

respect to  $\mathcal{L}$ . Using the same approach on a different model could yield a different functional form of  $q$ .

This was of course, just a small case constructed for demonstration purposes. For examples of real applications of variational learning with continuous variables in the context of deep learning, see [Goodfellow \*et al.\* \(2013d\)](#).

#### 19.4.4 Interactions between Learning and Inference

Using approximate inference as part of a learning algorithm affects the learning process, and this in turn affects the accuracy of the inference algorithm.

Specifically, the training algorithm tends to adapt the model in a way that makes the approximating assumptions underlying the approximate inference algorithm become more true. When training the parameters, variational learning increases

$$\mathbb{E}_{\mathbf{h} \sim q} \log p(\mathbf{v}, \mathbf{h}). \quad (19.68)$$

For a specific  $\mathbf{v}$ , this increases  $p(\mathbf{h} \mid \mathbf{v})$  for values of  $\mathbf{h}$  that have high probability under  $q(\mathbf{h} \mid \mathbf{v})$  and decreases  $p(\mathbf{h} \mid \mathbf{v})$  for values of  $\mathbf{h}$  that have low probability under  $q(\mathbf{h} \mid \mathbf{v})$ .

This behavior causes our approximating assumptions to become self-fulfilling prophecies. If we train the model with a unimodal approximate posterior, we will obtain a model with a true posterior that is far closer to unimodal than we would have obtained by training the model with exact inference.

Computing the true amount of harm imposed on a model by a variational approximation is thus very difficult. There exist several methods for estimating  $\log p(\mathbf{v})$ . We often estimate  $\log p(\mathbf{v}; \boldsymbol{\theta})$  after training the model, and find that the gap with  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q)$  is small. From this, we can conclude that our variational approximation is accurate for the specific value of  $\boldsymbol{\theta}$  that we obtained from the learning process. We should not conclude that our variational approximation is accurate in general or that the variational approximation did little harm to the learning process. To measure the true amount of harm induced by the variational approximation, we would need to know  $\boldsymbol{\theta}^* = \max_{\boldsymbol{\theta}} \log p(\mathbf{v}; \boldsymbol{\theta})$ . It is possible for  $\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) \approx \log p(\mathbf{v}; \boldsymbol{\theta})$  and  $\log p(\mathbf{v}; \boldsymbol{\theta}) \ll \log p(\mathbf{v}; \boldsymbol{\theta}^*)$  to hold simultaneously. If  $\max_q \mathcal{L}(\mathbf{v}, \boldsymbol{\theta}^*, q) \ll \log p(\mathbf{v}; \boldsymbol{\theta}^*)$ , because  $\boldsymbol{\theta}^*$  induces too complicated of a posterior distribution for our  $q$  family to capture, then the learning process will never approach  $\boldsymbol{\theta}^*$ . Such a problem is very difficult to detect, because we can only know for sure that it happened if we have a superior learning algorithm that can find  $\boldsymbol{\theta}^*$  for comparison.