

Boltzmann Machines

Sargur N. Srihari

srihari@cedar.buffalo.edu

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

1. Boltzmann Machines (BMs)

- Introduced as general connectionist approach for learning arbitrary probability distributions over binary vectors
- Variants of BMs include other kinds of variables
 - Surpassed popularity of the original
- First we consider binary Boltzmann machines and discuss their training and inference

Binary Boltzmann Machine

- We define a Boltzmann machine over a d -dimensional binary vector $\mathbf{x} \in \{0,1\}^d$
- Boltzmann machine is an energy-based model that defines the joint probability distribution

$$P(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z}$$

- Energy function $E(\mathbf{x})$ is defined by $E(\mathbf{x}) = \mathbf{x}^T U \mathbf{x} - \mathbf{b}^T \mathbf{x}$
 - where U is the weight matrix of model parameters and \mathbf{b} is the vector of bias parameters
- Z is the partition function that ensures $\sum_{\mathbf{x}} P(\mathbf{x}) = 1$

Boltzmann Machine as Linear Model

- Boltzmann machine is a joint probability distribution over observed variables

$$P(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z}$$
$$E(\mathbf{x}) = \mathbf{x}^T U \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

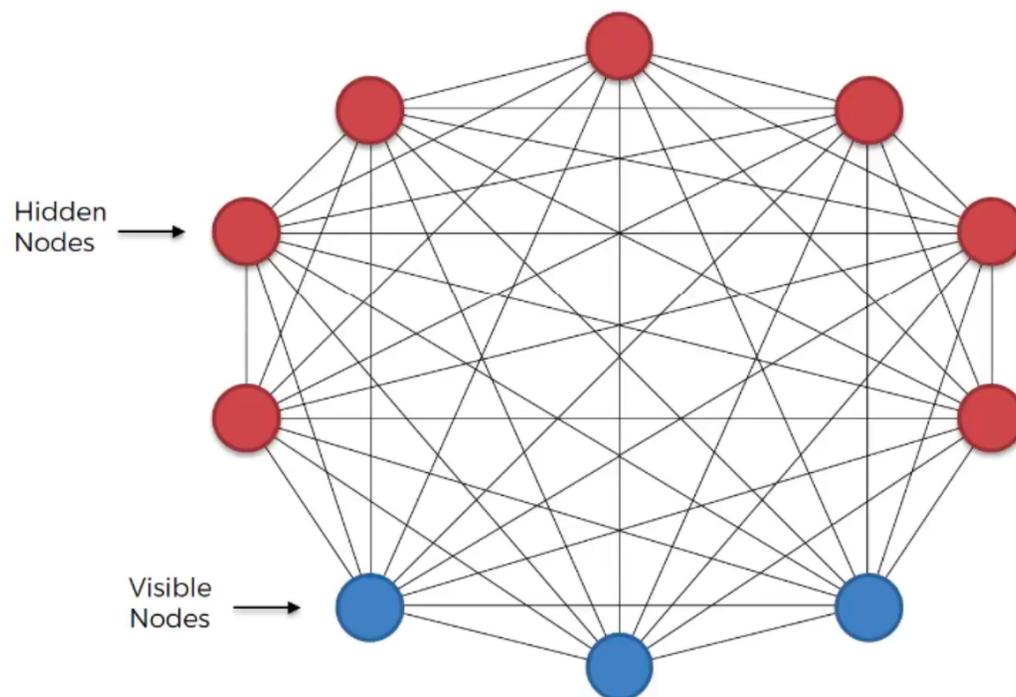
- It defines a distribution where the probability of a given unit being on is determined by a linear model (logistic regression) of the other variables
- In the general setting we are given training examples over all the variables

Boltzmann with hidden units

- Boltzmann becomes more powerful when not all variables are observed
- Just as hidden units convert logistic regression to MLP (a universal approximator of functions)
 - No longer limited to modeling linear relationships between variables
 - Boltzmann becomes universal approximator of probability mass functions over discrete variables
 - Units \mathbf{x} are decomposed into visible units \mathbf{v} and hidden units \mathbf{h} . Energy function becomes

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{v}^T R \mathbf{v} - \mathbf{v}^T W \mathbf{h} - \mathbf{h}^T S \mathbf{h} - \mathbf{b}^T V - \mathbf{c}^T H$$

Boltzmann with hidden nodes



Energy Model:

$$E(\mathbf{x})=E(\mathbf{v},\mathbf{h})=-\mathbf{v}^T R \mathbf{v} - \mathbf{v}^T W \mathbf{h} - \mathbf{h}^T S \mathbf{h} - \mathbf{b}^T V - \mathbf{c}^T H$$

Probability distribution:

$$P(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z}$$

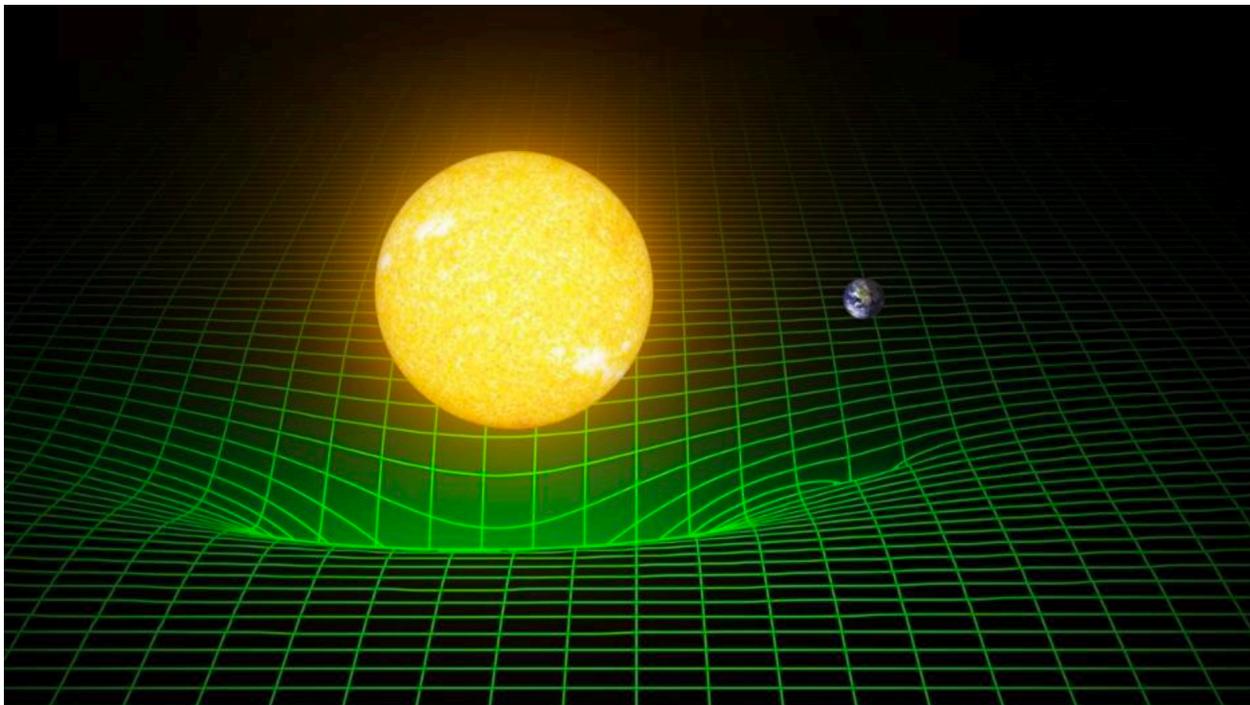
Characteristics of Boltzmann Machines

- Boltzmann Machines are neural networks
 - Where all neurons are connected to each other
 - Including *visible* to visible, and *hidden* to hidden
- Boltzmann Machine doesn't expect input data, it generates data
 - Neurons *generate information regardless they are hidden or visible.*
- Whole thing is a system and it generates the state of the system

Energy based model in physics

Gravitational energy of two body masses:

Describes the potential energy a body with mass has in relation to another massive object due to gravity.



Specifying energy model

- Joint distribution of \mathbf{v} and \mathbf{h}

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}$$

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

- It is difficult to calculate the partition function Z
- Much easier to calculate conditional probabilities

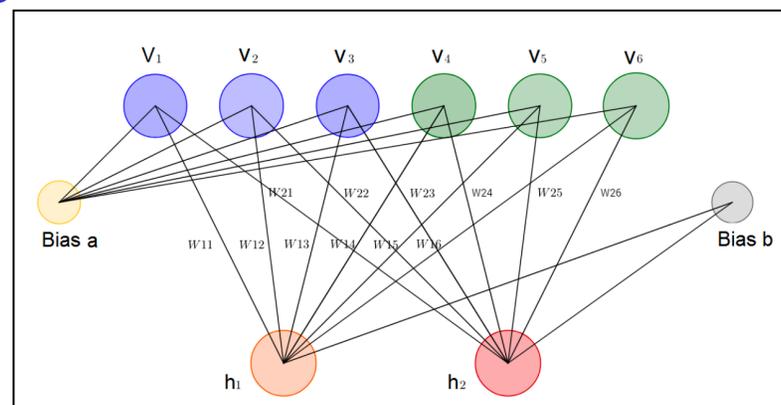
$$p(\mathbf{h}|\mathbf{v}) = \prod_i p(h_i|\mathbf{v})$$

$$p(\mathbf{v}|\mathbf{h}) = \prod_i p(v_i|\mathbf{h})$$

- Since each neuron is binary-valued, for an RBM:

$$p(h_j = 1|\mathbf{v}) = \frac{1}{1 + e^{-(b_j + \sum_i v_i w_{ij})}} = \sigma(b_j + \sum_i v_i w_{ij})$$

$$p(v_j = 1|\mathbf{h}) = \frac{1}{1 + e^{-(b_j + \sum_i v_i w_{ij})}} = \sigma(b_j + \sum_i v_i w_{ij})$$



Boltzmann Machine Learning

- Usually based on maximum likelihood
- All Boltzmann machines have an intractable partition function
 - So max. likelihood gradient has to be approximated
- Interesting property:
 - Update for a particular weight connecting two units depends only on the statistics of the two units collected under different distributions
 - $P_{\text{model}}(\mathbf{v})$, $P_{\text{data}}(\mathbf{v})$, $P_{\text{model}}(\mathbf{h}|\mathbf{v})$
 - The rest of model shapes those statistics
 - This means learning rule is “local” which makes Boltzmann learning biologically plausible

Biological Plausibility

- If each neuron were a random variable in a Boltzmann machine
 - Then axons and dendrites connecting two variables could learn only by observing the firing patterns of variables that they touch
 - Positive phase: two units that frequently fire together have their connection strengthen
 - This is an example of Hebbian learning rule
 - Oldest hypothesized biological learning