

Approximate Inference

Presenter: Ming-Chang Chiu

September 10, 2019

Many probabilistic models are difficult to train because it is difficult to perform inference in them. In deep learning, we usually have a set of visible variables \mathbf{v} and a set of latent variables \mathbf{h} and the difficulty often refers to computing $p(\mathbf{h}|\mathbf{v})$.

Most graphical models have many hidden layer and that makes the $p(\mathbf{h}|\mathbf{v})$ intractable, since exact inference requires an exponential amount of time in these models (See Ch16). In this chapter, we introduce several of the techniques for confronting these intractable inference problems.

1 Inference as Optimization

One approach to address the difficulty in inference is to cast the exact inference as an optimization problem, so that approximating the exact inference problem becomes approximating the underlying optimization problem. E.g. Assume we have a probabilistic model $p(\mathbf{v}, \mathbf{h})$, and we want to know log probability of observed data, $\log p(\mathbf{v}; \boldsymbol{\theta})$, but sometimes it is hard to marginalize out \mathbf{h} . Instead of computing it directly, we compute the lower bound on $L(\mathbf{v}, \boldsymbol{\theta}, q)$ on $\log p(\mathbf{v}; \boldsymbol{\theta})$, which is called *evidence lower bound* (ELBO), defined to be

$$L(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}; \boldsymbol{\theta})) \quad (1)$$

where q is an arbitrary probability distribution over \mathbf{h} .

Since KL divergence is always non-negative, we can see that L is always at most the same value as the resired probability, and so, to find the approximate value of $\log p(\mathbf{v}; \boldsymbol{\theta})$ becomes an optmization with objective to maximize L .

For an appropriate choice of q , L is tractable to compute. And algebraically, we can rearrange L into a more convenient form (see appendix):

$$L(\mathbf{v}, \boldsymbol{\theta}, q) = -\mathbf{E}_{\mathbf{h} \sim q}[\log q(\mathbf{h}|\mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})] \quad (2)$$

We can now think of inference as procedure for finding the q that maximizes L . We can make the optimization procedure less expensive but approximate by restricting the family of distributions q the optimization is allowed to search over or by using an imperfect optimization procedure that may not completely maximize L but merely increase it by a significant amount.

2 Expectation Maximization

The first algorithm to maximize the lower bound L is *expectation maximization* (EM) algorithm. Note that EM is not an approach to approximate inference, but rather an approach to learning with an approximate posterior.

The EM algorithm consists of alternating between two steps until convergence:

- E-step: Let $\theta^{(0)}$ denote the value of the parameters at the beginning of the step. Set $q(\mathbf{h}^{(i)}|\mathbf{v}) = p(\mathbf{h}^{(i)}|\mathbf{v}^{(i)}; \theta^{(0)})$ for all indices i of the training examples $\mathbf{v}^{(i)}$ we want to train on (both batch and minibatch variants are valid). By this we mean q is defined in terms of the current parameter value of $\theta^{(0)}$; if we vary θ then $p(\mathbf{h}|\mathbf{v}; \theta)$ will change but $q(\mathbf{h}|\mathbf{v})$ will remain equal to $p(\mathbf{h}|\mathbf{v}; \theta^{(0)})$.
- M-step: Completely or partially maximize

$$\sum_i L(\mathbf{v}^i, \theta, q) \tag{3}$$

with respect to using your optimization θ algorithm of choice.

This can be viewed as a coordinate ascent algorithm to maximize L . On one step, we maximize L with respect to q , and on the other, we maximize L with respect to θ .

One can use stochastic gradient ascent on latent variable models, which is seen as a special case of EM where M-step consists of taking a single gradient step. Other variants of the EM algorithm can make much larger steps.

Even though the E-step involves exact inference, we can think of the EM algorithm as using approximate inference in some sense. Specifically, the M-step assumes that the same value of q can be used for all values of θ . This will introduce a gap between L and the true $\log p(\mathbf{v})$ as the M-step moves further and further away from the value $\theta^{(0)}$ used in the E-step. Fortunately, the E-step reduces the gap to zero again as we enter the loop for the next time. (give example)

We can get an insight that there is the basic structure of the learning process, in which we update the model parameters to improve the likelihood of a completed dataset, where all missing variables have their values provided by an estimate of the posterior distribution.

3 MAP Inference and Sparse Coding

As before, we are interested in computing $p(\mathbf{h}|\mathbf{v})$. An alternative form of inference is to compute the **single** most likely value of the missing variables, rather than to infer the entire distribution over their possible values. In the context of latent variable models, this means computing

$$\mathbf{h}^* = \operatorname{argmax}_h p(\mathbf{h}|\mathbf{v}) \tag{4}$$

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

Recall from Eq.2,

$$L(\mathbf{v}, \boldsymbol{\theta}, q) = -\mathbf{E}_{\mathbf{h} \sim q}[\log q(\mathbf{h}|\mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})] \quad (5)$$

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:

$$q(\mathbf{h}|\mathbf{v}) = \delta(\mathbf{h} - \boldsymbol{\mu}) \quad (6)$$

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

$$\boldsymbol{\mu}^* = \underset{\boldsymbol{\mu}}{\operatorname{argmax}} \log p(\mathbf{h} = \boldsymbol{\mu}, \mathbf{v}) \quad (7)$$

which is equivalent to Eq.4.

As with EM, this is a form of coordinate ascent on L , where we alternate between using inference to optimize L with respect to q and using parameter updates to optimize L with respect to $\boldsymbol{\theta}$. MAP is primarily used for sparse coding models.

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

$$p(h_i) = \frac{\lambda}{2} e^{-\lambda|h_i|} \quad (8)$$

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

$$p(\mathbf{v}|\mathbf{h}) = \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h} + \mathbf{b}, \beta^{-1}\mathbf{I}) \quad (9)$$

Computing or even representing $p(\mathbf{h}|\mathbf{v})$ is difficult. Every pair of variables h_i and h_j are both parents of \mathbf{v} .

Because $p(\mathbf{h}|\mathbf{v})$ is intractable, 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

$$J(\mathbf{H}, \mathbf{W}) = \sum_{i,j} |H_{i,j}| + \sum_{i,j} (\mathbf{V} - \mathbf{H}\mathbf{W}^T)_{i,j}^2 \quad (10)$$

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.

4 Variational Inference and Learning

The core idea behind variational learning is that we can maximize L over a restricted family of distributions q . This family should be chosen so that it is easy to compute $\mathbf{E}_q \log p(\mathbf{h}, \mathbf{v})$. A typical way to do this is to introduce assumptions about how q factorizes.

A common approach to variational learning is to impose the restriction that q is a factorial distribution:

$$q(\mathbf{h}|\mathbf{v}) = \prod_i q(h_i|\mathbf{v}) \quad (11)$$

This is called the mean field approach. Generally speaking, we can impose any graphical model structure we choose on q , to flexibly determine how many interactions we want our approximation to capture. We do not need to specify a specific parametric form for q . We specify how it should factorize, but then the optimization problem determines the optimal probability distribution within those factorization constraints. For discrete latent variables, we use traditional optimization techniques; for continuous latent variables, we use calculus of variations to perform optimization, and actually determine which function should be used to represent q . Calculus of variations makes designer of the model must only specify how q factorizes, rather than q .

Recall from Eq.2, we can think of maximizing L w.r.t. q as minimizing $D_{KL}(q(\mathbf{h}|\mathbf{v})||p(\mathbf{h}|\mathbf{v}))$.

4.1 Discrete Latent Variables

Variational inference with discrete latent variables is relatively straightforward. We define a distribution q , typically one where each factor of q is just defined by a lookup table over discrete states. After determining how to represent q , we simply optimize its parameters. In principle this can be done with any optimization algorithm, such as gradient descent.

Because this optimization must occur in the inner loop of a learning algorithm, it must be very fast. To achieve this speed, we typically use special optimization algorithms that are designed to solve comparatively small and simple problems in very few iterations. A popular choice is to iterate fixed point equations, in other words, to solve

$$\frac{\partial}{\partial \hat{h}_i} L = 0 \quad (12)$$

for \hat{h}_i . We repeatedly update different elements of $\hat{\mathbf{h}}$ until we satisfy a convergence criterion.

4.2 Calculus of Variations

Calculus of Variations enables us to actually solve for a function $f(\mathbf{x})$, such as when we want to find the probability density function over some random variable. A function of a function f is known as a functional $J[f]$. We can take functional derivatives, also known as variational derivatives, of a functional $J[f]$ with respect to individual values of the function $f(\mathbf{x})$ at any specific value of \mathbf{x} . The functional derivative of the functional J with respect to the value of the function f at point

\mathbf{x} is denoted $\frac{\delta}{\delta f(\mathbf{x})}J$. For instance, for differentiable functions $f(\mathbf{x})$ and differentiable functions $g(y, \mathbf{x})$ with continuous derivatives, that

$$\frac{\delta}{\delta f(\mathbf{x})} \int g(f(\mathbf{x}), \mathbf{x}) d\mathbf{x} = \frac{\partial}{\partial y} g(f(\mathbf{x}), \mathbf{x}) \quad (13)$$

Similar to optimizing a function w.r.t. a vector, we can optimize a functional by solving for the function where the functional derivative at every point is equal to zero.

Example: consider the problem of finding the probability distribution function over $x \in \mathbb{R}$ that has maximal differential entropy. Recall that

$$H[p] = -\mathbf{E}_x \log p(x) = - \int p(x) \log p(x) dx \quad (14)$$

Adding constraints using Lagrange multipliers regarding mean, variance, and probability distribution.

$$L[p] = \lambda_1 \left(\int p(x) dx - 1 \right) + \lambda_2 (\mathbf{E}[x] - \mu) + \lambda_3 (\mathbf{E}[(x - \mu)^2] - \sigma^2) + H[p] \quad (15)$$

$$= \int (\lambda_1 p(x) + \lambda_2 p(x)x + \lambda_3 p(x)(x - \mu)^2 - p(x) \log p(x)) dx - \lambda_1 - \mu \lambda_2 - \sigma^2 \lambda_3 \quad (16)$$

And set functional derivatives equal to 0:

$$\forall x, \frac{\delta}{\delta p(x)} L = \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1 - \log p(x) = 0 - \sigma^2 \lambda_3 \quad (17)$$

And so,

$$p(x) = \exp(\lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2 - 1) \quad (18)$$

To finish this optimization problem, we are left with choosing the λ values, as long as the constraints are satisfied. In the book, the authors set $\lambda_1 = 1 - \log \sigma \sqrt{2\pi}$, $\lambda_2 = 0$, $\lambda_3 = -\frac{1}{2\sigma^2}$, to obtain $p(x) = \mathcal{N}(x; \mu, \sigma^2)$. This is one reason for using the normal distribution when we do not know the true distribution.

4.3 Continuous Latent Variables

When our graphical model contains continuous latent variables, we must now use calculus of variations when maximizing L with respect to $q(\mathbf{h}|\mathbf{v})$. There is a general equation for the mean field fixed point updates, so people need not solve calculus of variations themselves. If we make the mean field approximation

$$q(\mathbf{h}|\mathbf{v}) = \prod_i q(h_i|\mathbf{v}) \quad (19)$$

and fix $q(h_j|\mathbf{v})$ for all $j \neq i$, then the optimal $q(h_i|\mathbf{v})$ may be obtained by

$$\tilde{q}(\mathbf{h}|\mathbf{v}) = \exp(\mathbf{E}_{\mathbf{h}_{-i} \sim q(\mathbf{h}_{-i}|\mathbf{v})} \log \tilde{p}(\mathbf{v}, \mathbf{h})) \quad (20)$$

so long as $p(\mathbf{v}, \mathbf{h}) \neq 0$ at everywhere. Now Eq.20 is a fixed point equation, and can be iteratively for all i until convergence. We can also regard some of the values appear in it as parameters and optimize with any algorithm we like.

Example: Suppose that $p(\mathbf{h}) = \mathcal{N}(\mathbf{h}; 0, \mathbf{I})$ and $p(v|\mathbf{h}) = \mathcal{N}(v; \mathbf{w}^T \mathbf{h}; 1)$, with $\mathbf{h} \in \mathbb{R}^2$, and single v . The true posterior is given by

$$p(\mathbf{h}|\mathbf{v}) \tag{21}$$

$$\propto p(\mathbf{h}, \mathbf{v}) \tag{22}$$

$$= p(h_1)p(h_2)p(\mathbf{v}|\mathbf{h}) \tag{23}$$

$$\propto \exp(-\frac{1}{2}[h_1^2 + h_2^2 + (v - h_1w_1 - v - h_2w_2)^2]) \tag{24}$$

Applying Eq.20, we find that

$$\tilde{q}(h_1|\mathbf{v}) = \exp(\mathbf{E}_{h_2 \sim q(h_2|\mathbf{v})} \log \tilde{p}(\mathbf{v}, \mathbf{h})) \tag{25}$$

$$= \exp(-\frac{1}{2}\mathbf{E}_{h_2 \sim q(h_2|\mathbf{v})}[h_1^2 + h_2^2 + (v - h_1w_1 - v - h_2w_2)^2]) \tag{26}$$

$$= \exp(-\frac{1}{2}\mathbf{E}_{h_2 \sim q(h_2|\mathbf{v})}[h_1^2 + h_2^2 + v^2 + h_1^2w_1^2 + h_2^2w_2^2 - 2vh_1w_1 - 2vh_2w_2 + 2h_1w_1h_2w_2]) \tag{27}$$

$$= \exp(-\frac{1}{2}[h_1^2 + \langle h_2^2 \rangle + v^2 + h_1^2w_1^2 + \langle h_2^2 \rangle w_2^2 - 2vh_1w_1 - 2v\langle h_2 \rangle w_2 + 2h_1w_1\langle h_2 \rangle w_2]) \tag{28}$$

From this we can see that \tilde{q} has the functino form of gaussian, which is not assumed at the outset. Using the same approach on a different model could yield a different functional form of q .

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

$$E_{\mathbf{h} \sim q} \log p(\mathbf{v}, \mathbf{h}) \tag{29}$$

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

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.

5 Learned Approximate Inference

We can think of the optimization process as a function f that maps an input \mathbf{v} to an approximate distribution $q^* = \operatorname{argmax}_q L(\mathbf{v}, q)$. Other methods of approximate inference can be very expensive, such as fixed point equations or gradient-based methods. To avoid such, we can think of the multi-step iterative optimization process as just being a function, and approximate it with a neural network that implements an approximation $\hat{f}(\mathbf{v}; \boldsymbol{\theta})$.

5.1 Wake-Sleep

It is sometimes hard to train a model to infer \mathbf{h} from \mathbf{v} , since we do not have a supervised training set. The wake-sleep algorithm resolves this problem by drawing samples of both \mathbf{h} and \mathbf{v} from the model distribution.

Example: in a directed model, this can be done cheaply by performing ancestral sampling beginning at \mathbf{h} and ending at \mathbf{v} . The inference network can then be trained to perform the reverse mapping: predicting which \mathbf{h} caused the present \mathbf{v} .

5.2 Other Forms of Learned Inference

Other than wake-sleep, Salakhutdinov and Larochelle (2010) showed that a single pass in a learned inference network could yield faster inference than iterating the mean field fixed point equations in a DBM.

Variational Autoencoder, one of the recent approach to generative modeling, is also a form of learned approximate inference. In this elegant approach, there is no need to construct explicit targets for the inference network. Instead, the inference network is simply used to define L elegant approach.