



Figure 15.4: Example of a density over x that is a mixture over three components. The component identity is an underlying explanatory factor, y . Because the mixture components (e.g., natural object classes in image data) are statistically salient, just modeling $p(x)$ in an unsupervised way with no labeled example already reveals the factor y .

Next, let us see a simple example of how semi-supervised learning can succeed. Consider the situation where \mathbf{x} arises from a mixture, with one mixture component per value of \mathbf{y} , as illustrated in figure 15.4. If the mixture components are well-separated, then modeling $p(\mathbf{x})$ reveals precisely where each component is, and a single labeled example of each class will then be enough to perfectly learn $p(\mathbf{y} | \mathbf{x})$. But more generally, what could make $p(\mathbf{y} | \mathbf{x})$ and $p(\mathbf{x})$ be tied together?

If \mathbf{y} is closely associated with one of the causal factors of \mathbf{x} , then $p(\mathbf{x})$ and $p(\mathbf{y} | \mathbf{x})$ will be strongly tied, and unsupervised representation learning that tries to disentangle the underlying factors of variation is likely to be useful as a semi-supervised learning strategy.

Consider the assumption that \mathbf{y} is one of the causal factors of \mathbf{x} , and let \mathbf{h} represent all those factors. The true generative process can be conceived as structured according to this directed graphical model, with \mathbf{h} as the parent of \mathbf{x} :

$$p(\mathbf{h}, \mathbf{x}) = p(\mathbf{x} | \mathbf{h})p(\mathbf{h}). \quad (15.1)$$

As a consequence, the data has marginal probability

$$p(\mathbf{x}) = \mathbb{E}_{\mathbf{h}} p(\mathbf{x} | \mathbf{h}). \quad (15.2)$$

From this straightforward observation, we conclude that the best possible model of \mathbf{x} (from a generalization point of view) is the one that uncovers the above “true”