

Due Date : April 14th, 2018

Instructions

- For all questions, show your work !
- This part (theory) is to be done individually.
- Use a doc prep system such as LaTeX, or scan a **very neatly** hand written version.
- Submit your answers electronically via the course studium page.

1. (10 points) Reparameterization Trick of Variational Autoencoder

Consider a generative model that factorizes as follows $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})$, where $p(\mathbf{x} | \mathbf{z})$ is mapped through a neural net, i.e. $p(\mathbf{x} | \mathbf{z}) = p(\mathbf{x}; \mathbf{h}_\theta(\mathbf{z}))$, θ being the set of parameters for the generative network (i.e. decoder), a simple distribution parameterized by $h(\cdot)$ such as Gaussian or Bernoulli (i.e. $p(\mathbf{x} | \mathbf{z}) = \prod_j p(x_j | \mathbf{z})$). In the case of Gaussian, $\mathbf{h}_\theta(\mathbf{z})$ refers to the mean and variance, per dimension as it is fully factorized in the common setting. We have $\mathbf{z} \in \mathbb{R}^K$, which implies a continuous latent space model, and $p(\mathbf{z}) = \mathcal{N}(0, \mathbf{I}_K)$, where \mathbf{I} is the identity matrix. The framework of auto-encoding variational Bayes considers maximizing the variational lower bound on the log-likelihood $\mathcal{L}(\theta, \phi) \leq \log p(\mathbf{x})$, which is expressed as

$$\mathcal{L}(\theta, \phi) = \mathbf{E}_{q_\phi}[\log p(\mathbf{x} | \mathbf{z})] - \mathbf{KL}(q_\phi(\mathbf{z} | \mathbf{x}) \| p_\theta(\mathbf{z})) \quad (1)$$

where ϕ is the set of parameters used for the inference network (i.e. encoder). The reparameterization trick used in the original work rewrites the random variable in the variational distribution as

$$\mathbf{z} = \mu(\mathbf{x}) + \sigma(\mathbf{x}) \odot \epsilon \quad (2)$$

where $\epsilon \sim \mathcal{N}(\epsilon; 0, \mathbf{I})$, so that gradient can be backpropagated through the stochastic bottleneck.

- (a) Prove that the linearly transformed standard Gaussian noise (2) has the same mean and variance as $\mathcal{N}(\mathbf{z}; \mu(\mathbf{x}), \sigma^2(\mathbf{x}))$. What if we write $\mathbf{z} = \mu(\mathbf{x}) + S(\mathbf{x})\epsilon$, where $S(\mathbf{x}) \in \mathbb{R}^{K \times K}$ could be a reshaped K^2 dimensional output of a neural net? What is the new distribution this reparameterization induces.
- (b) If the traditional mean field variational method is used, i.e. if we factorize the variational distribution as a product of distributions : $q^{mf}(z_i) = \prod_j \mathcal{N}(z_{i,j} | m_{i,j}, \sigma_{i,j}^2)$ for each data instance x_i , and we maximize the lower bound with respect to the variational parameters and model parameters iteratively, can the inference network used in the variational autoencoder q_ϕ (2) outperform the mean field method? What is the advantage of using an encoder as in VAE?

2. (10 points) Importance Weighted Autoencoder

When training a variational autoencoder, the standard training objective is to maximize the evidence lower bound (ELBO). Here we consider another lower bound, called the Importance Weighted Lower Bound (IWLB), a tighter bound than ELBO, defined as

$$\mathcal{L}_k = \mathbf{E}_{\mathbf{z}_{1:k} \sim q(\mathbf{z} | \mathbf{x})} \left[\log \frac{1}{k} \sum_{j=1}^k \frac{p(\mathbf{x}, \mathbf{z}_j)}{q(\mathbf{z}_j | \mathbf{x})} \right] \quad (3)$$

for an observed variable \mathbf{x} and a latent variable \mathbf{z} , k being the number of importance samples. The model we are considering has joint that factorizes as $p(\mathbf{z}, \mathbf{x}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})$, \mathbf{x} and \mathbf{z}

being the observed and latent variables, respectively. In the following questions, one needs to make use of the Jensen's inequality :

$$f(\mathbf{E}[X]) \leq \mathbf{E}[f(X)] \quad (4)$$

for a convex function f .

- (a) Show that IWLB is a lower bound on the log likelihood $\log p(\mathbf{x})$.
- (b) Given a special case where $k = 2$, prove that \mathcal{L}_2 is a tighter bound than the ELBO (with $k = 1$).

3. (10 points) Maximum Likelihood for Generative Adversarial Networks

The original GAN objective is the following

$$\max_D \mathbf{E}_{P_{\text{data}}(x)} [\log D(x)] + \mathbf{E}_{P(z)} [\log(1 - D(G(z)))]; \quad \max_G \mathbf{E}_{P(z)} [\log D(G(z))]$$

This generator objective can be generalized by replacing the log with a general function f :

$$\max_G \mathbf{E}_{P(z)} [f(D(G(z)))]$$

Find a function f such that the objective corresponds to maximum likelihood, assuming the discriminator is optimal.

Hint : Use the optimal discriminator : $D^* = \frac{P_{\text{data}}}{(P_{\text{data}} + P_{\text{gen}})}$.