

Video Change Detection Project Update

Suya Wu

December 21, 2018

Week 12

An idea: Mixture Variational Autoencoders

Since the original VAEs cannot be directly applied to discrete variables with standard back-propagation through the reparametrisation trick, previous works either estimate the gradient of discrete variables or adjusting the architecture of the standard VAEs. Inspired by the framework concrete VAEs [Jang et al. \[2016\]](#), which enables backpropagation through discrete samples, we can train classes with original VAEs by choosing either gaussian mixture model (GMM) or bernoulli mixture model (BMM) as the prior distribution for continuous latent variables or discrete latent variables respectively.

The generative model

We consider the generative model $p_\theta(\mathbf{x}, \mathbf{z}, \mathbf{c}) = p_\theta(\mathbf{c})p_\theta(\mathbf{z}|\mathbf{c})p_\theta(\mathbf{x}|\mathbf{z}, \mathbf{c})$, where an observe sample \mathbf{x} is generated from a set of latent variables \mathbf{z} and \mathbf{z} follows the mixture distributions *w.r.t.* the classes variables \mathbf{c} . They are distributed as:

$$\mathbf{x} \sim \text{Bernoulli}(\boldsymbol{\mu}_x) \text{ or } \mathbf{x} \sim N(\boldsymbol{\mu}_x, \boldsymbol{\sigma}_x \mathbf{I}) \quad (1)$$

$$\mathbf{z} \sim \text{Bernoulli}(\boldsymbol{\mu}_z) \text{ or } \mathbf{z} \sim N(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z \mathbf{I}) \quad (2)$$

$$\mathbf{c} \sim \text{Categorical}(\boldsymbol{\pi}) \quad (3)$$

The parametric space θ therefore relevant to $\{\boldsymbol{\mu}_x, \boldsymbol{\mu}_z, \boldsymbol{\sigma}_x, \boldsymbol{\sigma}_z, \boldsymbol{\pi}\}$. We denote $q_\Phi(\mathbf{z}, \mathbf{c}|\mathbf{x})$ as the variational approximation to the intractable posterior $p_\theta(\mathbf{z}, \mathbf{c}|\mathbf{x})$. Then generated variables are distributed as:

$$\mathbf{c} \sim \text{Categorical}(\hat{\boldsymbol{\pi}}) \quad (4)$$

$$\mathbf{z}|\mathbf{c} \sim \text{Bernoulli}(\boldsymbol{\mu}_{z|\mathbf{c}}) \text{ or } \mathbf{z}|\mathbf{c} \sim N(\boldsymbol{\mu}_{z|\mathbf{c}}, \boldsymbol{\sigma}_{z|\mathbf{c}} \mathbf{I}) \quad (5)$$

$$\mathbf{x}|\mathbf{z}, \mathbf{c} \sim \text{Bernoulli}(\boldsymbol{\mu}_{x|\mathbf{z}}) \text{ or } \mathbf{x}|\mathbf{z}, \mathbf{c} \sim N(\boldsymbol{\mu}_{x|\mathbf{z}}, \boldsymbol{\sigma}_{x|\mathbf{z}} \mathbf{I}) \quad (6)$$

The observed sample \mathbf{x} can be generated from a neural network model parametrised by Φ , which is relevant to $\{\boldsymbol{\mu}_{x|\mathbf{z}}, \boldsymbol{\mu}_{z|\mathbf{c}}, \boldsymbol{\sigma}_{x|\mathbf{z}}, \boldsymbol{\sigma}_{z|\mathbf{c}}, \hat{\boldsymbol{\pi}}\}$. By the generative process, where \mathbf{x} and \mathbf{c} can be treated independently with conditioning on \mathbf{z} , we can approximate

$$q_\Phi(\mathbf{x}|\mathbf{z}, \mathbf{c}) = q_\Phi(\mathbf{x}|\mathbf{z})$$

The variational lower bound

The loglikelihood of the whole observed data \mathbf{X} , $\log p_\theta(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}) = \sum_{i=1}^N \log p_\theta(\mathbf{x}^{(i)})$. Each likelihood element can be written as:

$$\log p_\theta(\mathbf{x}^{(i)}) = KL(q_\Phi(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})) + \mathcal{L}(\theta, \Phi; \mathbf{x}^{(i)}) \quad (7)$$

The first RHS term is the KL-divergence of the approximate from the true posterior. Following the same argument in [Kingma & Welling \[2013\]](#) that the second RHS term is the lower bound we need to optimize since this KL-divergence is non-negative.

$$\mathcal{L}(\theta, \Phi; \mathbf{x}^{(i)}) = E_{q_{\Phi}(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}, \mathbf{z}, \mathbf{c}) - \log q_{\Phi}(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})] \quad (8)$$

which can also be written as

$$\begin{aligned} \mathcal{L}(\theta, \Phi; \mathbf{x}^{(i)}) &= E_{q_{\Phi}(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})} [\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}) + \log p_{\theta}(\mathbf{z}|\mathbf{c}) + \log p_{\theta}(\mathbf{c}) - \log q_{\Phi}(\mathbf{z}|\mathbf{x}^{(i)}) - \log q_{\Phi}(\mathbf{c}|\mathbf{z})] \\ &= E_{q_{\Phi}(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})} [\log p(\mathbf{x}^{(i)}|\mathbf{z})] + E_{q_{\Phi}(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})} \left[\log \frac{p(\mathbf{z}|\mathbf{c})}{q(\mathbf{z}|\mathbf{x}^{(i)})} \right] + E_{q_{\Phi}(\mathbf{z}, \mathbf{c}|\mathbf{x}^{(i)})} \left[\log \frac{p(\mathbf{c})}{q(\mathbf{c}|\mathbf{z})} \right] \\ &= \int_{\mathbf{z}} \sum_{\mathbf{c}} \log p(\mathbf{x}^{(i)}|\mathbf{z}) q(\mathbf{z}|\mathbf{x}^{(i)}) q(\mathbf{c}|\mathbf{z}) d\mathbf{z} + \int_{\mathbf{z}} \sum_{\mathbf{c}} q(\mathbf{z}|\mathbf{x}^{(i)}) q(\mathbf{c}|\mathbf{z}) \log \frac{p(\mathbf{z}|\mathbf{c})}{q(\mathbf{z}|\mathbf{x}^{(i)})} d\mathbf{z} \\ &\quad + \int_{\mathbf{z}} \sum_{\mathbf{c}} q(\mathbf{z}|\mathbf{x}^{(i)}) q(\mathbf{c}|\mathbf{z}) \log \frac{p(\mathbf{c})}{q(\mathbf{c}|\mathbf{z})} d\mathbf{z} \\ &\approx \frac{1}{L} \sum_{l=1}^L \sum_{\mathbf{c}} \log p(\mathbf{x}^{(i)}|\mathbf{z}^{(l)}) q(\mathbf{c}|\mathbf{z}^{(l)}) + \sum_{\mathbf{c}} KL(q(\mathbf{z}|\mathbf{x}^{(i)})||p(\mathbf{z}|\mathbf{c})) q(\mathbf{c}|\mathbf{z}) \\ &\quad + \frac{1}{L} \sum_{l=1}^L KL(q(\mathbf{c}|\mathbf{z}^{(l)})||p(\mathbf{c})) \end{aligned}$$

The last approximation may be confused, since we use only one sample ($L=1$) to estimate the integrate w.r.t $q_{\Phi}(\mathbf{z}|\mathbf{x}^{(i)})$ following the Stochastic Gradient Variational Bayes estimator (SGVB) [Kingma & Welling \[2013\]](#). Based on the one sample of \mathbf{z} , we also sample one \mathbf{c} .

1. The first term of loss is only depending on $p(\mathbf{x}^{(i)}|\mathbf{z}^{(1)})$, which can be learned by the decoder procedure, which is 2 layers linear neutral network.
2. The second term evaluate the distance between $q(\mathbf{z}|\mathbf{x}^{(i)})$ and $p(\mathbf{z}|\mathbf{c})$, where the first one is learned by the encoder procedure (2 layers linear neutral network) and the second one is assumed to be the true distribution of \mathbf{z} depending on different classes \mathbf{c} . We updated $p(\mathbf{z}|\mathbf{c})$ by

$$\begin{aligned} \mu_{\mathbf{z}|\mathbf{c}=k}^{update} &= \frac{\sum_{i=1}^n q(c=k|\mathbf{z}_i^{(l)}) \mathbf{z}_i^{(l)}}{\sum_{i=1}^n q(c=k|\mathbf{z}_i^{(l)})} \\ \Sigma_{\mathbf{z}|\mathbf{c}=k}^{update} &= \frac{\sum_{i=1}^n q(c=k|\mathbf{z}_i^{(l)}) (\mathbf{z}_i^{(l)} - \mu_{\mathbf{z}|\mathbf{c}=k}) (\mathbf{z}_i^{(l)} - \mu_{\mathbf{z}|\mathbf{c}=k})^T}{\sum_{i=1}^n q(c=k|\mathbf{z}_i^{(l)})} \end{aligned}$$

3. $q(\mathbf{c}|\mathbf{z}^{(1)})$ can be learned by 2 layers linear neutral network with gumbel-softmax distribution to estimate the discrete distribution of \mathbf{c} .

Then applied the Gumbel-Softmax estimator and SGVB to be capable of differentiating and optimizing the lower bound w.r.t. the variational parameters Φ and generative parameters θ . However, the third loss term is always converging to a constant and does not contribute to the whole loss function. It cannot distinguish different classes but only learn the first term (the reconstruction term).