

Boltzmann Machines for Real Data

Sargur N. Srihari

srihari@cedar.buffalo.edu

Topics

1. Boltzmann machines
 2. Restricted Boltzmann machines
 3. Deep Belief Networks
 4. Deep Boltzmann machines
 5. Boltzmann machines for continuous data
 6. Convolutional Boltzmann machines
 7. Boltzmann machines for structured and sequential outputs
 8. Other Boltzmann machines
9. Backpropagation through random operations
10. Directed generative nets
11. Drawing samples from autoencoders
12. Generative stochastic networks
13. Other generative schemes
14. Evaluating generative models
15. Conclusion

Topics in Boltzmann for Real Data

- Need for real-valued Boltzmann
- Gaussian-Bernoulli RBMs
- Undirected Models of Conditional Covariance

Need for Real-valued data

- While Boltzmann machines were originally developed for use with binary data, many applications such as image and audio modeling require ability to represent distributions over real values
- In some cases, it is possible to treat real-valued data in the interval $[0, 1]$ as representing the expectation of a binary variable.

Short-cut for real values

- Grayscale images in the training set as defining $[0, 1]$ probability values
 - Each pixel defines probability of a binary value being 1, binary pixels are sampled independently
 - A common procedure for evaluating binary models on grayscale image datasets
 - Not theoretically satisfying, and binary images sampled independently have a noisy appearance
- Here we present Boltzmann machines that define a probability density over real-valued data.

Gaussian-Bernoulli RBMs

- RBMs may be developed for many exponential family conditional distributions
- Of these, the most common is the RBM with:
 - Binary hidden units
 - Real-valued visible units
 - Conditional over visible units is Gaussian
 - whose mean is a function of the hidden units

Parameterizing Gaussian-Bernoulli

- We wish to have the conditional distribution

$$p(\mathbf{v} | \mathbf{h}) = N(\mathbf{v}; \mathbf{W}\mathbf{h}, \beta^{-1})$$

- Find terms to add to the energy function by expanding the unnormalized log conditional:

$$\log \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \beta^{-1}) = -\frac{1}{2} (\mathbf{v} - \mathbf{W}\mathbf{h})^\top \beta (\mathbf{v} - \mathbf{W}\mathbf{h}) + f(\beta)$$

- Here f encapsulates terms that are a function only of the parameters and not the random variables
- Discard f because its only role is to normalize
 - Partition function of whatever energy function we choose will carry out that role

Representing $p(\mathbf{v}|\mathbf{h})$

- If we include all of the terms (with their sign flipped) involving \mathbf{v} from equation

$$\log \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \beta^{-1}) = -\frac{1}{2} (\mathbf{v} - \mathbf{W}\mathbf{h})^\top \beta (\mathbf{v} - \mathbf{W}\mathbf{h}) + f(\beta)$$

- in our energy function and do not add any other terms involving \mathbf{v} , then our energy function will represent the desired conditional $p(\mathbf{v} | \mathbf{h})$

Representing $p(\mathbf{h}|\mathbf{v})$

- We have some freedom regarding $p(\mathbf{h} | \mathbf{v})$

$$\log \mathcal{N}(\mathbf{v}; \mathbf{W}\mathbf{h}, \beta^{-1}) = -\frac{1}{2} (\mathbf{v} - \mathbf{W}\mathbf{h})^\top \beta (\mathbf{v} - \mathbf{W}\mathbf{h}) + f(\beta) \quad \text{contains a term} \quad \frac{1}{2} \mathbf{h}^\top \mathbf{W}^\top \beta \mathbf{W} \mathbf{h}$$

- which cannot be included because of $h_i h_j$ terms,
 - i.e., edges between hidden units, which gives a linear factor model instead of RBM
 - Omitting $h_i h_j$ cross terms does not change $p(\mathbf{v} | \mathbf{h})$
- Whether to include terms involving only a single h_i
 - With diagonal precision, for each h_i we have a term

$$\frac{1}{2} h_i \sum_j \beta_j W_{j,i}^2$$

- Using $h_i^2 = h_i$ because $h_i \in \{0, 1\}$.
- If we include this term (with its sign flipped) in the energy function, then it will naturally bias h_i to be turned off when the weights for that unit are large and connected to visible units with high precision.