

of differentiable generator nets, the criteria are intractable because the data does not specify both the inputs \mathbf{z} and the outputs \mathbf{x} of the generator net. In the case of supervised learning, both the inputs \mathbf{x} and the outputs \mathbf{y} were given, and the optimization procedure needs only to learn how to produce the specified mapping. In the case of generative modeling, the learning procedure needs to determine how to arrange \mathbf{z} space in a useful way and additionally how to map from \mathbf{z} to \mathbf{x} .

Dosovitskiy *et al.* (2015) studied a simplified problem, where the correspondence between \mathbf{z} and \mathbf{x} is given. Specifically, the training data is computer-rendered imagery of chairs. The latent variables \mathbf{z} are parameters given to the rendering engine describing the choice of which chair model to use, the position of the chair, and other configuration details that affect the rendering of the image. Using this synthetically generated data, a convolutional network is able to learn to map \mathbf{z} descriptions of the content of an image to \mathbf{x} approximations of rendered images. This suggests that contemporary differentiable generator networks have sufficient model capacity to be good generative models, and that contemporary optimization algorithms have the ability to fit them. The difficulty lies in determining how to train generator networks when the value of \mathbf{z} for each \mathbf{x} is not fixed and known ahead of each time.

The following sections describe several approaches to training differentiable generator nets given only training samples of \mathbf{x} .

20.10.3 Variational Autoencoders

The **variational autoencoder** or VAE (Kingma, 2013; Rezende *et al.*, 2014) is a directed model that uses learned approximate inference and can be trained purely with gradient-based methods.

To generate a sample from the model, the VAE first draws a sample \mathbf{z} from the code distribution $p_{\text{model}}(\mathbf{z})$. The sample is then run through a differentiable generator network $g(\mathbf{z})$. Finally, \mathbf{x} is sampled from a distribution $p_{\text{model}}(\mathbf{x}; g(\mathbf{z})) = p_{\text{model}}(\mathbf{x} | \mathbf{z})$. However, during training, the approximate inference network (or encoder) $q(\mathbf{z} | \mathbf{x})$ is used to obtain \mathbf{z} and $p_{\text{model}}(\mathbf{x} | \mathbf{z})$ is then viewed as a decoder network.

The key insight behind variational autoencoders is that they may be trained by maximizing the variational lower bound $\mathcal{L}(q)$ associated with data point \mathbf{x} :

$$\mathcal{L}(q) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x}) + \mathcal{H}(q(\mathbf{z} | \mathbf{x})) \quad (20.76)$$

$$= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})} \log p_{\text{model}}(\mathbf{x} | \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) || p_{\text{model}}(\mathbf{z})) \quad (20.77)$$

$$\leq \log p_{\text{model}}(\mathbf{x}). \quad (20.78)$$