

Training sparse coding with maximum likelihood is intractable. Instead, the training alternates between encoding the data and training the decoder to better reconstruct the data given the encoding. This approach will be justified further as a principled approximation to maximum likelihood later, in section 19.3.

For models such as PCA, we have seen the use of a parametric encoder function that predicts  $\mathbf{h}$  and consists only of multiplication by a weight matrix. The encoder that we use with sparse coding is not a parametric encoder. Instead, the encoder is an optimization algorithm, that solves an optimization problem in which we seek the single most likely code value:

$$\mathbf{h}^* = f(\mathbf{x}) = \arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{x}). \quad (13.15)$$

When combined with equation 13.13 and equation 13.12, this yields the following optimization problem:

$$\arg \max_{\mathbf{h}} p(\mathbf{h} | \mathbf{x}) \quad (13.16)$$

$$= \arg \max_{\mathbf{h}} \log p(\mathbf{h} | \mathbf{x}) \quad (13.17)$$

$$= \arg \min_{\mathbf{h}} \lambda \|\mathbf{h}\|_1 + \beta \|\mathbf{x} - \mathbf{W}\mathbf{h}\|_2^2, \quad (13.18)$$

where we have dropped terms not depending on  $\mathbf{h}$  and divided by positive scaling factors to simplify the equation.

Due to the imposition of an  $L^1$  norm on  $\mathbf{h}$ , this procedure will yield a sparse  $\mathbf{h}^*$  (See section 7.1.2).

To train the model rather than just perform inference, we alternate between minimization with respect to  $\mathbf{h}$  and minimization with respect to  $\mathbf{W}$ . In this presentation, we treat  $\beta$  as a hyperparameter. Typically it is set to 1 because its role in this optimization problem is shared with  $\lambda$  and there is no need for both hyperparameters. In principle, we could also treat  $\beta$  as a parameter of the model and learn it. Our presentation here has discarded some terms that do not depend on  $\mathbf{h}$  but do depend on  $\beta$ . To learn  $\beta$ , these terms must be included, or  $\beta$  will collapse to 0.

Not all approaches to sparse coding explicitly build a  $p(\mathbf{h})$  and a  $p(\mathbf{x} | \mathbf{h})$ . Often we are just interested in learning a dictionary of features with activation values that will often be zero when extracted using this inference procedure.

If we sample  $\mathbf{h}$  from a Laplace prior, it is in fact a zero probability event for an element of  $\mathbf{h}$  to actually be zero. The generative model itself is not especially sparse, only the feature extractor is. Goodfellow *et al.* (2013d) describe approximate