_\theta(x, z)}{p_\theta(x)} \right] \quad (9)$$

$$= \mathbb{E}_{q_\phi(z|x)}[\log q_\phi(z | x)] - \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x, z)] + \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x)] \quad (10)$$

Since $\log p_\theta(x)$ is constant with respect to z , we can take it out of the expectation.

$$\mathbb{E}_{q_\phi(z|x)}[\log q_\phi(z | x)] - \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x, z)] + \log p_\theta(x) \quad (11)$$

Therefore maximizing the log-likelihood is equivalent to minimizing the KL-divergence.

$$\log p_\theta(x) = KL(q_\phi(z | x) || p_\theta(z | x)) + \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x, z)] - \mathbb{E}_{q_\phi(z|x)}[\log q_\phi(z | x)] \quad (12)$$

But again the KL divergence part is intractable due to $p_\theta(z | x)$ being intractable. Using the fact that the KL divergence is always greater than or equal to 0, we can drop the term and set a lower bound on the log likelihoods. This lower bound is called the *variational lower bound*.

$$\sum_{i=1}^N \log p_\theta(x^{(i)}) \geq \sum_{i=1}^N \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)}, z)] - \sum_{i=1}^N \mathbb{E}_{q_\phi(z|x^{(i)})}[\log q_\phi(z | x^{(i)})] \quad (13)$$

Definition 1.1 (Variational Lower Bound)

The **variational lower bound** of the dataset \mathcal{D} is defined

$$\text{ELBO}(\mathcal{D}, \phi, \theta) = \sum_{i=1}^N \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)}, z)] - \sum_{i=1}^N \mathbb{E}_{q_\phi(z|x^{(i)})}[\log q_\phi(z | x^{(i)})] \quad (14)$$

which can be decomposed into the sums of the variational lower bounds of the individual data points.

$$\text{ELBO}(\mathcal{D}, \phi, \theta) = \sum_i \text{ELBO}(x^{(i)}, \phi, \theta) \quad (15)$$

where

$$\text{ELBO}(x^{(i)}, \phi, \theta) = \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)}, z)] - \mathbb{E}_{q_\phi(z|x^{(i)})}[\log q_\phi(z | x^{(i)})] \quad (16)$$

Note that we can alternatively define ELBO using Jensen's inequality.

Definition 1.2 (Evidence Lower Bound)

To lower bound it, we can use Jensen's inequality^a with the concave function $f(x) = \log(x)$ over domain \mathbb{R}^+ and the following holds true for all θ and more importantly, for *any arbitrary density function* $q(z)$. Therefore, we have

$$\ell(\theta) = \log p_\theta(x) \quad (17)$$

$$= \log \int p_\theta(x, z) dz \quad (18)$$

$$= \log \int q_\phi(z) \frac{p_\theta(x, z)}{q_\phi(z)} dz \quad (19)$$

$$\geq \int q_\phi(z | x) \log \left(\frac{p_\theta(x, z)}{q_\phi(z)} \right) dz \quad (20)$$

$$= \text{ELBO}(x, q_\phi) \quad (21)$$

The lower bound is called the **evidence lower bound (ELBO)**, and the ELBO of the whole dataset is

$$\text{ELBO}(\mathcal{D}, \phi, \theta) = \sum_{i=1}^N \text{ELBO}(x^{(i)}, \phi, \theta) \quad (22)$$

^aGiven convex function $f : \mathbb{R} \rightarrow \mathbb{R}$, and random variable X , $\mathbb{E}[f(x)] \geq f(\mathbb{E}[X])$.

Note that this lower bound is with respect to *any* distribution q_ϕ , and it is because of this flexibility that we choose q_ϕ in the first place. Therefore, we can vary ϕ in hopes that the lower bound is maximized, and optimize with respect to this, hence the name *variational*. For more interpretability, look at the corollary.

Corollary 1.1 (Decomposition of ELBO)

The following decomposition of ELBO shows that maximizing the ELBO simultaneously attempts to keep q_ϕ close to p and concentrate $q_\phi(z | x)$ on those z that maximizes $\ln p_\theta(x | z)$. That is, the approximate posterior q_ϕ balances between staying close to the prior $p(z)$ and moving towards the maximum likelihood $\operatorname{argmax}_z \ln p_\theta(x | z)$.

$$\text{ELBO}(x^{(i)}, \phi, \theta) = \underbrace{\mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)} | z)]}_{\text{likelihood term (reconstruction part)}} - \underbrace{KL(q_\phi(z | x^{(i)}) || p(z))}_{\text{closeness of encoding to } p(z) \text{ (typically Gaussian)}} \quad (23)$$

Note the first expression is the likelihood term, which measures the reconstruction quality of the decoder $p_\theta(x^{(i)} | z)$ averaged over encodings sampled from $q_\phi(z | x^{(i)})$. The second term is the KL divergence between the encoder distribution $q_\phi(z | x^{(i)})$ and the prior $p(z)$, which acts as a regularizer by ensuring the encoded distributions remain close to the chosen prior, typically a standard normal distribution.

Proof.

Starting with the ELBO for a single data point:

$$\text{ELBO}(x^{(i)}, \phi, \theta) = \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)}, z)] - \mathbb{E}_{q_\phi(z|x^{(i)})}[\log q_\phi(z | x^{(i)})]$$

Using the chain rule of probability for the joint distribution:

$$p_\theta(x^{(i)}, z) = p_\theta(x^{(i)} | z)p(z)$$

Substituting this into our ELBO:

$$\begin{aligned} \text{ELBO}(x^{(i)}, \phi, \theta) &= \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)} | z) + \log p(z)] - \mathbb{E}_{q_\phi(z|x^{(i)})}[\log q_\phi(z | x^{(i)})] \\ &= \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)} | z)] + \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p(z)] - \mathbb{E}_{q_\phi(z|x^{(i)})}[\log q_\phi(z | x^{(i)})] \\ &= \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)} | z)] - \left(\mathbb{E}_{q_\phi(z|x^{(i)})}[\log q_\phi(z | x^{(i)})] - \mathbb{E}_{q_\phi(z|x^{(i)})}[\log p(z)] \right) \\ &= \underbrace{\mathbb{E}_{q_\phi(z|x^{(i)})}[\log p_\theta(x^{(i)} | z)]}_{\text{reconstruction term}} - \underbrace{KL(q_\phi(z | x^{(i)}) || p(z))}_{\text{KL divergence term}} \end{aligned}$$

Therefore, maximizing the ELBO will simultaneously allow us to obtain an accurate generative model $p_\theta(x | z) \approx p^*(x | z)$ and an accurate discriminative model $q_\phi(z | x) \approx p_\theta(z | x)$. The next step is to actually maximize the ELBO with respect to both θ and ϕ . To do this we need to compute the derivatives of ELBO w.r.t. to ϕ and θ .

$$\max_{\phi, \theta} \text{ELBO}(\mathcal{D}, \phi, \theta) \quad (24)$$

It turns out that this itself is a nonconvex optimization problem, and to make it doable we iterate between updating ϕ and θ . Remember that the ELBO is really an expectation, i.e. an integral, and to get a good estimate of its derivative we must try to change it from the derivative of an expectation to the expectation of a derivative. The gradient with respect to θ is very easy since from measure theory, we are deriving and integrating over different variables.

Lemma 1.1 (Gradient of ELBO w.r.t. θ)

For θ , its unbiased gradient is

$$\nabla_\theta \text{ELBO}(x, \theta, \phi) = \mathbb{E}_{q_\phi(z|x)}[\nabla_\theta \log p_\theta(x | z)] \quad (25)$$

and therefore we can approximate the gradient by sampling L points $p^{(1)}, \dots, p^{(L)}$ from $p(z)$ and computing the gradient of the log (since we know the closed form of the conditional distribution given z), and finally averaging them.

$$\nabla_{\theta} \text{ELBO}(x, \theta, \phi) \approx \frac{1}{L} \sum_{l=1}^L \nabla_{\theta} \log p_{\theta}(x | z^{(l)}) \quad (26)$$

which is guaranteed to converge by the law of large numbers, and furthermore, we can do this for any batch size L .

Proof.

Note that the KL divergence does not depend on θ and neither does the prior, so they can be removed

$$\nabla_{\theta} \text{ELBO}(x, \theta, \phi) = \nabla_{\theta} \{ \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z)] - \mathbb{E}_{q_{\phi}(z|x)} [\log q_{\phi}(z | x)] \} \quad (27)$$

$$= \nabla_{\theta} \{ \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z)] \} \quad (28)$$

$$= \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\theta} \{ \log p_{\theta}(x, z) \}] \quad (29)$$

$$= \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\theta} \{ \log p_{\theta}(x | z) - \log p(z) \}] \quad (30)$$

$$= \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\theta} \log p_{\theta}(x | z)] \quad (31)$$

However, taking the gradient w.r.t. ϕ is more complicated since we cannot put the gradient in the expectation, i.e. swap the derivative and integral (since we are deriving and integrating w.r.t. ϕ). Fortunately, we have a well-known mathematical identity often used in policy gradient algorithms in reinforcement learning. [Wil92]

Lemma 1.2 (Log-Derivative Trick)

The following identity holds.

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [f(z)] = \mathbb{E}_{q_{\phi}(z)} [f(z) \nabla_{\phi} \log q_{\phi}(z)] \quad (32)$$

Proof.

First, let's write out the left-hand side using the definition of expectation:

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [f(z)] = \nabla_{\phi} \int f(z) q_{\phi}(z) dz$$

Under suitable regularity conditions, we can exchange the gradient and integral operators:

$$= \int f(z) \nabla_{\phi} q_{\phi}(z) dz$$

Now, we multiply and divide by $q_{\phi}(z)$ inside the integral:

$$= \int f(z) q_{\phi}(z) \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)} dz$$

Recognize that $\nabla_{\phi} \log q_{\phi}(z) = \frac{\nabla_{\phi} q_{\phi}(z)}{q_{\phi}(z)}$ by the chain rule:

$$= \int f(z) q_{\phi}(z) \nabla_{\phi} \log q_{\phi}(z) dz$$

Finally, we can rewrite this back as an expectation:

$$= \mathbb{E}_{q_{\phi}(z)} [f(z) \nabla_{\phi} \log q_{\phi}(z)]$$

Example 1.1 (Gradient of Expection of $f(x) = x^2$ w.r.t. Gaussian)

Assume we have a normal distribution q that is parameterized by ϕ , specifically $q_\phi(x) = N(\phi, 1)$. We want to solve the below problem

$$\min_{\phi} \mathbb{E}_q[x^2] \quad (33)$$

This is of course a rather silly problem and the optimal $\phi = 0$ is obvious. One way to calculate $\nabla_{\phi} \mathbb{E}[x^2]$ is using the log-derivative trick as follows

$$\nabla_{\phi} \mathbb{E}_q[x^2] = \nabla_{\phi} \int q_{\phi}(x) x^2 dx \quad (34)$$

$$= \int x^2 \nabla_{\phi} q_{\phi}(x) \frac{q_{\phi}(x)}{q_{\phi}(x)} dx \quad (35)$$

$$= \int q_{\phi}(x) \nabla_{\phi} \log q_{\phi}(x) x^2 dx \quad (36)$$

$$= \mathbb{E}_q[x^2 \nabla_{\phi} \log q_{\phi}(x)] \quad (37)$$

For our example where $q_{\phi}(x) = N(\phi, 1)$, this method gives

$$\nabla_{\phi} \mathbb{E}[x^2] = \mathbb{E}_q[x^2(x - \phi)] \quad (38)$$

Using this on the gradient of ELBO w.r.t. ϕ gives the following form as the expectation of the gradient.

Lemma 1.3 ()

We can use the score function estimator.

$$\nabla_{\phi} \text{ELBO}(x, \theta, \phi) = \nabla_{\phi} \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)] \quad (39)$$

$$= \mathbb{E}_{q_{\phi}(z|x)} [\nabla_{\phi} \{ \log q_{\phi}(z|x) (\log p_{\theta}(x, z) - \log q_{\phi}(z|x)) \}] \quad (40)$$

Proof.

However, REINFORCE is known to have high variance, and so we need large batch sizes L for good convergence. Many methods such as [GBB01, PBJ12] were developed to reduce this. Later it was shown in [KW22] that the *reparameterization trick* beat everything else, allowing us to efficiently train neural-net-based non-linear latent variable models, e.g. the variational autoencoder. We will focus on the reparameterization trick in my deep learning notes and omit it here. Now that we have approximate closed form solutions for the gradients, we can optimize the two using coordinate ascent. Note that we have shown this for a single sample x , and ideally we would do this for a minibatch of samples $x^{(i)}$.

Algorithm 1.1 (Coordinate Ascent Variational Inference)

A common approach to maximize the ELBO is coordinate ascent, where we alternatively optimize with respect to ϕ and θ :

Algorithm 1 Coordinate Ascent Variational Inference (CAVI) with Reparameterization**Require:** Initial parameters $\theta^{[0]}$, $\phi^{[0]}$, batch size B, number of samples L

```

1: while not converged do
2:   // E-step: optimize variational parameters
3:   Sample minibatch  $\{x^{(1)}, \dots, x^{(B)}\}$  from dataset  $\mathcal{D}$ 
4:   Sample noise  $\{\epsilon^{(1)}, \dots, \epsilon^{(L)}\} \sim p(\epsilon)$  for reparameterization
5:   Transform noise to latent variables:  $z^{(l)} = g_{\phi^{[t]}}(\epsilon^{(l)}, x)$  for  $l = 1, \dots, L$ 
6:   // Approximate gradient using Monte Carlo samples
7:    $\hat{g}_\phi \leftarrow \frac{1}{BL} \sum_{i=1}^B \sum_{l=1}^L [\nabla_\phi \log p_{\theta^{[t]}}(x^{(i)} | z^{(l)}) - \nabla_\phi \log q_{\phi^{[t]}}(z^{(l)} | x^{(i)}) + \nabla_\phi \log p(z^{(l)})]$ 
8:    $\phi^{[t+1]} \leftarrow \phi^{[t]} + \eta_\phi \hat{g}_\phi$  ▷Update with learning rate  $\eta_\phi$ 
9:   // M-step: optimize model parameters
10:   $\hat{g}_\theta \leftarrow \frac{1}{BL} \sum_{i=1}^B \sum_{l=1}^L \nabla_\theta \log p_{\theta^{[t]}}(x^{(i)} | z^{(l)})$ 
11:   $\theta^{[t+1]} \leftarrow \theta^{[t]} + \eta_\theta \hat{g}_\theta$  ▷Update with learning rate  $\eta_\theta$ 
12: end while

```

Once we are done, we have our optimized encoder and decoders p_θ and q_ϕ .

1.2 EM Algorithm

Let's consider a slightly simpler sub-problem where we have covariates $x^{(i)} \sim X$ coming from distribution $p(x)$. We can again add latent random variables Z but rather than being fixed, the prior $p_\theta(z)$ is also parameterized by θ . Therefore, we would like to find

$$\operatorname{argmax}_{\theta} p_\theta(x) = \operatorname{argmax}_{\theta} \int p_\theta(x | z) p_\theta(z) dz \quad (41)$$

Even though this integral is not tractable, we will assume that $p_\theta(z | x)$ can be computed for a given θ . Let's try to redo our algorithm again with computable posterior assumptions. We have a training set $\mathcal{D} = \{x^{(i)}\}_{i=1}^n \in \mathbb{R}^d$, which we assume are generated by some latent distributions $p_\theta(z)$ followed by the generative component $p_\theta(x | z)$. Then, we bound the likelihood of each sample $x^{(i)}$ by an ELBO that varies for all distributions $q^{(i)}$ (we write q rather than q_ϕ since the ϕ will be irrelevant here).

$$\log p_\theta(x^{(i)}) \geq \text{ELBO}(x^{(i)}, q^{(i)}, \theta) = \mathbb{E}_{q^{(i)}(z|x^{(i)})}[\log p_\theta(x^{(i)}, z)] - \mathbb{E}_{q^{(i)}(z|x^{(i)})}[\log q^{(i)}(z | x^{(i)})] \quad (42)$$

Summing this all up gives the ELBO of our dataset, which is a lower bound for *all* collections of distributions $q^{(1)}, \dots, q^{(n)}$.

$$\sum_{i=1}^N \log p_\theta(x^{(i)}) \geq \text{ELBO}(\mathcal{D}, q^{(1)}, \dots, q^{(n)}, \theta) \quad (43)$$

$$= \sum_{i=1}^N \mathbb{E}_{q^{(i)}(z|x^{(i)})}[\log p_\theta(x^{(i)}, z)] - \sum_{i=1}^N \mathbb{E}_{q^{(i)}(z|x^{(i)})}[\log q^{(i)}(z | x^{(i)})] \quad (44)$$

We maximized the ELBO w.r.t. q and θ by using CAVI, but by invoking our assumption that the posterior $p_\theta(z | x)$ can be computed, we can immediately find a maximum.

Theorem 1.2 (Posterior Maximizes ELBO)

When we set $q^{(i)}(z | x) = p(z | x^{(i)})$, equality is achieved.

$$\sum_{i=1}^N \log p_{\theta}(x^{(i)}) = \text{ELBO}(\mathcal{D}, q^{(1)}, \dots, q^{(n)}, \theta) \quad (45)$$

$$= \sum_{i=1}^N \mathbb{E}_{q^{(i)}(z|x^{(i)})} [\log p_{\theta}(x^{(i)}, z)] - \sum_{i=1}^N \mathbb{E}_{q^{(i)}(z|x^{(i)})} [\log q^{(i)}(z | x^{(i)})] \quad (46)$$

Proof.

Let's start by examining the gap between $\log p_{\theta}(x^{(i)})$ and the ELBO. From our previous derivations, this gap is the KL divergence:

$$\log p_{\theta}(x^{(i)}) - \text{ELBO}(x^{(i)}, q^{(i)}, \theta) = KL(q^{(i)}(z|x^{(i)}) || p_{\theta}(z|x^{(i)})) \quad (47)$$

$$= \mathbb{E}_{q^{(i)}} [\log q^{(i)}(z|x^{(i)}) - \log p_{\theta}(z|x^{(i)})] \quad (48)$$

When we set $q^{(i)}(z|x^{(i)}) = p_{\theta}(z|x^{(i)})$:

$$KL(p_{\theta}(z|x^{(i)}) || p_{\theta}(z|x^{(i)})) = \mathbb{E}_{p_{\theta}} [\log p_{\theta}(z|x^{(i)}) - \log p_{\theta}(z|x^{(i)})] \quad (49)$$

$$= \mathbb{E}_{p_{\theta}} [0] = 0 \quad (50)$$

Therefore, when summing over all samples:

$$\sum_{i=1}^N \log p_{\theta}(x^{(i)}) - \text{ELBO}(\mathcal{D}, q^{(1)}, \dots, q^{(n)}, \theta) = \sum_{i=1}^N KL(q^{(i)}(z|x^{(i)}) || p_{\theta}(z|x^{(i)})) = 0 \quad (51)$$

Therefore, our CAVI algorithm has been decomposed into the following.

1. E-step. Maximizing ELBO over the variational parameters q_{θ} is really just setting all the $q^{(i)}$ to the posteriors. Note that this is with respect to a fixed θ only.
2. M-step. Maximizing ELBO over the model parameters θ with fixed q is the same by taking the gradient w.r.t. θ which is easy.

This results in the following algorithm.

Algorithm 1.2 (EM Algorithm)

The EM algorithm is described as such:

1. Initialize θ .
2. *E-Step*. Since $\log p_{\theta}(x)$ is bounded below for all $q^{(1)}, \dots, q^{(n)}$ as

$$\sum_{i=1}^N \log p_{\theta}(x^{(i)}) \geq \sum_{i=1}^N \text{ELBO}(x^{(i)}, q^{(i)}, \theta) \quad (52)$$

setting $q^{(i)}(z|x^{(i)}) = p_{\theta}(z|x^{(i)})$ for all $i = 1, \dots, N$ achieves equality. Note that this equality only holds for the current fixed value of θ .

3. *M-Step*. We maximize with respect to θ whilst fixing $q^{(i)}$.^a

$$\theta = \operatorname{argmax}_{\theta} \sum_{i=1}^N \text{ELBO}(x^{(i)}, q^{(i)}, \theta) \quad (53)$$

$$= \operatorname{argmax}_{\theta} \sum_{i=1}^N \mathbb{E}_{q^{(i)}(z|x^{(i)})} [\log p_{\theta}(x^{(i)}, z)] - \sum_{i=1}^N \mathbb{E}_{q^{(i)}(z|x^{(i)})} [\log q^{(i)}(z|x^{(i)})] \quad (54)$$

4. Repeat steps 2 and 3 until convergence. Step 2 brings improvements because changing θ creates a new sum of ELBO functions as a new lower bound.

^aFor specific models like GMM as we will see later, this maximization has closed-form solutions, e.g. ϕ = average of responsibilities μ_k =: weighted average of points, Σ_k = weighted covariance. For other distributions, this maximum must be found analytically or numerically.

The EM algorithm is a specific instance of ELBO optimization! The additional assumption that EM has is that we can calculate the posterior densities.

Corollary 1.2 (Connection to ELBO)

The EM algorithm can be viewed as coordinate ascent on the ELBO where:

- E-step: Sets $q(z) = p_{\theta^{[t]}}(z|x)$, maximizing ELBO over q
- M-step: Maximizes ELBO over θ with fixed q

Note that there is a duality between the true parameters θ and the latent variables z . If θ is known, then the values of z can be found by maximizing the log-likelihood over all possible values of z . Conversely, if we know the value of the latent variables z , then we can find an estimate of the parameters by grouping the data points into each value of z and optimizing $p_{\theta}(x|z)$, e.g. by averaging the values. This suggests an iterative algorithm in the case where both θ and z are unknown. We assume that we know θ and optimize z , then optimize θ , and so on, similar to k -means clustering.

We can formulate the algorithm alternatively yet equivalently.

Algorithm 1.3 (EM Algorithm)

The **Expectation-Maximization algorithm** optimizes the likelihood above with the following steps.

1. First initialize $\theta = \theta^{[0]}$ in some way.^a
2. *E-Step*. Define

$$Q(\theta | \theta^{[t]}) = \mathbb{E}_{p_{\theta}(z|x)} [\log p_{\theta}(x, z)] = \int p_{\theta^{[t]}}(z|x) \log p_{\theta}(x, z) dz \quad (55)$$

as the expected value of the log-likelihood with respect to the current conditional distribution of z , given x and $\theta^{[t]}$.

3. *M-Step*. Find the parameters that maximize this quantity.

$$\theta^{[t+1]} = \operatorname{argmax}_{\theta} Q(\theta | \theta^{[t]}) \quad (56)$$

^aNote that within this θ are the parameterizations of the initial multinomial density p_Z , which is our initial “guess” of the distribution of Z .

Theorem 1.3 (EM Monotonicity)

The EM algorithm monotonically increases the observed data log-likelihood:

$$\log p_{\theta^{[t+1]}}(x) \geq \log p_{\theta^{[t]}}(x) \quad (57)$$

Therefore, though there is no guarantee that this will hit the global maximum, it will hit a local maximum.

Proof.

Let's consider the difference in log-likelihoods between iterations:

$$\log p_{\theta^{[t+1]}}(x) - \log p_{\theta^{[t]}}(x) = \left[Q(\theta^{[t+1]}|\theta^{[t]}) - H(\theta^{[t+1]}|\theta^{[t]}) \right] \quad (58)$$

$$- \left[Q(\theta^{[t]}|\theta^{[t]}) - H(\theta^{[t]}|\theta^{[t]}) \right] \quad (59)$$

where $H(\theta|\theta^{[t]}) = \mathbb{E}_{z|x, \theta^{[t]}}[\log p_{\theta}(z|x)]$. By the M-step, we know $Q(\theta^{[t+1]}|\theta^{[t]}) \geq Q(\theta^{[t]}|\theta^{[t]})$. Also, by Jensen's inequality:

$$H(\theta^{[t+1]}|\theta^{[t]}) \leq H(\theta^{[t]}|\theta^{[t]}) \quad (60)$$

Therefore, the difference is non-negative.

For some intuition, we can visualize l as a function of θ . For the sake of visuals, we will assume that $\theta \in \mathbb{R}$ and $l : \mathbb{R} \rightarrow \mathbb{R}$. On the contrary to what a visual is supposed to do, we want to point out that we cannot just visualize l as a curve in $\mathbb{R} \times \mathbb{R}$. This can be misleading since then it implies that the optimal θ value is easy to find, as shown in the left. Rather, we have no clue what the whole curve of l looks like, but we can get little snippets (right).

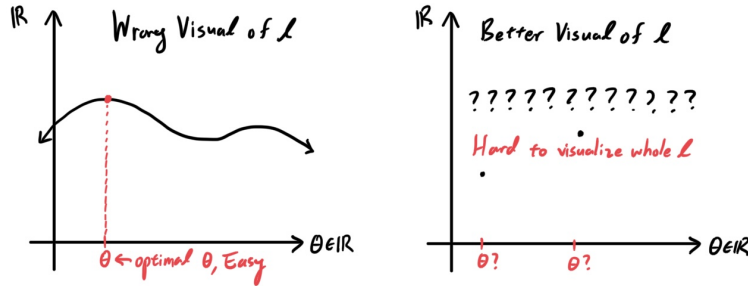


Figure 2

Rather, all we can do is hope to take whatever easier-to-visualize, lower-bound functions and maximize them as much as we can in hopes of converging onto l . Let us walk through the first two iterations of the EM algorithm. We first initialize θ to, say θ_0 . This immediately induces the lower-bound ELBO-sum function $\sum_i \text{ELBO}(x^{(i)}; p_Z^{*i}, \theta)$, which takes in multinomial density functions $p_Z^{*i} = p_1, p_2, \dots$ and outputs different functions of θ that are valid lower bounds. Two of these possible lower-bound functions are shown (in green) for when we input some arbitrary density p_1, p_2 . However, there exists a density $p_Z^{(i)}$ that produces not only the maximum possible lower-bound (called max ELBO, shown in red) but is equal to $l(\theta)$ for that density input $p_Z^{(i)}$. We maximize this function with respect to θ to get θ_1 as our next assignment of θ .

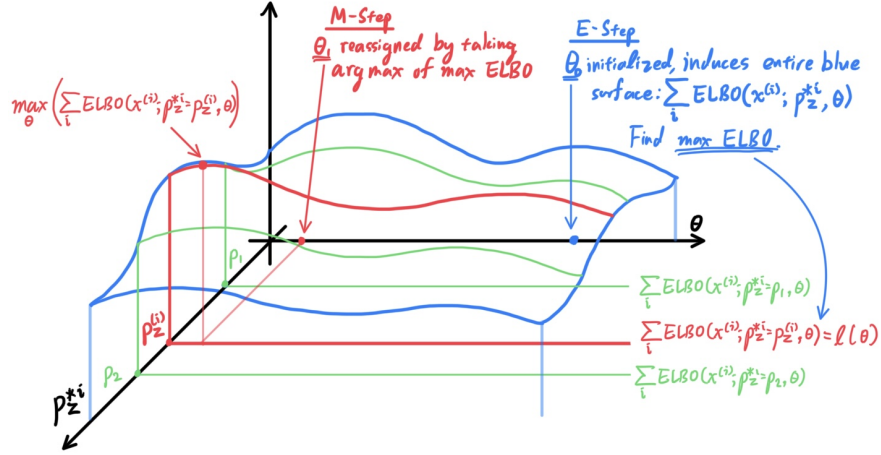


Figure 3

The next step is identical. Now that we have a new value of $\theta = \theta_1$, this induces the lower-bound ELBO-sum function $\sum_i \text{ELBO}(x^{(i)}; p_Z^{*i}, \theta)$ that also takes in multinomial densities p_Z^{*i} and outputs different functions of θ that are valid lower-bounds. Two possible lower bounds are shown (in green), but the maximum lower-bound (in blue) is produced when we input density $p_Z^{(i)}$. Since this max ELBO function is equal to θ for this fixed density input $p_Z^{(i)}$, we maximize this function with respect to θ to get θ_2 as our next assignment of θ .

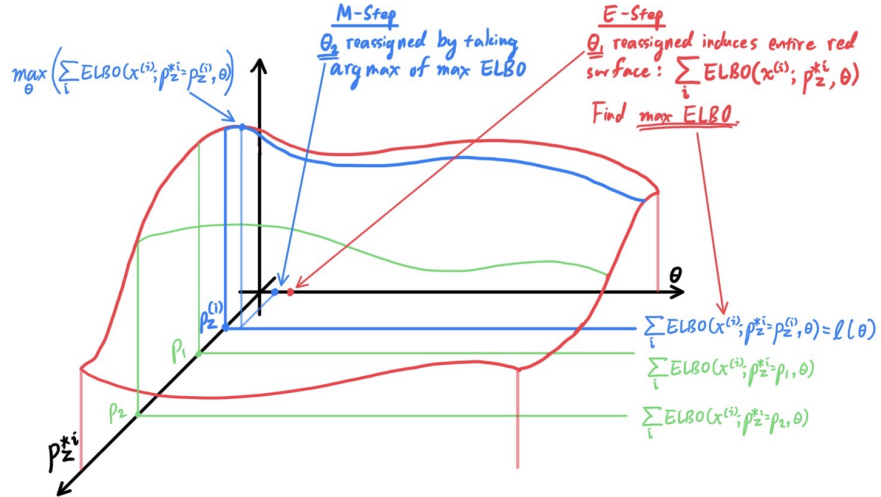


Figure 4

1.3 Gaussian Mixture Models

Given a training set $x_{i=1}^{(i)n}$ (without the y -labels and so in the unsupervised setting), there are some cases where it may seem like we can fit multiple Gaussian distributions in the input space \mathcal{X} . For example, the points below seem like they can be fitted well with 3 Gaussians.

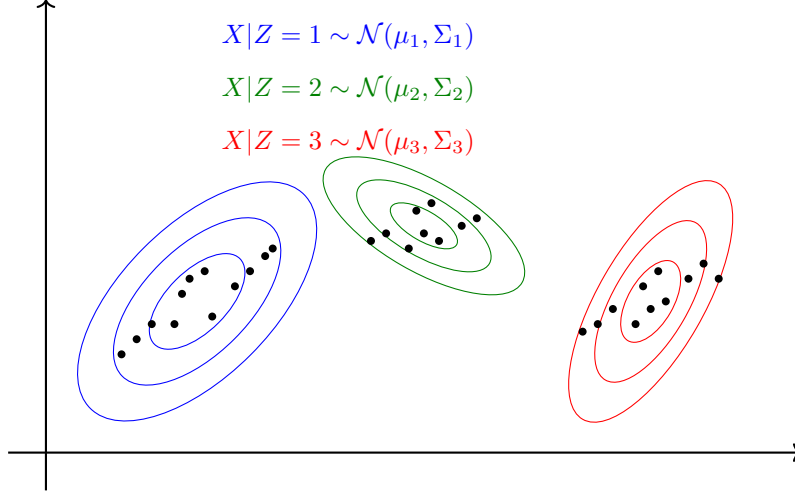


Figure 5: Example of data that can be fitted with 3 Gaussians

Therefore, we can construct a best-fit model as a composition of a multinomial distribution (to decide which one of the Gaussians x should follow) followed by a Gaussian.

Definition 1.3 (Gaussian Mixture Model)

The **Gaussian mixture model (GMM)** assumes that the covariates $x \sim X \in \mathbb{R}^d$ are generated by the following.^a The parameters are $\theta = \{\lambda, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k\}$.^b

1. A latent variable $z \sim \text{Multinomial}(\lambda)$, where $\lambda = (\lambda_1, \dots, \lambda_k)$ with PMF defined

$$p_\theta(z) = \lambda_z \quad (61)$$

2. The generative random variable $X | Z = z \sim \mathcal{N}(\mu_z, \Sigma_z)$ where $\mu_z \in \mathbb{R}^d, \Sigma_z \in \mathbb{R}^{d \times d}$ and PDF defined

$$p_\theta(x | z) = \frac{1}{(2\pi)^{d/2} |\Sigma_z|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_z)^\top \Sigma_z^{-1} (x - \mu_z) \right) \quad (62)$$

^aTherefore, our model says that each $x^{(i)}$ was generated by randomly choosing $z^{(i)}$ from $1, \dots, k$ according to some multinomial, and then the $x^{(i)}$ was drawn from one of the k Gaussians depending on $z^{(i)}$.

^bNote that λ really has $k - 1$ free parameters and Σ_i 's should be symmetric and positive-definite.

We can write down the log-likelihood of the given data $x^{(i)}$'s as a function of all the parameters above as

$$\sum_{i=1}^n \log p_\theta(x^{(i)}) = \sum_{i=1}^n \log \left(\sum_{z=1}^k p_\theta(x^{(i)} | z^{(i)}), p_\theta(z^{(i)}) \right) \quad (63)$$

Example 1.2 (Dual Nature of Latents and Parameters)

Note that since we only know that the *final* value of the i th sample is $x^{(i)}$ and not anything at all about which value $z^{(i)}$ the i th sample had, there is an extra unknown in this model. If we did know the values of the hidden variables $z^{(i)}$ (i.e. if we knew which of the k Gaussians each $x^{(i)}$ was generated from), then our log likelihood function would be much more simple since now, our givens will be both $x^{(i)}$ and $z^{(i)}$. Therefore, we don't have to condition on the $z^{(i)}$ and can directly calculate the log of the probability of us having sample values $(z^{(1)}, x^{(1)}), (z^{(2)}, x^{(2)}), \dots, (z^{(n)}, x^{(n)})$.

$$\sum_{i=1}^n \log p(x^{(i)}) = \sum_{i=1}^n \log p(x^{(i)}, z^{(i)}) = \sum_{i=1}^n \log p(x^{(i)} | z^{(i)}) p(z^{(i)}) \quad (64)$$

This model, with known $z^{(i)}$'s, is basically the GDA model, which is easy to calculate. That is, the maximum values of ϕ, μ, Σ are

$$\begin{aligned} \phi_j &= \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{z^{(i)}=j} \\ \mu_j &= \frac{\sum_{i=1}^n \mathbb{1}_{z^{(i)}=j} x^{(i)}}{\sum_{i=1}^n \mathbb{1}_{z^{(i)}=j}} \\ \Sigma_j &= \frac{1}{\sum_{i=1}^n \mathbb{1}_{z^{(i)}=j}} \sum_{i=1}^n \mathbb{1}_{z^{(i)}=j} (x^{(i)} - \mu_j)(x^{(i)} - \mu_j)^T \end{aligned}$$

But since we do *not* know the values of $z^{(i)}$, we first try to “guess” the values of the $z^{(i)}$'s and then update the parameters of our model assuming our guesses are correct.

Algorithm 1.4 (EM Algorithm on GMMs)

The EM Algorithm applied to GMMs has the following steps:

1. Randomly initialize $\theta^{[0]} = \{\lambda, \mu_1, \dots, \mu_k, \Sigma_1, \dots, \Sigma_k\}$.^a
2. **(E Step)** Calculate the posterior density $p(z | x)$ by applying Bayes rule to each sample keeping the parameter $\theta^{[t]}$ fixed.

$$p_{\theta^{[t]}}(z | x^{(i)}) = \frac{p_{\theta^{[t]}}(x^{(i)} | z) p_{\theta^{[t]}}(z)}{p(x)} = \frac{p_{\theta^{[t]}}(x^{(i)} | z) p_{\theta^{[t]}}(z)}{\sum_z p_{\theta^{[t]}}(x^{(i)} | z) p_{\theta^{[t]}}(z)} \quad (65)$$

We should have n different multinomial distribution parameters, each representing our best guess of what multinomial density $p(z | x^{(i)})$ each $x^{(i)}$ had followed in order to be at the given points. Let's label the updated parameters of the multinomial distribution of the i th sample to be $\lambda^{[t](i)}$ at the t th iteration.

3. **(M Step)** We update θ as such.

$$\lambda^{[t+1]} = \frac{1}{n} \sum_{i=1}^n \lambda^{[t](i)} \quad (66)$$

$$\mu_j^{[t+1]} = \frac{\sum_{i=1}^n \lambda_j^{[t](i)} x^{(i)}}{\sum_{i=1}^n \lambda^{[t](i)}} \quad (67)$$

$$\Sigma_j^{[t+1]} = \frac{1}{\sum_{i=1}^n \lambda_j^{[t](i)}} \sum_{i=1}^n \lambda_j^{[t](i)} (x^{(i)} - \mu_j^{[t+1]})(x^{(i)} - \mu_j^{[t+1]})^T \quad (68)$$

4. Repeat steps 2 and 3 until convergence.

^aThis might converge faster using K-means initialization.

Let us elaborate further on the intuition of this step. In the normal GDA with given values of $z^{(i)}$, we have $\lambda = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{z^{(i)} = j\} = \frac{1}{n}$ (Number of Samples in j th Gaussian), which is a sum of "hard" guesses, meaning that each $x^{(i)}$ is undoubtedly in cluster j or not, and so to find out our best guess for the true vector λ , all we have to do is find out the proportion of all examples in each of the k groups and we're done (without needing to iterate). However, in our EM model, we do not know the $z^{(i)}$'s, and so the best we can do is give the *probability* $\lambda_j^{(i)}$ that $x^{(i)}$ is in cluster j . So for each point $x^{(i)}$, the model has changed from it

being undoubtedly in group $z^{(i)} = j$ to it having a probability of being in $\lambda_j^{(i)}$ for $j = 1, \dots, k$.

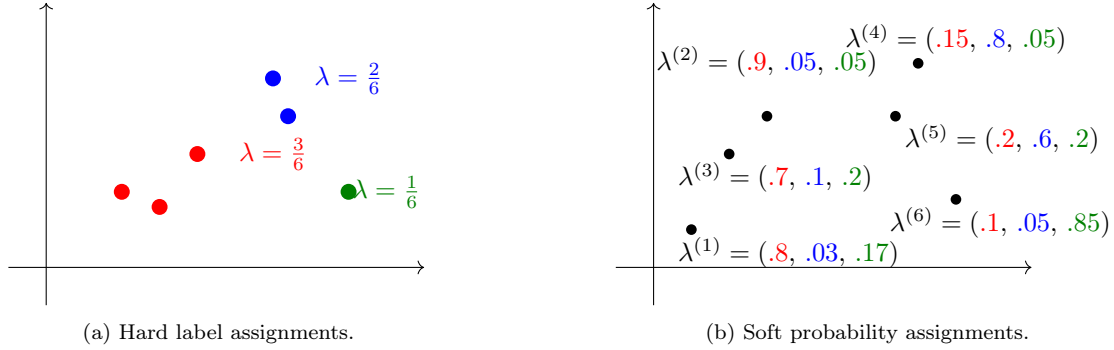


Figure 6: The superscript $[t]$ is omitted for clarity.

When we update the λ in the M-step, we can interpret the vectors $\lambda^{(i)}$ as tuples where $\lambda_j^{(i)}$ describes the expected "portion" of each sample $x^{(i)}$ to be in group j . So, we are adding up all the "portions" of the points that are expected to be in cluster j to get $\lambda = \sum_{i=1}^n \lambda^{(i)}$.

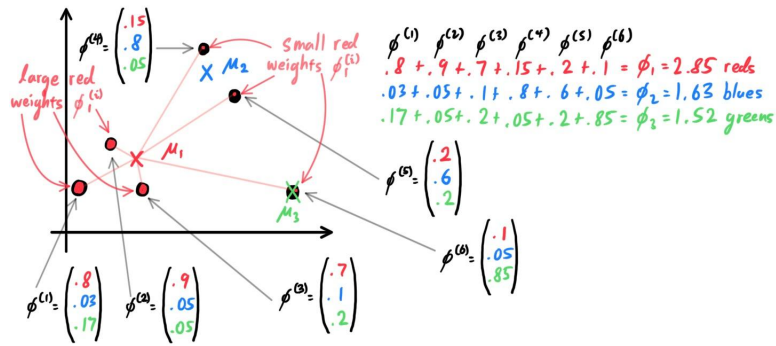


Figure 7

Now, given the j th Gaussian cluster, we would like to compute its mean μ_j . Since each $x^{(i)}$ has probability $\lambda_j^{(i)}$ of being in cluster j , we can weigh each of the n points by $\lambda_j^{(i)}$ (which determines how "relevant" $x^{(i)}$ is to cluster j) and average these (already weighted) points to get our "best-guess" of the mean μ_j . Given the MLE of the means, we can straightforwardly compute the MLE of the covariance matrices.

In summary, this entire algorithm results from modifying the "hard" data of each point $x^{(i)}$ being undoubtedly in one cluster to a model containing points $x^{(i)}$ that have been "smeared" around different clusters, with a probability $\lambda^{(i)}$ being in cluster j .

1.4 Nonlinear ICA

Bibliography

[GBB01] Evan Greensmith, Peter Bartlett, and Jonathan Baxter. Variance reduction techniques for gradient estimates in reinforcement learning. In T. Dietterich, S. Becker, and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems*, volume 14. MIT Press, 2001.

[KW22] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022.

- [PBJ12] John Paisley, David Blei, and Michael Jordan. Variational bayesian inference with stochastic search, 2012.
- [Wil92] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3):229–256, 1992.