

in one modality, a representation in the other, and the relationship (in general a joint distribution) between pairs  $(\mathbf{x}, \mathbf{y})$  consisting of one observation  $\mathbf{x}$  in one modality and another observation  $\mathbf{y}$  in the other modality (Srivastava and Salakhutdinov, 2012). By learning all three sets of parameters (from  $\mathbf{x}$  to its representation, from  $\mathbf{y}$  to its representation, and the relationship between the two representations), concepts in one representation are anchored in the other, and vice-versa, allowing one to meaningfully generalize to new pairs. The procedure is illustrated in figure 15.3.

## 15.3 Semi-Supervised Disentangling of Causal Factors

An important question about representation learning is “what makes one representation better than another?” One hypothesis is that an ideal representation is one in which the features within the representation correspond to the underlying causes of the observed data, with separate features or directions in feature space corresponding to different causes, so that the representation disentangles the causes from one another. This hypothesis motivates approaches in which we first seek a good representation for  $p(\mathbf{x})$ . Such a representation may also be a good representation for computing  $p(\mathbf{y} \mid \mathbf{x})$  if  $\mathbf{y}$  is among the most salient causes of  $\mathbf{x}$ . This idea has guided a large amount of deep learning research since at least the 1990s (Becker and Hinton, 1992; Hinton and Sejnowski, 1999), in more detail. For other arguments about when semi-supervised learning can outperform pure supervised learning, we refer the reader to section 1.2 of Chapelle *et al.* (2006).

In other approaches to representation learning, we have often been concerned with a representation that is easy to model—for example, one whose entries are sparse, or independent from each other. A representation that cleanly separates the underlying causal factors may not necessarily be one that is easy to model. However, a further part of the hypothesis motivating semi-supervised learning via unsupervised representation learning is that for many AI tasks, these two properties coincide: once we are able to obtain the underlying explanations for what we observe, it generally becomes easy to isolate individual attributes from the others. Specifically, if a representation  $\mathbf{h}$  represents many of the underlying causes of the observed  $\mathbf{x}$ , and the outputs  $\mathbf{y}$  are among the most salient causes, then it is easy to predict  $\mathbf{y}$  from  $\mathbf{h}$ .

First, let us see how semi-supervised learning can fail because unsupervised learning of  $p(\mathbf{x})$  is of no help to learn  $p(\mathbf{y} \mid \mathbf{x})$ . Consider for example the case where  $p(\mathbf{x})$  is uniformly distributed and we want to learn  $f(\mathbf{x}) = \mathbb{E}[\mathbf{y} \mid \mathbf{x}]$ . Clearly, observing a training set of  $\mathbf{x}$  values alone gives us no information about  $p(\mathbf{y} \mid \mathbf{x})$ .