

1 Introduction

1.1 Problem Setup

Suppose $X \in \mathbb{R}^n$ and $Z \in \mathbb{R}^m$ with $X = f(Z) + \varepsilon$, where $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ and $\varepsilon \sim N(0, \sigma^2 I)$. Assuming Z is normally distributed with independent marginals, this is equivalent to the following latent variable model (a special case of the well-known *nonlinear ICA* model):

$$\begin{aligned} Z &\sim N(0, I) \\ X | Z &\sim N(f(Z), \sigma^2 I). \end{aligned}$$

Let $\varphi(u; \mu, \Sigma)$ denote the density of a $N(\mu, \Sigma)$ random variable and $p_{\theta, \sigma^2}(x, z)$ denote the joint density under the model. It is easy to see that

$$\begin{aligned} p_{\theta, \sigma^2}(x, z) &= p_{\theta, \sigma^2}(x | z) p(z) = \varphi(x; f(z), \sigma^2 I) \varphi(z; 0, I) \\ L(\theta, \sigma^2; x) &= p_{\theta, \sigma^2}(x) = \int \varphi(x; f(z), \sigma^2 I) \varphi(z; 0, I) dz \end{aligned}$$

1.2 Objective Function

Now, suppose we let g_θ denote a family of deep neural network distributions parametrized by θ . To approximate the marginal density $p(x)$, we replace f with g_θ and try to find the choice of θ that maximizes the observed data likelihood. Given k observations $x^{(i)} \stackrel{i.i.d.}{\sim} p(x)$, we wish to solve the following maximum likelihood problem:

$$\max_{\theta, \sigma^2} \underbrace{\sum_{i=1}^k \log \int \varphi(x^{(i)}; g_\theta(z), \sigma^2 I) \varphi(z; 0, I) dz}_{:= \ell(\theta, \sigma^2)}$$

2 Direct MLE Method

In this section, we directly solve the MLE problem by computing gradients of $\ell(\theta, \sigma^2)$ w.r.t θ and σ^2 . This is, in general, intractable for arbitrary nonlinear ICA models but worst-case thinking does not apply to our special cases.

Gradient w.r.t θ

$$\begin{aligned} \nabla_\theta \ell(\theta, \sigma^2) &= \sum_{i=1}^k \frac{1}{L(\theta, \sigma^2; x^{(i)})} \int \nabla_\theta \varphi(x^{(i)}; g_\theta(z), \sigma^2 I) \varphi(z; 0, I) dz \\ &= \sum_{i=1}^k \frac{1}{L(\theta, \sigma^2; x^{(i)})} \int (2\pi\sigma^2)^{-n/2} \nabla_\theta \exp\left(-\frac{\|x^{(i)} - g_\theta(z)\|_2^2}{2\sigma^2}\right) \varphi(z; 0, I) dz \\ &= \sum_{i=1}^k \frac{1}{L(\theta, \sigma^2; x^{(i)})} \int \nabla_\theta \left(-\frac{\|x^{(i)} - g_\theta(z)\|_2^2}{2\sigma^2}\right) \varphi(x^{(i)}; g_\theta(z), \sigma^2 I) \varphi(z; 0, I) dz \\ &= \sum_{i=1}^k \frac{1}{L(\theta, \sigma^2; x^{(i)})} \int \left[\frac{1}{\sigma^2} \cdot (x^{(i)} - g_\theta(z))^T \nabla_\theta g_\theta(z)\right] \varphi(x^{(i)}; g_\theta(z), \sigma^2 I) \varphi(z; 0, I) dz \end{aligned}$$

Gradient w.r.t σ^2

$$\begin{aligned}
\nabla_{\sigma^2} \ell(\theta, \sigma^2) &= \sum_{i=1}^k \frac{1}{L(\theta, \sigma^2; x^{(i)})} \int \nabla_{\sigma^2} \varphi(x^{(i)}; g_\theta(z), \sigma^2 I) \varphi(z; 0, I) dz \\
&= \sum_{i=1}^k \frac{1}{L(\theta, \sigma^2; x^{(i)})} \int \nabla_{\sigma^2} (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{\|x^{(i)} - g_\theta(z)\|_2^2}{2\sigma^2}\right) \varphi(z; 0, I) dz \\
&= \sum_{i=1}^k \frac{1}{L(\theta, \sigma^2; x^{(i)})} \int \left[\frac{1}{2\sigma^2} \cdot \left(-n + \frac{\|x^{(i)} - g_\theta(z)\|_2^2}{\sigma^2}\right) \right] \varphi(x^{(i)}; g_\theta(z), \sigma^2 I) \varphi(z; 0, I) dz
\end{aligned}$$

From the two results above, we can iteratively update θ and σ^2 via gradient descent. The integrals can be approximated via numerical integration.

Algorithm 1 Direct MLE via Gradient Descent

- Initialise $\theta^{(0)}$ and $\sigma^{2(0)}$ and set $t = 0$
- Repeat until convergence
 - ▷ Compute the gradient $\nabla_{\theta} \ell(\theta^{(t)}, \sigma^{2(t)})$ and update the parameters

$$\theta^{(t+1)} = \theta^{(t)} - \eta_1 \nabla_{\theta} \ell(\theta^{(t)}, \sigma^{2(t)})$$

- ▷ Compute the gradient $\nabla_{\sigma^2} \ell(\theta^{(t+1)}, \sigma^{2(t)})$ and update the parameters

$$\sigma^{2(t+1)} = \sigma^{2(t)} - \eta_2 \nabla_{\sigma^2} \ell(\theta^{(t+1)}, \sigma^{2(t)})$$

- ▷ Set $t \leftarrow t + 1$
-

3 Variational Method

In this section, we consider variational inference and VAEs. We use the ELBO to obtain a lower bound on the likelihood $\ell(\theta, \sigma)$ and optimize the ELBO using SGD. The marginal likelihoods of individual datapoints can each be rewritten as

$$\log p_\theta(x^{(i)}) = D_{KL}(q_\phi(z | x^{(i)}) \| p_\theta(z | x^{(i)})) + \mathcal{L}(\theta, \phi; x^{(i)})$$

The second RHS term $\mathcal{L}(\theta, \phi; x^{(i)})$ is called the evidence lower bound on the marginal likelihood of datapoint i , and can be written as

$$\begin{aligned}
\log p_\theta(x^{(i)}) &\geq \mathcal{L}(\theta, \phi; x^{(i)}) = \mathbb{E}_{q_\phi(z|x)} [-\log q_\phi(z | x) + \log p_\theta(x, z)] \\
&= -D_{KL}(q_\phi(\mathbf{z} | \mathbf{x}^{(i)}) \| p_\theta(\mathbf{z})) + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_\theta(\mathbf{x}^{(i)} | \mathbf{z})]
\end{aligned}$$

We want to differentiate and optimize the lower bound $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i)})$ w.r.t. both the variational parameters ϕ and generative parameters θ .