

MAP inference is commonly used in deep learning as both a feature extractor and a learning mechanism. It is primarily used for sparse coding models.

Recall from section 13.4 that sparse coding is a linear factor model that imposes a sparsity-inducing prior on its hidden units. A common choice is a factorial Laplace prior, with

$$p(h_i) = \frac{\lambda}{2} e^{-\lambda|h_i|}. \quad (19.14)$$

The visible units are then generated by performing a linear transformation and adding noise:

$$p(\mathbf{x} \mid \mathbf{h}) = \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h} + \mathbf{b}, \beta^{-1} \mathbf{I}). \quad (19.15)$$

Computing or even representing $p(\mathbf{h} \mid \mathbf{v})$ is difficult. Every pair of variables h_i and h_j are both parents of \mathbf{v} . This means that when \mathbf{v} is observed, the graphical model contains an active path connecting h_i and h_j . All of the hidden units thus participate in one massive clique in $p(\mathbf{h} \mid \mathbf{v})$. If the model were Gaussian then these interactions could be modeled efficiently via the covariance matrix, but the sparse prior makes these interactions non-Gaussian.

Because $p(\mathbf{h} \mid \mathbf{v})$ is intractable, so is the computation of the log-likelihood and its gradient. We thus cannot use exact maximum likelihood learning. Instead, we use MAP inference and learn the parameters by maximizing the ELBO defined by the Dirac distribution around the MAP estimate of \mathbf{h} .

If we concatenate all of the \mathbf{h} vectors in the training set into a matrix \mathbf{H} , and concatenate all of the \mathbf{v} vectors into a matrix \mathbf{V} , then the sparse coding learning process consists of minimizing

$$J(\mathbf{H}, \mathbf{W}) = \sum_{i,j} |H_{i,j}| + \sum_{i,j} \left(\mathbf{V} - \mathbf{H}\mathbf{W}^\top \right)_{i,j}^2. \quad (19.16)$$

Most applications of sparse coding also involve weight decay or a constraint on the norms of the columns of \mathbf{W} , in order to prevent the pathological solution with extremely small \mathbf{H} and large \mathbf{W} .

We can minimize J by alternating between minimization with respect to \mathbf{H} and minimization with respect to \mathbf{W} . Both sub-problems are convex. In fact, the minimization with respect to \mathbf{W} is just a linear regression problem. However, minimization of J with respect to both arguments is usually not a convex problem.

Minimization with respect to \mathbf{H} requires specialized algorithms such as the feature-sign search algorithm (Lee *et al.*, 2007).