

maximum likelihood training of a generative model that has latent variables. Suppose we have a model with visible variables \mathbf{x} and latent variables \mathbf{h} , with an explicit joint distribution $p_{\text{model}}(\mathbf{x}, \mathbf{h}) = p_{\text{model}}(\mathbf{h})p_{\text{model}}(\mathbf{x} | \mathbf{h})$. We refer to $p_{\text{model}}(\mathbf{h})$ as the model's prior distribution over the latent variables, representing the model's beliefs prior to seeing \mathbf{x} . This is different from the way we have previously used the word "prior," to refer to the distribution $p(\boldsymbol{\theta})$ encoding our beliefs about the model's parameters before we have seen the training data. The log-likelihood can be decomposed as

$$\log p_{\text{model}}(\mathbf{x}) = \log \sum_{\mathbf{h}} p_{\text{model}}(\mathbf{h}, \mathbf{x}). \quad (14.3)$$

We can think of the autoencoder as approximating this sum with a point estimate for just one highly likely value for \mathbf{h} . This is similar to the sparse coding generative model (section 13.4), but with \mathbf{h} being the output of the parametric encoder rather than the result of an optimization that infers the most likely \mathbf{h} . From this point of view, with this chosen \mathbf{h} , we are maximizing

$$\log p_{\text{model}}(\mathbf{h}, \mathbf{x}) = \log p_{\text{model}}(\mathbf{h}) + \log p_{\text{model}}(\mathbf{x} | \mathbf{h}). \quad (14.4)$$

The $\log p_{\text{model}}(\mathbf{h})$ term can be sparsity-inducing. For example, the Laplace prior,

$$p_{\text{model}}(h_i) = \frac{\lambda}{2} e^{-\lambda|h_i|}, \quad (14.5)$$

corresponds to an absolute value sparsity penalty. Expressing the log-prior as an absolute value penalty, we obtain

$$\Omega(\mathbf{h}) = \lambda \sum_i |h_i| \quad (14.6)$$

$$-\log p_{\text{model}}(\mathbf{h}) = \sum_i \left(\lambda|h_i| - \log \frac{\lambda}{2} \right) = \Omega(\mathbf{h}) + \text{const} \quad (14.7)$$

where the constant term depends only on λ and not \mathbf{h} . We typically treat λ as a hyperparameter and discard the constant term since it does not affect the parameter learning. Other priors such as the Student- t prior can also induce sparsity. From this point of view of sparsity as resulting from the effect of $p_{\text{model}}(\mathbf{h})$ on approximate maximum likelihood learning, the sparsity penalty is not a regularization term at all. It is just a consequence of the model's distribution over its latent variables. This view provides a different motivation for training an autoencoder: it is a way of approximately training a generative model. It also provides a different reason for