

of latent variable models, this means computing

$$\mathbf{h}^* = \arg \max_{\mathbf{h}} p(\mathbf{h} \mid \mathbf{v}). \quad (19.9)$$

This is known as **maximum a posteriori** inference, abbreviated MAP inference.

MAP inference is usually not thought of as approximate inference—it does compute the exact most likely value of \mathbf{h}^* . However, if we wish to develop a learning process based on maximizing $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$, then it is helpful to think of MAP inference as a procedure that provides a value of q . In this sense, we can think of MAP inference as approximate inference, because it does not provide the optimal q .

Recall from section 19.1 that exact inference consists of maximizing

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \mathbb{E}_{\mathbf{h} \sim q} [\log p(\mathbf{h}, \mathbf{v})] + H(q) \quad (19.10)$$

with respect to q over an unrestricted family of probability distributions, using an exact optimization algorithm. We can derive MAP inference as a form of approximate inference by restricting the family of distributions q may be drawn from. Specifically, we require q to take on a Dirac distribution:

$$q(\mathbf{h} \mid \mathbf{v}) = \delta(\mathbf{h} - \boldsymbol{\mu}). \quad (19.11)$$

This means that we can now control q entirely via $\boldsymbol{\mu}$. Dropping terms of \mathcal{L} that do not vary with $\boldsymbol{\mu}$, we are left with the optimization problem

$$\boldsymbol{\mu}^* = \arg \max_{\boldsymbol{\mu}} \log p(\mathbf{h} = \boldsymbol{\mu}, \mathbf{v}), \quad (19.12)$$

which is equivalent to the MAP inference problem

$$\mathbf{h}^* = \arg \max_{\mathbf{h}} p(\mathbf{h} \mid \mathbf{v}). \quad (19.13)$$

We can thus justify a learning procedure similar to EM, in which we alternate between performing MAP inference to infer \mathbf{h}^* and then update $\boldsymbol{\theta}$ to increase $\log p(\mathbf{h}^*, \mathbf{v})$. As with EM, this is a form of coordinate ascent on \mathcal{L} , where we alternate between using inference to optimize \mathcal{L} with respect to q and using parameter updates to optimize \mathcal{L} with respect to $\boldsymbol{\theta}$. The procedure as a whole can be justified by the fact that \mathcal{L} is a lower bound on $\log p(\mathbf{v})$. In the case of MAP inference, this justification is rather vacuous, because the bound is infinitely loose, due to the Dirac distribution's differential entropy of negative infinity. However, adding noise to $\boldsymbol{\mu}$ would make the bound meaningful again.