



國立交通大學
National Chiao Tung University

Deep Learning

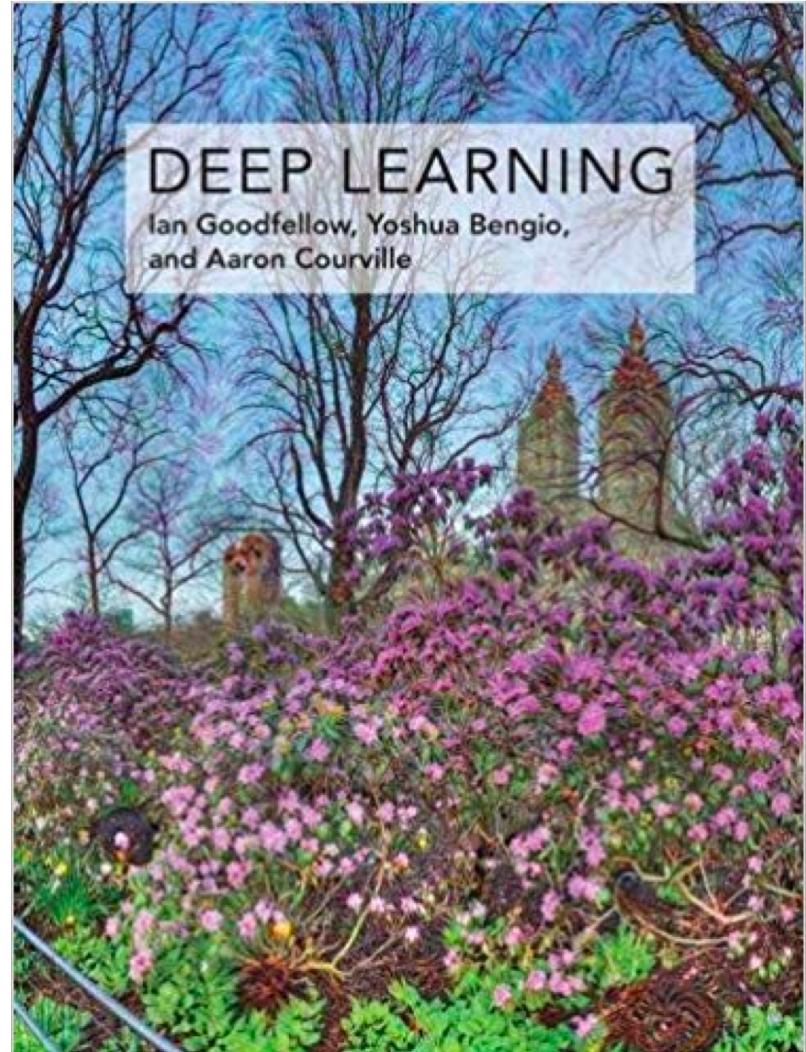
深度學習

Fall 2018

Boltzmann Machines

(Chapter 20)

Prof. Chia-Han Lee
李佳翰 副教授





Boltzmann machines

- Boltzmann machines were introduced to learning arbitrary probability distributions over binary vectors.
- We define the Boltzmann machine over a d -dimensional binary random vector $\mathbf{x} \in \{0,1\}^d$. The Boltzmann machine is an energy-based model, meaning we define the joint probability distribution using an energy function:

$$P(\mathbf{x}) = \frac{\exp(-E(\mathbf{x}))}{Z}, \quad (20.1)$$

where $E(\mathbf{x})$ is the energy function, and Z is the partition function that ensures that $\sum_{\mathbf{x}} P(\mathbf{x}) = 1$.

- The energy function of the Boltzmann machine is

$$E(\mathbf{x}) = -\mathbf{x}^\top \mathbf{U} \mathbf{x} - \mathbf{b}^\top \mathbf{x}, \quad (20.2)$$

where \mathbf{U} is the “weight” matrix of model parameters and \mathbf{b} is the vector of bias parameters.



Boltzmann machines

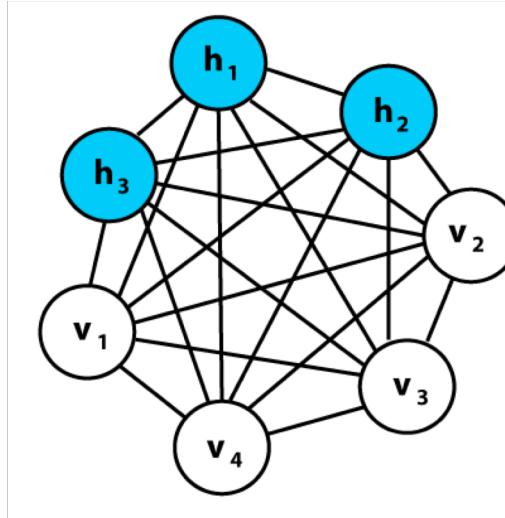
- The Boltzmann machine becomes more powerful when **not all the variables are observed**. In this case, the **latent variables** can act similarly to hidden units in a multilayer perceptron and **model higher-order interactions among the visible units**.
- Just as the addition of hidden units to convert logistic regression into an MLP results in the MLP being a universal approximator of functions, a Boltzmann machine with hidden units is no longer limited to modeling linear relationships between variables. Instead, **the Boltzmann machine becomes a universal approximator of probability mass functions over discrete variables**.



Boltzmann machines

- Formally, we decompose the units x into two subsets: the **visible units v** and the **latent (or hidden) units h** . The energy function becomes

$$E(v, h) = -v^\top R v - v^\top W h - h^\top S h - b^\top v - c^\top h. \quad (20.3)$$



- Learning algorithms for Boltzmann machines are usually based on **maximum likelihood**. All Boltzmann machines have an **intractable partition function**, so the maximum likelihood gradient must be approximated.



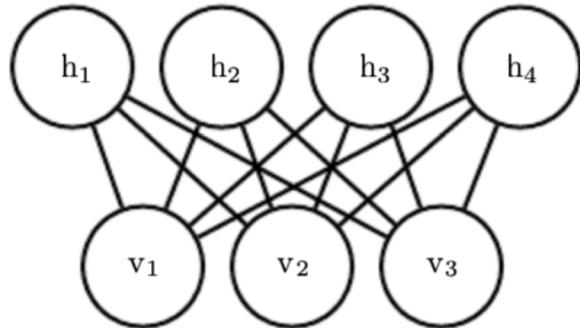
Restricted Boltzmann machines

- Restricted Boltzmann machines (RBMs) are undirected probabilistic graphical models containing a layer of observable variables and a single layer of latent variables. RBMs may be stacked (one on top of the other) to form deeper models.
- The graph of the RBM is a bipartite graph, with no connections permitted between any variables in the observed layer or between any units in the latent layer.
- We begin with the binary version of the restricted Boltzmann machine, but as we see later, there are extensions to other types of visible and hidden units.

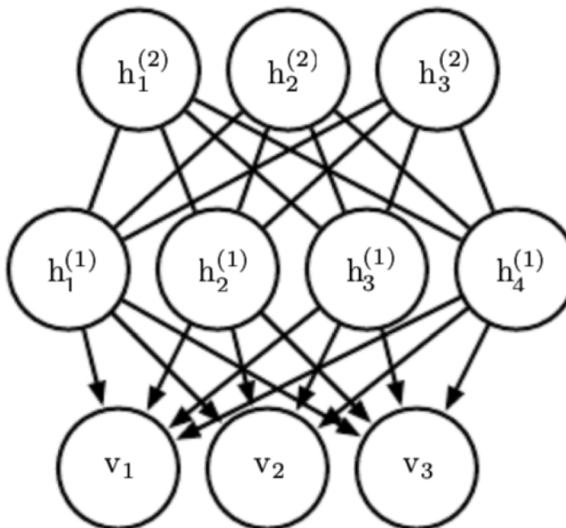


Restricted Boltzmann machines

RBM

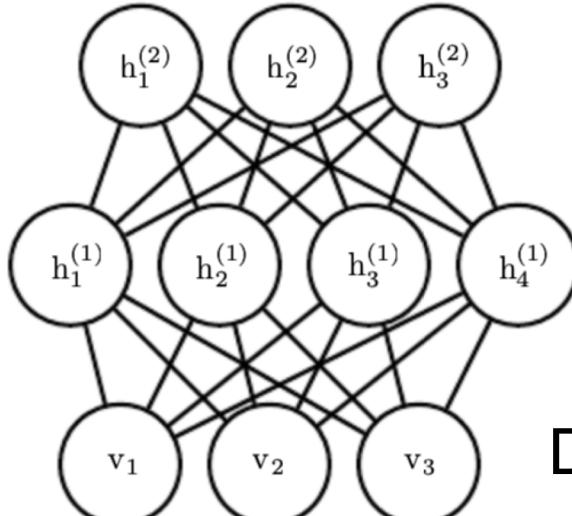


(a)



(b)

Deep belief networks



(c)

Deep Boltzmann machine



Restricted Boltzmann machines

- More formally, let the observed layer consist of a set of n_v binary random variables, which we refer to collectively with the vector \mathbf{v} . We refer to the latent, or hidden, layer of n_h binary random variables as \mathbf{h} .
- The restricted Boltzmann machine is an energy-based model with the **joint probability distribution** specified by its energy function:

$$P(\mathbf{v} = \mathbf{v}, \mathbf{h} = \mathbf{h}) = \frac{1}{Z} \exp(-E(\mathbf{v}, \mathbf{h})). \quad (20.4)$$

- The **energy function** for an RBM is given by

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h}, \quad (20.5)$$

and **Z** is the partition function:

$$Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp \{-E(\mathbf{v}, \mathbf{h})\}. \quad (20.6)$$



Restricted Boltzmann machines

- From the definition of the partition function Z , the naive method of computing Z (exhaustively summing over all states) could be computationally intractable.
- It was formally proved that for RBM the partition function Z is intractable. The intractable partition function Z implies that the normalized joint probability distribution $P(\mathbf{v})$ is also intractable to evaluate.
- Though $P(\mathbf{v})$ is intractable, the bipartite graph structure of the RBM has the special property of its conditional distributions $P(\mathbf{h}|\mathbf{v})$ and $P(\mathbf{v}|\mathbf{h})$ being factorial and relatively simple to compute and sample from.



Conditional distributions

- Deriving the conditional distributions from the joint distribution is straightforward:

$$P(\mathbf{h} \mid \mathbf{v}) = \frac{P(\mathbf{h}, \mathbf{v})}{P(\mathbf{v})} \quad (20.7)$$

$$= \frac{1}{P(\mathbf{v})} \frac{1}{Z} \exp \left\{ \mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \quad (20.8)$$

$$= \frac{1}{Z'} \exp \left\{ \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h} \right\} \quad (20.9)$$

$$= \frac{1}{Z'} \exp \left\{ \sum_{j=1}^{n_h} c_j h_j + \sum_{j=1}^{n_h} \mathbf{v}^\top \mathbf{W}_{:,j} \mathbf{h}_j \right\} \quad (20.10)$$

$$= \frac{1}{Z'} \prod_{j=1}^{n_h} \exp \left\{ c_j h_j + \mathbf{v}^\top \mathbf{W}_{:,j} \mathbf{h}_j \right\}. \quad (20.11)$$



Conditional distributions

- Since we are conditioning on the visible units \mathbf{v} , we can treat these as constant with respect to the distribution $P(\mathbf{h}|\mathbf{v})$. The factorial nature of the conditional $P(\mathbf{h}|\mathbf{v})$ follows from our ability to **write the joint probability over the vector \mathbf{h} as the product of (unnormalized) distributions over the individual elements, h_j** .
- It is now a simple matter of **normalizing the distributions over the individual binary h_j** .

$$P(h_j = 1 \mid \mathbf{v}) = \frac{\tilde{P}(h_j = 1 \mid \mathbf{v})}{\tilde{P}(h_j = 0 \mid \mathbf{v}) + \tilde{P}(h_j = 1 \mid \mathbf{v})} \quad (20.12)$$

$$= \frac{\exp \{c_j + \mathbf{v}^\top \mathbf{W}_{:,j}\}}{\exp \{0\} + \exp \{c_j + \mathbf{v}^\top \mathbf{W}_{:,j}\}} \quad (20.13)$$

$$= \sigma(c_j + \mathbf{v}^\top \mathbf{W}_{:,j}). \quad (20.14)$$



Conditional distributions

- We can now express the full conditional over the hidden layer as the factorial distribution:

$$P(\mathbf{h} \mid \mathbf{v}) = \prod_{j=1}^{n_h} \sigma \left((2\mathbf{h} - 1) \odot (\mathbf{c} + \mathbf{W}^\top \mathbf{v}) \right)_j. \quad (20.15)$$

- A similar derivation will show that the other condition of interest to us, $P(\mathbf{v}|\mathbf{h})$, is also a factorial distribution:

$$P(\mathbf{v} \mid \mathbf{h}) = \prod_{i=1}^{n_v} \sigma \left((2\mathbf{v} - 1) \odot (\mathbf{b} + \mathbf{W}\mathbf{h}) \right)_i. \quad (20.16)$$

- Compared to other undirected models used in deep learning, the RBM is relatively straightforward to train because we can compute $P(\mathbf{h}|\mathbf{v})$ exactly in closed form.



Deep belief networks

- Deep belief networks (DBNs) were one of the first nonconvolutional models to successfully admit training of deep architectures. The introduction of deep belief networks in 2006 began the current deep learning renaissance. Prior to the introduction of deep belief networks, deep models were considered too difficult to optimize. DBNs demonstrated that deep architectures can be successful by outperforming kernelized support vector machines on the MNIST dataset.
- Today, deep belief networks have mostly fallen out of favor and are rarely used, even compared to other unsupervised or generative learning algorithms, but they are still deservedly recognized for their important role in deep learning history.



Deep belief networks

- Deep belief networks are generative models with several layers of latent variables. The latent variables are typically binary, while the visible units may be binary or real. There are **no intralayer connections**. Usually, **every unit in each layer is connected to every unit in each neighboring layer**, though it is possible to construct more sparsely connected DBNs.
- The **connections between the top two layers are undirected**. The **connections between all other layers are directed**, with **the arrows pointed toward the layer that is closest to the data**.



Deep belief networks

- A DBN with l hidden layers contains l weight matrices: $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(l)}$. It also contains $l + 1$ bias vectors $\mathbf{b}^{(0)}, \dots, \mathbf{b}^{(l)}$, with $\mathbf{b}^{(0)}$ providing the biases for the visible layer.

- The probability distribution represented by the DBN is

$$P(\mathbf{h}^{(l)}, \mathbf{h}^{(l-1)}) \propto \exp \left(\mathbf{b}^{(l)\top} \mathbf{h}^{(l)} + \mathbf{b}^{(l-1)\top} \mathbf{h}^{(l-1)} + \mathbf{h}^{(l-1)\top} \mathbf{W}^{(l)} \mathbf{h}^{(l)} \right), \quad (20.17)$$

$$P(h_i^{(k)} = 1 \mid \mathbf{h}^{(k+1)}) = \sigma \left(b_i^{(k)} + \mathbf{W}_{:,i}^{(k+1)\top} \mathbf{h}^{(k+1)} \right) \forall i, \forall k \in 1, \dots, l-2, \quad (20.18)$$

$$P(v_i = 1 \mid \mathbf{h}^{(1)}) = \sigma \left(b_i^{(0)} + \mathbf{W}_{:,i}^{(1)\top} \mathbf{h}^{(1)} \right) \forall i. \quad (20.19)$$

- In the case of real-valued visible units, substitute

$$\mathbf{v} \sim \mathcal{N} \left(\mathbf{v}; \mathbf{b}^{(0)} + \mathbf{W}^{(1)\top} \mathbf{h}^{(1)}, \boldsymbol{\beta}^{-1} \right) \quad (20.20)$$

with $\boldsymbol{\beta}$ diagonal for tractability.



Deep belief networks

- To generate a sample from a DBN, we first run several steps of **Gibbs sampling** on the **top two hidden layers**. This stage is essentially drawing a sample from the RBM defined by the top two hidden layers. We can then use a **single pass of ancestral sampling through the rest of the model** to draw a sample from the visible units.
- Inference in a deep belief network is **intractable**. Evaluating or maximizing the log-likelihood requires confronting not just the problem of intractable inference to marginalize out the latent variables, but also the problem of an intractable partition function within the undirected model of the top two layers. Evaluating or maximizing the standard evidence lower bound on the log-likelihood is also intractable.



Deep belief networks

- To train a deep belief network, one begins by training an RBM to maximize $\mathbb{E}_{\mathbf{v} \sim p_{\text{data}}} \log p(\mathbf{v})$ using contrastive divergence or stochastic maximum likelihood. The parameters of the RBM then define the parameters of the first layer of the DBN.
- A second RBM is trained to approximately maximize

$$\mathbb{E}_{\mathbf{v} \sim p_{\text{data}}} \mathbb{E}_{\mathbf{h}^{(1)} \sim p^{(1)}(\mathbf{h}^{(1)} | \mathbf{v})} \log p^{(2)}(\mathbf{h}^{(1)}), \quad (20.21)$$

where $p^{(1)}$ is the probability distribution represented by the first RBM, and $p^{(2)}$ is the probability distribution represented by the second RBM. In other words, the second RBM is trained to model the distribution defined by sampling the hidden units of the first RBM, when the first RBM is driven by the data.



Deep belief networks

- This procedure can be repeated indefinitely, to add as many layers to the DBN as desired, with each new RBM modeling the samples of the previous one. Each RBM defines another layer of the DBN. This procedure can be justified as increasing a variational lower bound on the log-likelihood of the data under the DBN.
- In most applications, no effort is made to jointly train the DBN after the greedy layer-wise procedure is complete. However, it is possible to perform generative fine-tuning using the wake-sleep algorithm.



Deep belief networks

- The trained DBN may be used directly as a generative model, but most of the interest in DBNs arose from their ability to improve classification models. We can **take the weights from the DBN and use them to define an MLP**:

$$\mathbf{h}^{(1)} = \sigma \left(b^{(1)} + \mathbf{v}^\top \mathbf{W}^{(1)} \right), \quad (20.22)$$

$$\mathbf{h}^{(l)} = \sigma \left(b_i^{(l)} + \mathbf{h}^{(l-1)\top} \mathbf{W}^{(l)} \right) \forall l \in 2, \dots, m. \quad (20.23)$$

- After initializing this MLP with the weights and biases learned via generative training of the DBN, we can train the MLP to perform a classification task. This additional training of the MLP is an example of **discriminative fine-tuning**. This specific choice of MLP is somewhat arbitrary.



Deep Boltzmann machines

- A **deep Boltzmann machine**, or **DBM** is another kind of deep generative model. Unlike the DBN, it is an **entirely undirected model**. Unlike the RBM, the **DBM has several layers of latent variables** (RBMs have just one).
- Like the RBM, **within each layer, each of the variables are mutually independent, conditioned on the variables in the neighboring layers.**
- Like RBMs and DBNs, DBMs typically contain only binary units—as we assume for simplicity of our presentation of the model—but it is straightforward to include real-valued visible units.



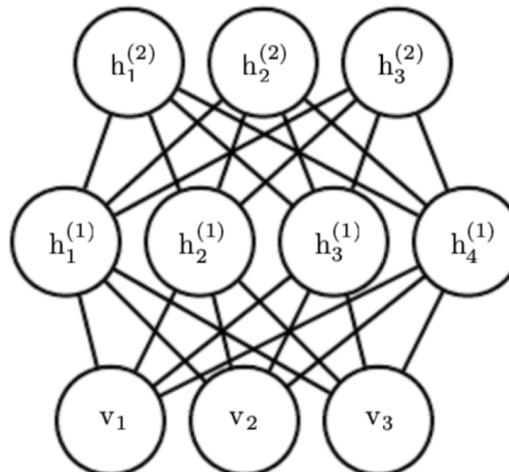
Deep Boltzmann machines

- In the case of a deep Boltzmann machine with one visible layer, v , and three hidden layers, $\mathbf{h}^{(1)}$, $\mathbf{h}^{(2)}$, and $\mathbf{h}^{(3)}$, the joint probability is given by

$$P(v, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}) = \frac{1}{Z(\theta)} \exp(-E(v, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \theta)). \quad (20.24)$$

- To simplify our presentation, we omit the bias parameters below. The DBM energy function is then defined as follows:

$$E(v, \mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \mathbf{h}^{(3)}; \theta) = -v^\top \mathbf{W}^{(1)} \mathbf{h}^{(1)} - \mathbf{h}^{(1)\top} \mathbf{W}^{(2)} \mathbf{h}^{(2)} - \mathbf{h}^{(2)\top} \mathbf{W}^{(3)} \mathbf{h}^{(3)}. \quad (20.25)$$





Deep Boltzmann machines

- In comparison to the RBM energy function, the DBM energy function includes connections between the hidden units (latent variables) in the form of the weight matrices ($W^{(2)}$ and $W^{(3)}$).
- In comparison to fully connected Boltzmann machines, the DBM offers some advantages that are similar to those offered by the RBM. **The DBM layers can be organized into a bipartite graph, with odd layers on one side and even layers on the other.** This immediately implies that **when we condition on the variables in the even (odd) layer, the variables in the odd (even) layers become conditionally independent.**



Deep Boltzmann machines

- The bipartite structure of the DBM means that we can apply the same equations we have previously used for the conditional distributions of an RBM to determine the conditional distributions in a DBM.
- The units within a layer are conditionally independent from each other given the values of the neighboring layers, so the distributions over binary variables can be fully described by the Bernoulli parameters, giving the probability of each unit being active. In our example with two hidden layers, the activation probabilities are

$$P(v_i = 1 \mid \mathbf{h}^{(1)}) = \sigma \left(\mathbf{W}_{i,:}^{(1)} \mathbf{h}^{(1)} \right), \quad (20.26)$$

$$P(h_i^{(1)} = 1 \mid \mathbf{v}, \mathbf{h}^{(2)}) = \sigma \left(\mathbf{v}^\top \mathbf{W}_{:,i}^{(1)} + \mathbf{W}_{i,:}^{(2)} \mathbf{h}^{(2)} \right), \quad (20.27)$$

$$P(h_k^{(2)} = 1 \mid \mathbf{h}^{(1)}) = \sigma \left(\mathbf{h}^{(1)\top} \mathbf{W}_{:,k}^{(2)} \right). \quad (20.28)$$



Deep Boltzmann machines

- The bipartite structure makes Gibbs sampling in a deep Boltzmann machine efficient. In RBMs, it is possible to update all the units in only two iterations. RBMs allow all the visible units to be updated in one block and all the hidden units to be updated in a second block.
- Gibbs sampling can be divided into two blocks of updates, one including all even layers (including the visible layer) and the other including all odd layers. Because of the bipartite DBM connection pattern, given the even (odd) layers, the distribution over the odd (even) layers is factorial and thus can be sampled simultaneously and independently as a block.