

This implicitly imposes a probability distribution over \mathbf{x} :

$$p_{\mathbf{x}}(\mathbf{x}) = \frac{p_{\mathbf{z}}(g^{-1}(\mathbf{x}))}{\left| \det\left(\frac{\partial g}{\partial \mathbf{z}}\right) \right|}. \quad (20.73)$$

Of course, this formula may be difficult to evaluate, depending on the choice of g , so we often use indirect means of learning g , rather than trying to maximize $\log p(\mathbf{x})$ directly.

In some cases, rather than using g to provide a sample of \mathbf{x} directly, we use g to define a conditional distribution over \mathbf{x} . For example, we could use a generator net whose final layer consists of sigmoid outputs to provide the mean parameters of Bernoulli distributions:

$$p(x_i = 1 \mid \mathbf{z}) = g(\mathbf{z})_i. \quad (20.74)$$

In this case, when we use g to define $p(\mathbf{x} \mid \mathbf{z})$, we impose a distribution over \mathbf{x} by marginalizing \mathbf{z} :

$$p(\mathbf{x}) = \mathbb{E}_{\mathbf{z}} p(\mathbf{x} \mid \mathbf{z}). \quad (20.75)$$

Both approaches define a distribution $p_g(\mathbf{x})$ and allow us to train various criteria of p_g using the reparametrization trick of section 20.9.

The two different approaches to formulating generator nets—emitting the parameters of a conditional distribution versus directly emitting samples—have complementary strengths and weaknesses. When the generator net defines a conditional distribution over \mathbf{x} , it is capable of generating discrete data as well as continuous data. When the generator net provides samples directly, it is capable of generating only continuous data (we could introduce discretization in the forward propagation, but doing so would mean the model could no longer be trained using back-propagation). The advantage to direct sampling is that we are no longer forced to use conditional distributions whose form can be easily written down and algebraically manipulated by a human designer.

Approaches based on differentiable generator networks are motivated by the success of gradient descent applied to differentiable feedforward networks for classification. In the context of supervised learning, deep feedforward networks trained with gradient-based learning seem practically guaranteed to succeed given enough hidden units and enough training data. Can this same recipe for success transfer to generative modeling?

Generative modeling seems to be more difficult than classification or regression because the learning process requires optimizing intractable criteria. In the context