

# Approximate Inference: MAP

Sargur N. Srihari  
[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

# Topics in Approximate Inference

- Task of Inference

- Intractability in Inference

1. Inference as Optimization

2. Expectation Maximization

- 3. MAP Inference and Sparse Coding**

4. Variational Inference and Learning

5. Learned Approximate Inference

# Definition of MAP inference

- Usual form of Inference
  - Compute distribution of a set of variables given another
  - With latent variables, our interest is  $p(\mathbf{h} | \mathbf{v})$
- Maximum a posteriori (MAP) inference:
  - Compute single most likely value of missing variables, rather than infer the entire distribution over their possible values
  - With latent variables, compute

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

# Role of MAP in Approximate Inference

- MAP inference is not approximate inference
  - it does compute the exact most likely value of  $h^*$
- However, if we wish to develop a learning process based on maximizing  $L(v, h, q)$ , then it is helpful to think of MAP inference as a procedure that provides a value of  $q$
- In this sense, MAP inference is approximate inference, because it does not provide the optimal  $q$

# Restricting the family $q$

- Recall that exact inference consists of maximizing

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

- wrt  $q$  over an unrestricted 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

# $q$ as a point estimate

- Specifically,  $q$  is a Dirac distribution:

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

- This means 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

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

- which is equivalent to the MAP inference problem

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

- We can use learning similar to EM

- alternate between MAP inference to infer  $\mathbf{h}^*$  and then update  $\theta$  to increase  $\log p(\mathbf{h}^*, \mathbf{v})$

# EM-like algorithm for MAP

- As with EM, this is coordinate ascent on  $\mathcal{L}$ 
  - Alternate between using inference to optimize  $\mathcal{L}$  wrt  $q$  and using parameter updates to optimize  $\mathcal{L}$  wrt  $\theta$ 
    - Justified by fact that  $\mathcal{L}$  is lower bound on  $\log p(v)$
    - With MAP inference, justification is vacuous, because bound is infinitely loose, due to the Dirac distribution's differential entropy of negative infinity
  - However, adding noise to  $\mu$  would make the bound meaningful again

# Use of MAP inference

- MAP inference is commonly used in deep learning as
  - Feature extractor
  - Learning mechanism
- It is primarily used for sparse coding models

# Sparse Coding

- Sparse coding is a linear factor model
  - It imposes a sparsity-inducing prior on hidden units
    - A common choice is a factorial Laplace prior
$$p(h_i) = \frac{\lambda}{2} e^{-\lambda|h_i|}$$
    - Visible units generated by a linear transformation + noise:
$$p(x | h) = \mathcal{N}(v; Wh + b, \beta^{-1} I)$$
- Even representing  $p(h | v)$  is difficult
  - Every pair  $h_i$  and  $h_j$  are both parents of  $v$ 
    - When  $v$  is observed, active path connects  $h_i$  and  $h_j$
    - All hidden units are in one massive clique in  $p(h | v)$
- MAP solution is used