

The above analysis shows that when we project the data  $\mathbf{x}$  to  $\mathbf{z}$ , via the linear transformation  $\mathbf{W}$ , the resulting representation has a diagonal covariance matrix (as given by  $\mathbf{\Sigma}^2$ ) which immediately implies that the individual elements of  $\mathbf{z}$  are mutually uncorrelated.

This ability of PCA to transform data into a representation where the elements are mutually uncorrelated is a very important property of PCA. It is a simple example of a representation that attempts to *disentangle the unknown factors of variation* underlying the data. In the case of PCA, this disentangling takes the form of finding a rotation of the input space (described by  $\mathbf{W}$ ) that aligns the principal axes of variance with the basis of the new representation space associated with  $\mathbf{z}$ .

While correlation is an important category of dependency between elements of the data, we are also interested in learning representations that disentangle more complicated forms of feature dependencies. For this, we will need more than what can be done with a simple linear transformation.

### 5.8.2 $k$ -means Clustering

Another example of a simple representation learning algorithm is  $k$ -means clustering. The  $k$ -means clustering algorithm divides the training set into  $k$  different clusters of examples that are near each other. We can thus think of the algorithm as providing a  $k$ -dimensional one-hot code vector  $\mathbf{h}$  representing an input  $\mathbf{x}$ . If  $\mathbf{x}$  belongs to cluster  $i$ , then  $h_i = 1$  and all other entries of the representation  $\mathbf{h}$  are zero.

The one-hot code provided by  $k$ -means clustering is an example of a sparse representation, because the majority of its entries are zero for every input. Later, we will develop other algorithms that learn more flexible sparse representations, where more than one entry can be non-zero for each input  $\mathbf{x}$ . One-hot codes are an extreme example of sparse representations that lose many of the benefits of a distributed representation. The one-hot code still confers some statistical advantages (it naturally conveys the idea that all examples in the same cluster are similar to each other) and it confers the computational advantage that the entire representation may be captured by a single integer.

The  $k$ -means algorithm works by initializing  $k$  different centroids  $\{\boldsymbol{\mu}^{(1)}, \dots, \boldsymbol{\mu}^{(k)}\}$  to different values, then alternating between two different steps until convergence. In one step, each training example is assigned to cluster  $i$ , where  $i$  is the index of the nearest centroid  $\boldsymbol{\mu}^{(i)}$ . In the other step, each centroid  $\boldsymbol{\mu}^{(i)}$  is updated to the mean of all training examples  $\mathbf{x}^{(j)}$  assigned to cluster  $i$ .