

Auto-Encoding Variational Bayes

Authors: Diederik P. Kingma & Max Welling

Sai Subramanyam Thoota
SPC Lab, Department of ECE
Indian Institute of Science

August 25, 2018

Table of contents

- 1 Problem
- 2 The Variational Bound
- 3 The SGVB estimator and AEVB algorithm
- 4 The Reparameterization Trick
- 5 Variational Auto-Encoder

Problem

- Consider a unlabelled dataset $\mathbf{X} = \{\mathbf{x}\}_{i=1}^N$ consisting of N i.i.d. samples of some continuous or discrete random variable \mathbf{x} with some unknown distribution
- Assume data are generated by a random process parameterized by θ
- A common approach in statistical inference is to consider a joint distribution involving an unobserved continuous r.v \mathbf{z} (latent variable in a feature space \mathcal{Z})
- Goal is to learn the posterior distribution of the latent variables given the dataset \mathbf{X}
- We impose a simple prior $p_\theta(\mathbf{z})$, and the likelihood function is $p_\theta(\mathbf{x}|\mathbf{z})$
- We propose a general algorithm that learns the posterior distribution $p_\theta(\mathbf{z}|\mathbf{x})$ and/or the marginal distribution $p_\theta(\mathbf{x})$
 - Learning the marginal distribution is computationally intractable, especially when the dataset is very large
- Approach:
 - Introduce a recognition model $q_\phi(\mathbf{z}|\mathbf{x})$ as an approximation to the true posterior $p_\theta(\mathbf{z}|\mathbf{x})$
 - Recognition model is not necessarily factorial and its parameters are not computed from closed form expectations
 - Learn the recognition model parameters ϕ jointly with the generative model parameters θ

- Inference using exact posterior is computationally intractable which necessitates variational approaches
- Variational Bayesian (VB) approach involves the optimization of an approximation to the intractable posterior
- Mean-field approach requires analytical solutions of expectations w.r.t. the approximate posterior, which are intractable in the general case
- A reparameterization of the variational lower bound yields a simple differentiable unbiased estimator of the lower bound
 - Stochastic Gradient Variational Bayes (SGVB) estimator
 - Optimized using standard stochastic gradient ascent techniques
- Auto-encoding VB (AEVB) algorithm is proposed for the case of an i.i.d. dataset and continuous latent variables
 - Inference and learning using the SGVB estimator to optimize a recognition model
 - Neural network used to learn the parameters of the approximate posterior, which leads to the variational auto-encoder (VAE)

The Variational Bound (Evidence Lower Bound)

- The marginal likelihood $\log p_{\theta}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)}) = \sum_{i=1}^M \log p_{\theta}(\mathbf{x}^{(i)})$ can be rewritten as:

$$\log p_{\theta}(\mathbf{x}^{(i)}) = \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) + D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) || p_{\theta}(\mathbf{z}|\mathbf{x}^{(i)})) \quad (1)$$

$$\geq \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) \quad (2)$$

- $\mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)})$ is the variational lower bound on the marginal likelihood of the data point i since KL divergence is non-negative

$$\mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) = \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[-\log q_{\phi}(\mathbf{z}|\mathbf{x}) + \log p_{\theta}(\mathbf{x}^{(i)}, \mathbf{z}) \right] \quad (3)$$

$$= -D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)}) || p_{\theta}(\mathbf{z})) + \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})} \left[\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) \right] \quad (4)$$

- We want to differentiate and optimize the lower bound w.r.t. both the variational parameters ϕ and generative parameters $\boldsymbol{\theta}$

The SGVB estimator and AEVB algorithm

- Under mild conditions for a chosen approximate posterior $q_\phi(\mathbf{z}|\mathbf{x})$, we can reparameterize the r.v $\tilde{\mathbf{z}} \sim q_\phi(\mathbf{z}|\mathbf{x})$ using a differentiable transformation $\mathbf{g}_\phi(\boldsymbol{\epsilon}, \mathbf{x})$ of an auxiliary noise variable $\boldsymbol{\epsilon}$:

$$\tilde{\mathbf{z}} = \mathbf{g}_\phi(\boldsymbol{\epsilon}, \mathbf{x}) \quad \text{with} \quad \boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon}) \quad (5)$$

- Monte Carlo estimates of expectations of some function $f(\mathbf{z})$ can be computed as follows:

$$\begin{aligned} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}^{(i)})} [f(\mathbf{z})] &= \mathbb{E}_{p(\boldsymbol{\epsilon})} \left[f \left(\mathbf{g}_\phi \left(\boldsymbol{\epsilon}, \mathbf{x}^{(i)} \right) \right) \right] \\ &\simeq \frac{1}{L} \sum_{l=1}^L f \left(\mathbf{g}_\phi \left(\boldsymbol{\epsilon}^{(l)}, \mathbf{x}^{(i)} \right) \right) \end{aligned} \quad (6)$$

$$\text{where } \boldsymbol{\epsilon}^{(l)} \sim p(\boldsymbol{\epsilon})$$

- Applying this technique to the variational lower bound, we get the SGVB estimator $\tilde{\mathcal{L}}^A(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) \simeq \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)})$ as follows:

$$\tilde{\mathcal{L}}^A(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) = \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)}, \mathbf{z}^{(i,l)}) - \log q_\phi(\mathbf{z}^{(i,l)}|\mathbf{x}^{(i)}) \quad (7)$$

$$\text{where } \mathbf{z}^{(i,l)} = \mathbf{g}_\phi(\boldsymbol{\epsilon}^{(l)}, \mathbf{x}^{(i)}) \quad \text{and} \quad \boldsymbol{\epsilon}^{(l)} \sim p(\boldsymbol{\epsilon})$$

- Second version of the SGVB estimator $\tilde{\mathcal{L}}^B(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) \simeq \mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)})$

$$\tilde{\mathcal{L}}^B(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) = -D_{KL} \left(q_\phi \left(\mathbf{z} | \mathbf{x}^{(i)} \right) \parallel p_\theta \left(\mathbf{z} \right) \right) + \frac{1}{L} \sum_{l=1}^L \log p_\theta \left(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)} \right) \quad (8)$$

- The KL divergence term can be integrated analytically, and hence only the expected reconstruction error requires estimation by sampling (see (8))
- Given a dataset \mathbf{X} with N datapoints, we can construct an estimator of the marginal likelihood lower bound of the full dataset based on minibatches:

$$\mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{X}) \simeq \tilde{\mathcal{L}}^M(\boldsymbol{\theta}, \phi; \mathbf{X}^M) = \frac{N}{M} \sum_{i=1}^M \tilde{\mathcal{L}}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) \quad (9)$$

The Reparameterization Trick

- Goal: To generate samples of a random variable $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x})$
- Possible to express \mathbf{z} as a deterministic variable $\mathbf{z} = \mathbf{g}_\phi(\boldsymbol{\epsilon}, \mathbf{x})$
 - $\boldsymbol{\epsilon}$ is an auxiliary variable with independent marginal p_ϵ
 - $\mathbf{g}_\epsilon(\cdot)$ is some vector-valued function parameterized by ϕ
- Proof: Given the mapping $\mathbf{z} = \mathbf{g}_\phi(\boldsymbol{\epsilon}, \mathbf{x})$, we know that

$$\begin{aligned}
 q_\phi(\mathbf{z}|\mathbf{x}) \prod_i dz_i &= p(\boldsymbol{\epsilon}) \prod_i d\epsilon_i \\
 \implies \int q_\phi(\mathbf{z}|\mathbf{x}) f(\mathbf{z}) d\mathbf{z} &= \int p(\boldsymbol{\epsilon}) f(\mathbf{g}_\phi(\boldsymbol{\epsilon}, \mathbf{x})) d\boldsymbol{\epsilon} \\
 &\simeq \frac{1}{L} \sum_{l=1}^L f(\mathbf{g}_\phi(\mathbf{x}, \boldsymbol{\epsilon})^{(l)})
 \end{aligned} \tag{10}$$

- This trick is used to obtain a differentiable estimator of the variational lower bound
- Example:
 - $z \sim \mathcal{N}(\mu, \sigma^2)$
 - Reparameterization $z = \mu + \sigma\epsilon$, where $\epsilon \sim \mathcal{N}(0, 1)$

Variational Auto-Encoder

Encoder Neural Network:

- Probabilistic encoder $q_\phi(\mathbf{z}|\mathbf{x})$
- Simple prior over the latent variables \mathbf{z} chosen as multivariate Gaussian $p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$
 - Note that the prior lacks parameters
- Let $p_\theta(\mathbf{z}|\mathbf{x})$ be a multivariate Gaussian or Bernoulli whose parameters are computed from \mathbf{z} using a fully connected MLP
- We choose the variational approximate posterior to be a multivariate Gaussian:

$$\log q_\phi(\mathbf{z}|\mathbf{x}^{(i)}) = \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}^{(i)}, \boldsymbol{\sigma}^{2(i)}\mathbf{I}) \quad (11)$$

where the mean and s.d. of the approximate posterior are the outputs of the encoding MLP, i.e., $\boldsymbol{\mu}^{(i)} = \mathbf{W}_4\mathbf{h}_e + \mathbf{b}_4$, $\log \boldsymbol{\sigma}^{2(i)} = \mathbf{W}_5\mathbf{h}_e + \mathbf{b}_5$, and the hidden layer output $\mathbf{h}_e = \tanh(\mathbf{W}_3\mathbf{x}^{(i)} + \mathbf{b}_3)$

- The parameter $\phi = \{\mathbf{W}_3, \mathbf{W}_4, \mathbf{W}_5, \mathbf{b}_3, \mathbf{b}_4, \mathbf{b}_5\}$ are the weights and biases, which are found by passing $\mathbf{x}^{(i)}$ to the MLP

Decoder Neural Network:

- Probabilistic decoder $p_\theta(\mathbf{x}|\mathbf{z})$
- Weights and biases found similar to the encoder NN

- Sample from the posterior $\mathbf{z}^{(i,l)} \sim q_\phi(\mathbf{z}|\mathbf{x}^{(i)})$ using $\mathbf{z}^{(i,l)} = g_\phi(\mathbf{x}^{(i)}, \epsilon^{(l)}) = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \epsilon^{(l)}$ where $\epsilon^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- The resulting estimator of the evidence lower bound for this model is given by

$$\mathcal{L}(\boldsymbol{\theta}, \phi; \mathbf{x}^{(i)}) \simeq \frac{1}{2} \sum_{j=1}^J \left(1 + \log \left((\sigma_j^{(i)})^2 \right) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2 \right) + \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)}) \quad (12)$$

where $\mathbf{z}^{(i,l)} = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \epsilon^{(l)}$ and $\epsilon^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

- The decoding term $p_\theta(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)})$ is a Bernoulli or Gaussian MLP depending on the type of data we are modelling

Algorithm 1 Minibatch version of the Auto-Encoding VB (AEVB) algorithm. Either of the two SGVB estimators in section 2.3 can be used. We use settings $M = 100$ and $L = 1$ in experiments.

$\boldsymbol{\theta}, \phi \leftarrow$ Initialize parameters

repeat

$\mathbf{X}^M \leftarrow$ Random minibatch of M datapoints (drawn from full dataset)

$\epsilon \leftarrow$ Random samples from noise distribution $p(\epsilon)$

$\mathbf{g} \leftarrow \nabla_{\boldsymbol{\theta}, \phi} \tilde{\mathcal{L}}^M(\boldsymbol{\theta}, \phi; \mathbf{X}^M, \epsilon)$ (Gradients of minibatch estimator [8])

$\boldsymbol{\theta}, \phi \leftarrow$ Update parameters using gradients \mathbf{g} (e.g. SGD or Adagrad [DHS10])

until convergence of parameters $(\boldsymbol{\theta}, \phi)$

return $\boldsymbol{\theta}, \phi$

THANK YOU!