

## Tutorial on Deep Generative Model

Loc Nguyen

ngphloc@gmail.com

Loc Nguyen's Academic Network, Vietnam

Doi: <https://doi.org/10.210985/hhv29f22>

Date of Submission: 18-12-2024	Date of Acceptance: 30-04-2025	Date of Publish: 30-04-2025
--------------------------------	--------------------------------	-----------------------------

### Abstract

Artificial intelligence (AI) is a current trend in computer science, which extends itself its amazing capacities to other technologies such as mechatronics and robotics. Going beyond technological applications, the philosophy behind AI is that there is a vague and potential convergence of artificial manufacture and natural world although the limiting approach may be still very far away, but why? The implicit problem is that Darwin theory of evolution focuses on natural world where breeding conservation is the cornerstone of the existence of creature world but there is no similar concept of breeding conservation in artificial world whose things are created by human. However, after developing for a long time until now, AI issues an interesting concept of generation in which artifacts created by computer science can derive their new generations inheriting their aspects / characteristics. Such generated artifacts make us look back on offsprings by the process of breeding conservation in natural world. Therefore, it is possible to think that AI generation, which is a recent subject of AI, is a significant development in computer science as well as high-tech domain. AI generation does not help us to reach near biological evolution even in the case that AI can combine with biological technology but, AI generation can help us to extend our viewpoint about Darwin theory of evolution as well as there may exist some uncertain relationship between man-made world and natural world. Anyhow AI generation is a current important subject in AI and there are two main generative models in computer science: 1) generative model that applies large language model into generating natural language texts understandable by human and 2) generative model that applies deep neural network into generating digital content such as sound, image, and video. This technical report focuses on deep generative model (DGM) for digital content generation, which is a short summary of approaches to implement DGMs. Researchers can read this work as an introduction to DGM with easily understandable explanations.

**Keywords:** generative artificial intelligence, deep neural network, deep generative model, data generation.

### 1. INTRODUCTION TO DEEP GENERATIVE MODEL (DGM)

By informal statement, generative artificial intelligence (GenAI) applications aim to reproduce original artifacts such as images, sounds, music, texts, and speeches into a new artifact with some changes. The problem is that reproduction or generation, which is not duplication, indeed derives a new piece of content which is large or small from whole content of the original artifacts. For example, given a smiling face of a specific person, GenAI application will generate a crying face of the same person. As a subdomain of GenAI, deep

generative model (DGM) applies deep neural network (DNN) into generating artifacts but many deep generative models (DGMs) are also relevant to applied statistics. Note, DNN is an artificial neural network having many hidden layers, besides input layer and output layer. Training deep neural network or learning deep neural network is known as deep learning. Given random variable vector  $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$  presenting any digital artifact or any digital data such as image and sound, let  $P(\mathbf{x})$  be probability density function (PDF) of  $\mathbf{x}$  but it is difficult to estimate such probabilistic distribution  $P(\mathbf{x})$  because data  $\mathbf{x}$  is complicated with suppose that  $\mathbf{x}$  belongs to the real field  $R^m$  where  $m$  is high dimension and so,  $P(\mathbf{x})$  is called intractable PDF. Suppose there is another random variable vector  $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$  belonging to the real field  $R^n$  where  $n$  is low dimension ( $n < m$ ) so that PDF of  $\mathbf{z}$  denoted  $P(\mathbf{z})$  is tractable and it is possible to understand  $P(\mathbf{z})$ . Moreover, it is most important that suppose there is a function  $g(\mathbf{z} | \Phi) = \mathbf{x}$  that maps tractable data  $\mathbf{z}$  to intractable data  $\mathbf{x}$  where  $\Phi$  is parameter of such mapping function. For some illustrations or examples in this report, random variable vector  $\mathbf{x}$  is flattened from two-dimension image. As a convention, the function  $g(\mathbf{z} | \Phi) = \mathbf{x}$  is called generator of  $\mathbf{x}$  and its parameter  $\Phi$  is called generator parameter (Ruthotto & Haber, 2021, p. 2).

$g: R^n \rightarrow R^m$  such that  $\mathbf{x} = g(\mathbf{z} | \Phi)$  where  $\mathbf{z} \in Z \subseteq R^n, \mathbf{x} \in X \subseteq R^m$

Where  $Z$  and  $X$  are domains of tractable data  $\mathbf{z}$  and intractable data  $\mathbf{x}$  with note that  $Z$  is called latent space and  $X$  is called sample space by convention. When  $Z$  is called latent space, tractable PDF  $P(\mathbf{z})$  is called latent distribution. Because  $g(\mathbf{z} | \Phi)$  is essentially vector-by-vector function whose input and output are vectors, it should have denoted as  $\mathbf{g}(\mathbf{z} | \Phi)$ . However, it is still denoted  $g(\mathbf{z} | \Phi)$  in context DNN because there are two reasons: 1)  $\mathbf{g}(\mathbf{z} | \Phi)$  is not bijection and 2) the output  $\mathbf{x}$  of  $\mathbf{g}(\mathbf{z} | \Phi)$  can be considered as scalar variable  $x$  corresponding to an output neuron of output layer in neural network. Therefore,  $g(\mathbf{z} | \Phi)$  also implies a vector-by-scalar function whose first-order derivative can be considered as gradient vector although the first-order derivative of vector-by-vector function  $\mathbf{g}(\mathbf{z} | \Phi)$  is Jacobian matrix. Note, in mathematical, the first-order derivative of scalar-by-vector function is called gradient vector and the first-order derivative of vector-by-vector function is called Jacobian matrix.

The ultimate purpose of any DGM is to determine parameter  $\Phi$  because generator  $g(\mathbf{z} | \Phi)$  is defined based on  $\Phi$ . In DGM, generator  $g(\mathbf{z} | \Phi)$  is constructed by a deep neural network (DNN) and its parameter  $\Phi$  is essentially weights of such DNN. When  $g(\mathbf{z} | \Phi)$  is constructed by DNN,  $g(\mathbf{z} | \Phi)$  is not totally equal to  $\mathbf{x}$  as  $g(\mathbf{z} | \Phi) = \mathbf{x}$  but it is expected that  $g(\mathbf{z} | \Phi)$  is approximated to  $\mathbf{x}$  in practice:

$$g(\mathbf{z} | \Phi) = \mathbf{x}' \cong \mathbf{x}$$

Note that there are many DGMs and some of them do not require explicit definition of the PDF  $P(\mathbf{z})$  of tractable data  $\mathbf{z}$  but how to estimate generator parameter  $\Phi$  for determining generator  $g(\mathbf{z} | \Phi) = \mathbf{x}$  is always concerned. When  $g(\mathbf{z} | \Phi)$  was determined, we can easily randomize some random tractable data  $\mathbf{z}'$  according to the known tractable PDF  $P(\mathbf{z})$  and then it is totally possible to generate new artifact  $\mathbf{x}'$  by  $\mathbf{x}' = g(\mathbf{z}' | \Phi)$  so that  $\mathbf{x}'$  is generated data / derived data of original intractable data  $\mathbf{x}$  with expectation that probability distribution of  $\mathbf{x}'$  is approximate to the true distribution  $P(\mathbf{x})$  of  $\mathbf{x}$ . The process to randomize  $\mathbf{z}'$  is called sampling tractable data ( $\mathbf{z}$ ). When  $g(\mathbf{z} | \Phi)$  is modeled by a DNN, how to estimate parameter  $\Phi$  is essentially to train such DNN and hence, the DNN is denoted as generator function  $g(\mathbf{z} | \Phi)$  for a convention, which is called generator DNN  $g(\mathbf{z} | \Phi)$ . Here we identify generator function with DNN.

Intractable PDF  $P(\mathbf{x})$  of  $\mathbf{x}$  is specified (Ruthotto & Haber, 2021, p. 3) based on law of total probability as follows:

$$P(\mathbf{x}) = \int_{\mathbf{z}} P_g(\mathbf{x} | \mathbf{z}) P(\mathbf{z}) d\mathbf{z}$$

Where  $P_g(\mathbf{x} | \mathbf{z})$  is conditional PDF of  $\mathbf{x}$  given  $\mathbf{z}$ , which implies that  $P_g(\mathbf{x} | \mathbf{z})$  depends on generator  $g(\mathbf{z} | \Phi)$  too because random variable  $\mathbf{x}$  inside condition PDF  $P_g(\mathbf{x} | \mathbf{z})$  is generated from  $\mathbf{z}$ . Note, notation  $P(\cdot)$  denotes probability distribution or probability density function (PDF) in this research. Therefore, it is possible to denote such conditional PDF as  $P(\mathbf{x} | \Phi, \mathbf{z})$ .

$$P_g(\mathbf{x} | \mathbf{z}) = P(\mathbf{x} | \Phi, \mathbf{z})$$

Such that:

$$P(\mathbf{x}) = \int_{\mathbf{z}} P(\mathbf{x} | \Phi, \mathbf{z}) P(\mathbf{z}) d\mathbf{z} \quad (1.1)$$

This implies intractable PDF  $P(\mathbf{x})$  can be known via tractable PDF  $P(\mathbf{z})$  and conditional PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$ ; however it is really difficult to compute  $P(\mathbf{x})$  due to complication of the integral but this difficulty is unimportant because the purpose of DGM is to estimate generator  $g(\mathbf{z} | \Phi)$ . As a convention, the conditional PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$  is called likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$ . Indeed,  $P(\mathbf{x} | \Phi, \mathbf{z})$  is likelihood function of intractable data  $\mathbf{x}$  given tractable data  $\mathbf{z}$ , which indicates how close generated data  $\mathbf{x}' = g(\mathbf{z} | \Phi)$  to  $\mathbf{x}$ . In practice tractable PDF  $P(\mathbf{z})$  is predefined and likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is determined based on generator DNN  $g(\mathbf{z} | \Phi)$ . For instance,  $P(\mathbf{x} | \Phi, \mathbf{z})$  is assumed to be normal distribution (Gaussian distribution) with mean  $\mu$  and variance  $\sigma^2$  in popular as follows (Ruthotto & Haber, 2021, p. 3):

$$P(\mathbf{x} | \Phi, \mathbf{z}) = (2\pi\sigma^2)^{-\frac{m}{2}} \exp\left(-\frac{\|g(\mathbf{z} | \Phi) - \mathbf{x} - \mu\|^2}{2\sigma^2}\right)$$

Let  $\mu = \mathbf{0}$  and  $\sigma^2=1$  for optimization, we have:

$$P(\mathbf{x} | \Phi, \mathbf{z}) = (2\pi)^{-\frac{m}{2}} \exp\left(-\frac{\|g(\mathbf{z} | \Phi) - \mathbf{x} - \mu\|^2}{2}\right)$$

Where notation  $\|\cdot\|$  denotes norm of vector. For instance, Euclidean norm of intractable data  $\mathbf{x}$  is:

$$\|\mathbf{x}\| = \sqrt{x_1^2 + x_2^2 + \dots + x_m^2}$$

That tractable PDF  $P(\mathbf{z})$  is predefined (constant with regard to  $\Phi$  and  $\mathbf{x}$ ) and likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is assumed to distribute normally indicates that intractable PDF  $P(\mathbf{x})$  is implied by the simpler conditional PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$  with support of generator DNN  $g(\mathbf{z} | \Phi)$ ; in other words,  $P(\mathbf{x} | \Phi, \mathbf{z})$  is probability distribution of  $\mathbf{x}$  from viewpoint of DNN  $g(\mathbf{z} | \Phi)$  indeed where  $\mathbf{z}$  is totally determined in latent space  $Z$  and  $P(\mathbf{x} | \Phi, \mathbf{z})$  is really simpler with support of DNN  $g(\mathbf{z} | \Phi)$ . Because how to determine generator  $g(\mathbf{z} | \Phi)$  is to estimate parameter  $\Phi$ , it is easy to calculate  $\Phi$  as maximizer of likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$ , which is the optimization problem as follows:

$$\Phi^* = \underset{\Phi}{\operatorname{argmax}} P(\mathbf{x} | \Phi, \mathbf{z}) \quad (1.2)$$

Taking natural logarithm of likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  aims to easily determine  $\Phi$  by maximizing the log-likelihood function  $\log(P(\mathbf{x} | \Phi, \mathbf{z}))$  as follows:

$$\Phi^* = \underset{\Phi}{\operatorname{argmax}} P(\mathbf{x} | \Phi, \mathbf{z}) = \underset{\Phi}{\operatorname{argmax}} (\log P(\mathbf{x} | \Phi, \mathbf{z})) = \underset{\Phi}{\operatorname{argmax}} \left(-\frac{1}{2} \|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2\right)$$

Note,

$$\log P(\mathbf{x} | \Phi, \mathbf{z}) = -\frac{m}{2} \log(2\pi\sigma^2) - \frac{\|g(\mathbf{z} | \Phi) - \mathbf{x} - \mu\|^2}{2\sigma^2}$$

Let  $\mu = \mathbf{0}$  and  $\sigma^2=1$  for optimization, we have:

$$\log P(\mathbf{x} | \Phi, \mathbf{z}) = -\frac{m}{2} \log(2\pi) - \frac{\|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2}{2}$$

This implies the minimization problem as follows:

$$\Phi^* = \underset{\Phi}{\operatorname{armin}} \frac{1}{2} \|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2$$

As a result, the estimation of generator parameter  $\Phi$  based on maximum likelihood estimation (MLE) with assumption of normal distribution of generator  $g(\mathbf{z} | \Phi)$  turns back minimization of error function  $\frac{1}{2}\|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2$  which is popular technique in learning DNN by backpropagation algorithm because  $\frac{1}{2}\|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2$  is, indeed, quadratic error function in neural network where  $g(\mathbf{z} | \Phi)$  and  $\mathbf{x}$  are output and real output of a neuron, respectively with note that the output  $g(\mathbf{z} | \Phi)$  is calculated by propagation rule and the real output  $\mathbf{x}$  is from training data. In other words, MLE is entry point to estimate generator parameter  $\Phi$  which is weights of DNN  $g(\mathbf{z} | \Phi)$  that is learned fully by backpropagation algorithm (Nguyen, 2023, pp. 8-20). Therefore, please pay attention to the association of MLE and backpropagation algorithm for determining totally generator  $g(\mathbf{z} | \Phi)$ , in which  $g(\mathbf{z} | \Phi)$  and  $\mathbf{x}$  are output at real output of neurons at the output layer of DNN so that backpropagation algorithm can be applied successively. Note, *error function* is also called *loss function*. Backpropagation algorithm is often associated with stochastic gradient descent (SGD) method to optimize loss function.

Let  $\nabla P(\mathbf{x} | \Phi, \mathbf{z})$  be gradient of *likelihood*  $P(\mathbf{x} | \Phi, \mathbf{z})$  which is first-order derivative of  $P(\mathbf{x} | \Phi, \mathbf{z})$  with regard to parameter  $\Phi$  where  $\mathbf{x}$  and  $\mathbf{z}$  are samples as follows:

$$\nabla P(\mathbf{x} | \Phi, \mathbf{z}) = \frac{dP(\mathbf{x} | \Phi, \mathbf{z})}{d\Phi}$$

*Stochastic gradient descent (SGD)* method (Nguyen, 2023, pp. 22-27) estimates  $\Phi$  by iterative process to update successively  $\Phi$  at every iteration as follows:

$$\Phi = \Phi + \gamma \nabla P(\mathbf{x} | \Phi, \mathbf{z})$$

Where  $0 < \gamma \leq 1$  is learning rate. SGD, which is an iterative process, pushes candidate solution at each iteration along the direction which is opposite to gradient of target function for minimization or has the same direction to gradient of target function for maximization with note that the step length is represented by learning rate. In practice, likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is replaced by its natural logarithm as follows:

$$\Phi = \Phi + \gamma \nabla \log P(\mathbf{x} | \Phi, \mathbf{z}) \quad (1.3)$$

Where  $\nabla P(\mathbf{x} | \Phi, \mathbf{z})$  is gradient of the log-likelihood function  $\log P(\mathbf{x} | \Phi, \mathbf{z})$  with regard to  $\Phi$ . Note the estimation equation above mentions maximization problem according to MLE method and hence, if error function denoted  $\varepsilon(\mathbf{x} | \Phi, \mathbf{z})$  which is the function related to  $g(\mathbf{z} | \Phi)$  and  $\mathbf{x}$  like  $\varepsilon(\mathbf{x} | \Phi, \mathbf{z}) = \frac{1}{2}\|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2$  aforementioned, then SGD modifies a little bit the estimation equation as follows:

$$\Phi = \Phi - \gamma \nabla \varepsilon(\mathbf{x} | \Phi, \mathbf{z}) \quad (1.4)$$

Where  $\nabla \varepsilon(\mathbf{x} | \Phi, \mathbf{z})$  is gradient of error  $\varepsilon(\mathbf{x} | \Phi, \mathbf{z})$ .

$$\nabla \varepsilon(\mathbf{x} | \Phi, \mathbf{z}) = \frac{d\varepsilon(\mathbf{x} | \Phi, \mathbf{z})}{d\Phi}$$

The main difference between maximizing likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  and minimizing error  $\varepsilon(\mathbf{x} | \Phi, \mathbf{z})$  is the changing from the sign “+” regarding maximization problem to the sign “-” regarding minimization problem, which is the essence of gradient descent method. In the example of assuming normal distribution, likelihood maximization is the same to error minimization but likelihood maximization gives broader applications to estimate generator parameter  $\Phi$  within context of DNN along with *backpropagation algorithm* to train DNN. Besides, it is possible to consider error function is the minus opposite of likelihood function:

$$\varepsilon(\mathbf{x} | \Phi, \mathbf{z}) = -P(\mathbf{x} | \Phi, \mathbf{z})$$

It is better that error is the minus opposite of log-likelihood function:

$$\varepsilon(\mathbf{x} | \Phi, \mathbf{z}) = -\log P(\mathbf{x} | \Phi, \mathbf{z})$$

Moreover, there many ways to define likelihood and error and so, the way to define them will contribute to form a concrete DGM, besides how to specify and design generator DNN  $g(\mathbf{z} | \Phi)$ . When  $\Phi$  is weight vector consisting of many weights of entire DNN, only elemental sub-

weights at the output layer are estimated by SGD which maximizes likelihood or minimizes error:

$$\Phi_{\text{output}} = \Phi_{\text{output}} + \gamma \nabla P(\mathbf{x}|\Phi, \mathbf{z})$$

Or

$$\Phi_{\text{output}} = \Phi_{\text{output}} - \gamma \nabla \varepsilon(\mathbf{x}|\Phi, \mathbf{z})$$

Then backpropagation algorithm continues to update remaining sub-weights at hidden layers based on such determined sub-weights at the output layer. Therefore, for convenience, we only focus on likelihood maximization or error minimization and parameter  $\Phi$  represents entire weights of DNN with assertion that backpropagation algorithm is always feasible. It means that there are two important equivalent estimation equations as follows:

$$\Phi = \Phi + \gamma \nabla P(\mathbf{x}|\Phi, \mathbf{z})$$

$$\Phi = \Phi - \gamma \nabla \varepsilon(\mathbf{x}|\Phi, \mathbf{z})$$

In similar to:

$$\Phi = \Phi + \gamma \nabla \log P(\mathbf{x}|\Phi, \mathbf{z})$$

When DGM is trained with big data, training data is fed to DGM at very time point  $i$  as a pair  $\mathbf{d}^{(i)} = (\mathbf{x}^{(i)}, \mathbf{z}^{(i)})$  and therefore, a set of pairs over  $N$  time points is called epoch. As a convention, epoch of size  $N$  is denoted as  $\mathbf{D} = (\mathbf{d}^{(1)} = (\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), \mathbf{d}^{(2)} = (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, \mathbf{d}^{(N)} = (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$ . An interesting result from SGD is that DGM can be learned with epoch  $\mathbf{D}$  without significant change as follows:

$$\Phi = \Phi + \gamma \nabla \frac{1}{N} \sum_{i=1}^N P(\mathbf{x}^{(i)}|\Phi, \mathbf{z}^{(i)}) = \Phi + \gamma \frac{1}{N} \sum_{i=1}^N \nabla P(\mathbf{x}^{(i)}|\Phi, \mathbf{z}^{(i)})$$

Therefore, training data is counted according to every epoch  $\mathbf{D}$  instead of every pair  $(\mathbf{x}, \mathbf{z})$  so that  $\mathbf{D}$  is fed to SGD at every time point  $k$ . Moreover, it is essential that SGD aims to update current parameter at current iteration based on previous parameter at previous iteration. Exactly, let  $\Phi^{(k+1)}$  be generator parameter at the  $(k+1)^{\text{th}}$  iteration, then  $\Phi^{(k+1)}$  is calculated based on previous generator parameter  $\Phi^{(k)}$  at the  $k^{\text{th}}$  iteration as follows:

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma \nabla \frac{1}{N} \sum_{i=1}^N P(\mathbf{x}^{(i)}|\Phi^{(k)}, \mathbf{z}^{(i)})$$

The equation above is the most precise equation for parameter estimation with SGD, which is called epoch estimation with note that SGD is an iterative process. It can also be replaced by following equations:

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma \nabla \frac{1}{N} \sum_{i=1}^N \log P(\mathbf{x}^{(i)}|\Phi^{(k)}, \mathbf{z}^{(i)})$$

$$\Phi^{(k+1)} = \Phi^{(k)} - \gamma \nabla \frac{1}{N} \sum_{i=1}^N \varepsilon(\mathbf{x}^{(i)}|\Phi^{(k)}, \mathbf{z}^{(i)})$$

The first equation

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma \nabla \frac{1}{N} \sum_{i=1}^N \log P(\mathbf{x}^{(i)}|\Phi^{(k)}, \mathbf{z}^{(i)}) \quad (1.5)$$

Which is most popular. Moreover, the index  $k$  indicates time point as well as iteration of SGD. If learning rate  $\gamma$  is varied at every iteration as  $\gamma^{(k)}$ , we have:

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma^{(k)} \nabla \frac{1}{N} \sum_{i=1}^N \log P(\mathbf{x}^{(i)}|\Phi^{(k)}, \mathbf{z}^{(i)})$$

There are two problems related to construct a DGM: 1) how to define likelihood or error to train generator DNN  $g(\mathbf{z} | \Phi)$  and 2) how to define tractable PDF  $P(\mathbf{z})$  which implies the way



to randomize  $\mathbf{z}$ . The second problem relates to assert qualification of random data  $\mathbf{z}'$  and hence, the second problem is stated as qualification problem of how to qualify random data. Therefore, the two problems of constructing DGM are 1) how to train generator DNN  $g(\mathbf{z} | \Phi)$  and 2) how to qualify such training task which often relates to another optimization task or another training task. Some basic principles related to DGM are introduced in this section but the two problems cannot be mentioned because there are many specific DGMs which have own specifications. Anyhow generator likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  based on definition of generator  $g(\mathbf{z} | \Phi)$  is always important regardless that if it is specified explicitly and thus, suppose it was defined, then SGD is favorite method to optimize it. As an example aforementioned, suppose  $P(\mathbf{x} | \Phi, \mathbf{z})$  distributes normally with mean  $\boldsymbol{\mu}$  and variance  $\sigma^2$  in some DGM as follows (Ruthotto & Haber, 2021, p. 3):

$$P(\mathbf{x} | \Phi, \mathbf{z}) = (2\pi\sigma^2)^{-\frac{m}{2}} \exp\left(-\frac{\|g(\mathbf{z} | \Phi) - \mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}\right)$$

Generator log-likelihood is natural logarithm of generator likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$ :

$$\log P(\mathbf{x} | \Phi, \mathbf{z}) = -\frac{m}{2} \log(2\pi\sigma^2) - \frac{\|g(\mathbf{z} | \Phi) - \mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}$$

Gradient of this log-likelihood with regard to  $\Phi$  is:

$$\nabla \log P(\mathbf{x} | \Phi, \mathbf{z}) = \frac{d \log P(\mathbf{x} | \Phi, \mathbf{z})}{d\Phi} = -\frac{(g(\mathbf{z} | \Phi) - \mathbf{x} - \boldsymbol{\mu})}{\sigma^2} \frac{dg(\mathbf{z} | \Phi)}{d\Phi}$$

Where  $dg(\mathbf{z} | \Phi) / d\Phi$  is differential of  $g(\mathbf{z} | \Phi)$  with regard to  $\Phi$ . Let  $\boldsymbol{\mu} = \mathbf{0}$  and  $\sigma^2=1$  for optimization, we have:

$$\nabla \log P(\mathbf{x} | \Phi, \mathbf{z}) = \frac{d \log P(\mathbf{x} | \Phi, \mathbf{z})}{d\Phi} = -(g(\mathbf{z} | \Phi) - \mathbf{x}) \frac{dg(\mathbf{z} | \Phi)}{d\Phi}$$

As usual, estimation equation resulted from SGD is:

$$\Phi = \Phi + \gamma \nabla \log P(\mathbf{x} | \Phi, \mathbf{z}) = \Phi + \gamma (\mathbf{x} - g(\mathbf{z} | \Phi)) \frac{dg(\mathbf{z} | \Phi)}{d\Phi}$$

There is a question that how to calculate the differential  $dg(\mathbf{z} | \Phi) / d\Phi$ . Indeed, it is not difficult to calculate it in context of neural network associated with backpropagation algorithm so that the last output layer as well as last neuron  $o$  of DNN  $g(\mathbf{z} | \Phi)$  is acted by activation function  $a(\cdot)$  as follows:

$$a(o) = a(\mathbf{w}^T \mathbf{i})$$

$$o = \mathbf{w}^T \mathbf{i}$$

Where  $\mathbf{i}$  is input of the last layer  $o$  and weight parameter  $\mathbf{w}$  is a part of entire parameter  $\Phi$  and hence, we need to focus on calculating differential  $da(o) / d\mathbf{w}$  which is equivalent to differential  $dg(\mathbf{z} | \Phi) / d\Phi$  so that backpropagation algorithm will solve remaining parts of entire parameter  $\Phi$ .

$$\frac{da(o)}{d\mathbf{w}} \cong \frac{dg(\mathbf{z} | \Phi)}{d\Phi}$$

Indeed, we have:

$$\frac{da(o)}{d\mathbf{w}} = \frac{da(\mathbf{w}^T \mathbf{i})}{d\mathbf{w}} = \frac{da(o)}{do} \mathbf{i}^T = a'(o) \mathbf{i}^T$$

Note, the subscript “ $T$ ” denotes transposition operator of vector and matrix in which row vector becomes column vector and vice versa. It is easy to calculate the derivative  $a'(o)$  when activation function was specified, for instance, if  $a(o)$  is sigmoid function, we have:

$$a(o) = y = \frac{1}{1 + e^{-o}}$$

$$a'(o) = \frac{1}{1 + e^{-o}} \left(1 - \frac{1}{1 + e^{-o}}\right) = a(o)(1 - a(o)) = y(1 - y)$$

In practice,  $y$  is replaced by  $a(y)$  in order to prevent  $o$  from being out of space:

$$a'(o) \cong a(y)(1 - a(y)) = a(a(o))(1 - a(a(o)))$$

As a result, we have:

$$\frac{dg(\mathbf{z}|\Phi)}{d\Phi} \cong a(a(o))(1 - a(a(o))) \mathbf{i}^T$$

For fast computation, it is possible to set the derivative  $a'(o)$  to be small enough constants like 1 such that  $dg(\mathbf{z}|\Phi) / d\Phi = \mathbf{i}^T$ .

Suppose some other DGM assumes that  $\mathbf{x}$  is binary ( $\mathbf{x} = 0$  or  $\mathbf{x} = 1$ ) and follows Bernoulli (Ruthotto & Haber, 2021, p. 3) and so, its generator DNN  $g(\mathbf{z}|\Phi)$  derives values in interval  $[0, 1]$ . In other words, image of  $g(\mathbf{z}|\Phi)$  is the real number interval  $[0, 1]$ , which leads to a specification that  $g(\mathbf{z}|\Phi)$  is probability of the event  $\mathbf{x}=1$  with note that  $\mathbf{x}$  is scalar variable ( $x$ ) for convenience:

$$g(\mathbf{z}|\Phi) = P(\mathbf{x} = 1)$$

Because  $g(\mathbf{z}|\Phi)$  becomes a (scalar) random variable whose value is probability, it is possible to identify  $g(\mathbf{z}|\Phi)$  with its parameter  $\Phi$  as a convention:

$$\Phi = g(\mathbf{z}|\Phi) = P(\mathbf{x} = 1)$$

Given  $N$  trials with binary values of  $\mathbf{x}$ , let  $N(\mathbf{x})$  be the number of event  $\mathbf{x}=1$  among  $N$  trials, then generator likelihood  $P(\mathbf{x}|\Phi, \mathbf{z})$  is specified according to Bernoulli distribution as follows:

$$P(\mathbf{x}|\Phi, \mathbf{z}) = \Phi^{N(\mathbf{x})} (1 - \Phi)^{N - N(\mathbf{x})}$$

The generator log-likelihood is:

$$\log P(\mathbf{x}|\Phi, \mathbf{z}) = N(\mathbf{x}) \log \Phi + (N - N(\mathbf{x})) \log(1 - \Phi)$$

Gradient of the generator log-likelihood with regard to  $\Phi$  is:

$$\nabla \log P(\mathbf{x}|\Phi, \mathbf{z}) = \frac{d \log P(\mathbf{x}|\Phi, \mathbf{z})}{d\Phi} = \frac{N(\mathbf{x})}{\Phi} - \frac{N - N(\mathbf{x})}{1 - \Phi} = \frac{N(\mathbf{x}) - N\Phi}{\Phi(1 - \Phi)}$$

As a result, estimation equation resulted from SGD is:

$$\Phi = \Phi + \gamma \nabla \log P(\mathbf{x}|\Phi, \mathbf{z}) = \Phi + \gamma \frac{N(\mathbf{x}) - N\Phi}{\Phi(1 - \Phi)}$$

Although normal distribution and Bernoulli distribution are two popular distributions to specify generator likelihood  $P(\mathbf{x}|\Phi, \mathbf{z})$ , there are other specifications which depend on specific DGM.

Given epoch  $\mathbf{D} = (\mathbf{d}^{(1)} = (\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), \mathbf{d}^{(2)} = (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, \mathbf{d}^{(N)} = (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$  implies that the epoch is created or sent by equilateral distribution  $1/N$  but in general case,  $\mathbf{D}$  can follow an arbitrary distribution denoted by PDF  $P(\mathbf{d})$ , which makes the optimization problem and the SGD estimation changed a little bit by theoretical expectation given distribution  $P(\mathbf{d})$ .

$$\Phi^* = \underset{\Phi}{\operatorname{armax}} E(P(\mathbf{x}|\Phi, \mathbf{z})|P(\mathbf{d}))$$

$$\Phi = \Phi + \gamma \nabla E(P(\mathbf{x}|\Phi, \mathbf{z})|P(\mathbf{d}))$$

Where,

$$E(P(\mathbf{x}|\Phi, \mathbf{z})|P(\mathbf{d})) = \int_{\mathbf{d}} P(\mathbf{x}|\Phi, \mathbf{z}) P(\mathbf{d}) d\mathbf{d}$$

However, there is no significant change in aforementioned practical technique to estimate parameters.

Turning back the assumption that generator likelihood  $P(\mathbf{x}|\Phi, \mathbf{z})$  distributes normally with mean  $\boldsymbol{\mu}$  and variance  $\sigma^2$  in some DGM as follows (Ruthotto & Haber, 2021, p. 3):

$$P(\mathbf{x}|\Phi, \mathbf{z}) = (2\pi\sigma^2)^{-\frac{m}{2}} \exp\left(-\frac{\|g(\mathbf{z}|\Phi) - \mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}\right)$$

This assumption is not totally exact because the distribution above mentions the error  $g(\mathbf{z}|\Phi) - \mathbf{x}$  between generated data  $g(\mathbf{z}|\Phi)$  and real data  $\mathbf{x}$ . Exactly, generator likelihood  $P(\mathbf{x}|\Phi, \mathbf{z})$

is defined as distribution of the error  $\varepsilon = g(\mathbf{z} | \Phi) - \mathbf{x}$  and such error distribution is assumed to follow normal distribution with mean  $\boldsymbol{\mu}$  and variance  $\sigma^2$ .

$$P(\varepsilon | \boldsymbol{\mu}, \sigma^2) = P(\mathbf{x} | \Phi, \mathbf{z}) = (2\pi\sigma^2)^{-\frac{m}{2}} \exp\left(-\frac{\|g(\mathbf{z} | \Phi) - \mathbf{x} - \boldsymbol{\mu}\|^2}{2\sigma^2}\right)$$

Therefore, setting error mean and error variance to be zero and one as  $\boldsymbol{\mu} = \mathbf{0}$ ,  $\sigma^2 = 1$  is for best optimization because of smallest error mean  $\mathbf{0}$  but the setting is not totally diverse in data generation.

$$P(\varepsilon | \boldsymbol{\mu}, \sigma^2) = P(\mathbf{x} | \Phi, \mathbf{z}) = (2\pi)^{-\frac{m}{2}} \exp\left(-\frac{\|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2}{2}\right)$$

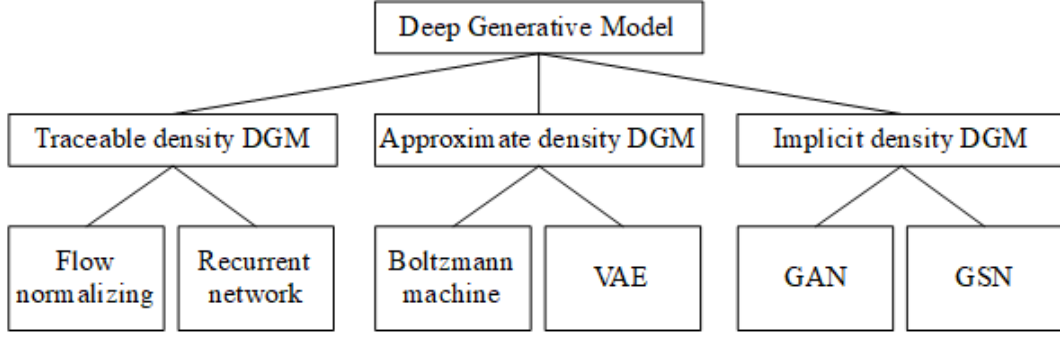
When learning generator DNN by backpropagation algorithm associated with SGD, it is possible to estimate dynamically  $\boldsymbol{\mu}$  and  $\sigma^2$  by maximum likelihood estimation (MLE) method. Given epoch  $\mathbf{D} = (\mathbf{d}^{(1)} = (\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), \mathbf{d}^{(2)} = (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, \mathbf{d}^{(N)} = (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$ , error mean and error variance are estimated as follows:

$$\begin{aligned} \hat{\boldsymbol{\mu}} &= \frac{1}{N} \sum_{i=1}^N (g(\mathbf{z}^{(i)} | \Phi) - \mathbf{x}^{(i)}) \\ \hat{\sigma}^2 &= \frac{1}{N} \sum_{i=1}^N \|g(\mathbf{z}^{(i)} | \Phi) - \mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}}\|^2 \end{aligned} \tag{1.6}$$

When error mean and error variance are dynamically estimated instead of fixing them by zero and unit, generator DNN  $g(\mathbf{z} | \Phi)$  may produce new data in high diversity, which is similar to add noises to generated data. In other words, estimation of error mean and error variance based on epoch makes the data generation more diverse because  $\mathbf{z}$  may be randomized in interval  $[0, 1]$  although DGMs try to diversify  $\mathbf{z}$  or  $\mathbf{x}$  like Variational Autoencoders (VAE) and Generative Adversarial Network (GAN). Note, if  $\mathbf{z}$  is randomized only in interval  $[0, 1]$ , generated data  $\mathbf{x}' = g(\mathbf{z} | \Phi)$  may not be different much from sample  $\mathbf{x}$  in epoch in case that error mean  $\boldsymbol{\mu}$  and error variance  $\sigma^2$  are fixed by  $\mathbf{0}$  and 1. However, quality of data generation is the best with zero error mean  $\mathbf{0}$ .

Recall that the two problems of constructing DGM are 1) how to train generator DNN  $g(\mathbf{z} | \Phi)$  and 2) how to qualify such training task which often relates to another optimization task or another training task. The first problem relates to how to establish generator likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  which is the probability density function (PDF) of intractable  $\mathbf{x}$  given tractable data  $\mathbf{z}$  and this establishment is based on generator DNN  $g(\mathbf{z} | \Phi)$ . However, there are some DGMs do not specify explicitly the density function  $P(\mathbf{x} | \Phi, \mathbf{z})$ , which is cause of the fact that there are two DGM approaches: 1) DGM specifies explicitly generator PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$  and 2) vice versa. In group of explicit PDF approach, there are two built-in approaches: 1) tractable density DGM specifies explicitly well-known distributions for generator likelihood and 2) approximate density DGM tries to estimate approximately generator PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$  or derive other PDF that is similar to  $P(\mathbf{x} | \Phi, \mathbf{z})$ . In general, there are three main approaches for constructing DGM such as tractable density DGM, approximate density DGM, and implicit density DGM which are mentioned in next sections. Following figure depicts taxonomy of DGM (Oussidi & Elhassouny, 2018, p. 7) by Goodfellow.





**Figure 1.1.** Taxonomy of DGM

Especially, if data is image, there is another categorization way that there are two main approaches: 1) pixel density approach tries to model pixel distribution and 2) block density approach tries to model entire image distribution as any data distribution. In other words, likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is defined based on probabilistic distribution of pixels where  $\mathbf{x}$  is considered as set of pixels according to pixel density approach. On the other hand, block density approach considers likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is PDF of a block or entire image (unified big block) where  $\mathbf{x}$  is considered as any arbitrary data. As usual, pixel density approach belongs to tractable density approach of the first categorization.

## 2. TRACTABLE DENSITY DGM

According to tractable density approach, DGMs specify explicitly generator PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$  with note that PDF is abbreviation of probability density function. Recall that the two problems of constructing DGM are 1) how to train generator DNN  $g(\mathbf{z} | \Phi)$  and 2) how to qualify such training task which often relates to another optimization task or another training task. However, the two problems are merged into the first problem which is to train  $g(\mathbf{z} | \Phi)$  according to normalizing flow technique in which  $g(\mathbf{z} | \Phi)$  is invertible given tractable data  $\mathbf{z}$  and intractable data  $\mathbf{x}$  have the same dimension  $n$ . Therefore, latent space  $Z$  and sample space  $X$  are the same with dimension  $n$ . As a convention,  $g^{-1}(\mathbf{x} | \Phi)$  is called inversed generator.

$$\begin{aligned} \mathbf{x} &= g(\mathbf{z} | \Phi) \\ \mathbf{z} &= g^{-1}(\mathbf{x} | \Phi) \end{aligned} \quad (2.1)$$

Where,

$$\mathbf{z} = (z_1, z_2, \dots, z_n)^T, \mathbf{x} = (x_1, x_2, \dots, x_n)^T$$

$$\mathbf{x} \in X \subseteq \mathbb{R}^n$$

$$\mathbf{z} \in Z \subseteq \mathbb{R}^n$$

$$X = Z$$

Note, the subscript “ $T$ ” denotes transposition operator of vector and matrix in which row vector becomes column vector and vice versa. Because  $g(\mathbf{z} | \Phi)$  is essentially vector-by-vector function whose input and output are vectors, it should have denoted as  $\mathbf{g}(\mathbf{z} | \Phi)$ , especially, when  $\mathbf{g}(\mathbf{z} | \Phi)$  here is bijection. However, it is still denoted  $g(\mathbf{z} | \Phi)$  for convenience. Therefore, the first-order derivative of vector-by-vector function  $g(\mathbf{z} | \Phi)$  here is Jacobian matrix but is still called gradient. Note, in mathematical, the first-order derivative of scalar-by-vector function is called gradient vector and the first-order derivative of vector-by-vector function is called Jacobian matrix.

As a result, **normalizing flow (NL)** technique focuses on maximizing intractable PDF  $P(\mathbf{x})$  now called sample PDF or sample likelihood rather than maximizing generator likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  because  $P(\mathbf{x})$  is now proportional to tractable PDF  $P(\mathbf{z})$  and  $P(\mathbf{x})$  is stronger than  $P(\mathbf{x} | \Phi, \mathbf{z})$ . When  $P(\mathbf{x})$  has generator parameter  $\Phi$ , it is denoted as  $P(\mathbf{x} | \Phi)$ . According to applied statistics literature, sample likelihood  $P(\mathbf{x} | \Phi)$  is determined based on tractable PDF  $P(\mathbf{z})$  and generator  $g(\mathbf{z} | \Phi)$  as follows (Ruthotto & Haber, 2021, p. 7):

$$P(\mathbf{x}|\Phi) = P(\mathbf{z})|\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)| = P(g^{-1}(\mathbf{x}|\Phi))|\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)| \quad (2.2)$$

Where  $|\cdot|$  or  $\det(\cdot)$  denotes determinant of square matrix with note that the gradient  $\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)$  of the inverse  $g^{-1}(\mathbf{x}|\Phi)$  is Jacobian matrix which is the first-order derivative of  $g^{-1}(\mathbf{x}|\Phi)$  with regard to  $\mathbf{x}$ . As a convention,  $\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)$  is called inversed gradient because  $\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)$  is the first-order derivative of inversed generator  $g^{-1}(\mathbf{x}|\Phi)$  with regard to  $\mathbf{x}$ .

$$\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi) = \frac{dg^{-1}(\mathbf{x}|\Phi)}{d\mathbf{x}}$$

The equation of sample likelihood  $P(\mathbf{x}|\Phi)$  is much more definite than the integral formulation of  $P(\mathbf{x})$  as aforementioned

$$P(\mathbf{x}) = \int_{\mathbf{z}} P(\mathbf{x}|\Phi, \mathbf{z})P(\mathbf{z}) d\mathbf{z}$$

It is explained from the equation of sample likelihood  $P(\mathbf{x}|\Phi)$  that given source and target with a function from source to target, target distribution is calculated by multiplying source distribution with determinant of gradient of inversed function.

$$P(\mathbf{x}|\Phi) = P(g^{-1}(\mathbf{x}|\Phi))|\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)|$$

For optimization,  $P(\mathbf{z})$  is assumed to follow standard normal distribution with mean 0 and variance 1:

$$P(\mathbf{z}) = (2\pi)^{-\frac{n}{2}} \exp\left(-\frac{1}{2}\|\mathbf{z}\|^2\right)$$

Such that:

$$P(\mathbf{x}|\Phi) = (2\pi)^{-\frac{n}{2}} \exp\left(-\frac{1}{2}\|g^{-1}(\mathbf{x}|\Phi)\|^2\right)|\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)|$$

Where notation  $\|\cdot\|$  denotes norm of vector. Exactly,  $P(\mathbf{z})$  follows standard normal distribution with mean vector  $\mathbf{0}$  and identity covariance matrix  $I$ . Sample log-likelihood is derived by taking natural logarithm of sample likelihood:

$$\log P(\mathbf{x}|\Phi) = -\frac{n}{2} \log 2\pi - \frac{1}{2}\|g^{-1}(\mathbf{x}|\Phi)\|^2 + \log|\nabla_{\mathbf{x}}g^{-1}(\mathbf{x}|\Phi)|$$

NL aims to maximize sample log-likelihood so as to estimate generator parameter  $\Phi$ :

$$\Phi^* = \underset{\Phi}{\operatorname{argmax}} \log P(\mathbf{x}|\Phi)$$

Stochastic gradient descent (SGD) method is used to estimate  $\Phi$  by iterative process to update successively  $\Phi$  at every iteration as follows:

$$\Phi = \Phi + \gamma \nabla \log P(\mathbf{x}|\Phi)$$

Where  $\nabla \log P(\mathbf{x}|\Phi)$  which is called sample log-likelihood gradient is gradient of sample log-likelihood  $\log P(\mathbf{x}|\Phi)$  and  $\gamma$  ( $0 < \gamma \leq 1$ ) is learning rate. Note that SGD, which is an iterative process, pushes candidate solution at each iteration along the direction which is opposite to gradient of target function for minimization or has the same direction to gradient of target function for maximization with note that the step length is represented by learning rate. Given epoch of size  $N$  is denoted as  $\mathbf{D} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$ , the estimation equation of  $\Phi$  is extended exactly as epoch estimation at every iteration of SGD:

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma \nabla \frac{1}{N} \sum_{i=1}^N \log P(\mathbf{x}^{(i)}|\Phi^{(k)})$$

It is necessary to determine sample log-likelihood gradient  $\nabla \log P(\mathbf{x}|\Phi)$  with regard to parameter  $\Phi$ . Due to (Nguyen, Matrix Analysis and Calculus, 2015, pp. 45-46):

$$\begin{aligned}
\frac{d \log(|\nabla_x g^{-1}(\mathbf{x}|\Phi)|)}{d\Phi} &= \frac{d \log(|\nabla_x g^{-1}(\mathbf{x}|\Phi)|)}{d\nabla_x g^{-1}(\mathbf{x}|\Phi)} \frac{d\nabla_x g^{-1}(\mathbf{x}|\Phi)}{d\Phi} \\
&= \frac{d|\nabla_x g^{-1}(\mathbf{x}|\Phi)|}{d\nabla_x g^{-1}(\mathbf{x}|\Phi)} \frac{1}{|\nabla_x g^{-1}(\mathbf{x}|\Phi)|} \frac{d\nabla_x g^{-1}(\mathbf{x}|\Phi)}{d\Phi} \\
&= \frac{|\nabla_x g^{-1}(\mathbf{x}|\Phi)| (\nabla_x g^{-1}(\mathbf{x}|\Phi))^{-1}}{|\nabla_x g^{-1}(\mathbf{x}|\Phi)|} \frac{d\nabla_x g^{-1}(\mathbf{x}|\Phi)}{d\Phi} \\
&= (\nabla_x g^{-1}(\mathbf{x}|\Phi))^{-1} \frac{d\nabla_x g^{-1}(\mathbf{x}|\Phi)}{d\Phi}
\end{aligned}$$

And

$$\frac{d \frac{1}{2} \|g^{-1}(\mathbf{x}|\Phi)\|^2}{d\Phi} = g^{-1}(\mathbf{x}|\Phi) \nabla_x g^{-1}(\mathbf{x}|\Phi)$$

We have following equation to calculate log-likelihood gradient  $\nabla \log P(\mathbf{x} | \Phi)$ :

$$\nabla \log P(\mathbf{x}|\Phi) = \frac{d \log P(\mathbf{x}|\Phi)}{d\Phi} = (\nabla_x g^{-1}(\mathbf{x}|\Phi))^{-1} \frac{d\nabla_x g^{-1}(\mathbf{x}|\Phi)}{d\Phi} - g^{-1}(\mathbf{x}|\Phi) \nabla_x g^{-1}(\mathbf{x}|\Phi)$$

The notation  $(\nabla_x g^{-1}(\mathbf{x} | \Phi))^{-1}$  denotes inverse of matrix  $\nabla_x g^{-1}(\mathbf{x} | \Phi)$ . Because  $\nabla_x g^{-1}(\mathbf{x} | \Phi)$  called Jacobian matrix is a square matrix, the derivative  $d\nabla_x g^{-1}(\mathbf{x} | \Phi) / d\Phi$  is determined by taking first-order derivative for every element of  $\nabla_x g^{-1}(\mathbf{x} | \Phi)$  with regard to  $\Phi$ , which produces a tensor. Therefore,  $d\nabla_x g^{-1}(\mathbf{x} | \Phi) / d\Phi$  is the second-order derivative of inversed generator  $g^{-1}(\mathbf{x} | \Phi)$  with regard to  $\mathbf{x}$  and  $\Phi$ . Let

$$\nabla_{\Phi} \nabla_x g^{-1}(\mathbf{x}|\Phi) = \frac{d\nabla_x g^{-1}(\mathbf{x}|\Phi)}{d\Phi}$$

It is possible to calculate this second-order derivative if inversed gradient  $\nabla_x g^{-1}(\mathbf{x} | \Phi)$  is determined. Log-likelihood gradient  $\nabla \log P(\mathbf{x} | \Phi)$  is rewritten:

$$\nabla \log P(\mathbf{x}|\Phi) = (\nabla_x g^{-1}(\mathbf{x}|\Phi))^{-1} \nabla_{\Phi} \nabla_x g^{-1}(\mathbf{x}|\Phi) - g^{-1}(\mathbf{x}|\Phi) \nabla_x g^{-1}(\mathbf{x}|\Phi) \quad (2.3)$$

Where:

$$\begin{aligned}
\nabla_x g^{-1}(\mathbf{x}|\Phi) &= \frac{d g^{-1}(\mathbf{x}|\Phi)}{d\mathbf{x}} \\
\nabla_{\Phi} \nabla_x g^{-1}(\mathbf{x}|\Phi) &= \frac{d \nabla_x g^{-1}(\mathbf{x}|\Phi)}{d\Phi}
\end{aligned}$$

According to traditional neural network, let  $\varphi_i$  be the  $i^{\text{th}}$  row vector of matrix  $\Phi$ , then generator  $g(\mathbf{x} | \varphi)$  is linear composition as follows:

$$g(\mathbf{z}|\varphi) = a(\varphi_i^T \mathbf{z} + \delta_i) = x_i$$

Where  $\delta_i$  is the  $i^{\text{th}}$  bias parameter associated with each  $x_i$ . Note,  $x_i$  is the  $i^{\text{th}}$  elemental variable in  $\mathbf{x}$  whereas activation  $a(\cdot)$  is invertible, whose inverse is  $a^{-1}(\cdot)$ . In traditional neural network,  $x_i$  represents a neuron or unit. Due to:

$$a^{-1}(x_i) = \varphi_i^T \mathbf{z} + \delta_i = \sum_{t=1}^n \varphi_{it} z_t + \delta_i = \varphi_{ij} z_j + \sum_{t=1, t \neq j}^n \varphi_{it} z_t + \delta_i$$

When  $\varphi_i = (\varphi_{i1}, \varphi_{i2}, \dots, \varphi_{in})^T$  and  $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$ , without loss of generality, given  $\varphi_{ij}$  and  $z_j$  are the  $j^{\text{th}}$  elements of  $\varphi_i$  and  $\mathbf{z}$ , respectively we have fine-tuned inversed generator  $g^{-1}(x_i | \varphi_{ij})$ :

$$g^{-1}(x_i | \varphi_{ij}) = z_j = \frac{a^{-1}(x_i) - b_{ij} - \delta_i}{\varphi_{ij}}$$

Where,

$$b_{ij} = \sum_{t=1, t \neq j}^n \varphi_{it} z_t$$

It is easy to calculate the inversed gradient:

$$\nabla_{x_i} g^{-1}(x_i | \varphi_{ij}) = \frac{d(a^{-1}(x_i) - b_{ij} - \delta_i)/dx_i}{\varphi_{ij}} = \frac{(a^{-1}(x_i))'}{\varphi_{ij}}$$

Where  $a^{-1}(x_i)$  is the first-order derivative of  $a^{-1}(\cdot)$  at  $x_i$ . The second-order derivative  $\nabla_{\varphi_{ij}} \nabla_{x_i} g^{-1}(x_i | \varphi_{ij})$  is determined as follows:

$$\nabla_{\varphi_{ij}} \nabla_{x_i} g^{-1}(x_i | \varphi_{ij}) = \frac{d \nabla_{x_i} g^{-1}(x_i | \varphi_{ij})}{d \varphi_{ij}} = -\frac{(a^{-1}(x_i))'}{\varphi_{ij}^2}$$

Log-likelihood gradient with regard to  $\varphi_{ij}$  is fine-tuned as  $\nabla \log P(x_i | \varphi_{ij})$  is expended again:

$$\begin{aligned} \nabla \log P(x_i | \varphi_{ij}) &= -\frac{\varphi_{ij}}{(a^{-1}(x_i))'} \frac{(a^{-1}(x_i))'}{\varphi_{ij}^2} - \frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}} \\ &= -\frac{1}{\varphi_{ij}} - \frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}^2} = -\frac{\varphi_{ij} + a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}^2} \end{aligned}$$

Because  $\delta_i$  is the  $i^{\text{th}}$  bias parameter, the second-order derivative  $\nabla_{\delta_i} \nabla_{x_i} g^{-1}(x_i | \varphi_{ij})$  is determined as follows:

$$\nabla_{\delta_i} \nabla_{x_i} g^{-1}(x_i | \delta_i) = \frac{d \nabla_{x_i} g^{-1}(x_i | \delta_i)}{d \delta_i} = 0$$

Where,

$$\begin{aligned} \nabla_{x_i} g^{-1}(x_i | \delta_i) &= \frac{d(a^{-1}(x_i) - b_{ij} - \delta_i)/dx_i}{\varphi_{ij}} = \frac{(a^{-1}(x_i))'}{\varphi_{ij}} \\ g^{-1}(x_i | \delta_i) &= z_j = \frac{a^{-1}(x_i) - b_{ij} - \delta_i}{\varphi_{ij}} \end{aligned}$$

Log-likelihood gradient with regard to  $\delta_i$  is fine-tuned as  $\nabla \log P(x_i | \delta_i)$  is expended again:

$$\nabla \log P(x_i | \delta_i) = -\frac{\varphi_{ij}}{(a^{-1}(x_i))'} 0 - \frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}} = -\frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}^2}$$

In general, log-likelihood gradient  $\nabla \log P(x_i | \varphi_{ij}, \delta_i)$  is specified as follows:

$$\begin{aligned} \nabla \log P(x_i | \varphi_{ij}) &= -\frac{1}{\varphi_{ij}} \left( 1 + \frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}} \right), \forall i, j \\ \nabla \log P(x_i | \delta_i) &= -\frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}^2}, \forall i \end{aligned} \tag{2.4}$$

Where  $a(\cdot)$  and  $a^{-1}(\cdot)$  are invertible activation function and its inverse and,

$$x_i \in \mathbf{x} = (x_1, x_2, \dots, x_n)^T$$

$$\varphi_{ij} \in \varphi_i = (\varphi_{i1}, \varphi_{i2}, \dots, \varphi_{in})^T$$

$$\varphi_i \in \Phi$$

$$x_i = a(\varphi_i^T \mathbf{z} + \delta_i)$$

SGD estimation is fine-tuned as follows:

$$\begin{aligned} \varphi_{ij} &= \varphi_{ij} + \gamma \nabla \log P(x_i | \varphi_{ij}) = \varphi_{ij} - \gamma \frac{1}{\varphi_{ij}} \left( 1 + \frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}} \right), \forall i, j \\ \delta_i &= \delta_i + \gamma \nabla \log P(x_i | \delta_i) = \delta_i - \gamma \frac{a^{-1}(x_i)(a^{-1}(x_i))'}{\varphi_{ij}^2}, \forall i \end{aligned}$$

Given epoch of size  $N$  is denoted as  $\mathbf{D} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$ , the estimation equation of  $\varphi_{ij}$  and  $\delta_i$  is extended exactly as epoch estimation at every  $u^{\text{th}}$  iteration of SGD with regard to log-likelihood gradient  $\nabla \log P(x_i | \varphi_{ij}, \delta_i)$ .

$$\begin{aligned}\varphi_{ij}^{(u+1)} &= \varphi_{ij}^{(u)} + \gamma \frac{1}{N} \nabla \sum_{v=1}^N \log P(x_i^{(v)} | \varphi_{ij}^{(u)}) \\ &= \varphi_{ij}^{(u)} - \gamma \frac{1}{\varphi_{ij}^{(u)}} \left( 1 + \frac{1}{N \varphi_{ij}^{(u)}} \sum_{v=1}^N a^{-1}(x_i^{(v)}) (a^{-1}(x_i^{(v)}))' \right), \forall i, j\end{aligned}$$

$$\begin{aligned}\delta_i^{(u+1)} &= \delta_i^{(u)} + \gamma \frac{1}{N} \nabla \sum_{v=1}^N \log P(x_i^{(v)} | \delta_i^{(u)}) \\ &= \delta_i^{(u)} - \gamma \frac{1}{N (\varphi_{ij}^{(u)})^2} \sum_{v=1}^N a^{-1}(x_i^{(v)}) (a^{-1}(x_i^{(v)}))', \forall i\end{aligned}$$

Where  $x_i^{(v)}$  is the  $i^{\text{th}}$  element of  $\mathbf{x}^{(v)}$  in epoch. As a result, NL trained with SGD is specified as follows:

Initialize all  $\varphi_{ij}$ ,  $\delta_i$  and set  $u = 0$ .

Repeat

Sampling epoch  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$  or receiving epoch  $\mathbf{X}$  from big data / data stream.

$$\begin{aligned}\varphi_{ij}^{(u+1)} &= \varphi_{ij}^{(u)} - \gamma \frac{1}{\varphi_{ij}^{(u)}} \left( 1 + \frac{1}{N \varphi_{ij}^{(u)}} \sum_{v=1}^N a^{-1}(x_i^{(v)}) (a^{-1}(x_i^{(v)}))' \right), \forall i, j \\ \delta_i^{(u+1)} &= \delta_i^{(u)} - \gamma \frac{1}{N (\varphi_{ij}^{(u)})^2} \sum_{v=1}^N a^{-1}(x_i^{(v)}) (a^{-1}(x_i^{(v)}))', \forall i\end{aligned}$$

Increase  $u = u + 1$ .

Until some terminating conditions are met.

Note, a terminating condition is customized, for example, parameters  $\varphi_{ij}$  and  $\delta_i$  are not changed significantly or there is no more coming epoch  $\mathbf{X}$ . Moreover, the index  $u$  indicates time point as well as iteration of SGD. After finite NL is trained, it can generate new data  $\mathbf{x}'$  by generator  $g(\mathbf{z} | \Phi) = \mathbf{x}'$  with any  $\mathbf{z}$  randomized from standard normal distribution with mean  $\mathbf{0}$  and variance 1.

It is interesting that log-likelihood gradient  $\nabla \log P(x_i | \varphi_{ij})$  is determined based on inversed gradient  $\nabla_{x_i} g^{-1}(x_i | \varphi_{ij})$ . Therefore, how to estimate generator parameter  $\Phi$  by SGD estimation focuses on calculating inversed gradient  $\nabla_{x_i} g^{-1}(x_i | \varphi_{ij})$  which is central point of normalizing flow (NL) technique. Moreover, how to calculate  $\nabla_{x_i} g^{-1}(x_i | \varphi_{ij})$  is based on how to determine inversed generator  $g^{-1}(\mathbf{x} | \Phi)$ . In other words, the main problem of NL is how to determine inversed generator  $g^{-1}(\mathbf{x} | \Phi)$  because it is easy to calculate gradient of function  $f(\mathbf{x}) = g^{-1}(\mathbf{x} | \Phi)$  with regard to  $\mathbf{x}$ . Especially, when generator  $g(\mathbf{x} | \Phi)$  is implemented by DNN, NL will have some special techniques so that determining its inverse  $g^{-1}(\mathbf{x} | \Phi)$  is easier. One of these technique is *finite normalizing flow* (finite NL) in which generator  $g(\mathbf{x} | \Phi)$  is implemented by a DNN having  $K$  layers from layer 1 to layer  $K$  where layer 0 is input layer with note that each layer is represented by partial generator function  $f_k$  (Ruthotto & Haber, 2021, p. 8):

$$g(\mathbf{z} | \Phi) = f_K \circ f_{K-1} \circ \dots \circ f_2 \circ f_1(\mathbf{z}) \quad (2.5)$$

Note, all layers  $f_k$  have the same number of neurons which is the dimension  $n$ . Because  $f_k$  is essentially vector-by-vector function whose input and output are vectors, it should have denoted as  $\mathbf{f}_k$ , especially, when  $f_k$  here is bijection. However, it is still denoted  $f_k$  for convenience. Let  $\mathbf{z}^{(k+1)}$  be output of partial generator  $f_k$ , we have:



$$\begin{aligned} \mathbf{z}^{(k+1)} &= f_k(\mathbf{z}^{(k)}) \\ \mathbf{z} &= \mathbf{z}^{(1)} = \mathbf{z}^{(0)} = f_0(\mathbf{z}^{(0)}) \equiv f_0 \\ \mathbf{x} &= \mathbf{z}^{(K+1)} = f_K(\mathbf{z}^{(K)}) \end{aligned}$$

Inversed generator  $g^{-1}(\mathbf{x} | \Phi)$  representing inversed DNN is determined:

$$g^{-1}(\mathbf{x} | \Phi) = f_1^{-1} \circ f_2^{-1} \circ \dots \circ f_{K-1}^{-1} \circ f_K^{-1}(\mathbf{x}) \quad (2.6)$$

Each  $f_k^{-1}$  is called inversed partial generator function which is the inverse of partial generator function  $f_k$ . Let  $\mathbf{x}^{(k-1)}$  be output of partial generator  $f_k^{-1}$ , we have:

$$\begin{aligned} \mathbf{x}^{(k-1)} &= f_k^{-1}(\mathbf{x}^{(k)}) \\ \mathbf{x} &= \mathbf{x}^{(K+1)} = \mathbf{x}^{(K)} = f_{K+1}^{-1}(\mathbf{x}^{(K+1)}) \equiv f_{K+1}^{-1} \\ \mathbf{z} &= \mathbf{x}^{(0)} = f_2^{-1}(\mathbf{x}^{(1)}) \end{aligned}$$

“Input layer” of “inversed DNN” is  $f_{K+1}^{-1}$ . The inversed generator DNN may be pseudo in case that only one generator DNN is designed so that inversed generator function  $f_k^{-1}$  is existent. An interesting result of the design of finite NL is that inversed gradient  $\nabla g^{-1}(\mathbf{x} | \Phi)$  is product of gradients of inversed partial generator  $f_k^{-1}$ .

$$\begin{aligned} \nabla_{\mathbf{x}} g^{-1}(\mathbf{x} | \Phi) &= \prod_{k=1}^K \nabla_{\mathbf{x}^{(k)}} f_k^{-1}(\mathbf{x}^{(k)} | \Phi^{(k)}) \\ \log |\nabla_{\mathbf{x}} g^{-1}(\mathbf{x} | \Phi)| &= \sum_{k=1}^K |\nabla_{\mathbf{x}^{(k)}} f_k^{-1}(\mathbf{x}^{(k)} | \Phi^{(k)})| \end{aligned}$$

Where  $\Phi^{(k)}$  is parameter of  $f_k$ . It is now necessary to determine fine-tuned partial inversed gradient  $\nabla_{\mathbf{x}_i^{(k)}} f_k^{-1}(\mathbf{x}_i^{(k)} | \varphi_{ij}^{(k)})$  in order to determine fine-tuned partial log-likelihood gradient  $\nabla \log P(x_i^{(k)} | \varphi_{ij}^{(k)})$  where  $x_i^{(k)}$  is an elemental variable in  $\mathbf{x}^{(k)} = (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})^T$  and  $\varphi_{ij}^{(k)}$  is the  $j^{\text{th}}$  element in  $\varphi_i^{(k)} = (\varphi_{i1}^{(k)}, \varphi_{i2}^{(k)}, \dots, \varphi_{in}^{(k)})^T$  with note that  $\varphi_i^{(k)}$  is the  $i^{\text{th}}$  row vector of matrix  $\Phi^{(k)}$ . Moreover, let  $\delta_i^{(k)}$  be the  $i^{\text{th}}$  bias parameter associated with each  $x_i^{(k)}$ . Without loss of generality, given  $\varphi_j^{(k)}$ ,  $\delta_i^{(k)}$ , and  $z_j^{(k)}$  along with invertible activation  $a(\cdot)$ , we have fine-tuned inversed generator  $g^{-1}(x_i^{(k)} | \varphi_{ij}^{(k)}, \delta_i^{(k)})$ .

$$f_k^{-1}(\mathbf{x}_i^{(k)} | \varphi_{ij}^{(k)}) = z_j^{(k)} = \frac{a^{-1}(\mathbf{x}_i^{(k)}) - b_{ij}^{(k)} - \delta_i^{(k)}}{\varphi_{ij}^{(k)}}$$

Where,

$$b_{ij}^{(k)} = \sum_{t=1, t \neq j}^n \varphi_{it}^{(k)} z_t^{(k)}$$

Where  $z_i^{(k)}$  is the  $i^{\text{th}}$  elemental variable in  $\mathbf{z}^{(k)} = (z_1^{(k)}, z_2^{(k)}, \dots, z_n^{(k)})^T$ . By similar way aforementioned, log-likelihood gradient  $\nabla \log P(x_i^{(k)} | \varphi_{ij}^{(k)}, \delta_i^{(k)})$  is specified as follows:

$$\begin{aligned} \nabla \log P(\mathbf{x}_i^{(k)} | \varphi_{ij}^{(k)}) &= -\frac{1}{\varphi_{ij}^{(k)}} \left( 1 + \frac{a^{-1}(\mathbf{x}_i^{(k)}) (a^{-1}(\mathbf{x}_i^{(k)}))'}{\varphi_{ij}^{(k)}} \right), \forall i, j \\ \nabla \log P(\mathbf{x}_i^{(k)} | \delta_i^{(k)}) &= -\frac{a^{-1}(\mathbf{x}_i^{(k)}) (a^{-1}(\mathbf{x}_i^{(k)}))'}{(\varphi_{ij}^{(k)})^2}, \forall i \end{aligned} \quad (2.7)$$

Where  $a(\cdot)$  and  $a^{-1}(\cdot)$  are invertible activation function and its inverse and,

$$\begin{aligned} \mathbf{x}_i^{(k)} \in \mathbf{x}^{(k)} &= (x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)})^T \\ \varphi_{ij}^{(k)} \in \varphi_i^{(k)} &= (\varphi_{i1}^{(k)}, \varphi_{i2}^{(k)}, \dots, \varphi_{in}^{(k)})^T \end{aligned}$$

$$\varphi_i^{(k)} \in \Phi^{(k)}$$

$$x_i^{(k)} = a\left(\left(\varphi_i^{(k)}\right)^T \mathbf{z} + \delta_i^{(k)}\right)$$

SGD estimation is fine-tuned as follows:

$$\varphi_{ij}^{(k)} = \varphi_{ij}^{(k)} + \gamma \nabla \log P\left(x_i^{(k)} \middle| \varphi_{ij}^{(k)}\right) = \varphi_{ij}^{(k)} - \gamma \frac{1}{\varphi_{ij}^{(k)}} \left(1 + \frac{a^{-1}\left(x_i^{(k)}\right)\left(a^{-1}\left(x_i^{(k)}\right)\right)'}{\varphi_{ij}^{(k)}}\right), \forall i, j$$

$$\delta_i^{(k)} = \delta_i^{(k)} + \gamma \nabla \log P\left(x_i^{(k)} \middle| \delta_i^{(k)}\right) = \delta_i^{(k)} - \gamma \frac{a^{-1}\left(x_i^{(k)}\right)\left(a^{-1}\left(x_i^{(k)}\right)\right)'}{\left(\varphi_{ij}^{(k)}\right)^2}, \forall i$$

Given epoch of size  $N$  is denoted as  $\mathbf{D} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$ , the estimation equation of  $\varphi_{ij}^{(k)}$  and  $\delta_i^{(k)}$  is extended exactly as epoch estimation at every  $u^{\text{th}}$  iteration of SGD with regard to log-likelihood gradient  $\nabla \log P(x_i^{(k)} | \varphi_{ij}^{(k)}, \delta_i^{(k)})$ :

$$\begin{aligned} \left(\varphi_{ij}^{(k)}\right)^{(u+1)} &= \left(\varphi_{ij}^{(k)}\right)^{(u)} + \gamma \frac{1}{N} \nabla \sum_{v=1}^N \log P\left(\left(x_i^{(k)}\right)^{(v)} \middle| \left(\varphi_{ij}^{(k)}\right)^{(u)}\right) \\ &= \left(\varphi_{ij}^{(k)}\right)^{(u)} \\ &\quad - \gamma \frac{1}{\left(\varphi_{ij}^{(k)}\right)^{(u)}} \left(1 + \frac{1}{N \left(\varphi_{ij}^{(k)}\right)^{(u)} \sum_{v=1}^N a^{-1}\left(\left(x_i^{(k)}\right)^{(v)}\right)\left(a^{-1}\left(\left(x_i^{(k)}\right)^{(v)}\right)\right)'}\right), \forall i, j \end{aligned}$$

$$\begin{aligned} \left(\delta_i^{(k)}\right)^{(u+1)} &= \left(\delta_i^{(k)}\right)^{(u)} + \gamma \frac{1}{N} \nabla \sum_{v=1}^N \log P\left(\left(x_i^{(k)}\right)^{(v)} \middle| \left(\delta_i^{(k)}\right)^{(u)}\right) \\ &= \left(\delta_i^{(k)}\right)^{(u)} - \gamma \frac{1}{N \left(\varphi_{ij}^{(k)}\right)^{(u)2} \sum_{v=1}^N a^{-1}\left(\left(x_i^{(k)}\right)^{(v)}\right)\left(a^{-1}\left(\left(x_i^{(k)}\right)^{(v)}\right)\right)'} \sum_{v=1}^N a^{-1}\left(\left(x_i^{(k)}\right)^{(v)}\right)\left(a^{-1}\left(\left(x_i^{(k)}\right)^{(v)}\right)\right)', \forall i \end{aligned}$$

Where  $(x_i^{(k)})^{(v)}$  is the  $i^{\text{th}}$  element of  $\mathbf{x}^{(v)}$  in epoch with regard to inversed generator  $f_k^{-1}$ . As a result, finite NL trained with SGD is specified as follows:

Initialize all  $\varphi_{ij}^{(k)}$ ,  $\delta_i^{(k)}$  and set  $u = 0$ .

Repeat

Sampling epoch  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$  or receiving epoch  $\mathbf{X}$  from big data / data stream.

$$\begin{aligned}
\left(\varphi_{ij}^{(k)}\right)^{(u+1)} &= \left(\varphi_{ij}^{(k)}\right)^{(u)} \\
&\quad - \gamma \frac{1}{\left(\varphi_{ij}^{(k)}\right)^{(u)}} \left( 1 \right. \\
&\quad \left. + \frac{1}{N \left(\varphi_{ij}^{(k)}\right)^{(u)}} \sum_{v=1}^N a^{-1} \left( \left(x_i^{(k)}\right)^{(v)} \right) \left( a^{-1} \left( \left(x_i^{(k)}\right)^{(v)} \right) \right)' \right), \forall i, j \\
\left(\delta_i^{(k)}\right)^{(u+1)} &= \left(\delta_i^{(k)}\right)^{(u)} - \gamma \frac{1}{N \left(\left(\varphi_{ij}^{(k)}\right)^{(u)}\right)^2} \sum_{v=1}^N a^{-1} \left( \left(x_i^{(k)}\right)^{(v)} \right) \left( a^{-1} \left( \left(x_i^{(k)}\right)^{(v)} \right) \right)', \forall i
\end{aligned}$$

Increase  $u = u + 1$ .

Until some terminating conditions are met.

Note, a terminating condition is customized, for example, parameters  $\varphi_{ij}^{(k)}$  and  $\delta_i^{(k)}$  are not changed significantly or there is no more coming epoch  $X$ . Moreover, the index  $u$  indicates time point as well as iteration of SGD. After finite NL is trained, it can generate new data  $\mathbf{x}'$  by generator  $g(\mathbf{z} \mid \Phi) = \mathbf{x}'$  with any  $\mathbf{z}$  randomized from standard normal distribution with mean  $\mathbf{0}$  and variance 1.

Because it is not easy to calculate inversed gradient  $\nabla_{\mathbf{x}} g^{-1}(\mathbf{x} \mid \Phi)$  as well as its determinant  $|\nabla_{\mathbf{x}} g^{-1}(\mathbf{x} \mid \Phi)|$  according to finite NL except the decomposition technique above of entire matrix parameter  $\Phi$  into partial vector parameters  $\varphi_i^{(k)}$ , there is technique called real RVP (Ruthotto & Haber, 2021, p. 9) which defines each layer or partial generator  $f_k(\mathbf{z}^{(k)})$  by special way where  $\mathbf{z}^{(k)}$  is split into two parts such as  $\mathbf{z}_1^{(k)}$  and  $\mathbf{z}_2^{(k)}$  so that:

$$f_k(\mathbf{z}^{(k)}) = f_k \begin{pmatrix} \mathbf{z}_1^{(k)} \\ \mathbf{z}_2^{(k)} \end{pmatrix} = \begin{pmatrix} \mathbf{z}_1^{(k)} \\ \mathbf{z}_2^{(k)} \otimes \exp(s_k(\mathbf{z}_1^{(k)})) + \mathbf{t}_k(\mathbf{z}_1^{(k)}) \end{pmatrix} \quad (2.8)$$

Of course, we have:

$$f_k(\mathbf{z}^{(k)}) = f_k \begin{pmatrix} \mathbf{z}_1^{(k)} \\ \mathbf{z}_2^{(k)} \end{pmatrix} = \begin{pmatrix} \mathbf{z}_1^{(k)} \\ \mathbf{z}_2^{(k)} \otimes \exp(s_k(\mathbf{z}_1^{(k)})) + \mathbf{t}_k(\mathbf{z}_1^{(k)}) \end{pmatrix} = \begin{pmatrix} \mathbf{z}_1^{(k+1)} \\ \mathbf{z}_2^{(k+1)} \end{pmatrix} = \mathbf{z}^{(k+1)}$$

Where  $s_k$  and  $\mathbf{t}_k$  are two neural networks for scaling and translation, whose inputs and outputs have the same dimension. The operator  $\otimes$  denotes component-wise multiplication of two vectors where every pair of two corresponding elements of the two vectors are multiplied together, for instance, given two arbitrary vectors  $\mathbf{u} = (u_1, u_2, \dots, u_n)^T$  and  $\mathbf{v} = (v_1, v_2, \dots, v_n)^T$ , we have  $\mathbf{u} \otimes \mathbf{v} = (u_1 v_1, u_2 v_2, \dots, u_n v_n)^T$ . Moreover, the exponential function  $\exp(\cdot)$  above whose input is vector produces a vector by taking exponential function over every element of input vector. Inversed generator  $f_k^{-1}$  is specified from generator  $f_k$ .

$$f_k^{-1}(\mathbf{x}^{(k)}) = f_k^{-1} \begin{pmatrix} \mathbf{x}_1^{(k)} \\ \mathbf{x}_2^{(k)} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^{(k)} \\ \left( \mathbf{x}_2^{(k)} - \mathbf{t}_k(\mathbf{x}_1^{(k)}) \right) \otimes \exp(-s_k(\mathbf{x}_1^{(k)})) \end{pmatrix} \quad (2.9)$$

Of course, we have:

$$f_k^{-1}(\mathbf{x}^{(k)}) = f_k^{-1} \begin{pmatrix} \mathbf{x}_1^{(k)} \\ \mathbf{x}_2^{(k)} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^{(k)} \\ \left( \mathbf{x}_2^{(k)} - \mathbf{t}_k(\mathbf{x}_1^{(k)}) \right) \otimes \exp(-s_k(\mathbf{x}_1^{(k)})) \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1^{(k-1)} \\ \mathbf{x}_2^{(k-1)} \end{pmatrix} = \mathbf{x}^{(k-1)}$$

Inversed gradient  $\nabla_{\mathbf{x}^{(k)}} f_k^{-1}(\mathbf{x}^{(k)})$  is the 2x2 Jacobian matrix determined as follows:

$$\nabla_{\mathbf{x}^{(k)}} f_k^{-1}(\mathbf{x}^{(k)}) = \begin{pmatrix} 1 & 0 \\ \frac{d(\mathbf{x}_2^{(k)} - \mathbf{t}_k(\mathbf{x}_1^{(k)})) \otimes \exp(-\mathbf{s}_k(\mathbf{x}_1^{(k)}))}{d\mathbf{x}_1^{(k)}} & \exp(-\mathbf{s}_k(\mathbf{x}_1^{(k)})) \end{pmatrix}$$

It is interesting that taking determinant of inversed gradient  $\nabla_{\mathbf{x}^{(k)}} f_k^{-1}(\mathbf{x}^{(k)})$  becomes simple:

$$|\nabla_{\mathbf{x}^{(k)}} f_k^{-1}(\mathbf{x}^{(k)})| = \exp(-\mathbf{s}_k(\mathbf{x}_1^{(k)})) \quad (2.10)$$

When this determinant is determined, it is possible to maximize log-likelihood  $\log P(\mathbf{x} | \Phi)$  to estimate  $\Phi$  where  $\Phi$  here are weights of scaling neural network  $\mathbf{s}_k$  and translation neural network  $\mathbf{t}_k$ . Log-likelihood  $\log P(\mathbf{x} | \Phi)$  is written:

$$\begin{aligned} \log P(\mathbf{x} | \Phi) &= -\frac{n}{2} \log 2\pi - \frac{1}{2} \|\mathbf{g}^{-1}(\mathbf{x} | \Phi)\|^2 + \log |\nabla_{\mathbf{x}} \mathbf{g}^{-1}(\mathbf{x} | \Phi)| \\ &= \frac{n}{2} \log 2\pi - \frac{1}{2} \|\mathbf{g}^{-1}(\mathbf{x} | \Phi)\|^2 + \sum_{k=1}^K \log |\nabla_{\mathbf{x}^{(k)}} f_k^{-1}(\mathbf{x}^{(k)})| \\ &= \frac{n}{2} \log 2\pi - \frac{1}{2} \|\mathbf{g}^{-1}(\mathbf{x} | \Phi)\|^2 - \sum_{k=1}^K \mathbf{s}_k(\mathbf{x}_1^{(k)}) \end{aligned}$$

Because parameter  $\Phi$  is now only weights of scaling neural network  $\mathbf{s}_k$  and translation neural network  $\mathbf{t}_k$ , maximizing log-likelihood  $\log P(\mathbf{x} | \Phi)$  is now to optimize (train)  $\mathbf{s}_k$  and  $\mathbf{t}_k$  by some algorithms like backpropagation algorithm.

$$(\Phi_s^{(k)})^* = \underset{\Phi_s^{(k)}}{\text{optimize}} \mathbf{s}_k(\mathbf{x}_1^{(k)} | \Phi_s^{(k)})$$

$$(\Phi_t^{(k)})^* = \underset{\Phi_t^{(k)}}{\text{optimize}} \mathbf{t}_k(\mathbf{x}_1^{(k)} | \Phi_t^{(k)})$$

Beside finite NL there is another NL technique called continuous NL but it is not mentioned here because continuous NL is relevant to hazard problem of differential equation which is not main subject of DNN.

Recall that there are three main approaches for constructing DGM such as tractable density DGM, approximate density DGM, and implicit density DGM. However, if data is image, there is another categorization way that there are two main approaches: 1) pixel density approach tries to model pixel distribution and 2) block density approach tries to model entire image distribution as any data distribution. In other words, likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is defined based on probabilistic distribution of pixels where  $\mathbf{x}$  is considered as set of pixels according to pixel density approach. On the other hand, block density approach considers likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is PDF of a block or entire image (unified big block) where  $\mathbf{x}$  is considered as any arbitrary data. For instance, NL belongs to both tractable density DGM and block density approach. It is interesting that pixel density approach also belongs to tractable density approach because its PDF is defined obviously. Moreover, pixel density approach merges the two problems of training generator  $g(\mathbf{z} | \Phi)$  and qualifying such training task into the first problem which is to train  $g(\mathbf{z} | \Phi)$  by learning sample PDF  $P(\mathbf{x})$  because  $P(\mathbf{x})$  or  $P(\mathbf{x} | \Phi)$  now replaces  $P(\mathbf{x} | \Phi, \mathbf{z})$ .

Shortly, **pixel density (PD)** approach defines  $P(\mathbf{x})$  as product of all pixel distribution. Concretely, let  $\mathbf{x} = (x_1, x_2, \dots, x_{n^2})^T$  denote an image whose every  $i^{\text{th}}$  pixel is represented by elemental variable  $x_i$  and  $P(\mathbf{x})$  called image PDF is defined according to joint probability rule as follows:

$$P(\mathbf{x}) = \prod_{i=1}^{n^2} P(x_i | x_{i-1}, x_{i-2}, \dots, x_1) \quad (2.11)$$

Where  $n$  implies image width with suppose that image width and image height are equal for convenience, and

$$P(x_1|x_0) = P(x_1)$$

In other words, image PDF  $P(\mathbf{x})$  is product of all conditional PDFs  $P(x_i | x_{i-1}, x_{i-2}, \dots, x_1)$  where every  $P(x_i | x_{i-1}, x_{i-2}, \dots, x_1)$  is called conditional pixel PDF or pixel PDF in short. There is neither tractable data  $\mathbf{z}$  nor explicit generator  $g(\mathbf{z} | \Phi)$  for generating new data in PD because generation task is based on the entire PDF  $P(\mathbf{x})$ . For instance, without loss of generality, if we randomize  $k$  first pixels  $x_1, x_2, \dots, x_k$ , we can generate  $n^2-k$  remaining pixels by the recurrent process: determining  $P(x_{k+1} | x_k, x_{k-1}, \dots, x_1)$  based on  $x_1$  to  $x_k$ , generating  $x_{k+1}$  according to  $P(x_{k+1} | x_k, x_{k-1}, \dots, x_1)$  and determining  $P(x_{k+2} | x_{k+1}, x_k, \dots, x_1)$  based on  $x_1$  to  $x_{k+1}$ , generating  $x_{k+2}$  according to  $P(x_{k+2} | x_{k+1}, x_k, \dots, x_1)$  and determining  $P(x_{k+3} | x_{k+2}, x_{k+1}, \dots, x_1)$  based on  $x_1$  to  $x_{k+2}$ , ..., generating  $x_{n^2-1}$  according to  $P(x_{n^2-1} | x_{n^2-2}, x_{n^2-3}, \dots, x_1)$  and determining  $P(x_{n^2} | x_{n^2-1}, x_{n^2-2}, \dots, x_1)$  based on  $x_1$  to  $x_{n^2-1}$ , generating the last  $x_{n^2}$  according to  $P(x_{n^2} | x_{n^2-1}, x_{n^2-2}, \dots, x_1)$ . By another viewpoint, the joint probability of  $n^2-k$  remaining pixels denoted  $P(x_k, x_{k+1}, \dots, x_{n^2})$  is determined and then,  $n^2-k$  remaining pixels are generated according to this joint probability. Indeed, the joint probability  $P(x_k, x_{k+1}, \dots, x_{n^2})$  of  $n^2-k$  remaining pixels is totally determined when  $P(\mathbf{x})$  and  $k$  probabilities  $P(x_i | x_{i-1}, x_{i-2}, \dots, x_1)$  are determined where  $i$  is from 1 to  $k$ .

$$P(x_{k+1}, x_{k+2}, \dots, x_{n^2}) = \prod_{i=k+1}^{n^2} P(x_i | x_{i-1}, x_{i-2}, \dots, x_1) = \frac{P(\mathbf{x})}{\prod_{i=1}^k P(x_i | x_{i-1}, x_{i-2}, \dots, x_1)}$$

Because there are a large number of pixels in a large image which produces a large number of pixel PDFs as well as every pixel PDF  $P(x_i | x_{i-1}, x_{i-2}, \dots, x_1)$  of a given pixel  $x_i$  is itself also complicated with a lot of its previous pixels  $x_{i-1}, x_{i-2}, \dots, x_1$ , there are many techniques proposed to PD in order to decrease complexity and increase computation effectiveness. Anyhow, the equation of image PDF  $P(\mathbf{x})$  above is important one in theory. One of PD techniques is to apply long short-term memory (LSTM) (Theis & Bethge, 2015) into modeling and learning sample PDF  $P(\mathbf{x})$ .

The default artificial neural network is feedforward neural network where data is fed to input layer which, in turn, is evaluated and passed across hidden layers to output layer in one-way direction, finally. However, there is an extension of neural network, which is called recurrent neural work (RNN), where an output can be turned back in order to feed on network as input. In other words, RNN has circle, which allow that output can become input. For convenience and easy explanation, given  $T$  time points  $t = 1, 2, \dots, T$ , current state of a RNN at time point  $t$  is represented by three layers such as input layer  $\mathbf{x}_t$ , hidden layer  $\mathbf{h}_t$ , and output layer  $\mathbf{o}_t$  without loss of generality with note that  $\mathbf{h}_t$  can represent many hidden layers when RNN is a DNN too. Obviously, RNN is an extension of neural network because every triple  $(\mathbf{x}_t, \mathbf{h}_t, \mathbf{o}_t)$  is, essentially, a feedforward neural network, even a DNN. Hidden layer  $\mathbf{h}_t$  as well as output layer  $\mathbf{o}_t$  at current state  $t$  is calculated based on both current input layer  $\mathbf{x}_t$  and previous hidden layer  $\mathbf{h}_{t-1}$  of previous state at time point  $t-1$ . Without loss of generality, input layer, hidden layer, and output layer are considered as input neuron, hidden neuron, and output neuron for convenience (Wikipedia, Recurrent neural network, 2005).

$$\begin{aligned} \mathbf{h}_t &= \sigma_h(\mathbf{W}_h \mathbf{x}_t + \mathbf{U}_h \mathbf{h}_{t-1} + \mathbf{b}_h) \\ \mathbf{o}_t &= \sigma_o(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o) \end{aligned} \quad (2.12)$$

Where  $\mathbf{W}_h$  is weight matrix of current hidden neuron  $\mathbf{h}_t$  regarding current input neuron  $\mathbf{x}_t$ ,  $\mathbf{U}_h$  is weight matrix of current hidden neuron  $\mathbf{h}_t$  regarding previous hidden neuron  $\mathbf{h}_{t-1}$ , and  $\mathbf{b}_h$  is bias vector of current hidden neuron  $\mathbf{h}_t$  whereas  $\mathbf{W}_o$  is weight matrix of current output neuron  $\mathbf{o}_t$  regarding current hidden neuron  $\mathbf{h}_t$  and  $\mathbf{b}_o$  is bias vector of current output neuron  $\mathbf{o}_t$ . Moreover,  $\sigma_h(\cdot)$  and  $\sigma_o(\cdot)$  are activation functions of  $\mathbf{h}_t$  and  $\mathbf{o}_t$ , respectively with note that  $\sigma_h(\cdot)$  and  $\sigma_o(\cdot)$  are vector-by-vector functions. Backpropagation algorithm can be applied into



learning RNN as usual. It is interesting that structure of RNN defined by the triple  $(\mathbf{x}_t, \mathbf{h}_t, \mathbf{o}_t)$  is not changed but its parameters  $\mathbf{W}_h, \mathbf{U}_h, \mathbf{b}_h, \mathbf{W}_o$ , and  $\mathbf{b}_o$  are changed by backpropagation algorithm when RNN is learned. Of course, values of the triple  $(\mathbf{x}_t, \mathbf{h}_t, \mathbf{o}_t)$  are changed over time points. Note,  $\mathbf{W}_h, \mathbf{U}_h$ , and  $\mathbf{W}_o$  are matrix parameters and  $\mathbf{b}_h$  and  $\mathbf{b}_o$  are vector parameters whereas  $\mathbf{x}_t, \mathbf{h}_t$ , and  $\mathbf{o}_t$  are vector variables.

Long short-term memory (LSTM) is an extension of RNN, which implies that RNN is used to implement short-term memory so that the short-term memory can last for a longer time through  $T$  time points  $t = 1, 2, \dots, T$  built in RNN. Consequently, the short-term memory is represented by a so-called cell associated with three gates such as input gate, forget gate, and output gate. Cell represents information piece stored in memory at current time (Wikipedia, 2007). Input gate controls which new information to be put to cell, forget gate decides which information to be discarded, and output gate controls which information to be sent to next state (Wikipedia, 2007). As a convention, the cell at current state  $t$  is represented by the pair  $(\mathbf{c}_t, \mathbf{h}_t)$  whereas forget gate, input gate, and output gate are represented by vector variables  $\mathbf{f}_t, \mathbf{i}_t$ , and  $\mathbf{o}_t$ , respectively. Note, let  $\mathbf{g}_t$  and  $\mathbf{c}_t$  be cell input activation variable and cell state variable where cell input activation variable  $\mathbf{g}_t$  represents the activated input part of a cell, which is the important input part being different from the forgotten part, whereas cell state variable  $\mathbf{c}_t$  represents real information stored in cell which is, exactly, the short memory at current state. In literature,  $\mathbf{g}_t$  is also called cell gate. Some LSTM variants merge  $\mathbf{g}_t$  and  $\mathbf{c}_t$  into the same cell state variable. Although output gate  $\mathbf{o}_t$  represents which information to be sent to next state, it is consolidated with current cell memory  $\mathbf{c}_t$  in order to produce the real output information  $\mathbf{h}_t$  which represents bright and clear-cut memory. In other words, given cell  $(\mathbf{c}_t, \mathbf{h}_t)$ , then  $\mathbf{c}_t$  represents the real information stored in memory and  $\mathbf{h}_t$  represents the clear-cut memory which displays brightly at the outside for next state. It is possible to consider that  $\mathbf{c}_t$  is evaluated value of cell  $t$  and  $\mathbf{h}_t$  is predictive value of cell  $t$ . Following equations specify LSTM based on specification of RNN (Wikipedia, 2007), which indicates how to calculate cell and gates.

$$\begin{aligned} \mathbf{f}_t &= \sigma_f(\mathbf{W}_f \mathbf{x}_t + \mathbf{U}_f \mathbf{h}_{t-1} + \mathbf{b}_f) \\ \mathbf{i}_t &= \sigma_i(\mathbf{W}_i \mathbf{x}_t + \mathbf{U}_i \mathbf{h}_{t-1} + \mathbf{b}_i) \\ \mathbf{o}_t &= \sigma_o(\mathbf{W}_o \mathbf{x}_t + \mathbf{U}_o \mathbf{h}_{t-1} + \mathbf{b}_o) \\ \mathbf{g}_t &= \sigma_g(\mathbf{W}_g \mathbf{x}_t + \mathbf{U}_g \mathbf{h}_{t-1} + \mathbf{b}_g) \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \mathbf{g}_t \\ \mathbf{h}_t &= \mathbf{o}_t \otimes \sigma_h(\mathbf{c}_t) \end{aligned} \tag{2.13}$$

Note, weight matrix  $\mathbf{W}_i$ , weight matrix  $\mathbf{U}_i$ , and bias vector  $\mathbf{b}_i$  are parameters of input gate  $\mathbf{i}_t$ . Weight matrix  $\mathbf{W}_f$ , weight matrix  $\mathbf{U}_f$ , and bias vector  $\mathbf{b}_f$  are parameters of forget gate  $\mathbf{f}_t$ . Weight matrix  $\mathbf{W}_o$ , weight matrix  $\mathbf{U}_o$ , and bias vector  $\mathbf{b}_o$  are parameters of output gate  $\mathbf{o}_t$ . Weight matrix  $\mathbf{W}_g$ , weight matrix  $\mathbf{U}_g$ , and bias vector  $\mathbf{b}_g$  are parameters of cell gate  $\mathbf{g}_t$ . Vector variables  $\mathbf{i}_t, \mathbf{f}_t$ , and  $\mathbf{o}_t$  are often in range  $[0, 1]$  whereas vector variables  $\mathbf{c}_t$  and  $\mathbf{h}_t$  are often in range  $[-1, 1]$ . Activation functions  $\sigma_i(\cdot)$ ,  $\sigma_f(\cdot)$ , and  $\sigma_o(\cdot)$  are often sigmoid (logistic) functions whereas activation functions  $\sigma_g(\cdot)$  and  $\sigma_h(\cdot)$  are hyperbolic tangent functions. The operator  $\otimes$  denotes component-wise multiplication of two vectors where every pair of two corresponding elements of the two vectors are multiplied together, for instance, given two arbitrary vectors  $\mathbf{u} = (u_1, u_2, \dots, u_n)^T$  and  $\mathbf{v} = (v_1, v_2, \dots, v_n)^T$ , we have  $\mathbf{u} \otimes \mathbf{v} = (u_1 v_1, u_2 v_2, \dots, u_n v_n)^T$ . Note, backpropagation algorithm can be applied into learning LSTM as usual.

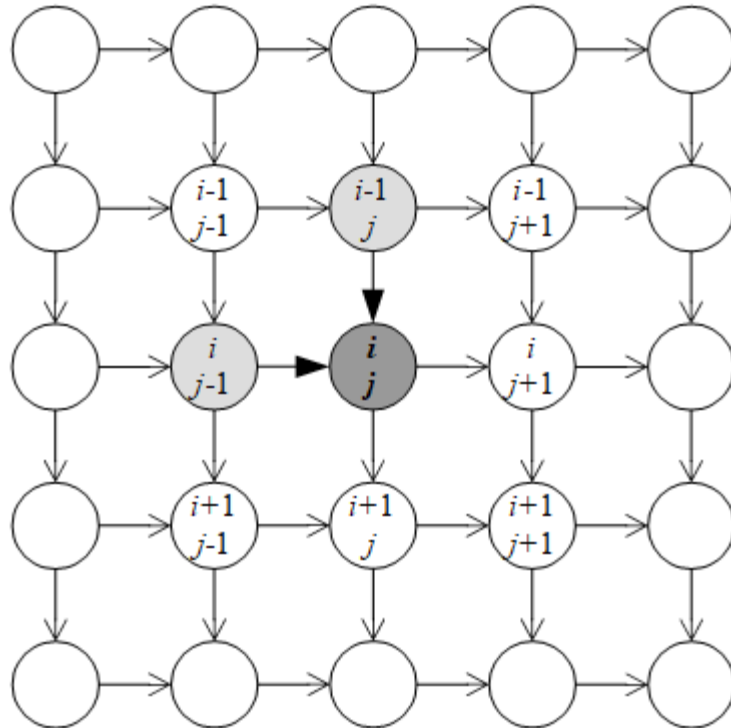
By applying LSTM into pixel density (PD) approach for modeling DGM, each pixel  $x_i$  is represented by cell  $\mathbf{c}_i$  when pixel index  $i$  is considered as time point  $t$ . Because each cell  $\mathbf{c}_i$  is dependent on its one right previous cell  $\mathbf{c}_{i-1}$  whereas conditional pixel PDF  $P(\mathbf{x}_i | \mathbf{x}_{i-1}, \mathbf{x}_{i-2}, \dots, \mathbf{x}_1)$  of pixel  $\mathbf{x}_i$  is dependent on  $i-1$  previous pixels  $\mathbf{x}_{i-1}, \mathbf{x}_{i-2}, \dots, \mathbf{x}_1$ , Markov property is applied so that conditional pixel PDF of pixel  $\mathbf{x}_i$  depends on only one previous pixel  $\mathbf{x}_{i-1}$ .

$$P(\mathbf{x}_i | \mathbf{x}_{i-1}, \mathbf{x}_{i-2}, \dots, \mathbf{x}_1) = P(\mathbf{x}_i | \mathbf{x}_{i-1})$$

It is now possible to apply LSTM to model PD by matching each pixel  $\mathbf{x}_i$  with each cell  $\mathbf{c}_i$  so that cell  $\mathbf{c}_i$  is considered as evaluated value of pixel  $\mathbf{x}_i$  as well as each  $\mathbf{h}_i$  is predictive value of pixel  $\mathbf{x}_i$ . Because image is two-dimension array, each pixel  $\mathbf{x}_{ij}$  or each cell  $\mathbf{c}_{ij}$  is indexed by two indices  $i$  and  $j$  following image height and image width. The event that cell  $\mathbf{c}_{ij}$  or  $\mathbf{c}_{i,j}$  indexed by two indices  $i$  and  $j$  makes LSTM extended into two-dimension LSTM as follows:

$$\begin{aligned} f_{i,j} &= \sigma_f(\mathbf{W}_f \mathbf{x}_{i,j} + \mathbf{U}_f \mathbf{h}_{i,j-1} + \mathbf{V}_f \mathbf{h}_{i-1,j} + \mathbf{b}_f) \\ i_{i,j} &= \sigma_i(\mathbf{W}_i \mathbf{x}_{i,j} + \mathbf{U}_i \mathbf{h}_{i,j-1} + \mathbf{V}_i \mathbf{h}_{i-1,j} + \mathbf{b}_i) \\ o_{i,j} &= \sigma_o(\mathbf{W}_o \mathbf{x}_{i,j} + \mathbf{U}_o \mathbf{h}_{i,j-1} + \mathbf{V}_o \mathbf{h}_{i-1,j} + \mathbf{b}_o) \\ g_{i,j} &= \sigma_g(\mathbf{W}_g \mathbf{x}_{i,j} + \mathbf{U}_g \mathbf{h}_{i,j-1} + \mathbf{V}_g \mathbf{h}_{i-1,j} + \mathbf{b}_g) \\ \mathbf{c}_{i,j} &= f_{i,j} \otimes (\mathbf{c}_{i,j-1} + \mathbf{c}_{i-1,j}) + i_{i,j} \otimes g_{i,j} \\ \mathbf{h}_{i,j} &= o_{i,j} \otimes \sigma_h(\mathbf{c}_{i,j}) \end{aligned} \quad (2.14)$$

The equations above specify core ideology of PD associated with two-dimension LSTM where the contextual meaning of weight and bias parameters  $\mathbf{W}_{(\cdot)}$ ,  $\mathbf{U}_{(\cdot)}$ ,  $\mathbf{V}_{(\cdot)}$ , and  $\mathbf{b}_{(\cdot)}$  is not changed with note that  $\mathbf{W}_{(\cdot)}$ ,  $\mathbf{U}_{(\cdot)}$ , and  $\mathbf{V}_{(\cdot)}$  are weight matrices regarding current pixel, previous pixel  $(i, j-1)$ , and previous pixel  $(i-1, j)$ , respectively. In literature, such PD is called *PixelRNN* associated with diagonal two-dimension LSTM (Oord, Kalchbrenner, & Kavukcuoglu, 2016, pp. 3-4). According to diagonal two-dimension LSTM each pixel  $(i, j)$  at  $i^{\text{th}}$  row and  $j^{\text{th}}$  column has two previous neighbors such as previous left pixel  $(i, j-1)$  and previous upper pixel  $(i-1, j)$ . For extension, each pixel  $(i, j)$  can have up four previous neighbors such as pixel  $(i, j-1)$ , pixel  $(i-1, j-1)$ , pixel  $(i-1, j)$ , and pixel  $(i-1, j+1)$ . Following figure depicts PixelRNN with diagonal two-dimension LSTM (Oord, Kalchbrenner, & Kavukcuoglu, 2016, p. 4).



**Figure 2.1.** PixelRNN with diagonal two-dimension LSTM

It is easy to add more weight parameters to these extensive cases. For example, cell gates and cell state with regard to the four previous neighbors are specified as follows:

$$\begin{aligned} f_{i,j} &= \sigma_f(\mathbf{W}_f \mathbf{x}_{i,j} + \mathbf{U}_f \mathbf{h}_{i,j-1} + \mathbf{V}_f \mathbf{h}_{i-1,j} + \mathbf{R}_f \mathbf{h}_{i-1,j-1} + \mathbf{S}_f \mathbf{h}_{i-1,j+1} + \mathbf{b}_f) \\ i_{i,j} &= \sigma_i(\mathbf{W}_i \mathbf{x}_{i,j} + \mathbf{U}_i \mathbf{h}_{i,j-1} + \mathbf{V}_i \mathbf{h}_{i-1,j} + \mathbf{R}_i \mathbf{h}_{i-1,j-1} + \mathbf{S}_i \mathbf{h}_{i-1,j+1} + \mathbf{b}_i) \end{aligned}$$

$$\begin{aligned}
o_{i,j} &= \sigma_o(W_o x_{i,j} + U_o h_{i,j-1} + V_o h_{i-1,j} + R_o h_{i-1,j-1} + S_o h_{i-1,j+1} + b_i) \\
g_{i,j} &= \sigma_g(W_g x_{i,j} + U_g h_{i,j-1} + V_g h_{i-1,j} + R_g h_{i-1,j-1} + S_g h_{i-1,j+1} + b_g) \\
c_{i,j} &= f_{i,j} \otimes (c_{i,j-1} + c_{i-1,j} + c_{i-1,j-1} + c_{i-1,j+1}) + i_{i,j} \otimes g_{i,j} \\
h_{i,j} &= o_{i,j} \otimes \sigma_h(c_{i,j})
\end{aligned}$$

Where matrices  $R_{(\cdot)}$  and  $S_{(\cdot)}$  are additional weight parameters regarding two new neighbor pixels such as pixel  $(i-1, j-1)$  and pixel  $(i-1, j+1)$ .

Recall that  $c_{i,j}$  is considered as evaluated value of pixel  $x_{i,j}$  and  $h_{i,j}$  is predictive value of pixel  $x_{i,j}$ . It is interesting that  $h_{i,j}$  is generated pixel within the aforementioned generation process by PD. Turning back the generation process, without loss of generality, given  $k$  randomized pixels  $x_{i-1,1}, x_{i-1,2}, \dots, x_{i-1,j+1}, \dots, x_{i,1}, x_{i,2}, \dots, x_{i,j}$ , we will generate the next pixel  $x_{i,j+1}$ . Firstly, PD model must be trained by some dataset as a set of images. Secondly,  $k$  randomized pixels  $x_{i-1,1}, x_{i-1,2}, \dots, x_{i-1,j+1}, \dots, x_{i,1}, x_{i,2}, \dots, x_{i,j}$  are fed to PD again so as to update  $k$  sets of parameters  $W_{(\cdot)}, U_{(\cdot)}$ , and  $b_{(\cdot)}$  as well as compute  $k$  predictive values  $h_{i-1,1}, h_{i-1,2}, \dots, h_{i-1,j+1}, \dots, h_{i,1}, h_{i,2}, \dots, h_{i,j}$ . Finally, it is possible to determine the predictive value  $h_{i,j+1}$  of the next pixel  $(i, j+1)$  given  $x_{i,j+1}, h_{i,j}$ , and  $h_{i-1,j+1}$  along with learned parameters of two-dimension LSTM PD. It is important to note that  $x_{i,j+1}$  is randomized arbitrarily whereas  $h_{i,j}$  and  $h_{i-1,j+1}$  were computed previously. Obviously, it is easy to generate next predictive values  $h_{i,j+2}, h_{i,j+3}, \dots, h_{i+1,j}, h_{i+1,j+1}$ , etc. by the similar process. Note, backpropagation algorithm can be applied into learning two-dimension LSTM as usual. Note, backpropagation algorithm is often associated with stochastic gradient descent (SGD) method and so, please pay attention to SGD.

### 3. APPROXIMATE DENSITY DGM

According to approximate density approach, DGMs try to estimate approximately generator PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$  or derive other PDF that is similar to  $P(\mathbf{x} | \Phi, \mathbf{z})$  with note that PDF is abbreviation of probability density function.

Recall that there are two problems related to construct a DGM: 1) how to define likelihood or error to train generator DNN  $g(\mathbf{z} | \Phi)$  and 2) how to define tractable PDF  $P(\mathbf{z})$  which implies the way to randomize  $\mathbf{z}$ . The second problem relates to assert qualification of random data  $\mathbf{z}'$  and hence, the second problem is stated as qualification problem of how to qualify random data. According to implicit density approach, a discrimination DNN is used to qualify randomized data  $\mathbf{z}$  instead of defining tractable PDF  $P(\mathbf{z})$  by Generative Adversarial Network (GAN) which is a typical method belonging to implicit density approach. In different way belonging to this approximate density approach, **Variational Autoencoders (VAE)** method developed by Kingma and Welling (Kingma & Welling, 2022) proposed another DNN called encoder  $f(\mathbf{x} | \Theta)$  to expectedly convert intractable data  $\mathbf{x}$  into tractable data  $\mathbf{z}$ . In other words, encoder  $f(\mathbf{x} | \Theta)$  approximates tractable data  $\mathbf{z}$  by encoded data  $\mathbf{z}'$ .

$$f(\mathbf{x} | \Theta) = \mathbf{z}' \cong \mathbf{z} \quad (3.1)$$

It is easy to recognize that encoder  $f(\mathbf{x} | \Theta)$  is an approximation of the inverse of generator  $g(\mathbf{z} | \Phi)$  when  $g(\mathbf{z} | \Phi)$  is invertible where  $\mathbf{x}$ -dimension  $m$  is larger than  $\mathbf{z}$ -dimension  $n$  ( $m > n$ ), which is the reason that generator  $g(\mathbf{z} | \Phi)$  is called decoder  $g(\mathbf{z} | \Phi)$  in VAE. Like decoder  $g(\mathbf{z} | \Phi)$ , encoder  $f(\mathbf{x} | \Theta)$  is modeled by a so-called encoder DNN whose weights are parameter  $\Theta$  called encoder parameter and so parameter  $\Phi$  is called decoder parameter in VAE. By following the fact that encoder  $f(\mathbf{x} | \Theta)$  approximates tractable data  $\mathbf{z}$  by encoded data  $\mathbf{z}'$ , tractable PDF  $P(\mathbf{z})$  is approximated by a so-called encoder PDF  $P_f(\mathbf{z}')$ .

$$P_f(\mathbf{z}') \cong P(\mathbf{z})$$

Because encoder  $f(\mathbf{x} | \Theta)$  depends on its parameter  $\Theta$ , we can denote:

$$P_f(\mathbf{z}') = P(\mathbf{z}' | \Theta, \mathbf{x})$$

Essential, encoder PDF  $P(\mathbf{z}' | \Theta, \mathbf{x})$  is likelihood function of  $\mathbf{z}'$  given  $\mathbf{x}$  which is conditional PDF of  $\mathbf{z}'$  given  $\mathbf{x}$  and hence,  $P(\mathbf{z}' | \Theta, \mathbf{x})$  is called encoder likelihood which depends on encoder  $f(\mathbf{x} | \Theta)$ , of course. On the other hand,  $P(\mathbf{z}' | \Theta, \mathbf{x})$  is posterior PDF of tractable data given tractable data  $\mathbf{x}$  where  $P(\mathbf{z})$  is prior PDF of tractable data. In practice,  $\mathbf{z}'$  is assumed to conform multivariate normal distribution and therefore, let  $\mu(\mathbf{x})$  and  $\Sigma(\mathbf{x})$  be mean vector and covariance matrix of  $\mathbf{z}'$ . Encoder likelihood  $P(\mathbf{z}' | \Theta, \mathbf{x})$  becomes  $P(\mathbf{z}' | \Theta, \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  so that output of encoder DNN  $f(\mathbf{x} | \Theta)$  is mean  $\mu(\mathbf{x} | \Theta)$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta)$  while its input is  $\mathbf{x}$  and its weights are  $\Theta$ , of course.

$$f(\mathbf{x} | \Theta) = \begin{pmatrix} \mu(\mathbf{x} | \Theta) \\ \Sigma(\mathbf{x} | \Theta) \end{pmatrix} \quad (3.2)$$

$$P(\mathbf{z}' | \Theta, \mathbf{x}) = P(\mathbf{z}' | \Theta, \mu(\mathbf{x}), \Sigma(\mathbf{x})) = \mathcal{N}(\mathbf{z}' | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$$

Note,  $\mathcal{N}(\cdot)$  denotes normal distribution and thus,  $\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  represents encoder likelihood. That  $\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  is encoder likelihood is an important improvement in developing VAE because encoder DNN  $f(\mathbf{x} | \Theta)$  is learned by minimizing a so-called encoder error which is represented by the difference between encoder likelihood and predefined tractable PDF  $P(\mathbf{z})$ . Let  $\text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$  be Kullback-Leibler divergence of encoder likelihood  $\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  and predefined tractable PDF  $P(\mathbf{z})$ . As a result,  $\text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$  becomes an ideal encoder error, which is called encoder KL divergence. The smaller the encoder KL divergence is, the closer the encoder likelihood  $\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  is to tractable PDF  $P(\mathbf{z})$ , the better the encoder DNN  $f(\mathbf{x} | \Theta)$  is.

$$\varepsilon(\mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta) | \Theta) = \text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z})) \quad (3.3)$$

Therefore, encoder KL divergence  $\text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$  is minimized by stochastic gradient descent (SGD) method in order to estimate decoder parameter  $\Theta$  for training encoder DNN  $f(\mathbf{x} | \Theta)$  as follows:

$$\Theta^* = \underset{\Theta}{\text{armin}} \text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$$

Which results estimation equation according to SGD:

$$\Theta = \Theta - \gamma \nabla \text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$$

Where  $\nabla \text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$  is gradient of encoder KL divergence  $\text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$  with regard to  $\mu(\mathbf{x} | \Theta)$  and  $\Sigma(\mathbf{x} | \Theta)$  while  $\gamma$  is learning rate. Recall that SGD, which is an iterative process, pushes candidate solution at each iteration along the direction which is opposite to gradient of target function for minimization or has the same direction to gradient of target function for maximization with note that the step length is represented by learning rate. We have:

$$\nabla \text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z})) = \frac{d \text{KL}(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))}{d \Theta}$$

There can be no change in estimating decoder parameter  $\Phi$  within VAE so that decoder error  $\varepsilon(\mathbf{x} | \Phi, \mathbf{z}) = \frac{1}{2} \|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2$  is minimized to produce optimal  $\Phi$ .

$$\Phi^* = \underset{\Phi}{\text{armin}} \frac{1}{2} \|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2$$

Which results estimation equation according to SGD:

$$\Phi = \Phi - \gamma \nabla \frac{1}{2} \|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2$$

Recall that generator  $g(\mathbf{z} | \Phi)$  is called decoder  $g(\mathbf{z} | \Phi)$  in VAE. As a result, encoder parameter  $\Theta$  and decoder parameter  $\Phi$  are estimated as follows:

$$\Theta = \Theta - \gamma \nabla \text{KL} \left( \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| P(\mathbf{z}) \right)$$

$$\Phi = \Phi - \gamma \nabla \frac{1}{2} \|g(\mathbf{z} | \Phi) - \mathbf{x}\|^2 = \Phi - \gamma (g(\mathbf{z} | \Phi) - \mathbf{x}) \frac{dg(\mathbf{z} | \Phi)}{d\Phi}$$

Where  $dg(\mathbf{z} | \Phi) / d\Phi$  is differential of  $g(\mathbf{z} | \Phi)$  with regard to  $\Phi$  while  $0 < \gamma \leq 1$  is learning rate and tractable PDF  $P(\mathbf{z})$  is predefined with note that VAE replaces tractable PDF  $P(\mathbf{z})$  by likelihood  $P(\mathbf{z} | \Theta, \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  with fixed  $P(\mathbf{z})$ . As usual,  $P(\mathbf{z})$  is assumed to conform standard normal distribution with mean  $\mathbf{0}$  and covariance matrix  $I$ .

$$P(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, I)$$

This implies:

$$\Theta = \Theta - \gamma \nabla \text{KL} \left( \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| \mathcal{N}(\mathbf{z} | \mathbf{0}, I) \right)$$

Where  $I$  is identity matrix:

$$I = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

It is easier to determine gradient of encoder KL divergence  $\nabla \text{KL}(\mathcal{N}(\mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mathbf{z} | \mathbf{0}, I))$  with regard to  $\Theta$  between the multivariate normal distribution  $\mathcal{N}(\mu(\mathbf{x}), \Sigma(\mathbf{x}) | \Theta)$  and the standard multivariate normal distribution  $\mathcal{N}(\mathbf{z} | \mathbf{0}, I)$ . We have following equation to calculate such gradient (Kingma & Welling, 2022, p. 5), (Doersch, 2016, p. 9), (Nguyen, 2015, p. 43):

$$\begin{aligned} \nabla \text{KL} \left( \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| \mathcal{N}(\mathbf{z} | \mathbf{0}, I) \right) &= \frac{d \text{KL} \left( \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| \mathcal{N}(\mathbf{z} | \mathbf{0}, I) \right)}{d\Theta} \\ &= \left( \mu(\mathbf{x} | \Theta) \frac{d\mu(\mathbf{x} | \Theta)}{d\Theta} \right)^T \\ &\quad \left( \frac{1}{2} \left( I - (\Sigma(\mathbf{x} | \Theta))^{-1} \right) \frac{d\Sigma(\mathbf{x} | \Theta)}{d\Theta} \right) \\ &= \left( \mu(\mathbf{x} | \Theta) \frac{d\mu(\mathbf{x} | \Theta)}{d\Theta}, \frac{1}{2} \left( I - (\Sigma(\mathbf{x} | \Theta))^{-1} \right) \frac{d\Sigma(\mathbf{x} | \Theta)}{d\Theta} \right) \end{aligned}$$

Where  $(\Sigma(\mathbf{x} | \Theta))^{-1}$  is inverse of covariance matrix  $\Sigma(\mathbf{x} | \Theta)$  and the subscript “ $T$ ” denotes transposition operator of matrix and vector whereas  $d\mu(\mathbf{x} | \Theta) / d\Theta$  and  $d\Sigma(\mathbf{x} | \Theta) / d\Theta$  are differentials of  $\mu(\mathbf{x} | \Theta)$  and  $\Sigma(\mathbf{x} | \Theta)$  with regard to  $\Theta$ , respectively. As a result, encoder parameter  $\Theta$  and decoder parameter  $\Phi$  are totally estimated according to SGD as follows:

$$\begin{aligned} \Theta &= \Theta - \gamma \left( \mu(\mathbf{x} | \Theta) \frac{d\mu(\mathbf{x} | \Theta)}{d\Theta} \right. \\ &\quad \left. \frac{1}{2} \left( I - (\Sigma(\mathbf{x} | \Theta))^{-1} \right) \frac{d\Sigma(\mathbf{x} | \Theta)}{d\Theta} \right) \\ \Phi &= \Phi - \gamma (g(\mathbf{z} | \Phi) - \mathbf{x}) \frac{dg(\mathbf{z} | \Phi)}{d\Phi} \end{aligned}$$

The estimation equations above are simple explanation of VAE but its formal construction is more complicated. We begin the aforementioned intractable PDF  $P(\mathbf{x})$  specified by law of total probability:

$$P(\mathbf{x}) = \int_{\mathbf{z}} P(\mathbf{x} | \Phi, \mathbf{z}) P(\mathbf{z}) d\mathbf{z}$$

However,  $P(\mathbf{x})$  is interpreted by another way which is based on Bayes’ rule within VAE:

$$P(\mathbf{x}) = \frac{P(\mathbf{x}, \mathbf{z})}{P(\mathbf{z} | \mathbf{x})}$$



Because the conditional probability  $P(\mathbf{z} | \mathbf{x})$  is arbitrary without formal specification, it should be approximated by another PDF denoted  $Q(\mathbf{z} | \mathbf{x})$  with assumption that the PDF  $Q(\mathbf{z} | \mathbf{x})$  has formal specification like normal distribution.

$$Q(\mathbf{z} | \mathbf{x}) \cong P(\mathbf{z} | \mathbf{x})$$

Logarithm of intractable PDF  $P(\mathbf{x})$  is specified as follows (Ruthotto & Haber, 2021, p. 13):

$$\log P(\mathbf{x}) = \log \frac{P(\mathbf{x}, \mathbf{z})}{P(\mathbf{z} | \mathbf{x})} = \log \left( \frac{P(\mathbf{x}, \mathbf{z}) Q(\mathbf{z} | \mathbf{x})}{Q(\mathbf{z} | \mathbf{x}) P(\mathbf{z} | \mathbf{x})} \right)$$

This implies:

$$\log P(\mathbf{x}) = \log \frac{P(\mathbf{x}, \mathbf{z})}{Q(\mathbf{z} | \mathbf{x})} + \log \frac{Q(\mathbf{z} | \mathbf{x})}{P(\mathbf{z} | \mathbf{x})} \quad (3.4)$$

The second term  $\log(Q(\mathbf{z} | \mathbf{x}) / P(\mathbf{z} | \mathbf{x}))$  is not variant because  $Q(\mathbf{z} | \mathbf{x})$  is approximated to  $P(\mathbf{z} | \mathbf{x})$ . Therefore, the first term  $\log(P(\mathbf{x}, \mathbf{z}) / Q(\mathbf{z} | \mathbf{x}))$  is called variation lower bound or *evidence lower bound* because it is variant. Let  $l(\mathbf{x}, \mathbf{z})$  be loss function or error function on VAE which is defined as the minus opposite of expectation of the evidence lower bound  $\log(P(\mathbf{x}, \mathbf{z}) / Q(\mathbf{z} | \mathbf{x}))$  given PDF  $Q(\mathbf{z} | \mathbf{x})$  with note that  $Q(\mathbf{z} | \mathbf{x})$  has formal probabilistic distribution.

$$l(\mathbf{x}, \mathbf{z}) = -E_{Q(\mathbf{z} | \mathbf{x})} \left( \log \frac{P(\mathbf{x}, \mathbf{z})}{Q(\mathbf{z} | \mathbf{x})} \right) = - \int_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}) \log \frac{P(\mathbf{x} | \mathbf{z}) P(\mathbf{z})}{Q(\mathbf{z} | \mathbf{x})} d\mathbf{z}$$

Loss function  $l(\mathbf{x}, \mathbf{z})$  is expended as follows:

$$\begin{aligned} l(\mathbf{x}, \mathbf{z}) &= -\log P(\mathbf{x} | \mathbf{z}) \int_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}) d\mathbf{z} + \int_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}) \left( \log \frac{Q(\mathbf{z} | \mathbf{x})}{P(\mathbf{z})} \right) d\mathbf{z} \\ &= -\log P(\mathbf{x} | \mathbf{z}) + \int_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}) \left( \log \frac{Q(\mathbf{z} | \mathbf{x})}{P(\mathbf{z})} \right) d\mathbf{z} \end{aligned}$$

Because  $Q(\mathbf{z} | \mathbf{x})$  and  $P(\mathbf{x} | \mathbf{z})$  depend on encoder  $f(\mathbf{x} | \Theta)$  and decoder  $g(\mathbf{z} | \Phi)$ , respectively, their parameters are  $\Theta$  and  $\Phi$ , respectively.

$$Q(\mathbf{z} | \mathbf{x}) = Q(\mathbf{z} | \Theta, \mathbf{x})$$

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

Exactly,  $Q(\mathbf{z} | \Theta, \mathbf{x})$  is encoder likelihood which is the same to the aforementioned  $P(\mathbf{z}' | \Theta, \mathbf{x})$  except that it is focused that  $Q(\mathbf{z} | \Theta, \mathbf{x})$  has formal probabilistic specification like normal distribution. Loss function  $l(\Theta, \Phi | \mathbf{x}, \mathbf{z})$ , which is now function of encoder parameter  $\Theta$  and decoder parameter  $\Phi$ , is written as follows (Ruthotto & Haber, 2021, p. 14):

$$l(\Theta, \Phi | \mathbf{x}, \mathbf{z}) = -\log P(\mathbf{x} | \Phi, \mathbf{z}) + \int_{\mathbf{z}} Q(\mathbf{z} | \Theta, \mathbf{x}) \left( \log \frac{Q(\mathbf{z} | \Theta, \mathbf{x})}{P(\mathbf{z})} \right) d\mathbf{z}$$

Firstly, please pay attention to the first term loss function  $l(\Theta, \Phi | \mathbf{x}, \mathbf{z})$  where  $P(\mathbf{x} | \Phi, \mathbf{z})$  depends only on  $\Phi$  although it can be considered as a conditional PDF of  $\mathbf{x}$  given  $\mathbf{z}$  because  $P(\mathbf{x} | \Phi, \mathbf{z})$  is defined for output layer containing only  $\mathbf{x}$  of decoder DNN  $g(\mathbf{x} | \Phi)$  whose input is  $\mathbf{x}$ . Therefore, we had the following assertion:

$$\int_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}) \log P(\mathbf{x} | \Phi, \mathbf{z}) d\mathbf{z} = \log P(\mathbf{x} | \Phi, \mathbf{z}) \int_{\mathbf{z}} Q(\mathbf{z} | \mathbf{x}) d\mathbf{z} = \log P(\mathbf{x} | \Phi, \mathbf{z}) = \log P(\mathbf{x} | \Phi)$$

Secondly, the second term in loss function  $l(\Theta, \Phi | \mathbf{x}, \mathbf{z})$  is, actually, Kullback-Leibler divergence of encoder likelihood  $Q(\mathbf{z} | \Theta, \mathbf{x})$  and predefined tractable PDF  $P(\mathbf{z})$ , which measure the difference between  $Q(\mathbf{z} | \Theta, \mathbf{x})$  and  $P(\mathbf{z})$ . As a convention, this Kullback-Leibler divergence is called encoder KL divergence which is an ideal encoder error.

$$\text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z})) = E_{Q(\mathbf{z} | \mathbf{x})} \left( \log \frac{Q(\mathbf{z} | \Theta, \mathbf{x})}{P(\mathbf{z})} \right) = \int_{\mathbf{z}} Q(\mathbf{z} | \Theta, \mathbf{x}) \left( \log \frac{Q(\mathbf{z} | \Theta, \mathbf{x})}{P(\mathbf{z})} \right) d\mathbf{z}$$

The smaller the encoder KL divergence is, the closer the encoder likelihood  $Q(\mathbf{z} | \Theta, \mathbf{x})$  is to tractable PDF  $P(\mathbf{z})$ , the better the encoder DNN  $f(\mathbf{x} | \Theta)$  is. Loss function is rewritten again:

$$l(\Theta, \Phi | \mathbf{x}, \mathbf{z}) = -\log P(\mathbf{x} | \Phi, \mathbf{z}) + \text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z}))$$

Or,

$$l(\Theta, \Phi | \mathbf{x}, \mathbf{z}) = -\log P(\mathbf{x} | \Phi) + \text{KL}(Q(\mathbf{z} | \Theta) | P(\mathbf{z})) \quad (3.5)$$

According to the two problem of construct a DGM, the first term  $-\log(P(\mathbf{x} | \Phi, \mathbf{z}))$  in loss function indicates the first problem of how to train decoder DNN  $g(\mathbf{z} | \Phi)$  which is called reconstruction error in literature and the second term  $\text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z}))$  in loss function indicates the second problem of how to qualify training task for training encoder DNN  $f(\mathbf{x} | \Theta)$  which is called regularity in literature. Loss function  $l(\Theta, \Phi | \mathbf{x}, \mathbf{z})$  is minimized to estimate  $\Theta$  and  $\Phi$  as follows:

$$\begin{aligned} \Theta^* &= \underset{\Theta}{\text{armin}} l(\Theta, \Phi | \mathbf{x}, \mathbf{z}) \\ \Phi^* &= \underset{\Phi}{\text{armin}} l(\Theta, \Phi | \mathbf{x}, \mathbf{z}) \end{aligned} \quad (3.6)$$

Because  $P(\mathbf{x} | \Theta, \mathbf{z})$  depends only on  $\Theta$  and encoder KL divergence  $\text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z}))$  depends only on  $\Phi$ , the optimization problem is specified as follows:

$$\begin{aligned} \Theta^* &= \underset{\Theta}{\text{armin}} \text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z})) \\ \Phi^* &= \underset{\Phi}{\text{armin}} (-\log P(\mathbf{x} | \Phi, \mathbf{z})) = \underset{\Phi}{\text{armax}} (\log P(\mathbf{x} | \Phi, \mathbf{z})) \end{aligned}$$

Which results estimation equations according to SGD:

$$\begin{aligned} \Theta &= \Theta - \gamma \nabla \text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z})) \\ \Phi &= \Phi + \gamma \nabla \log P(\mathbf{x} | \Phi, \mathbf{z}) \end{aligned} \quad (3.7)$$

Where  $\nabla \text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z}))$  is gradient of encoder KL divergence  $\text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z}))$  with regard to encoder parameter  $\Theta$ . Note that tractable PDF  $P(\mathbf{z})$  is predefined (fixed). While  $Q(\mathbf{z} | \Theta, \mathbf{x})$  is called encoder likelihood,  $P(\mathbf{x} | \Phi, \mathbf{z})$  is called decoder likelihood. On the other hand, while  $P(\mathbf{z})$  is prior PDF of intractable data  $\mathbf{z}$ , then  $Q(\mathbf{z} | \Theta, \mathbf{x})$  is approximated posterior PDF of  $\mathbf{z}$  given  $\mathbf{x}$  where both  $P(\mathbf{z})$  and  $Q(\mathbf{z} | \Theta, \mathbf{x})$  have formal probabilistic specifications and moreover,  $P(\mathbf{z})$  is fixed (predefined).

$$Q(\mathbf{z} | \Theta, \mathbf{x}) = P(\mathbf{z}' | \Theta, \mathbf{x}) \cong P(\mathbf{z} | \Theta, \mathbf{x})$$

Both  $P(\mathbf{z} | \Theta, \mathbf{x})$  and  $Q(\mathbf{z} | \Theta, \mathbf{x})$  are encoder likelihood as well as posterior PDF of tractable data  $\mathbf{z}$  but  $Q(\mathbf{z} | \Theta, \mathbf{x})$  is approximated one whose probabilistic distribution is specified formally. Therefore (Ruthotto & Haber, 2021, p. 16), randomized data  $\mathbf{z}'$  in latent space  $Z$  is sampled from approximated distribution  $Q(\mathbf{z} | \Theta, \mathbf{x})$  instead of sampling from true distribution  $P(\mathbf{z} | \Theta, \mathbf{x})$ .

Given epoch of size  $N$  is denoted as  $\mathbf{D} = (\mathbf{d}^{(1)} = (\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), \mathbf{d}^{(2)} = (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, \mathbf{d}^{(N)} = (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$ , the estimation equations of  $\Theta$  and  $\Phi$  are extended exactly as epoch estimation at every iteration of SGD:

$$\begin{aligned} \Theta^{(k+1)} &= \Theta^{(k)} - \gamma \nabla \frac{1}{N} \sum_{i=1}^N \text{KL}(Q(\mathbf{z}^{(i)} | \Theta^{(k)}, \mathbf{x}^{(i)}) | P(\mathbf{z})) \\ \Phi^{(k+1)} &= \Phi^{(k)} + \gamma \nabla \frac{1}{N} \sum_{i=1}^N \log P(\mathbf{x}^{(i)} | \Phi, \mathbf{z}^{(i)}) \end{aligned}$$

Please distinguish that the tractable data  $\mathbf{z}^{(i)}$  in the first equation above follows distribution  $P(\mathbf{z})$  but the tractable data  $\mathbf{z}^{(i)}$  in the second equation above follows distribution  $Q(\mathbf{z} | \Theta, \mathbf{x})$ . As a result, VAE trained with SGD is specified as follows:

Initialize  $\Theta$  and  $\Phi$  and set  $k = 0$ .

Repeat

Sampling epoch  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$  or receiving epoch  $\mathbf{X}$  from big data / data stream.

Randomize random epoch  $\mathbf{Z} = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$  in which each  $\mathbf{z}^{(i)}$  is randomized from distribution  $Q(\mathbf{z} | \Theta^{(k)}, \mathbf{x}^{(i)})$ .

$$\Theta^{(k+1)} = \Theta^{(k)} - \gamma \nabla \frac{1}{N} \sum_{i=1}^N \text{KL} \left( Q(\mathbf{z}^{(i)} | \Theta^{(k)}, \mathbf{x}^{(i)}) \middle| P(\mathbf{z}) \right)$$

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma \nabla \frac{1}{N} \sum_{i=1}^N \log P(\mathbf{x}^{(i)} | \Phi, \mathbf{z}^{(i)})$$

Increase  $k = k + 1$ .

Until some terminating conditions are met.

Note, a terminating condition is customized, for example, parameters  $\Theta$  and  $\Phi$  are not changed significantly or there is no more coming epoch  $\mathbf{X}$ . Moreover, the index  $k$  indicates time point as well as iteration of SGD. Because PDF  $P(\mathbf{z})$  is predefined, it is easy to calculate encoder KL divergence  $\text{KL}(Q(\mathbf{z}^{(i)} | \Theta^{(k)}, \mathbf{x}^{(i)}) | P(\mathbf{z}))$  but it is necessary to define  $P(\mathbf{x})$  by well-known distribution. However, randomizing random epoch  $\mathbf{Z} = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$  from distribution  $Q(\mathbf{z} | \Theta^{(k)}, \mathbf{x}^{(i)})$  is not easy and so, VAE trained with SGD will be fine-tuned. It is interesting that when  $Q(\mathbf{z} | \Theta^{(k)}, \mathbf{x}^{(i)})$  is posterior PDF of  $\mathbf{z}$  and  $P(\mathbf{z})$  is prior PDF of  $\mathbf{z}$ , the event that  $\mathbf{z}$  is randomized from the posterior PDF  $Q(\mathbf{z} | \Theta^{(k)}, \mathbf{x}^{(i)})$  and  $Q(\mathbf{z} | \Theta^{(k)}, \mathbf{x}^{(i)})$  itself is updated continuously based on its previous evidence  $\mathbf{x}^{(i)}$  over SGD iterations implies that VAE conforms Bayesian statistics in estimation. Moreover,  $P(\mathbf{z})$  is an alignment that  $Q(\mathbf{z} | \Theta^{(k)}, \mathbf{x}^{(i)})$  adjusts itself with support of encoder KL divergence  $\text{KL}(Q(\mathbf{z}^{(i)} | \Theta^{(k)}, \mathbf{x}^{(i)}) | P(\mathbf{z}))$ .

Because encoder likelihood  $Q(\mathbf{z} | \Theta, \mathbf{x})$  must always have formal probabilistic distribution, it is assumed to follow multivariate normal distribution in practice. Therefore, let  $\mu(\mathbf{x} | \Theta)$  and  $\Sigma(\mathbf{x} | \Theta)$  be mean vector and covariance matrix of  $\mathbf{z}$ , then encoder likelihood  $Q(\mathbf{z} | \Theta, \mathbf{x})$  becomes  $Q(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  so that output of encoder DNN  $f(\mathbf{x} | \Theta)$  is mean  $\mu(\mathbf{x} | \Theta)$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta)$  while its input is  $\mathbf{x}$  and its weights are  $\Theta$ , of course. Please pay attention to the fact that output of encoder DNN  $f(\mathbf{x} | \Theta)$  is now  $\mu(\mathbf{x} | \Theta)$  and  $\Sigma(\mathbf{x} | \Theta)$  which are corresponding to  $\mathbf{z}$ . Moreover,  $\mu(\mathbf{x} | \Theta)$  and  $\Sigma(\mathbf{x} | \Theta)$  are functions of  $\mathbf{x}$ , whose parameter is  $\Theta$ .

$$f(\mathbf{x} | \Theta) = \begin{pmatrix} \mu(\mathbf{x} | \Theta) \\ \Sigma(\mathbf{x} | \Theta) \end{pmatrix} \sim \mathbf{z} \quad (3.8)$$

$$Q(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) = \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$$

Note,  $\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  denotes multivariate normal distribution with mean  $\mu(\mathbf{x} | \Theta)$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta)$ .

$$\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) = (2\pi)^{-\frac{n}{2}} |\Sigma(\mathbf{x} | \Theta)|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (\mathbf{z} - \mu(\mathbf{x} | \Theta))^T (\Sigma(\mathbf{x} | \Theta))^{-1} (\mathbf{z} - \mu(\mathbf{x} | \Theta)) \right)$$

Note, dimension of tractable data  $\mathbf{z}$  is  $n$ . Moreover, notation  $|\cdot|$  or notation  $\det(\cdot)$  denotes determinant of matrix whereas  $(\Sigma(\mathbf{x} | \Theta))^{-1}$  is inverse of covariance matrix  $\Sigma(\mathbf{x} | \Theta)$  and the subscript “ $T$ ” denotes transposition operator of matrix and vector. It is easy to recognize that  $\mathbf{z}'$  is approximation of  $\mathbf{z}$ . When tractable PDF  $P(\mathbf{z})$  is fixed, it is often assumed to follow multivariate normal distribution with predefined mean  $\mu_0$  and predefined covariance matrix  $\Sigma_0$  as follows:

$$P(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mu_0, \Sigma_0) = \mathcal{N}(\mu_0, \Sigma_0) = (2\pi)^{-\frac{n}{2}} |\Sigma_0|^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (\mathbf{z} - \mu_0)^T \Sigma_0^{-1} (\mathbf{z} - \mu_0) \right)$$

Encoder KL divergence  $\text{KL}(Q(\mathbf{z} | \Theta, \mathbf{x}) | P(\mathbf{z}))$  between  $Q(\mathbf{z} | \Theta, \mathbf{x})$  and  $P(\mathbf{z})$  becomes encoder KL divergence  $\text{KL}(Q(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | P(\mathbf{z}))$  between  $Q(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta))$  and  $P(\mathbf{z})$  as follows:

$$\text{KL} \left( Q(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| P(\mathbf{z}) \right) = \text{KL} \left( \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| \mathcal{N}(\mathbf{z} | \mu_0, \Sigma_0) \right)$$

Which is, essentially, encoder KL divergence between two normal distributions,  $KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0))$ . As a convention, this divergence is called encoder KL divergence which is determined in literature as follows (Doersch, 2016, p. 9):

$$KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) \\ = \frac{1}{2} \left( \text{tr}(\Sigma_0^{-1} \Sigma(\mathbf{x} | \Theta)) + (\mu(\mathbf{x} | \Theta) - \mu_0)^T \Sigma_0^{-1} (\mu(\mathbf{x} | \Theta) - \mu_0) + \log \left( \frac{|\Sigma_0|}{|\Sigma(\mathbf{x} | \Theta)|} \right) - n \right)$$

Where  $\text{tr}(\cdot)$  denotes trace operator of square matrix which is sum of elements on main diagonal, for instance, given  $n \times n$  matrix  $A$ , then  $\text{tr}(A) = a_{11} + a_{22} + \dots + a_{nn}$  with note that  $a_{ij}$  is the element at row  $i$  and column  $j$ . Moreover, notation  $|\cdot|$  or notation  $\det(\cdot)$  denotes determinant of matrix. Gradient of encoder KL divergence consists of two elemental gradients with regard to mean  $\mu(\mathbf{x} | \Theta)$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta)$ .

$$\nabla KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) \\ = \begin{pmatrix} \nabla_{\mu} KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) \\ \nabla_{\Sigma} KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) \end{pmatrix}^T$$

Where,

$$\nabla_{\mu} KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) \\ = \frac{dKL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0))}{d\mu(\mathbf{x} | \Theta)} \frac{d\mu(\mathbf{x} | \Theta)}{d\Theta} \\ \nabla_{\Sigma} KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) \\ = \frac{dKL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0))}{d\Sigma(\mathbf{x} | \Theta)} \frac{d\Sigma(\mathbf{x} | \Theta)}{d\Theta}$$

Where  $d\mu(\mathbf{x} | \Theta) / d\Theta$  and  $d\Sigma(\mathbf{x} | \Theta) / d\Theta$  are differentials of  $\mu(\mathbf{x} | \Theta)$  and  $\Sigma(\mathbf{x} | \Theta)$  with regard to  $\Theta$ , respectively. It is not difficult to calculate KL gradient  $\nabla_{\mu}$ :

$$\nabla_{\mu} KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) = (\mu(\mathbf{x} | \Theta) - \mu_0)^T \Sigma_0^{-1} \frac{d\mu(\mathbf{x} | \Theta)}{d\Theta} \quad (3.9)$$

Due to (Nguyen, Matrix Analysis and Calculus, 2015, p. 35):

$$\frac{d(\mu(\mathbf{x} | \Theta) - \mu_0)^T \Sigma_0^{-1} (\mu(\mathbf{x} | \Theta) - \mu_0)}{d\mu(\mathbf{x} | \Theta)} = 2(\mu(\mathbf{x} | \Theta) - \mu_0)^T \Sigma_0^{-1}$$

It is not difficult to calculate KL gradient  $\nabla_{\Sigma}$  too:

$$\nabla_{\Sigma} KL(\mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) | \mathcal{N}(\mu_0, \Sigma_0)) = \frac{1}{2} (\Sigma_0^{-1} - (\Sigma(\mathbf{x} | \Theta))^{-1}) \frac{d\Sigma(\mathbf{x} | \Theta)}{d\Theta} \quad (3.10)$$

Due to (Nguyen, Matrix Analysis and Calculus, 2015, pp. 45-46):

$$\frac{d\text{tr}(\Sigma_0^{-1} \Sigma(\mathbf{x} | \Theta))}{d\Sigma(\mathbf{x} | \Theta)} = \Sigma_0^{-1} \\ \frac{d \log(|\Sigma(\mathbf{x} | \Theta)|)}{d\Sigma(\mathbf{x} | \Theta)} = \frac{d|\Sigma(\mathbf{x} | \Theta)|}{d\Sigma(\mathbf{x} | \Theta)} \frac{1}{|\Sigma(\mathbf{x} | \Theta)|} = \frac{|\Sigma(\mathbf{x} | \Theta)| (\Sigma(\mathbf{x} | \Theta))^{-1}}{|\Sigma(\mathbf{x} | \Theta)|} = (\Sigma(\mathbf{x} | \Theta))^{-1}$$

As a result, encoder parameter  $\Theta$  consists of two elemental parameters according to with regard to mean  $\mu(\mathbf{x} | \Theta)$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta)$  as follows:

$$\Theta = \begin{pmatrix} \Theta_{\mu} \\ \Theta_{\Sigma} \end{pmatrix}$$

Where,

$$\Theta_\mu = (\mu_1, \mu_2, \dots, \mu_n)^T$$

$$\Theta_\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix}$$

Note, given random vector  $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$  whose elements  $z_i$  are random variables too,  $\sigma_{ij}$  where  $i \neq j$  is covariance between two random variables  $z_i$  and  $z_j$  and  $\sigma_i^2$  is variance of random variable  $z_i$ . It is easy to calculate encoder parameters  $\Theta_\mu$  and  $\Theta_\Sigma$  by SGD estimation:

$$\begin{aligned} \Theta_\mu &= \Theta_\mu - \gamma \nabla_{\Theta_\mu} \text{KL} \left( \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| \mathcal{N}(\mathbf{z} | \mu_0, \Sigma_0) \right) \\ &= \Theta_\mu - \gamma \Sigma_0^{-1} (\mu(\mathbf{x} | \Theta_\mu) - \mu_0) \frac{d\mu(\mathbf{x} | \Theta_\mu)}{d\Theta_\mu} \end{aligned}$$

$$\begin{aligned} \Theta_\Sigma &= \Theta_\Sigma - \gamma \nabla_{\Theta_\Sigma} \text{KL} \left( \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta), \Sigma(\mathbf{x} | \Theta)) \middle| \mathcal{N}(\mathbf{z} | \mu_0, \Sigma_0) \right) \\ &= \Theta_\Sigma - \gamma \frac{1}{2} \left( \Sigma_0^{-1} - (\Sigma(\mathbf{x} | \Theta_\Sigma))^{-1} \right) \frac{d\Sigma(\mathbf{x} | \Theta_\Sigma)}{d\Theta_\Sigma} \end{aligned}$$

Where  $d\mu(\mathbf{x} | \Theta_\mu) / d\Theta_\mu$  and  $d\Sigma(\mathbf{x} | \Theta_\Sigma) / d\Theta_\Sigma$  are differentials of  $\mu(\mathbf{x} | \Theta_\mu)$  and  $\Sigma(\mathbf{x} | \Theta_\Sigma)$  with regard to  $\Theta_\mu$  and  $\Theta_\Sigma$ , respectively. In practice,  $P(\mathbf{z})$  is assumed to conform standard normal distribution with zero mean  $\mu_0 = \mathbf{0}$  and identity covariance matrix  $\Sigma_0 = I$  where  $I$  is identity matrix so that encoder parameters  $\Theta_\mu$  and  $\Theta_\Sigma$  are computed effectively.

$$\begin{aligned} \Theta_\mu &= \Theta_\mu - \gamma \mu(\mathbf{x} | \Theta_\mu) \frac{d\mu(\mathbf{x} | \Theta_\mu)}{d\Theta_\mu} \\ \Theta_\Sigma &= \Theta_\Sigma - \gamma \frac{1}{2} \left( I - (\Sigma(\mathbf{x} | \Theta_\Sigma))^{-1} \right) \frac{d\Sigma(\mathbf{x} | \Theta_\Sigma)}{d\Theta_\Sigma} \end{aligned} \tag{3.11}$$

In order to improve more computational effectiveness, it is possible to suppose that elemental variables  $z_i$  in  $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$  within context  $P(\mathbf{z})$  are mutually independent so that covariance  $\sigma_{ij}$  between two variables  $z_i$  and  $z_j$  where  $i \neq j$  is 0, which results that there only exist variances  $\sigma_i^2$  of  $z_i$ . Covariance matrix  $\Sigma(\mathbf{x} | \Theta)$  becomes diagonal matrix:

$$\Sigma(\mathbf{x} | \Theta_\Sigma) = (\sigma^2(\mathbf{x} | \Theta_\Sigma))_{n \times n} = \begin{pmatrix} \sigma_1^2(\mathbf{x} | \Theta_\Sigma) & 0 & \dots & 0 \\ 0 & \sigma_2^2(\mathbf{x} | \Theta_\Sigma) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2(\mathbf{x} | \Theta_\Sigma) \end{pmatrix}$$

Note,

$$\sigma^2(\mathbf{x} | \Theta_\Sigma) = \begin{pmatrix} \sigma_1^2(\mathbf{x} | \Theta_\Sigma) \\ \sigma_2^2(\mathbf{x} | \Theta_\Sigma) \\ \vdots \\ \sigma_n^2(\mathbf{x} | \Theta_\Sigma) \end{pmatrix}$$

Where  $\sigma_i^2(\mathbf{x} | \Theta)$  is variance of elemental variable  $x_i$  in  $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$  given  $\mathbf{x}$  according to encoder  $f(\mathbf{x} | \Theta)$ . As a result, encoder parameter  $\Theta_\Sigma$ , which is now diagonal matrix represented by its diagonal vector  $\Theta_{\sigma^2}$ , is computed easier.

$$\Theta_{\sigma^2} = \Theta_{\sigma^2} - \gamma \frac{1}{2} \left( \mathbf{1} - (\sigma^2(\mathbf{x} | \Theta_{\sigma^2}))^{-1} \right) \frac{d\sigma^2(\mathbf{x} | \Theta_{\sigma^2})}{d\Theta_{\sigma^2}}$$

Where,



$$\Theta_{\Sigma} = (\Theta_{\sigma^2})_{n \times n} = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{pmatrix}$$

$$\Theta_{\sigma^2} = (\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)^T$$

$$\sigma^2(\mathbf{x}|\Theta_{\sigma^2}) = \begin{pmatrix} \sigma_1^2(\mathbf{x}|\Theta_{\sigma^2}) \\ \sigma_2^2(\mathbf{x}|\Theta_{\sigma^2}) \\ \vdots \\ \sigma_n^2(\mathbf{x}|\Theta_{\sigma^2}) \end{pmatrix}$$

$$(\sigma^2(\mathbf{x}|\Theta_{\sigma^2}))^{-1} = \begin{pmatrix} 1/\sigma_1^2(\mathbf{x}|\Theta_{\sigma^2}) \\ 1/\sigma_2^2(\mathbf{x}|\Theta_{\sigma^2}) \\ \vdots \\ 1/\sigma_n^2(\mathbf{x}|\Theta_{\sigma^2}) \end{pmatrix}$$

$$\mathbf{1} = (1, 1, \dots, 1)^T$$

In general, estimation equations for encoder parameter  $\Theta = (\Theta_{\mu}, \Theta_{\sigma^2})^T$  are specified as follows:

$$\begin{aligned} \Theta_{\mu} &= \Theta_{\mu} - \gamma \mu(\mathbf{x}|\Theta_{\mu}) \frac{d\mu(\mathbf{x}|\Theta_{\mu})}{d\Theta_{\mu}} \\ \Theta_{\sigma^2} &= \Theta_{\sigma^2} - \gamma \frac{1}{2} \left( \mathbf{1} - (\sigma^2(\mathbf{x}|\Theta_{\sigma^2}))^{-1} \right) \frac{d\sigma^2(\mathbf{x}|\Theta_{\sigma^2})}{d\Theta_{\sigma^2}} \end{aligned} \quad (3.12)$$

Where  $d\sigma^2(\mathbf{x} | \Theta_{\sigma^2}) / d\Theta_{\sigma^2}$  is differential of  $\sigma^2(\mathbf{x} | \Theta_{\sigma^2})$  with regard to  $\Theta_{\sigma^2}$ .

There can be no change in estimating decoder parameter  $\Phi$  within VAE so that decoder log-likelihood  $\log(P(\mathbf{x} | \Phi, \mathbf{z}))$  is maximized.

$$\Phi^* = \underset{\Phi}{\operatorname{argmin}} (-\log P(\mathbf{x}|\Phi, \mathbf{z})) = \underset{\Phi}{\operatorname{argmax}} (\log P(\mathbf{x}|\Phi, \mathbf{z}))$$

As usual, decoder likelihood  $P(\mathbf{x} | \Phi, \mathbf{z})$  is assumed to distribute normally with mean  $\delta$  and variance  $\sigma^2$ .

$$P(\mathbf{x}|\Phi, \mathbf{z}) = (2\pi\sigma^2)^{-\frac{m}{2}} \exp\left(-\frac{\|g(\mathbf{z}|\Phi) - \mathbf{x} - \delta\|^2}{2\sigma^2}\right)$$

Which implies decoder log-likelihood  $\log(P(\mathbf{x} | \Phi, \mathbf{z}))$  as follows:

$$\log P(\mathbf{x}|\Phi, \mathbf{z}) = -\frac{m}{2} \log(2\pi\sigma^2) - \frac{\|g(\mathbf{z}|\Phi) - \mathbf{x} - \delta\|^2}{2\sigma^2}$$

Where  $\|\cdot\|$  denotes Euclidean norm of vector. Gradient of decoder log-likelihood is:

$$\nabla \log P(\mathbf{x}|\Phi, \mathbf{z}) = \frac{d \log P(\mathbf{x}|\Phi, \mathbf{z})}{d\Phi} = -\frac{(g(\mathbf{z}|\Phi) - \mathbf{x} - \delta)}{\sigma^2} \frac{dg(\mathbf{z}|\Phi)}{d\Phi}$$

Where  $dg(\mathbf{z} | \Phi) / d\Phi$  is differential of  $g(\mathbf{z} | \Phi)$  with regard to  $\Phi$ . Let  $\delta = \mathbf{0}$  and  $\sigma^2=1$  optimization, we have:

$$\nabla \log P(\mathbf{x}|\Phi, \mathbf{z}) = -(g(\mathbf{z}|\Phi) - \mathbf{x}) \frac{dg(\mathbf{z}|\Phi)}{d\Phi}$$

Which implies estimation equation for decoder parameter  $\Phi$  by SGD as follows:

$$\Phi = \Phi + \gamma \nabla \log P(\mathbf{x}|\Phi, \mathbf{z}) = \Phi - \gamma (g(\mathbf{z}|\Phi) - \mathbf{x}) \frac{dg(\mathbf{z}|\Phi)}{d\Phi}$$

Because data  $\mathbf{z}$  in the decoder estimation equation above follows encoder likelihood  $Q(\mathbf{z} | \Theta, \mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma})) = \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma}))$  rather than tractable PDF  $P(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mu_0, \Sigma_0)$ , it is denoted as  $\mathbf{z}'$  such that:

$$\Phi = \Phi + \gamma \nabla \log P(\mathbf{x}|\Phi, \mathbf{z}') = \Phi - \gamma (g(\mathbf{z}'|\Phi) - \mathbf{x}) \frac{dg(\mathbf{z}'|\Phi)}{d\Phi}$$

Given epoch of size  $N$  is denoted as  $\mathbf{D} = (\mathbf{d}^{(1)} = (\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), \mathbf{d}^{(2)} = (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, \mathbf{d}^{(N)} = (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$ , the estimation equations of  $\Theta$  and  $\Phi$  are extended exactly as epoch estimation at every iteration of SGD:

$$\begin{aligned}\Theta_{\mu}^{(k+1)} &= \Theta_{\mu}^{(k)} - \gamma \frac{1}{N} \sum_{i=1}^N \mu(\mathbf{x}^{(i)} | \Theta_{\mu}^{(k)}) \frac{d\mu(\mathbf{x}^{(i)} | \Theta_{\mu}^{(k)})}{d\Theta_{\mu}} \\ \Theta_{\sigma^2}^{(k+1)} &= \Theta_{\sigma^2}^{(k)} - \gamma \frac{1}{2N} \sum_{i=1}^N \left( 1 - \left( \sigma^2(\mathbf{x}^{(i)} | \Theta_{\sigma^2}^{(k)}) \right)^{-1} \right) \frac{d\sigma^2(\mathbf{x}^{(i)} | \Theta_{\sigma^2}^{(k)})}{d\Theta_{\sigma^2}} \\ \Phi^{(k+1)} &= \Phi^{(k)} - \gamma \frac{1}{N} \sum_{i=1}^N (g(\mathbf{z}'^{(i)} | \Phi^{(k)}) - \mathbf{x}^{(i)}) \frac{dg(\mathbf{z}'^{(i)} | \Phi^{(k)})}{d\Phi}\end{aligned}\tag{3.13}$$

As a result, VAE trained with SGD is specified as follows:

Initialize  $\Theta = (\Theta_{\mu}, \Theta_{\sigma^2})^T$  and  $\Phi$  and set  $k = 0$ .

Repeat

Sampling epoch  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$  or receiving epoch  $\mathbf{X}$  from big data / data stream.

$$\begin{aligned}\Theta_{\mu}^{(k+1)} &= \Theta_{\mu}^{(k)} - \gamma \frac{1}{N} \sum_{i=1}^N \mu(\mathbf{x}^{(i)} | \Theta_{\mu}^{(k)}) \frac{d\mu(\mathbf{x}^{(i)} | \Theta_{\mu}^{(k)})}{d\Theta_{\mu}} \\ \Theta_{\sigma^2}^{(k+1)} &= \Theta_{\sigma^2}^{(k)} - \gamma \frac{1}{2N} \sum_{i=1}^N \left( 1 - \left( \sigma^2(\mathbf{x}^{(i)} | \Theta_{\sigma^2}^{(k)}) \right)^{-1} \right) \frac{d\sigma^2(\mathbf{x}^{(i)} | \Theta_{\sigma^2}^{(k)})}{d\Theta_{\sigma^2}}\end{aligned}$$

Randomize random epoch  $\mathbf{Z} = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$  from standard normal distribution  $P(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I)$  with mean  $\mathbf{0}$  and identity covariance matrix  $I$ . For each randomized data  $\mathbf{z}^{(i)}$ , let  $\mathbf{z}'^{(i)}$  be calculated based on  $\mathbf{z}^{(i)}$  so that  $\mathbf{z}'^{(i)}$  follows multivariate normal distribution  $Q(\mathbf{z}' | \mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma})) = \mathcal{N}(\mathbf{z}' | \mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma}))$  with mean  $\mu(\mathbf{x} | \Theta_{\mu})$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta_{\Sigma})$  with note that  $\Theta_{\Sigma} = (\Theta_{\sigma^2})_{n \times n}$  is diagonal matrix.

$$\begin{aligned}\mathbf{z}'^{(i)} &= \mu(\mathbf{x}^{(i)} | \Theta_{\mu}^{(k)}) + \left( \Sigma(\mathbf{x}^{(i)} | \Theta_{\Sigma}^{(k)}) \right)^{\frac{1}{2}} \mathbf{z}^{(i)} \\ \Phi^{(k+1)} &= \Phi^{(k)} - \gamma \frac{1}{N} \sum_{i=1}^N (g(\mathbf{z}'^{(i)} | \Phi^{(k)}) - \mathbf{x}^{(i)}) \frac{dg(\mathbf{z}'^{(i)} | \Phi^{(k)})}{d\Phi}\end{aligned}$$

Increase  $k = k + 1$ .

Until some terminating conditions are met.

Note, a terminating condition is customized, for example, parameters  $\Theta$  and  $\Phi$  are not changed significantly or there is no more coming epoch  $\mathbf{X}$ . Moreover, the index  $k$  indicates time point as well as iteration of SGD. Because it is not easy to randomize  $\mathbf{z}$  according to normal distribution  $Q(\mathbf{z} | \mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma})) = \mathcal{N}(\mathbf{z} | \mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma}))$  with mean  $\mu(\mathbf{x} | \Theta_{\mu})$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta_{\Sigma})$ , there is a trick that simple data  $\mathbf{z}$  is randomized firstly by simple normal distribution  $P(\mathbf{z}) = \mathcal{N}(\mathbf{0}, I)$  with mean  $\mathbf{0}$  and identity covariance matrix  $I$  and, then random data  $\mathbf{z}'$  is calculated based on  $\mathbf{z}$  and  $\mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma})$  as follows:

$$\mathbf{z}' = \mu(\mathbf{x} | \Theta_{\mu}) + \left( \Sigma(\mathbf{x} | \Theta_{\Sigma}) \right)^{\frac{1}{2}} \mathbf{z}\tag{3.14}$$

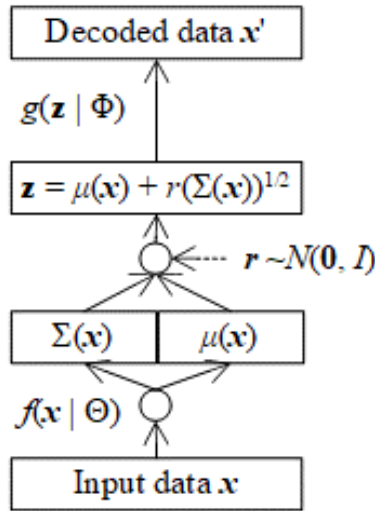
Such that  $\mathbf{z}'$  follows normal distribution  $\mathcal{N}(\mathbf{z}' | \mu(\mathbf{x} | \Theta_{\mu}), \Sigma(\mathbf{x} | \Theta_{\Sigma}))$  with mean  $\mu(\mathbf{x} | \Theta_{\mu})$  and covariance matrix  $\Sigma(\mathbf{x} | \Theta_{\Sigma})$  according to some rule of normal distribution in applied statistics (Hardle & Simar, 2013, p. 157). The notation  $A = \Sigma(\mathbf{x} | \Theta_{\Sigma})^{1/2}$  implies  $AA = \Sigma(\mathbf{x} | \Theta_{\Sigma})$  and so, we can consider it as square root of  $\Sigma(\mathbf{x} | \Theta_{\Sigma})$ . Calculating this square root is not so easy because of complexity of singular decomposition for calculating it. Fortunately, it is easier to calculate the square root when  $\Theta_{\Sigma}$  was simplified by diagonal elements  $(\sigma^2(\mathbf{x} | \Theta_{\Sigma}))_{n \times n}$ . Indeed, we have:

$$(\Sigma(\mathbf{x}|\Theta_\Sigma))^{\frac{1}{2}} = ((\sigma^2(\mathbf{x}|\Theta_\Sigma))_{n \times n})^{\frac{1}{2}} = \begin{pmatrix} \sqrt{\sigma_1^2(\mathbf{x}|\Theta_\Sigma)} & 0 & \cdots & 0 \\ 0 & \sqrt{\sigma_2^2(\mathbf{x}|\Theta_\Sigma)} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sqrt{\sigma_n^2(\mathbf{x}|\Theta_\Sigma)} \end{pmatrix}$$

Where,

$$\Theta_\Sigma = (\Theta_{\sigma^2})_{n \times n} = \begin{pmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{pmatrix}$$

Following figure depicts VAE.



**Figure 3.1.** Variational Autoencoders (VAE)

There is a question that how to calculate the differentials  $d\mu(\mathbf{x} | \Theta_\mu) / d\Theta_\mu$ ,  $d\sigma^2(\mathbf{x} | \Theta_{\sigma^2}) / d\Theta_{\sigma^2}$ , and  $dg(\mathbf{z}' | \Phi) / d\Phi$ . Indeed, it is not difficult to calculate them in context of neural network associated with backpropagation algorithm so that the last output layer as well as last neuron  $o$  of any DNN  $f(\mathbf{x} | \Theta)$  or  $g(\mathbf{z} | \Phi)$  is acted by activation function  $a(\cdot)$  as follows:

$$a(o) = a(\mathbf{w}^T \mathbf{i})$$

$$o = \mathbf{w}^T \mathbf{i}$$

Where  $\mathbf{i}$  is input of the last layer  $o$  and weight parameter  $\mathbf{w}$  is a part of entire parameter  $\Theta$  or  $\Phi$  and hence, we need to focus on calculating differential  $da(o) / d\mathbf{w}$  which is equivalent to any differential  $d\mu(\mathbf{x} | \Theta_\mu) / d\Theta_\mu$ ,  $d\sigma^2(\mathbf{x} | \Theta_{\sigma^2}) / d\Theta_{\sigma^2}$ , or  $dg(\mathbf{z}' | \Phi) / d\Phi$  so that backpropagation algorithm will solve remaining parts of entire parameter  $\Theta$  or  $\Phi$ .

$$\frac{da(o)}{d\mathbf{w}} \cong \text{any differential}$$

Indeed, we have:

$$\frac{da(o)}{d\mathbf{w}} = \frac{da(\mathbf{w}^T \mathbf{i})}{d\mathbf{w}} = \frac{da(o)}{do} \mathbf{i}^T = a'(o) \mathbf{i}^T$$

Note, the subscript “ $T$ ” denotes transposition operator of vector and matrix in which row vector becomes column vector and vice versa. It is easy to calculate the derivative  $a'(o)$  when activation function was specified, for instance, if  $a(o)$  is sigmoid function, we have:

$$a(o) = y = \frac{1}{1 + e^{-o}}$$

$$a'(o) = \frac{1}{1 + e^{-o}} \left( 1 - \frac{1}{1 + e^{-o}} \right) = a(o)(1 - a(o)) = y(1 - y)$$

In practice,  $y$  is replaced by  $a(y)$  in order to prevent  $o$  from being out of space:

$$a'(o) \cong a(y)(1 - a(y)) = a(a(o))(1 - a(a(o)))$$

As a result, we have:

$$\text{any differential} \cong a(a(o))(1 - a(a(o))) \mathbf{i}^T$$

For fast computation, it is possible to set the derivative  $a'(o)$  to be small enough constants like 1 such that any differential is  $\mathbf{i}^T$ .

Given epoch  $\mathbf{D} = (\mathbf{d}^{(1)} = (\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), \mathbf{d}^{(2)} = (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, \mathbf{d}^{(N)} = (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$  implies that the epoch is created or sent by equilateral distribution  $1/N$  but in general case,  $\mathbf{D}$  can follow an arbitrary distribution denoted by PDF  $P(\mathbf{d})$ , which makes the optimization problem and the SGD estimation changed a little bit by theoretical expectation given distribution  $P(\mathbf{d})$ .

$$\Theta^* = \underset{\Theta}{\operatorname{armin}} E \left( \text{KL}(Q(\mathbf{z}|\Theta, \mathbf{x})|P(\mathbf{z})) \middle| P(\mathbf{d}) \right)$$

$$\Phi^* = \underset{\Phi}{\operatorname{armax}} E(\log P(\mathbf{x}|\Phi, \mathbf{z}) | P(\mathbf{d}))$$

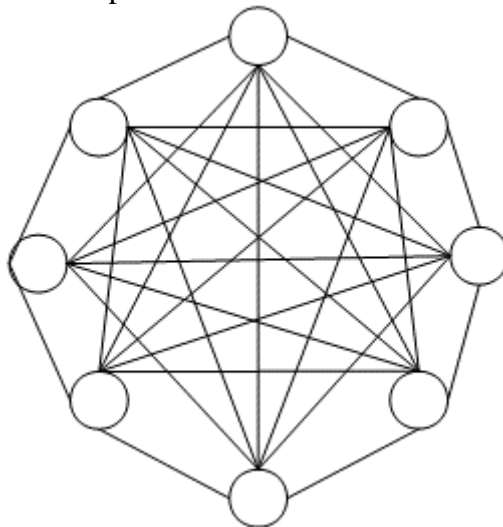
Where,

$$E \left( \text{KL}(Q(\mathbf{z}|\Theta, \mathbf{x})|P(\mathbf{d})) \middle| P(\mathbf{d}) \right) = \int_{\mathbf{d}} \text{KL}(Q(\mathbf{z}|\Theta, \mathbf{x})|P(\mathbf{d})) P(\mathbf{d}) d\mathbf{d}$$

$$E(\log P(\mathbf{x}|\Phi, \mathbf{z}) | P(\mathbf{d})) = \int_{\mathbf{d}} \log P(\mathbf{x}|\Phi, \mathbf{z}) P(\mathbf{d}) d\mathbf{d}$$

However, there is no significant change in aforementioned practical technique to estimate parameters.

Recall that the default artificial neural network is feedforward neural network where data is fed to input layer which, in turn, is evaluated and passed across hidden layers to output layer in one-way direction, finally. However, there is an extension of neural network, which is called recurrent neural work (RNN), where an output can be turned back in order to feed on network as input. In other words, RNN has circle, which allow that output can become input. There are many kinds of RNN, for instance, long short-term memory is a case of RNN aforementioned. Boltzmann machine (Wikipedia, Boltzmann machine, 2004) is another variant of RNN, in which there is no separation of inputs from outputs. Like Hopfield network (Wikipedia, Hopfield network, 2004), every neuron (unit) in Boltzmann machine connects to all remaining neurons. In other words, Boltzmann machine applies an interesting aspect that all input neurons are output neurons too.



**Figure 3.2.** Topology of Hopfield network and Boltzmann machine

Boltzmann machine named by the name of Austrian physicist Ludwig Eduard Boltzmann, also called Sherrington-Kirkpatrick model with external field or stochastic Ising-Lenz-Little

model, is a stochastic spin-glass model with an external field and classified as a Markov random field too. For easy explanation, Boltzmann machine simulates spinning glass process or annealing metal process, in which melt glass or melt metal will be frozen or get stable at some energy and some temperature where such energy and temperature are called stable energy and stable temperature at stable state of glass. The annealing process aims to reach the stable state of metal (glass) at which time the metal is frozen. Given concrete temperature, the smaller the energy is, the more stable the metal state is. Similarly, given concrete energy, the smaller the temperature is, the more stable the metal state is. Therefore, annealing process is cooling process where probability of metal state, which is proportional to energy and temperature, follows the so-called Boltzmann distribution specified as follows:

$$P(s) = \frac{\exp\left(-\frac{E(s)}{\kappa T}\right)}{\sum_{i=1}^M \exp\left(-\frac{E(s)}{\kappa T}\right)} \quad (3.15)$$

Where  $P(s)$  is probability of current state  $s$  and  $E(s)$  is energy applied to metal at state  $s$  given temperature  $T$  while  $\kappa$  is Boltzmann constant and  $M$  is the number of states. Note,  $T$  can be considered as a parameter. If the denominator is constant, Boltzmann probability is approximated as follows:

$$P(s) \cong \exp\left(-\frac{E(s)}{\kappa T}\right)$$

In annealing process, if next energy is concerned by observing current energy because of successive annealing process, energy deviation or energy difference  $\Delta E(s, s_{new})$  between current energy  $E(s)$  and next energy  $E(s_{new})$  is concerned so that Boltzmann probability derives a so-called acceptance probability  $P(s, s_{new}, T)$  as follows:

$$P(s, s_{new}, T) = \begin{cases} 1 & \text{if } \Delta E(s, s_{new}) < 0 \\ \exp\left(-\frac{\Delta E(s, s_{new})}{\kappa T}\right) & \text{otherwise} \end{cases}$$

Where,

$$\Delta E(s, s_{new}) = E(s_{new}) - E(s)$$

Given a certain temperature  $T$ , the larger the acceptance probability is, the higher likely the annealing process stops, the higher the likelihood of stability is. In other words, acceptance probability  $P(s, s_{new}, T)$  decides whether or not the new state  $s_{new}$  is moved next in annealing process. When applied into solving optimization problem as well as learning problem, simulated annealing (SA) algorithm codes candidate solution as states. Indeed, SA is iterative process including many enough iterations where SA decreases temperature  $T$  at each iteration and then, randomize a new state  $s_{new}$  and calculates energy  $E(s_{new})$  of the new state. Whether or not the new state (new candidate solution)  $s_{new}$  is based on the acceptance probability  $P(s, s_{new}, T)$  based on current state  $s$ , new state  $s_{new}$ , current temperature  $T$ . If the new candidate solution  $s_{new}$  is selected as current solution, SA will decrease temperature in the next iteration. Following is pseudo code of SA:

Initialize current temperature  $T$  by highest temperature  $T_0$  as  $T = T_0$ .

Repeat

    Decrease current temperature, for example,  $T = \text{decrease}(T)$ .

    Select a random neighbor of current state as  $s_{new} = \text{neighbor}(s)$ .

    If  $P(s, s_{new}, T)$  is larger than a predefined threshold then

$s = s_{new}$

    End if

Until terminating conditions are met.



The terminating conditions can be that best state (solution) is reached, the current state  $s$  is good enough, the current temperature  $T$  is low enough, or the number of iterations is large enough. As usual, given a maximum iteration number  $K$  and the current iteration number  $k$ , the temperature decreasing function can be defined as follows:

$$\text{decrease}(T) = T - \frac{k}{K}T$$

It is easy to infer that it is possible to set the initial temperature to be the maximum number of iterations as  $T_0 = K$  in practice. There is no significant change when applying SA into training Boltzmann machine where the most important problem is how to specify energy of Boltzmann machine. Fortunately, global energy of Boltzmann machine inherits from global energy of Hopfield network because Boltzmann machine is a type of Hopfield network which in turn is a variant of RNN. Suppose an entire Boltzmann machine is represented by a vector  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  in which each  $x_i$  is a neuron or unit. It is exact that a certain state of Boltzmann machine is represented by  $\mathbf{x}$  which is evaluated at certain time point. It is possible to denote current state of Boltzmann machine as  $\mathbf{x}$  instead. For convenience, the next state of Boltzmann machine is denoted  $\mathbf{x}'$ . Energy  $E(\mathbf{x})$  of Boltzmann machine at state  $\mathbf{x}$  is defined based on global energy of Hopfield network as follows (Hinton, 2007, p. 2):

$$E(\mathbf{x}) = - \sum_{i=1}^{n-1} \sum_{j=i+1}^n w_{ij} x_i x_j - \sum_{i=1}^n b_i x_i$$

Note,  $w_{ij}$  is weight between neuron  $x_i$  and neuron  $x_j$  whereas  $b_i$  is bias of neuron  $x_i$ . As usual, biases  $b_i$  are considered as parameters like weights  $w_{ij}$ . Because there are  $n(n-1)/2$  connections as well as  $n(n-1)/2$  weights, the equation of energy is rewritten for convenience as follows (Wikipedia, Boltzmann machine, 2004):

$$E(\mathbf{x}) = - \sum_{i < j} w_{ij} x_i x_j - \sum_i b_i x_i \quad (3.16)$$

All weights  $w_{ij}$  compose weight matrix  $\mathbf{W} = (w_{ij})_{n \times n}$  whose elements on diagonal are zero. Note,  $\mathbf{W}$  is  $n \times n$  symmetric matrix.

$$w_{ij} = w_{ji}, \forall i, j$$

$$w_{ii} = 0, \forall i \text{ by convention}$$

Every neuron  $x_i$  is evaluated by propagation rule:

$$x_i = \sum_{j=1, j \neq i}^n w_{ij} x_j$$

Neurons in traditional Boltzmann machine are binary variables such that  $x_i$  belongs to  $\{0, 1\}$  but it is extended to allow neurons  $x_i$  to belong to arbitrary real interval and so, suppose every  $x_i$  ranges in interval  $[0, 1]$  without loss of generality. Rectified Linear Unit (ReLU) function is used to ramp  $x_i$  in interval  $[0, 1]$  so as to modify the propagation rule a little bit but learning algorithm mentioned later is not changed because the first-order derivative of ReLU function within valid domain  $[0, 1]$  is 1.

$$x_i = \text{ReLU} \left( \sum_{j=1, j \neq i}^n w_{ij} x_j \right)$$

Where

$$\text{ReLU}(x) = \begin{cases} x & \text{if } 0 \leq x \leq 1 \\ 0 & \text{if } x < 0 \\ 1 & \text{if } x > 1 \end{cases}$$

It implies:

$$x_i = \text{ReLU}(x_i)$$

So that the propagation rule is not changed in theory:

$$x_i = \sum_{j=1, j \neq i}^n w_{ij} x_j$$

Based on definition of global energy, Boltzmann probability density function (PDF) of Boltzmann machine is determined as follows:

$$P(\mathbf{x}|\mathbf{W}, \mathbf{b}) = \frac{\exp\left(-\frac{E(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\kappa T}\right)}{\int_{\mathbf{x}} \exp\left(-\frac{E(\mathbf{x}|\mathbf{W})}{\kappa T}\right) d\mathbf{x}}$$

Recall that:

$$E(\mathbf{x}|\mathbf{W}, \mathbf{b}) = -\sum_{i < j} w_{ij} x_i x_j - \sum_i b_i x_i$$

$$\mathbf{b} = (b_1, b_2, \dots, b_n)^T$$

Within context of DGM, such PDF is generator likelihood whose parameter is  $\Phi = (\mathbf{W}, \mathbf{b})$ .

$$P(\mathbf{x}|\Phi) = P(\mathbf{x}|\mathbf{W}, \mathbf{b}) = \frac{\exp\left(-\frac{E(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\kappa T}\right)}{\int_{\mathbf{x}} \exp\left(-\frac{E(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\kappa T}\right) d\mathbf{x}}$$

Because the denominator is constant with regard to  $\mathbf{W}$  and  $\mathbf{b}$ , Boltzmann PDF is approximated as follows:

$$P(\mathbf{x}|\mathbf{W}, \mathbf{b}) \cong \exp\left(-\frac{E(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\kappa T}\right) \quad (3.17)$$

For learning Boltzmann, maximum likelihood estimation (MLE) method (Goodfellow, Bengio, & Courville, Deep Learning, 2016, p. 655) is applied into estimating weight parameter  $\mathbf{W}$  and bias parameter  $\mathbf{b}$  by maximizing Boltzmann PDF with regard to  $w_{ij}$  and  $b_i$ .

$$w_{ij}^* = \max_{w_{ij}} P(\mathbf{x}|\mathbf{W}, \mathbf{b})$$

$$b_i^* = \max_{b_i} P(\mathbf{x}|\mathbf{W}, \mathbf{b})$$

By taking natural logarithm of Boltzmann PDF, the optimization becomes easier to be solved.

$$w_{ij}^* = \max_{w_{ij}} \log P(\mathbf{x}|\mathbf{W}, \mathbf{b})$$

$$b_i^* = \max_{b_i} \log P(\mathbf{x}|\mathbf{W}, \mathbf{b})$$

Where  $\log P(\mathbf{x} | \mathbf{W}, \mathbf{b})$  is called Boltzmann log-likelihood or Boltzmann log-PDF.

$$\log P(\mathbf{x}|\mathbf{W}, \mathbf{b}) \cong -\frac{E(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\kappa T}$$

The first-order partial derivatives of Boltzmann log-likelihood are:

$$\nabla_{w_{ij}} \log P(\mathbf{x}|\mathbf{W}, \mathbf{b}) \cong \frac{\partial \log P(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\partial w_{ij}} = \frac{\partial \log E(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\partial w_{ij}} = \frac{x_i x_j}{\kappa T}$$

$$\nabla_{b_i} \log P(\mathbf{x}|\mathbf{W}, \mathbf{b}) \cong \frac{\partial \log P(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\partial b_i} = \frac{\partial \log E(\mathbf{x}|\mathbf{W}, \mathbf{b})}{\partial b_i} = \frac{x_i}{\kappa T}$$

As a convention, these first-order partial derivatives are called (partial) gradients. By applying stochastic gradient descent (SGD) method into estimating  $w_{ij}$  and  $b_i$  given Boltzmann log-likelihood, we have:

$$\begin{aligned} w_{ij} &= w_{ij} + \gamma \nabla_{w_{ij}} \log P(\mathbf{x}|\mathbf{W}, \mathbf{b}) = w_{ij} + \gamma \frac{x_i x_j}{\kappa T} \\ b_i &= b_i + \gamma \nabla_{b_i} \log P(\mathbf{x}|\mathbf{W}, \mathbf{b}) = b_i + \gamma \frac{x_i}{\kappa T} \end{aligned} \quad (3.18)$$

Where  $0 < \gamma \leq 1$  is learning rate. It is easy to recognize that the estimation equations above confirm Hebbian learning rule in which the strength of connection represented by weight is consolidated by agreement of two nodes to which the connection is attached. As a result, Boltzmann machine trained with SGD is specified as follows:

Initialize  $\mathbf{W}$  and set  $k = 0$ .

Repeat

Data (state)  $\mathbf{x}$  is received from some real sample, or it can be kept intact.

$$w_{ij}^{(k+1)} = w_{ij}^{(k)} + \gamma \frac{x_i x_j}{\kappa T}$$

$$b_i^{(k+1)} = b_i^{(k)} + \gamma \frac{x_i}{\kappa T}$$

Increase  $k = k + 1$ .

Until some terminating conditions are met.

Note, a terminating condition is customized, for example, parameters  $\mathbf{W}$  and  $\mathbf{b}$  are not changed significantly, the maximum number of iterations is reached, or Boltzmann machine gets stable. The terminating condition that Boltzmann machine gets stable receives more concerns because stability is important property of spinning glass process or annealing process that Boltzmann machine. However, checking the stability in which global energy  $E(\mathbf{x})$  is not changed may consume a lot of iterations. Fortunately, SA can be incorporated into SGD so as to derive a more effective estimation. Boltzmann machine trained with SGD and SA is specified as follows:

Initialize current temperature  $T$  by highest temperature  $T_0$  as  $T = T_0$ .

Repeat

Data (state)  $\mathbf{x}$  is received from some real sample, or it can be kept intact.

$$w_{ij}^{(k+1)} = w_{ij}^{(k)} + \gamma \frac{x_i x_j}{\kappa T}$$

$$b_i^{(k+1)} = b_i^{(k)} + \gamma \frac{x_i}{\kappa T}$$

Evaluate Boltzmann machine given current parameter  $\mathbf{W}^{(k+1)}$  and  $\mathbf{b}^{(k+1)}$  so as to produce a new state  $\mathbf{x}'$ :

$$x'_i = \text{ReLU} \left( \sum_{j=1, j \neq i}^n w_{ij}^{(k+1)} x_j \right), \text{ for all } x'_i \text{ not from real sample}$$

If  $P(\mathbf{x}, \mathbf{x}', T | \mathbf{W}^{(k+1)}, \mathbf{b}^{(k+1)})$  is larger than a predefined threshold then  
 $\mathbf{x} = \mathbf{x}'$

Decrease current temperature, for example,  $T = \text{decrease}(T)$ .

End if

Increase  $k = k + 1$ .

Until terminating conditions are met.

The terminating conditions can be that best state ( $\mathbf{x}'$ ) is reached, the current state  $\mathbf{x}$  is good enough, or the current temperature  $T$  is low enough. These terminating conditions reflect the stable state of Boltzmann machine. As usual, given a maximum iteration number  $K$  and the current iteration number  $k$ , the temperature decreasing function can be defined as follows:

$$\text{decrease}(T) = T - \frac{k}{K} T$$

Of course, the acceptance probability is:

$$P(\mathbf{x}, \mathbf{x}', T | \mathbf{W}, \mathbf{b}) = \begin{cases} 1 & \text{if } \Delta E(\mathbf{x}, \mathbf{x}' | \mathbf{W}, \mathbf{b}) < 0 \\ \exp\left(-\frac{\Delta E(\mathbf{x}, \mathbf{x}' | \mathbf{W}, \mathbf{b})}{\kappa T}\right) & \text{otherwise} \end{cases} \quad (3.19)$$

Where,

$$\Delta E(\mathbf{x}, \mathbf{x}' | \mathbf{W}, \mathbf{b}) = E(\mathbf{x}' | \mathbf{W}, \mathbf{b}) - E(\mathbf{x} | \mathbf{W}, \mathbf{b})$$

There is a so-called restricted Boltzmann machine (RBM) in which neurons are separated into two groups such as input group denoted  $\mathbf{x}_z$  and hidden group denoted  $\mathbf{x}_x$ .

$$\mathbf{x} = (\mathbf{x}_z, \mathbf{x}_x)^T$$

The training algorithm by incorporation of SA and SGD is not changed except that there is neither connection between input neurons and input neurons nor connection between hidden neurons and hidden neurons. In other words, all connections are made between input group and hidden group, for instance, suppose cardinality of input group is  $k$  then, the number of connections is  $k(n - k)$ . Therefore, the two groups are considered layers such as input layer and hidden layer. Of course, both layers are output layers because connections in Boltzmann machine are two-direction connections whereas feed-forward neural network accepts only one-direction connections. RBM is trained faster than traditional Boltzmann machine because its number of connections is smaller. Moreover, it is clear to apply RBM into DGM because generator function in DGM  $\mathbf{x} = g(\mathbf{z} | \Phi)$  is modeled by RBM whose input is input group  $\mathbf{x}_z$  and whose output is output group  $\mathbf{x}_x$  such as  $\mathbf{x}_x = g(\mathbf{x}_z | \mathbf{W}, \mathbf{b})$  where  $\mathbf{x}_x$  is calculated by evaluating RBM given input  $\mathbf{x}_z$ .

$$\mathbf{x}_x = g(\mathbf{x}_z | \mathbf{W}, \mathbf{b}) \stackrel{\text{def}}{=} \text{evaluate RBM at } \mathbf{x}_z \text{ to produce } \mathbf{x}_x$$

The reason that the RBM approach for DGM is classified into approximate density DGM is that generator likelihood  $P(\mathbf{x} | \mathbf{W}, \mathbf{b})$  is defined indirectly based on the energy  $E(\mathbf{x} | \mathbf{W}, \mathbf{b})$ . Of course,  $\mathbf{x}_z$  is randomized such that  $\mathbf{x}_x$  is generated data.

#### 4. IMPLICIT DENSITY DGM

According to implicit density approach, DGMs do not specify explicitly generator PDF  $P(\mathbf{x} | \Phi, \mathbf{z})$ , which does not mean that such PDF is not existent but it is simple that such PDF is not concerned. Note, PDF is abbreviation of probability density function.

Recall that there are two problems related to construct a DGM: 1) how to define likelihood or error to train generator DNN  $g(\mathbf{z} | \Phi)$  and 2) how to define tractable PDF  $P(\mathbf{z})$  which implies the way to randomize  $\mathbf{z}$ . The second problem relates to assert qualification of random data  $\mathbf{z}'$  and hence, the second problem is stated as qualification problem of how to qualify random data. However, it is essential that the qualification problem aims to improve generator DNN  $g(\mathbf{z} | \Phi)$  because  $g(\mathbf{z} | \Phi)$  translate intractable  $\mathbf{z}$  into tractable  $\mathbf{x}$ . **Generative Adversarial Network (GAN)** developed by Goodfellow et al. (Goodfellow, et al., 2014) aims to reinforcing quality of generator DNN  $g(\mathbf{z} | \Phi) = \mathbf{x}' \approx \mathbf{x}$  by adding a so-called discriminator which is a discrimination function  $d(\mathbf{x} | \Psi): \mathbf{x} \rightarrow [0, 1]$  from concerned data  $\mathbf{x}$  or  $\mathbf{x}'$  to range  $[0, 1]$  in which  $d(\mathbf{x} | \Psi)$  can distinguish fake data from real data. In other words, the larger result the discriminator  $d(\mathbf{x}' | \Psi)$  derives, the more realistic the generated data  $\mathbf{x}'$  is. Obviously, discriminator  $d(\mathbf{x} | \Psi)$  is implemented by a DNN whose weights are  $\Psi$  called discriminator parameter with note that this discriminator DNN has only one output neuron denoted  $d_0$ .

$$d: R^m \rightarrow [0, 1] \text{ such that } d_0 = d(\mathbf{x} | \Psi) \text{ where } d_0 \in [0, 1], \mathbf{x} \in X \subseteq R^m \quad (4.1)$$

Actually, the task of discriminator  $d(\mathbf{x} | \Psi)$  is classification task with regard to class  $d_0$  belonging to interval  $[0, 1]$ . GAN does not establish explicitly PDFs of generator  $g(\mathbf{z} | \Phi)$  and discriminator  $d(\mathbf{x} | \Psi)$  such as  $P(\mathbf{x} | \Phi, \mathbf{z})$  and  $P(d_0 | \Psi, \mathbf{x})$  and hence, GAN does not define explicitly and separately likelihoods / errors of  $g(\mathbf{z} | \Phi)$  and  $d(\mathbf{x} | \Psi)$  too. Indeed, GAN instead unifies optimization constraints of  $g(\mathbf{z} | \Phi)$  and  $d(\mathbf{x} | \Psi)$  into a target function  $l(\Phi, \Psi | \mathbf{x}, \mathbf{z})$ .

$$l(\Phi, \Psi | \mathbf{x}, \mathbf{z}) = \log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi) | \Psi)) \quad (4.2)$$

Indeed, target function  $l(\Phi, \Psi | \mathbf{x}, \mathbf{z})$  is error function and so it is called loss function in literature. As a result, GAN tries to optimize dually generator parameter  $\Phi$  and discriminator parameter  $\Psi$  so that optimal estimate  $\Phi^*$  and optimal estimate  $\Psi^*$  are minimizer and maximizer of loss function  $l(\Phi, \Psi | \mathbf{x}, \mathbf{z})$  with expectation that Nash equilibrium will be achieved at the saddle point  $(\Phi^*, \Psi^*)$  with note that loss function  $l(\Phi, \Psi | \mathbf{x}, \mathbf{z})$  is function of  $\Phi$  and  $\Psi$  given data  $\mathbf{x}$  and  $\mathbf{z}$ .

$$\begin{aligned} \Phi^* &= \underset{\Phi}{\operatorname{argmin}} l(\Phi, \Psi^* | \mathbf{x}, \mathbf{z}) = \underset{\Phi}{\operatorname{argmin}} \log(1 - d(g(\mathbf{z} | \Phi) | \Psi^*)) \\ \Psi^* &= \underset{\Psi}{\operatorname{argmax}} l(\Phi^*, \Psi | \mathbf{x}, \mathbf{z}) = \underset{\Psi}{\operatorname{argmax}} (\log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi^*) | \Psi))) \end{aligned} \quad (4.3)$$

This is min-max problem in game theory (Goodfellow, et al., 2014):

$$\min_{\Phi} \max_{\Psi} l(\Phi, \Psi | \mathbf{x}, \mathbf{z})$$

Which results estimation equation according to stochastic gradient descent (SGD) method:

$$\Phi = \Phi - \gamma \nabla_{\Phi} l(\Phi, \Psi^* | \mathbf{x}, \mathbf{z})$$

$$\Psi = \Psi + \gamma \nabla_{\Psi} l(\Phi^*, \Psi | \mathbf{x}, \mathbf{z})$$

Where  $\gamma$  is learning rate. Recall that SGD, which is an iterative process, pushes candidate solution at each iteration along the direction which is opposite to gradient of target function for minimization or has the same direction to gradient of target function for maximization with note that the step length is represented by learning rate. Note,  $\nabla_{\Phi} l(\Phi, \Psi^* | \mathbf{x}, \mathbf{z})$  is gradient of loss function  $l(\Phi, \Psi^* | \mathbf{x}, \mathbf{z})$  fixed  $\Psi^*$  with regard to generator parameter  $\Phi$  and  $\nabla_{\Psi} l(\Phi^*, \Psi | \mathbf{x}, \mathbf{z})$  is gradient of loss function  $l(\Phi^*, \Psi | \mathbf{x}, \mathbf{z})$  fixed  $\Phi^*$  with regard to discriminator parameter  $\Psi$  as follows:

$$\begin{aligned} \nabla_{\Phi} l(\Phi, \Psi^* | \mathbf{x}, \mathbf{z}) &= \frac{dl(\Phi, \Psi^* | \mathbf{x}, \mathbf{z})}{d\Phi} = \frac{\partial \log(1 - d(g(\mathbf{z} | \Phi) | \Psi^*))}{\partial \Phi} \\ \nabla_{\Psi} l(\Phi^*, \Psi | \mathbf{x}, \mathbf{z}) &= \frac{dl(\Phi^*, \Psi | \mathbf{x}, \mathbf{z})}{d\Psi} = \frac{\partial (\log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi^*) | \Psi)))}{\partial \Psi} \end{aligned}$$

Therefore, the estimation equation is rewritten as follows:

$$\begin{aligned} \Phi &= \Phi - \gamma \nabla_{\Phi} \log(1 - d(g(\mathbf{z} | \Phi) | \Psi^*)) \\ \Psi &= \Psi + \gamma \nabla_{\Psi} (\log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi^*) | \Psi))) \end{aligned} \quad (4.4)$$

According to equations above, real data  $\mathbf{x}$  aims to maximize discriminator  $d(\mathbf{x} | \Psi)$  and in opposite, generated data  $\mathbf{x}' = g(\mathbf{z} | \Phi)$  aims to minimize discriminator  $d(\mathbf{x}' | \Psi)$ . Although both GAN and VAE use two DNNs for data generation but the underlying theory of GAN is slightly more succinct than VAE because there is no requirement of specifying probabilistic distribution  $P(\mathbf{z})$  of tractable  $\mathbf{z}$ . As a convention, the gradient  $\nabla_{\Phi}(\log(1 - d(g(\mathbf{z} | \Phi) | \Psi^*)))$  related to generator parameter  $\Theta$  is called generator gradient and the gradient  $\nabla_{\Psi}(\log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi^*) | \Psi)))$  related to discriminator parameter  $\Theta$  is called discriminator gradient.

Given epoch of size  $N$  is denoted as  $\mathbf{D} = ((\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$ , the estimation equations of  $\Phi$  and  $\Psi$  are extended exactly as epoch estimation at every iteration of SGD:

$$\begin{aligned} \Phi^{(k+1)} &= \Phi^{(k)} - \gamma \nabla_{\Phi} \frac{1}{N} \sum_{i=1}^N \log(1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)})) \\ \Psi^{(k+1)} &= \Psi^{(k)} + \gamma \nabla_{\Psi} \frac{1}{N} \sum_{i=1}^N (\log(d(\mathbf{x}^{(i)} | \Psi^{(k)})) + \log(1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)}))) \end{aligned}$$

As a result, GAN trained with SGD is specified as follows:

Initialize  $\Phi$  and  $\Psi$  and set  $k = 0$ .

Repeat



Sampling epoch  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$  or receiving epoch  $\mathbf{X}$  from big data / data stream.  
Randomize random epoch  $\mathbf{Z} = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$  from standard normal distribution  $P(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, I)$  with mean  $\mathbf{0}$  and identity covariance matrix  $I$ .

$$\Phi^{(k+1)} = \Phi^{(k)} - \gamma \nabla_{\Phi} \frac{1}{N} \sum_{i=1}^N \log \left( 1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)}) \right)$$

$$\Psi^{(k+1)} = \Psi^{(k)} + \gamma \nabla_{\Psi} \frac{1}{N} \sum_{i=1}^N \left( \log(d(\mathbf{x}^{(i)} | \Psi^{(k)})) + \log \left( 1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)}) \right) \right)$$

Increase  $k = k + 1$ .

Until some terminating conditions are met.

Note, a terminating condition is customized, for example, parameters  $\Phi$  and  $\Psi$  are not changed significantly or there is no more coming epoch  $\mathbf{X}$ . Moreover, the index  $k$  indicates time point as well as iteration of SGD.

Recall that the estimation equations of generator parameter  $\Phi$  and discriminator parameter  $\Psi$  are:

$$\Phi = \Phi - \gamma \nabla_{\Phi} \log(1 - d(g(\mathbf{z} | \Phi) | \Psi^*))$$

$$\Psi = \Psi + \gamma \nabla_{\Psi} (\log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi^*) | \Psi)))$$

It is necessary to calculate generator gradient and discriminator gradient. Indeed, we have:

$$\nabla_{\Phi} \log(1 - d(g(\mathbf{z} | \Phi) | \Psi^*)) = \frac{\partial \log(1 - d(g(\mathbf{z} | \Phi) | \Psi^*))}{\partial \Phi} = - \frac{\partial d(g(\mathbf{z} | \Phi) | \Psi^*) / \partial \Phi}{1 - d(g(\mathbf{z} | \Phi) | \Psi^*)}$$

$$\nabla_{\Psi} (\log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi^*) | \Psi)))$$

$$= \frac{\partial (\log(d(\mathbf{x} | \Psi)) + \log(1 - d(g(\mathbf{z} | \Phi^*) | \Psi)))}{\partial \Psi}$$

$$= \frac{\partial d(\mathbf{x} | \Psi) / \partial \Psi}{d(\mathbf{x} | \Psi)} - \frac{\partial d(g(\mathbf{z} | \Phi^*) | \Psi) / \partial \Psi}{1 - d(g(\mathbf{z} | \Phi^*) | \Psi)}$$

Where  $\partial d(.) / \partial \Phi$  and  $\partial d(.) / \partial \Psi$  denotes differentials of discriminator function with regard to  $\Phi$  and  $\Psi$ , respectively.

Given epoch of size  $N$  is denoted as  $\mathbf{D} = ((\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$ , the estimation equations of  $\Phi$  and  $\Psi$  are extended exactly as epoch estimation at every iteration of SGD:

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma \frac{1}{N} \sum_{i=1}^N \frac{\partial d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)}) / \partial \Phi}{1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)})}$$

$$\Psi^{(k+1)} = \Psi^{(k)} + \gamma \frac{1}{N} \sum_{i=1}^N \left( \frac{\partial d(\mathbf{x}^{(i)} | \Psi^{(k)}) / \partial \Psi}{d(\mathbf{x}^{(i)} | \Psi^{(k)})} - \frac{\partial d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)}) / \partial \Psi}{1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)})} \right) \quad (4.5)$$

As a result, GAN trained with SGD is specified as follows:

Initialize  $\Phi$  and  $\Psi$  and set  $k = 0$ .

Repeat

Sampling epoch  $\mathbf{X} = (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})$  or receiving epoch  $\mathbf{X}$  from big data / data stream.

Randomize random epoch  $\mathbf{Z} = (\mathbf{z}^{(1)}, \mathbf{z}^{(2)}, \dots, \mathbf{z}^{(N)})$  from standard normal distribution  $P(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, I)$  with mean  $\mathbf{0}$  and identity covariance matrix  $I$ .

$$\Phi^{(k+1)} = \Phi^{(k)} + \gamma \frac{1}{N} \sum_{i=1}^N \frac{\partial d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)}) / \partial \Phi}{1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)})}$$

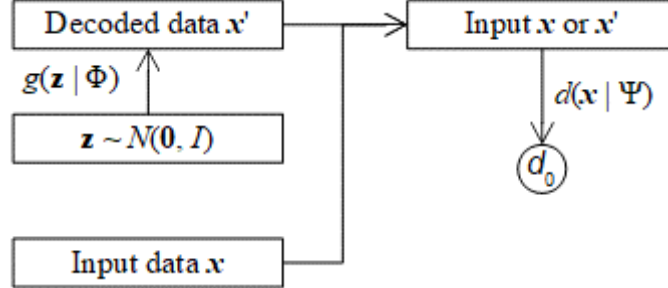
$$\Psi^{(k+1)} = \Psi^{(k)} + \gamma \frac{1}{N} \sum_{i=1}^N \left( \frac{\partial d(\mathbf{x}^{(i)} | \Psi^{(k)}) / \partial \Psi}{d(\mathbf{x}^{(i)} | \Psi^{(k)})} - \frac{\partial d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)}) / \partial \Psi}{1 - d(g(\mathbf{z}^{(i)} | \Phi^{(k)}) | \Psi^{(k)})} \right)$$

Increase  $k = k + 1$ .

Until some terminating conditions are met.

Note, a terminating condition is customized, for example, parameters  $\Phi$  and  $\Psi$  are not changed significantly or there is no more coming epoch  $X$ . Moreover, the index  $k$  indicates time point as well as iteration of SGD.

Following figure depicts GAN.



**Figure 4.1.** Generative Adversarial Network (GAN)

There is a question that how to calculate the differentials  $\partial d(.) / \partial \Phi$  and  $\partial d(.) / \partial \Psi$ . Indeed, it is not difficult to calculate them in context of neural network associated with backpropagation algorithm so that the last output layer as well as last neuron  $o$  of any DNN  $f(x | \Theta)$  or  $g(z | \Phi)$  is acted by activation function  $a(.)$  as follows:

$$a(o) = a(\mathbf{w}^T \mathbf{i})$$

$$o = \mathbf{w}^T \mathbf{i}$$

Where  $\mathbf{i}$  is input of the last layer  $o$  and weight parameter  $\mathbf{w}$  is a part of entire parameter  $\Phi$  or  $\Psi$  and hence, we need to focus on calculating differential  $da(o) / d\mathbf{w}$  which is equivalent to any differential  $\partial d(.) / \partial \Phi$  or  $\partial d(.) / \partial \Psi$  so that backpropagation algorithm will solve remaining parts of entire parameter  $\Phi$  or  $\Psi$ .

$$\frac{da(o)}{d\mathbf{w}} \cong \text{any differential}$$

Indeed, we have:

$$\frac{da(o)}{d\mathbf{w}} = \frac{da(\mathbf{w}^T \mathbf{i})}{d\mathbf{w}} = \frac{da(o)}{do} \mathbf{i}^T = a'(o) \mathbf{i}^T$$

Note, the subscript “ $T$ ” denotes transposition operator of vector and matrix in which row vector becomes column vector and vice versa. It is easy to calculate the derivative  $a'(o)$  when activation function was specified, for instance, if  $a(o)$  is sigmoid function, we have:

$$a(o) = y = \frac{1}{1 + e^{-o}}$$

$$a'(o) = \frac{1}{1 + e^{-o}} \left( 1 - \frac{1}{1 + e^{-o}} \right) = a(o)(1 - a(o)) = y(1 - y)$$

In practice,  $y$  is replaced by  $a(y)$  in order to prevent  $o$  from being out of space:

$$a'(o) \cong a(y)(1 - a(y)) = a(a(o))(1 - a(a(o)))$$

As a result, we have:

$$\text{any differential} \cong a(a(o))(1 - a(a(o))) \mathbf{i}^T$$

For fast computation, it is possible to set the derivative  $a'(o)$  to be small enough constants like 1 such that any differential is  $\mathbf{i}^T$ .

Given epoch  $\mathbf{D} = (\mathbf{d}^{(1)} = (\mathbf{x}^{(1)}, \mathbf{z}^{(1)}), \mathbf{d}^{(2)} = (\mathbf{x}^{(2)}, \mathbf{z}^{(2)}), \dots, \mathbf{d}^{(N)} = (\mathbf{x}^{(N)}, \mathbf{z}^{(N)}))$  implies that the epoch is created or sent by equilateral distribution  $1/N$  but in general case,  $\mathbf{D}$  can follow an arbitrary distribution denoted by PDF  $P(\mathbf{d})$ , which makes loss function  $l(\Phi, \Psi)$  changed a little bit by theoretical expectation given distribution  $P(\mathbf{d})$ .

$$l(\Phi, \Psi) = \int_d (\log(d(\mathbf{x}|\Psi)) + \log(1 - d(g(\mathbf{z}|\Phi)|\Psi))) d\mathbf{d}$$

Suppose  $\mathbf{x}$  and  $\mathbf{z}$  distribute separately as  $P(\mathbf{x})$  and  $Q(\mathbf{z})$  such that  $P(\mathbf{x})$  is called original data PDF and  $Q(\mathbf{z})$  is called generated data PDF, we have:

$$l(\Phi, \Psi) = \int_{\mathbf{x}} P(\mathbf{x}) \log(d(\mathbf{x}|\Psi)) d\mathbf{x} + \int_{\mathbf{z}} Q(\mathbf{z}) \log(1 - d(g(\mathbf{z}|\Phi)|\Psi)) d\mathbf{z} \quad (4.6)$$

Although there is no significant change in aforementioned practical technique to estimate parameters, it is necessary to research original data PDF  $P(\mathbf{x})$  and generated data PDF  $Q(\mathbf{z})$  as well as expectation form of loss function  $l(\Phi, \Psi)$  so as to prove convergence of GAN. Recall that the min-max problem is:

$$\min_{\Phi} \max_{\Psi} l(\Phi, \Psi)$$

That is:

$$\Phi^* = \operatorname{argmin}_{\Phi} l(\Phi, \Psi^*)$$

$$\Psi^* = \operatorname{argmax}_{\Psi} l(\Phi^*, \Psi)$$

The convergence of GAN is equivalent to the convergence of this min-max problem. In other words, Goodfellow et al. (Goodfellow, et al., 2014) proved the existence of global optimal value  $l^*$  such that min-max problem approach  $l^*$  as follows:

$$l^* = l(\Phi^*, \Psi^*) = \min_{\Phi} l(\Phi, \Psi^*) = \min_{\Phi} \max_{\Psi} l(\Phi, \Psi) \quad (4.7)$$

Because  $\mathbf{z}$  is generated by distribution  $Q(\mathbf{z})$  and  $g(\mathbf{z} | \Phi)$  is valued as  $\mathbf{x}$  as  $g(\mathbf{z} | \Phi) = \mathbf{x}$ , loss function  $l(\Phi, \Psi)$  is rewritten by changing variable (Goodfellow, et al., 2014, p. 5).

$$\begin{aligned} l(\Phi, \Psi) &= \int_{\mathbf{x}} P(\mathbf{x}) \log(d(\mathbf{x}|\Psi)) d\mathbf{x} + \int_{\mathbf{x}} Q(\mathbf{x}) \log(1 - d(\mathbf{x}|\Psi)) d\mathbf{x} \\ &= \int_{\mathbf{x}} (P(\mathbf{x}) \log(d(\mathbf{x}|\Psi)) + Q(\mathbf{x}) \log(1 - d(\mathbf{x}|\Psi))) d\mathbf{x} \end{aligned}$$

In mathematical literature, function  $a \log(y) + b \log(1-y)$  gets maximal at maximizer  $y = a/(a+b)$  such that:

$$d(\mathbf{x}|\Psi^*) = \frac{P(\mathbf{x})}{P(\mathbf{x}) + Q(\mathbf{x})}$$

Therefore, we have (Goodfellow, et al., 2014, p. 5):

$$\begin{aligned} \max_{\Psi} l(\Phi, \Psi) &= l(\Phi, \Psi^*) = \int_{\mathbf{x}} \left( P(\mathbf{x}) \log\left(\frac{P(\mathbf{x})}{P(\mathbf{x}) + Q(\mathbf{x})}\right) + Q(\mathbf{x}) \log\left(\frac{Q(\mathbf{x})}{P(\mathbf{x}) + Q(\mathbf{x})}\right) \right) d\mathbf{x} = \\ &= \int_{\mathbf{x}} P(\mathbf{x}) \log\left(\frac{2P(\mathbf{x})}{P(\mathbf{x}) + Q(\mathbf{x})}\right) d\mathbf{x} + \int_{\mathbf{x}} Q(\mathbf{x}) \log\left(\frac{2Q(\mathbf{x})}{P(\mathbf{x}) + Q(\mathbf{x})}\right) d\mathbf{x} - \log 4 \\ &= \text{KL}\left(P(\mathbf{x}) \middle| \frac{P(\mathbf{x}) + Q(\mathbf{x})}{2}\right) + \text{KL}\left(Q(\mathbf{x}) \middle| \frac{P(\mathbf{x}) + Q(\mathbf{x})}{2}\right) - \log 4 \end{aligned}$$

Where  $\text{KL}(\cdot)$  denotes Kullback-Leibler divergence of two distributions. The sum of two KL divergences above is a so-called Jensen-Shannon divergence of original data distribution  $P(\mathbf{x})$  and generated data distribution  $Q(\mathbf{z})$ , denoted  $\text{JS}(P(\mathbf{x}) | Q(\mathbf{x}))$ . Therefore, we have (Goodfellow, et al., 2014, p. 5):

$$\max_{\Psi} l(\Phi, \Psi) = l(\Phi, \Psi^*) = \text{JS}(P(\mathbf{x}) | Q(\mathbf{x})) - \log 4$$

Because Jensen-Shannon divergence is always nonnegative, we have (Goodfellow, et al., 2014, p. 5):

$$\max_{\Psi} l(\Phi, \Psi) = l(\Phi, \Psi^*) \geq -\log 4$$

The sign “=” occurs because Jensen-Shannon divergence is zero if the two distributions are equal, for instance,  $P(x) = Q(x)$ . Therefore,  $l(\Phi, \Psi^*)$  has maximal value  $-\log 4$ . In other words, we have:

$$l^* = l(\Phi^*, \Psi^*) = \min_{\Phi} \max_{\Psi} l(\Phi, \Psi) = -\log 4 \quad (4.8)$$

Due to the existence of global optimal value  $l^* = -\log 4$ , the convergence of GAN is asserted.

## 5. CONCLUSIONS

Recall that there are three main approaches for constructing deep generative model (DGM) such as tractable density DGM, approximate density DGM, and implicit density DGM. When skimming these approaches, it is easy to recognize that applied statistical problems such as probability distribution and parameter estimation are often mentioned but the effectiveness of a deep generative model is also dependent on how to structure the deep neural network (DNN) and how to train such network. Anyhow data generation function called generator is always defined by DNN in DGM. Backpropagation (BP) algorithm associated with stochastic gradient descent (SGD) method is focused as typical example in this research but there are some more effective training algorithms. Essential, training DNN generator is unsupervised learning task because there is no data class in DGM although generating data distribution (distribution of tractable data  $\mathbf{z}$ ) is often assumed to follow normal distribution whereas BP belongs to supervised learning algorithm. This is the reason that the two problems of constructing DGM are 1) how to train generator DNN  $g(\mathbf{z} | \Phi)$  and 2) how to qualify such training task which often relates to another optimization task or another training task so that the qualification task tries to attach supervised learning BP to unsupervised learning mechanism. For instance, PixelRNN allows output data becomes input data by recurrent neural network, VAE applies Kullback-Leibler divergence into forming data distribution as implicit data class, and GAN issues target function with expectation of Nash equilibrium. Essentially, these mechanisms make the exchange or transformation between supervised learning and unsupervised learning, which plays the role of a hinge for creating the DNN generator.

One of problems issued by BP is the zero derivative problem when SGD cannot improve parameters after some large enough iteration because the gradient (derivative) approaches zero at that time. In other words, SGD may not converge even though there is a large enough number of iterations. Moreover, basic DGM approaches mentioned here require a continuous data provision for training DNN, which consumes more resources than reinforcement learning.

## APPENDICES

### A1. Backpropagation algorithm

Because *backpropagation* (BP) algorithm is often associated with stochastic gradient descent (SGD) method to optimize loss function, it is necessary to describe a little bit BP and SGD, especially, in case of DGM where artificial neural network is deep neural network (DNN) with many hidden layers so that learning DNN (training DNN) is essential to deep learning. A neural network often has one input layer, one output layer, and hidden layers. The simplest neural network has one input layer and one output layer without a hidden layer. A DNN is a neural network often having many enough hidden layers. Each layer has a list of units called neurons and there are full connections of neurons between two successive layers. Feed-forward network which is the neural network whose connections are one-way from input layer to hidden layers to output layer is focused here. BP is a reverse process that begins from output layer back to input layer. Without loss of generality, input neurons, hidden neurons, and output neurons are concerned rather than input layer, hidden layers, and output layers,

respectively because BP processes layers by processing neurons of layers. As a convention, let  $i$ ,  $h$ , and  $o$  denote indices of input neurons, hidden neurons, and output neurons, respectively. Let  $x$  and  $y$  denote input variable and output variable of neurons. For instance,  $x_i$  and  $y_i$  are input and output of an input neuron and  $x_h$  and  $y_h$  are input and output of a hidden neuron whereas  $x_o$  and  $y_o$  are input and output of an output neuron, respectively. Because BP is a reverse process that begins from output layer back to input layer, output neuron is concerned firstly by starting its propagation rule as follows:

$$x_o = \sum_h w_{ho} y_h + \theta_o$$

$$y_o = f(x_o)$$

Where  $w_{ho}$  and  $\theta_o$  are weights of connections between previous hidden neurons  $h$  and current output neurons  $o$  while  $\theta_o$  is bias of current output neurons  $o$ . Moreover,  $f$  denotes activation function that squashes input into valid range, which is often sigmoid (logistic) function or hyperbolic tangent function. Some literature documents use letter  $b$  to denote such bias. BP aims to learn the parameters such as connection weights  $w_{(.)}$  and biases  $\theta_{(.)}$  from sample data. Note, propagation rule is cornerstone of evaluating neural network. Let  $v_o$  is real value of output neuron  $o$ , error function  $\varepsilon(y_o)$  of output neuron is half the square deviation between  $y_o$  and  $v_o$ .

$$\varepsilon(y_o) = \frac{1}{2} (v_o - y_o)^2 = \frac{1}{2} (y_o - v_o)^2$$

Note,  $y_o$  is variable and  $v_o$  represents sample data. Weight parameter and bias parameter are estimated by minimizing output error function  $\varepsilon(y_o)$  according to BP.

$$w_{ho}^* = \underset{w_{ho}}{\operatorname{argmin}} \frac{1}{2} (y_o - v_o)^2$$

$$\theta_o^* = \underset{\theta_o}{\operatorname{argmin}} \frac{1}{2} (y_o - v_o)^2$$

Minimizing output error function  $\varepsilon(y_o)$  is equivalent to maximizing likelihood (PDF) of random variable  $y_o$ . Indeed, output likelihood  $P(y_o)$  is specified as follows:

$$P(y_o) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} (y_o - v_o)^2\right)$$

Where,

$$y_o = f(x_o) = f\left(\sum_h w_{ho} y_h + \theta_o\right)$$

Note, mean and variance of output likelihood  $P(y_o)$  are 0 and 1, respectively. Exactly,  $P(y_o)$  is probability density function (PDF) of the error  $y_o - v_o$ . By maximum likelihood estimation (MLE) method, maximizing  $P(y_o)$  is to maximize its natural logarithm  $\log P(y_o)$ :

$$\log P(y_o) = -\frac{1}{2} \pi - \frac{1}{2} (y_o - v_o)^2$$

Because  $\pi$  is constant, it is obtained:

$$w_{ho}^* = \underset{w_{ho}}{\operatorname{argmax}} \log P(y_o) = \underset{w_{ho}}{\operatorname{argmax}} \left(-\frac{1}{2} (y_o - v_o)^2\right)$$

$$\theta_o^* = \underset{\theta_o}{\operatorname{argmax}} \log P(y_o) = \underset{\theta_o}{\operatorname{argmax}} \left(-\frac{1}{2} (y_o - v_o)^2\right)$$

This implies:



$$w_{ho}^* = \underset{w_{ho}}{\operatorname{argmin}} \frac{1}{2} (y_o - v_o)^2$$

$$\theta_o^* = \underset{\theta_o}{\operatorname{argmin}} \frac{1}{2} (y_o - v_o)^2$$

Which confirm the equivalence of error minimization and likelihood maximization. Therefore, BP becomes more potential because MLE can be extended with more techniques so that error function is defined indirectly and more flexibly via likelihood function. It is interesting that the minimization problem can be exchanged with the maximization problem. Shortly, the optimization is ultimate purpose.

The problem now is how to minimize error function  $\varepsilon(y_o)$ , which is the optimization problem. Fortunately, *stochastic gradient descent* (SGD) method is applied into solving this optimization when square error function is Lipschitz continuous and bounded. Given target function has both variable and parameter where parameter is the subject of optimization, SGD pushes parameter candidate-point along with the same direction (for maximization) or the opposite direction (for minimization) of gradient of target function. There are two important aspects of SGD: 1) the gradient is the first-order derivative of target function with regard to parameter, and 2) the variable is considered as input data which is fed by stochastic process or random way. Moreover, candidate point considered as candidate solution or candidate minimizer / maximizer of optimal parameter is pushed with step length which is coded as learning rate  $\gamma$ . It is proved that SGD will be converged to the optimal solution (good enough minimizer / maximizer) after many enough iterations and many enough data. Shortly, SGD updates weight parameter  $w_{ho}$  and bias parameter  $\theta_o$  of output error function  $\varepsilon(y_o)$  at every iteration within context of BP as follows:

$$w_{ho} = w_{ho} - \gamma \nabla_{w_{ho}} \varepsilon(y_o)$$

$$\theta_o = \theta_o - \gamma \nabla_{\theta_o} \varepsilon(y_o)$$

Where  $\nabla_{w_{ho}} \varepsilon(y_o)$  and  $\nabla_{\theta_o} \varepsilon(y_o)$  are gradients of output error function  $\varepsilon(y_o)$  with regard to weight parameter  $w_{ho}$  and bias parameter  $\theta_o$ , respectively:

$$\nabla_{w_{ho}} \varepsilon(y_o) = \frac{d\varepsilon(y_o)}{dw_{ho}}$$

$$\nabla_{\theta_o} \varepsilon(y_o) = \frac{d\varepsilon(y_o)}{d\theta_o}$$

It is not difficult to calculate these gradients. Due to chain rule of derivative and propagation rule

$$y_o = f\left(x_o = \sum_h w_{ho} y_h + \theta_o\right)$$

We have:

$$\nabla_{w_{ho}} \varepsilon(y_o) = \frac{d\varepsilon(y_o)}{dw_{ho}} = \frac{d\varepsilon(y_o)}{dy_o} \frac{dy_o}{dx_o} \frac{dx_o}{dw_{ho}} = (y_o - v_o) f'(x_o) y_h$$

$$\nabla_{\theta_o} \varepsilon(y_o) = \frac{d\varepsilon(y_o)}{d\theta_o} = \frac{d\varepsilon(y_o)}{dy_o} \frac{dy_o}{dx_o} \frac{dx_o}{d\theta_o} = (y_o - v_o) f'(x_o)$$

Note,

$$\frac{d\varepsilon(y_o)}{dy_o} = \frac{d \frac{1}{2} (y_o - v_o)^2}{dy_o} = (y_o - v_o)$$

$$\frac{dy_o}{dx_o} = \frac{df(x_o)}{dx_o} = f'(x_o)$$

$$\frac{dx_o}{dw_{ho}} = \frac{d(\sum_h w_{ho} y_h + \theta_o)}{dw_{ho}} = y_h$$

$$\frac{dx_o}{d\theta_o} = \frac{d(\sum_h w_{ho} y_h + \theta_o)}{d\theta_o} = 1$$

And due to:

$$w_{ho} = w_{ho} - \gamma \nabla_{w_{ho}} \varepsilon(y_o)$$

$$\theta_o = \theta_o - \gamma \nabla_{\theta_o} \varepsilon(y_o)$$

As a result, parameter  $w_{ho}$  and bias parameter  $\theta_o$  at output neuron is updated at every iteration of SGD within context of BP as follows:

$$w_{ho} = w_{ho} + \gamma y_h (v_o - y_o) f'(x_o)$$

$$\theta_o = \theta_o + \gamma (v_o - y_o) f'(x_o) \quad (A1.1)$$

Note,  $f'(x_o)$  is the first derivative of activation function at  $x_o$  and  $\gamma$  is learning rate ( $0 < \gamma \leq 1$ ). It is expected that weight parameter  $w_{ho}$  and bias parameter  $\theta_o$  will converge the aforementioned optimizers  $w_{ho}^*$  and  $\theta_o^*$ . Let  $\delta_o$  denote output error, we have:

$$w_{ho} = w_{ho} + \gamma y_h \delta_o$$

$$\theta_o = \theta_o + \gamma \delta_o$$

Where,

$$\delta_o = (v_o - y_o) f'(x_o)$$

BP continues to estimate weight and bias parameters of previous neurons which are hidden neurons. Without loss of generality, given hidden neuron  $h$  whose error  $\varepsilon(y_h)$  is sum of errors of all output neuron  $o$  to which such hidden neuron connects.

$$\varepsilon(y_h) = \sum_o \varepsilon(y_o)$$

Note,  $y_h$  is calculated by propagation rule as usual:

$$y_h = f\left(x_h = \sum_j w_{jh} y_j + \theta_h\right)$$

Note,  $y_j$  is output value of previous hidden neuron  $j$  which connects to current hidden neuron  $h$ . Because SGD is continuously applied into estimating weight parameters  $w_{jh}$  and bias parameter  $\theta_h$  of hidden neuron  $h$ , gradients of hidden error function  $\varepsilon(y_h)$  with regard to  $w_{jh}$  and  $\theta_h$  need to be determined, respectively. According to chain rule of derivative, we have:

$$\nabla_{w_{jh}} \varepsilon(y_h) = \frac{d\varepsilon(y_h)}{dw_{jh}} = \frac{d\varepsilon(y_h)}{dy_h} \frac{dy_h}{dx_h} \frac{dx_h}{dw_{jh}}$$

$$\nabla_{\theta_h} \varepsilon(y_h) = \frac{d\varepsilon(y_h)}{d\theta_h} = \frac{d\varepsilon(y_h)}{dy_h} \frac{dy_h}{dx_h} \frac{dx_h}{d\theta_h}$$

Due to:

$$\frac{dy_h}{dx_h} = \frac{df(x_h)}{dx_h} = f'(x_h)$$

$$\frac{dx_h}{dw_{jh}} = \frac{d(\sum_j w_{jh} y_j + \theta_h)}{dw_{jh}} = y_j$$

$$\frac{dx_h}{d\theta_h} = \frac{d(\sum_j w_{jh} y_j + \theta_h)}{d\theta_h} = 1$$

Equations of gradients  $\nabla_{w_{jh}} \varepsilon(y_h)$  and  $\nabla_{\theta_h} \varepsilon(y_h)$  are written:

$$\nabla_{w_{jh}} \varepsilon(y_h) = \frac{d\varepsilon(y_h)}{dy_h} y_j f'(x_h)$$

$$\nabla_{\theta_h} \varepsilon(y_h) = \frac{d\varepsilon(y_h)}{d\theta_h} = \frac{d\varepsilon(y_h)}{dy_h} f'(x_h)$$

It is now necessary to calculate the derivative  $d\varepsilon(y_h)/dy_h$  of hidden error function  $\varepsilon(y_h)$  with regard to  $y_h$ . Indeed, we have:

$$\frac{d\varepsilon(y_h)}{dy_h} = \frac{d \sum_o \varepsilon(y_o)}{dy_h} = \sum_n \frac{d\varepsilon(y_o)}{dy_h} = \sum_n \frac{d\varepsilon(y_o)}{dy_o} \frac{dy_o}{dx_o} \frac{dx_o}{dy_h}$$

Due to:

$$\frac{d\varepsilon(y_o)}{dy_o} \frac{dy_o}{dx_o} = \frac{d \frac{1}{2} (y_o - v_o)^2}{dy_o} \frac{df(x_o)}{dx_o} = (y_o - v_o) f'(x_o) = -\delta_o$$

This produces:

$$\frac{d\varepsilon(y_h)}{dy_h} = - \sum_n \frac{dx_o}{dy_h} \delta_o = - \sum_n \frac{d(\sum_h w_{ho} y_h + \theta_o)}{dy_h} \delta_o = - \sum_n w_{ho} \delta_o$$

Gradients  $\nabla_{w_{jh}} \varepsilon(y_h)$  and  $\nabla_{\theta_h} \varepsilon(y_h)$  are totally determined:

$$\nabla_{w_{jh}} \varepsilon(y_h) = -y_j \left( \sum_n w_{ho} \delta_o \right) f'(x_h)$$

$$\nabla_{\theta_h} \varepsilon(y_h) = - \left( \sum_n w_{ho} \delta_o \right) f'(x_h)$$

Due to SGD estimation:

$$w_{jh} = w_{jh} - \gamma \nabla_{w_{jh}} \varepsilon(y_h)$$

$$\theta_h = \theta_h - \gamma \nabla_{\theta_h} \varepsilon(y_h)$$

Parameter  $w_{jh}$  and bias parameter  $\theta_h$  at output neuron is updated at every iteration of SGD within context of BP as follows:

$$w_{jh} = w_{jh} + \gamma y_j \left( \sum_o w_{ho} \delta_o \right) f'(x_h)$$

$$\theta_h = \theta_h + \gamma \left( \sum_o w_{ho} \delta_o \right) f'(x_h)$$
(A1.2)

Let  $\delta_h$  denote hidden error:

$$\delta_h = \left( \sum_n w_{ho} \delta_o \right) f'(x_h)$$

SGD estimation equations of hidden neurons become more succinct:

$$w_{jh} = w_{jh} + \gamma y_j \delta_h$$

$$\theta_h = \theta_h + \gamma \delta_h$$

The reverse process of BP recurrently continues to estimate other parameters of previous hidden neurons, which is an interesting and effective aspect of BP. Finally, following system of estimation equations is the summary of association of BP and SGD, in which weight parameters and bias parameters of feed-forward neural network are updated at every iteration whenever data sample  $v$  is received.

$$w_{ho} = w_{ho} + \gamma y_h \delta_o$$

$$\theta_o = \theta_o + \gamma \delta_o$$

$$w_{jh} = w_{jh} + \gamma y_j \delta_h$$

$$\theta_h = \theta_h + \gamma \delta_h$$
(A1.3)

Where output error  $\delta_o$  and hidden error  $\delta_h$  are calculated as follows:

$$\begin{aligned}\delta_o &= (v_o - y_o)f'(x_o) \\ \delta_h &= \left( \sum_n w_{ho} \delta_o \right) f'(x_h)\end{aligned}\tag{A1.4}$$

For easy explanation, according to the BP recurrent process, weight parameters  $w_{kj}$  and bias parameter  $\theta_j$  of previous hidden neuron  $j$  of hidden neuron  $h$  are calculated as follows:

$$\begin{aligned}w_{kj} &= w_{kj} + \gamma y_k \delta_j \\ \theta_j &= \theta_j + \gamma \delta_j\end{aligned}$$

Where error  $\delta_j$  of hidden neuron  $j$  is:

$$\delta_j = \left( \sum_h w_{jh} \delta_h \right) f'(x_j)$$

The back recurrent process continues until reaching input layer (exclusive). It is easy to recognize that the entry point of BP is the output error  $\delta_o$  which relates to derivative of error function of output neuron at sample point  $v$ . Recall that such error function can be replaced by likelihood function instead. Therefore,  $\delta_o$  is the opposite of gradient of output error function if the error function is applied to estimation within context of minimization. Otherwise,  $\delta_o$  is gradient of output likelihood function if the likelihood function is applied to estimation within context of maximization. The interesting result allows us to extend BP applications by defining error function or likelihood function at output layer without changing BP recurrent process.

It is necessary to consider activation function  $f(x)$  and its derivative  $f'(x)$  which are evaluated at output neuron and hidden neuron as  $f(x_o)$ ,  $f(x_h)$ ,  $f'(x_o)$ , and  $f'(x_h)$ . For instance, if  $f(x)$  is sigmoid function (logistic function), we have:

$$\begin{aligned}f(x) &= y = \frac{1}{1 + e^{-x}} \\ f'(x) &= \frac{1}{1 + e^{-x}} \left( 1 - \frac{1}{1 + e^{-x}} \right) = f(x)(1 - f(x)) = y(1 - y)\end{aligned}$$

In practice,  $y$  is replaced by  $f(y)$  in order to prevent  $o$  from being out of space:

$$f'(x) \cong f(y)(1 - f(y)) = f(f(x))(1 - f(f(x)))$$

It is possible to fix the derivative by 1 as  $f'(x) = 1$  for all  $x$  for fast computation but this approximation is not optimal.

## A2. Kullback-Leibler divergence

In information theory, entropy is the quantity that indicates uncertainty of a random variable. Exactly, given random variable  $x$  and its probability density function (PDF)  $P(x)$ , entropy of  $x$  denoted  $H(x)$  is the metric which is minus expected value of natural logarithm of  $x$  given distribution  $P(x)$ .

$$H(x) = - \int_x P(x) \ln P(x) dx$$

If random variable  $x$  is discrete, its entropy becomes simpler:

$$H(x) = - \sum_x P(x) \ln P(x)$$

It is easy to recognize that entropy  $H(x)$  measures the level of uncertainty or the level of surprise for random variable  $x$  and such level is average level. As the opposite of probability, such uncertainty which measures the variation of random variable is also called information content, self-information, surprise, Shannon information, or information, in short:

$$I(x) = I(x|P(x)) = - \ln P(x)$$

It is interesting that information of  $x$  is defined based on probability of  $x$  although it is opposite to the probability. The larger the information of  $x$  is, the more the uncertainty of  $x$  is, the less the probability of  $x$  is. Entropy  $H(x)$  is the expectation of information  $I(x)$  given distribution  $P(x)$ .

$$H(x) = \int_x P(x) I(x) dx$$

In thermodynamics, entropy represents the level of chaos in movement of particles. *Kullback-Leibler (KL) divergence* which is defined based on concept of entropy measures the difference of two distributions. For instance, given distributions  $P(x)$  and  $Q(x)$ , Kullback-Leibler divergence of  $P(x)$  given  $Q(x)$  denoted  $D_{KL}(P(x) | Q(x))$  measures how much the distribution  $P(x)$  is different from the distribution  $Q(x)$  (Wikipedia, Kullback-Leibler divergence, 2004).

$$D_{KL}(P(x) || Q(x)) = \int_x P(x) \ln \frac{P(x)}{Q(x)} dx$$

For convenience, we denote:

$$KL(P(x) | Q(x)) = \int_x P(x) \ln \frac{P(x)}{Q(x)} dx \quad (A2.1)$$

The larger the KL divergence  $KL(P(x) | Q(x))$  is, the more different focused distribution  $P(x)$  is from distribution  $Q(x)$ . However, such difference does not represent distance metric between  $P(x)$  and  $Q(x)$  because KL divergence is not symmetric:

$$KL(P(x) | Q(x)) \neq KL(Q(x) | P(x))$$

KL divergence does not satisfy triangle inequality too (Wikipedia, Kullback-Leibler divergence, 2004). KL divergence is expended:

$$\begin{aligned} KL(P(x) | Q(x)) &= \int_x P(x) \ln \frac{P(x)}{Q(x)} dx = \int_x P(x) \ln P(x) dx - \int_x P(x) \ln Q(x) dx \\ &= - \int_x P(x) \ln Q(x) dx - H(x) \end{aligned}$$

Let  $H(x | Q(x))$  denote entropy of  $x$  such that  $x$  is quantified by distribution  $Q(x)$ .

$$H(x | Q(x)) = - \int_x P(x) \ln Q(x) dx = \int_x P(x) I(x | Q(x)) dx$$

Where,

$$I(x | Q(x)) = - \ln Q(x)$$

We obtain:

$$KL(P(x) | Q(x)) = H(x | Q(x)) - H(x) = \int_x P(x) (I(x | Q(x)) - I(x | P(x))) dx$$

This implies KL divergence measures the expected value of uncertainty when focused distribution  $P(x)$  is replaced by distribution  $Q(x)$  for quantifying random variable  $x$ . This is the reason that KL divergence is also called relative entropy. KL divergence is always nonnegative due to:

$$\int_x P(x) \ln P(x) dx \geq \int_x P(x) \ln Q(x) dx$$

Two distributions  $P(x)$  and  $Q(x)$  are identical if KL divergence  $KL(P(x) | Q(x))$  is zero. Moreover, if  $P(x)$  is 0 for all  $x$  then  $Q(x)$  must be 0 for all  $x$ . As usual,  $P(x)$  represents data distribution and  $Q(x)$  represents theoretical distribution so that it is possible to compare or fit observational model with hypothesis model (Wikipedia, Kullback-Leibler divergence, 2004).



While Kullback-Leibler divergence is not a metric because it satisfies neither symmetry nor triangle inequality, *Jensen-Shannon (JS) divergence* is a real metric that measures the distance between two distributions although JS divergence is defined based on KL divergence. Given the two distributions  $P(x)$  and  $Q(x)$ , its JS divergence is the following average KL divergence (Wikipedia, Jensen-Shannon divergence, 2006):

$$D_{JS}(P(x)||Q(x)) = \frac{1}{2} D_{KL}(P(x)||M(x)) + \frac{1}{2} D_{KL}(Q(x)||M(x))$$

Where  $M(x)$  is mixture distribution or mean distribution of  $P(x)$  and  $Q(x)$ :

$$M(x) = \frac{P(x) + Q(x)}{2}$$

Note, JS divergence satisfies both symmetry and triangle inequality.

$$D_{JS}(P(x)||Q(x)) = D_{JS}(Q(x)||P(x))$$

$$D_{JS}(P(x)||R(x)) < D_{JS}(P(x)||Q(x)) + D_{JS}(Q(x)||R(x))$$

For convenience, we denote:

$$JS(P(x)|Q(x)) = \frac{1}{2} KL(P(x)|M(x)) + \frac{1}{2} KL(Q(x)|M(x)) \quad (A2.2)$$

JS divergence is bounded in interval  $[0, 1]$  such that  $0 \leq JS(P(x) | Q(x)) \leq 1$ . Square root of JS divergence is called JS distance of two distributions.

## REFERENCES

- Doersch, C. (2016, January 3). Tutorial on Variational Autoencoders. *arXiv preprint*. Retrieved from <https://arxiv.org/abs/1606.05908>
- Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. Cambridge, Massachusetts, US: The MIT Press. Retrieved from <https://www.deeplearningbook.org/>
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., . . . Bengio, Y. (2014). Generative Adversarial Nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, & K. Weinberger (Ed.), *Advances in Neural Information Processing Systems 27 (NIPS 2014)*. 27. Montreal: NeurIPS. Retrieved from [https://proceedings.neurips.cc/paper\\_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf)
- Hardle, W., & Simar, L. (2013). *Applied Multivariate Statistical Analysis*. Berlin, Germany: Research Data Center, School of Business and Economics, Humboldt University.
- Hinton, G. E. (2007). Boltzmann machine. *Brain Corporation*, 2(5), 1668. doi:10.4249/scholarpedia.1668
- Kingma, D. P., & Welling, M. (2022, December 10). Auto-Encoding Variational Bayes. *arXiv Preprint*, 1-14. doi:10.48550/arXiv.1312.6114
- Nguyen, L. (2015). *Matrix Analysis and Calculus* (1st ed.). (C. Evans, Ed.) Hanoi, Vietnam: Lambert Academic Publishing. Retrieved March 3, 2014, from <https://www.shuyuan.sg/store/gb/book/matrix-analysis-and-calculus/isbn/978-3-659-69400-4>
- Nguyen, L. (2023). *Tutorial on artificial neural network*. Loc Nguyen's Academic Network. Open Science Framework (OSF). doi:<https://osf.io/k8syc>
- Oord, A. v., Kalchbrenner, N., & Kavukcuoglu, K. (2016, August 19). Pixel Recurrent Neural Networks. *arXiv preprint*, 1-11. doi:10.48550/arXiv.1601.06759
- Oussidi, A., & Elhassouny, A. (2018). Deep generative models: Survey. *2018 International Conference on Intelligent Systems and Computer Vision (ISCV)*. Fez, Morocco: IEEE. doi:10.1109/ISACV.2018.8354080
- Ruthotto, L., & Haber, E. (2021, April 12). An Introduction to Deep Generative Modeling. *arXiv preprint*. doi:10.48550/arXiv.2103.05180

- Theis, L., & Bethge, M. (2015, September 18). Generative Image Modeling Using Spatial LSTMs. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, & R. Garnett (Ed.), *Advances in Neural Information Processing Systems 28 (NIPS 2015)*. 28, pp. 1-9. Montreal: NeurIPS. doi:10.5555/2969442.2969455
- Wikipedia. (2004, November 13). *Boltzmann machine*. (Wikimedia Foundation) Retrieved from Wikipedia website: [https://en.wikipedia.org/wiki/Boltzmann\\_machine](https://en.wikipedia.org/wiki/Boltzmann_machine)
- Wikipedia. (2004, November 15). *Hopfield network*. (Wikimedia Foundation) Retrieved from Wikipedia website: [https://en.wikipedia.org/wiki/Hopfield\\_network](https://en.wikipedia.org/wiki/Hopfield_network)
- Wikipedia. (2004, February 13). *Kullback-Leibler divergence*. (Wikimedia Foundation) Retrieved from Wikipedia website: [https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler\\_divergence](https://en.wikipedia.org/wiki/Kullback%E2%80%93Leibler_divergence)
- Wikipedia. (2005, April 7). *Recurrent neural network*. (Wikimedia Foundation) Retrieved from Wikipedia website: [https://en.wikipedia.org/wiki/Recurrent\\_neural\\_network](https://en.wikipedia.org/wiki/Recurrent_neural_network)
- Wikipedia. (2006, February 9). *Jensen-Shannon divergence*. (Wikimedia Foundation) Retrieved from Wikipedia: [https://en.wikipedia.org/wiki/Jensen%E2%80%93Shannon\\_divergence](https://en.wikipedia.org/wiki/Jensen%E2%80%93Shannon_divergence)
- Wikipedia. (2007, April 16). *Long short-term memory*. (Wikimedia Foundation) Retrieved from Wikipedia website: [https://en.wikipedia.org/wiki/Long\\_short-term\\_memory](https://en.wikipedia.org/wiki/Long_short-term_memory)