

be expected via learning a generative model that attempts to recover the causal factors  $\mathbf{h}$  and  $p(\mathbf{x} \mid \mathbf{h})$ .

## 15.4 Distributed Representation

Distributed representations of concepts—representations composed of many elements that can be set separately from each other—are one of the most important tools for representation learning. Distributed representations are powerful because they can use  $n$  features with  $k$  values to describe  $k^n$  different concepts. As we have seen throughout this book, both neural networks with multiple hidden units and probabilistic models with multiple latent variables make use of the strategy of distributed representation. We now introduce an additional observation. Many deep learning algorithms are motivated by the assumption that the hidden units can learn to represent the underlying causal factors that explain the data, as discussed in section 15.3. Distributed representations are natural for this approach, because each direction in representation space can correspond to the value of a different underlying configuration variable.

An example of a distributed representation is a vector of  $n$  binary features, which can take  $2^n$  configurations, each potentially corresponding to a different region in input space, as illustrated in figure 15.7. This can be compared with a *symbolic representation*, where the input is associated with a single symbol or category. If there are  $n$  symbols in the dictionary, one can imagine  $n$  feature detectors, each corresponding to the detection of the presence of the associated category. In that case only  $n$  different configurations of the representation space are possible, carving  $n$  different regions in input space, as illustrated in figure 15.8. Such a symbolic representation is also called a one-hot representation, since it can be captured by a binary vector with  $n$  bits that are mutually exclusive (only one of them can be active). A symbolic representation is a specific example of the broader class of non-distributed representations, which are representations that may contain many entries but without significant meaningful separate control over each entry.

Examples of learning algorithms based on non-distributed representations include:

- Clustering methods, including the  $k$ -means algorithm: each input point is assigned to exactly one cluster.
- $k$ -nearest neighbors algorithms: one or a few templates or prototype examples are associated with a given input. In the case of  $k > 1$ , there are multiple