

some hidden units to learn to model the input using features that are relevant to the class but also to learn extra hidden units that explain nuisance details that are necessary for the samples of  $\mathbf{v}$  to be realistic but do not determine the class of the example. Another use of higher-order interactions is to gate some features. [Sohn \*et al.\* \(2013\)](#) introduced a Boltzmann machine with third-order interactions with binary mask variables associated with each visible unit. When these masking variables are set to zero, they remove the influence of a visible unit on the hidden units. This allows visible units that are not relevant to the classification problem to be removed from the inference pathway that estimates the class.

More generally, the Boltzmann machine framework is a rich space of models permitting many more model structures than have been explored so far. Developing a new form of Boltzmann machine requires some more care and creativity than developing a new neural network layer, because it is often difficult to find an energy function that maintains tractability of all of the different conditional distributions needed to use the Boltzmann machine, but despite this required effort the field remains open to innovation.

## 20.9 Back-Propagation through Random Operations

Traditional neural networks implement a deterministic transformation of some input variables  $\mathbf{x}$ . When developing generative models, we often wish to extend neural networks to implement stochastic transformations of  $\mathbf{x}$ . One straightforward way to do this is to augment the neural network with extra inputs  $\mathbf{z}$  that are sampled from some simple probability distribution, such as a uniform or Gaussian distribution. The neural network can then continue to perform deterministic computation internally, but the function  $f(\mathbf{x}, \mathbf{z})$  will appear stochastic to an observer who does not have access to  $\mathbf{z}$ . Provided that  $f$  is continuous and differentiable, we can then compute the gradients necessary for training using back-propagation as usual.

As an example, let us consider the operation consisting of drawing samples  $y$  from a Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ :

$$y \sim \mathcal{N}(\mu, \sigma^2). \quad (20.54)$$

Because an individual sample of  $y$  is not produced by a function, but rather by a sampling process whose output changes every time we query it, it may seem counterintuitive to take the derivatives of  $y$  with respect to the parameters of its distribution,  $\mu$  and  $\sigma^2$ . However, we can rewrite the sampling process as