

Restricted Boltzmann Machine and Deep Belief Network

Deep Learning (ENGG*6600*01)

School of Engineering,
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghogh
Summer 2023

Introduction

Introduction



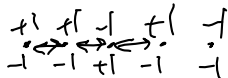
- Centuries ago, the Boltzmann distribution (1868) [1], also called the Gibbs distribution (1902) [2], was proposed.
- This energy-based distribution was found to be useful for modeling the physical systems statistically [3].
- One of these systems was the Ising model which modeled interacting particles with binary spins [4, 5].
- Later, the Ising model was found to be able to be a neural network [6]. Hence, Hopfield network was proposed which modeled an Ising model in a network for modeling memory (1982) [7].
- Inspired by the Hopfield network [6, 7], which was itself inspired by the physical Ising model [4, 5], Hinton et. al. proposed Boltzmann Machine (BM) and Restricted Boltzmann Machine (RBM) (1983-1985) [8, 9].
- These models are energy-based models [10] and the names come from the Boltzmann distribution [1, 2] used in these models.

Introduction

- During the winter of neural networks, Hinton tried to save neural networks from being forgotten in the history of machine learning. So, he returned to his previously proposed RBM and proposed a learning method for RBM with the help of some other researchers including Max Welling (2002-2004) [11, 12]. They proposed training the weights of BM and RBM using maximum likelihood estimation.
- BM and RBM can be seen as generative models where new values for neurons can be generated using Gibbs sampling [13].
- Hinton noticed RBM because he knew that the set of weights between every two layers of a neural network is an RBM.
- It was in the year 2006 [14, 15] that he thought it is possible to train a network in a greedy way [16] where the weights of every layer of network is trained using RBM training.
- This stack of RBM models with a greedy algorithm for training was named Deep Belief Network (DBN) [15, 17].
- DBN allowed the networks to become deep by preparing a good initialization of weights (using RBM training) for backpropagation. This good starting point for backpropagation optimization did not face the problem of vanishing gradients anymore.
- Since the breakthrough in 2006 [14], the winter of neural networks started to end gradually because the networks could get deep to become more nonlinear and handle more nonlinear data.
- Two important techniques were proposed, which were the ReLU activation function (2004) [18] and the dropout technique (2014) [19]. These two regularization methods prevented overfitting [20] and resolved vanishing gradients even without RBM pre-training.

**Statistical Physics,
Ising Model, and
Hopfield Network**

Boltzmann (Gibbs) Distribution



- Assume we have several particles $\{x_i\}_{i=1}^d$ in statistical physics.
- These particles can be seen as random variables which can randomly have a state. For example, if the particles are electrons, they can have states +1 and -1 for counterclockwise and clockwise spins, respectively.
- The Boltzmann distribution (1868) [1], also called the Gibbs distribution (1902) [2], can show the probability that a physical system can have a specific state. i.e., every of the particles has a specific state. The probability mass function of this distribution is [3]:

$$\mathbb{P}(x) = \frac{e^{-\beta E(x)}}{Z}, \quad (1)$$

where $E(x)$ is the energy of variable x and Z is the normalization constant so that the probabilities sum to one.

- This normalization constant is called the partition function which is hard to compute as it sums over all possible configurations of states (values) that the particles can have. If we define $\mathbb{R}^d \ni \mathbf{x} := [x_1, \dots, x_d]^\top$, we have:

$$Z := \sum_{\mathbf{x} \in \mathbb{R}^d} e^{-\beta E(\mathbf{x})}. \quad (2)$$



Boltzmann (Gibbs) Distribution

- We had:

$$\mathbb{P}(x) = \frac{e^{-\beta E(x)}}{Z}. \quad \leftarrow$$

- The coefficient $\beta \geq 0$ is defined as:

$$\boxed{\beta := \frac{1}{k_\beta T}} \propto \left(\frac{1}{T}\right), \quad (3)$$

where k_β is the Boltzmann constant and $T \geq 0$ is the absolute thermodynamic temperature in Kelvins.

- If the temperature tends to absolute zero, $T \rightarrow 0$, we have $\beta \rightarrow \infty$ and $\mathbb{P}(x) \rightarrow 0$, meaning that the absolute zero temperature occurs extremely rarely in the universe.

Boltzmann (Gibbs) Distribution

- Recall Eqs. (1) and (2):

$$\begin{cases} \mathbb{P}(x) = \frac{e^{-\beta E(x)}}{Z}, \\ Z := \sum_{x \in \mathbb{R}^d} e^{-\beta E(x)}. \end{cases}$$

- The free energy is defined as:

$$F(\beta) := \frac{-1}{\beta} \ln(Z), \quad (4)$$

where $\ln(\cdot)$ is the natural logarithm.

- The internal energy is defined as:

$$U(\beta) := \frac{\partial}{\partial \beta} (\beta F(\beta)). \quad (5)$$

- Therefore, we have:

$$U(\beta) = \frac{\partial}{\partial \beta} (-\ln(Z)) = \frac{-1}{Z} \frac{\partial Z}{\partial \beta} \stackrel{(2)}{=} \sum_{x \in \mathbb{R}^d} E(x) \frac{e^{-\beta E(x)}}{Z} \stackrel{(1)}{=} \sum_{x \in \mathbb{R}^d} \mathbb{P}(x) E(x). \quad (6)$$

Boltzmann (Gibbs) Distribution

- Recall Eqs. (1) and (4) and (6):

$$\begin{aligned}\star \mathbb{P}(x) &= \frac{e^{-\beta E(x)}}{Z}, \\ \star F(\beta) &:= \frac{-1}{\beta} \ln(Z), \\ \star U(\beta) &= \sum_{x \in \mathbb{R}^d} \mathbb{P}(x) E(x).\end{aligned}$$

- The **entropy** is defined as:

$$\begin{aligned}H(\beta) &:= - \sum_{x \in \mathbb{R}^d} \mathbb{P}(x) \ln(\mathbb{P}(x)) \stackrel{(1)}{=} - \sum_{x \in \mathbb{R}^d} \mathbb{P}(x) (-\beta E(x) - \ln(Z)) \\ &= \beta \sum_{x \in \mathbb{R}^d} \mathbb{P}(x) E(x) + \ln(Z) \underbrace{\sum_{x \in \mathbb{R}^d} \mathbb{P}(x)}_{=1} \stackrel{(a)}{=} -\beta F(\beta) + \beta U(\beta),\end{aligned}\tag{7}$$

where (a) is because of Eqs. (6) and (4).

Boltzmann (Gibbs) Distribution

Lemma ★

A physical system prefers to be in low energy; hence, the system always loses energy to have less energy.

Proof.

On the one hand, according to the second law of thermodynamics, entropy of a physical system always increases by passing time [21]. Entropy is a measure of randomness and disorder in system. On the other hand, when a system loses energy to its surrounding, it becomes less ordered. Hence, by passing time, the energy of system decreases to have more entropy.
Q.E.D. □

Corollary

According to Eq. (1):

$$\mathbb{P}(x) = \frac{e^{-\beta \langle E(x) \rangle}}{Z},$$

$$E \downarrow \Rightarrow e^{-\beta E} \uparrow \Rightarrow P \uparrow$$

and Lemma 1, the probability $\mathbb{P}(x)$ of states in a system tend to increase by passing time.

Boltzmann (Gibbs) Distribution

- This corollary makes sense because systems tend to become more probable. This idea is also used in simulated annealing¹ [22] where the temperature of system is cooled down gradually.
- Simulated annealing and temperature-based learning have been used in BM models [23, 24, 25].

¹ Simulated annealing is a metaheuristic optimization algorithm in which a temperature parameter controls the amount of global search versus local search. It reduces the temperature gradually to decrease the exploration and increase the exploitation of the search space, gradually.

Ising Model

- Recall the Boltzmann distribution [3]:

$$\mathbb{P}(x) = \frac{e^{-\beta E(x)}}{Z}.$$

- The Ising model [4, 5], also known as the Lenz-Ising model, is a model in which the particles can have -1 or +1 spins [26]. Therefore, $x_i \in \{-1, +1\}, \forall i \in \{1, \dots, d\}$.
- The Ising model uses the Boltzmann distribution, Eq. (1), where the energy function is defined as:

$$E(x) := \mathcal{H}(x) = - \sum_{(i,j)} J_{ij} x_i x_j, \quad (8)$$

where $\mathcal{H}(x)$ is called the Hamiltonian, $J_{ij} \in \mathbb{R}$ is the coupling parameter, and the summation is over particles which interact with each other.

- Note that as energy is proportional to the reciprocal of squared distance, nearby particles are only assumed to be interacting. Therefore, usually the interaction graph of particles is a chain (one dimensional grid), mesh grid (lattice), closed chain (loop), or torus (multi-dimensional loop).

Ising Model

- Based on the characteristic of model, the coupling parameter has different values. If for all interacting i and j , we have $J_{ij} \geq 0$ or $J_{ij} < 0$, the model is named **ferromagnetic** and **anti-ferromagnetic**, respectively. If J_{ij} can be both positive and negative, the model is called a **spin glass**. If the coupling parameters are all constant, the model is **homogeneous**.
- The BM and RBM are Ising models whose coupling parameters are considered as weights and these weights are learned using maximum likelihood estimation [27]. Hence, we can say that BM and RBM are energy-based learning methods [10].

Hopfield Network

- It was proposed in (1974) [6] to use the Ising model in a neural network structure. Hopfield extended this idea to model the memory by a neural network. The resulting network was the Hopfield network (1982) [7].
- This network has some units or neurons denoted by $\{x_i\}_{i=1}^d$. The states or outputs of units are all binary $x_i \in \{-1, +1\}, \forall i$. Let w_{ij} denote the weight of link connecting unit i to unit j .
- The weights of Hopfield network are learned using the Hebbian learning (Hebb's law of association) [28]:

$$w_{ij} := \begin{cases} x_i \times x_j & \text{if } i \neq j, \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

- After training, the outputs of units can be determined for an input if the weighted summation of inputs to unit passes a threshold θ :

$$x_i := \begin{cases} +1 & \text{if } \sum_{j=1}^d w_{ij} x_j \geq \theta, \\ -1 & \text{otherwise.} \end{cases} \quad (10)$$

- In the original paper of Hopfield network [7], the binary states are $x_i \in \{0, 1\}, \forall i$ so the Hebbian learning is $w_{ij} := (2x_i - 1) \times (2x_j - 1), \forall i \neq j$.

Hopfield Network

- Hopfield network is an Ising model so it uses Eq. (8):

$$\underbrace{E(x) := \mathcal{H}(x)} = - \sum_{(i,j)} J_{ij} x_i x_j,$$

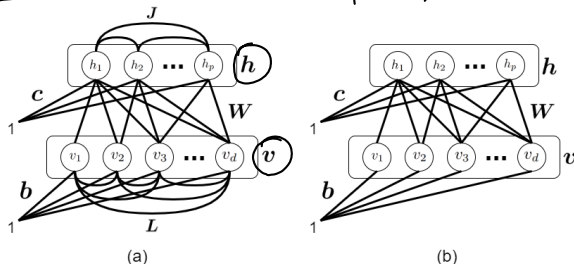
as its energy. This energy is also used in the Boltzmann distribution which is Eq. (1).

- It is noteworthy that there are also Hopfield networks with continuous states (1984) [29]. Modern Hopfield networks, such as (2020) [30], are often based on dense associative memories [31]. Some other recent works on associative memories are [32, 33].
- The BM and RBM models are Hopfield networks whose weights are learned using maximum likelihood estimation and not Hebbian learning.

Structure of Restricted Boltzmann Machine

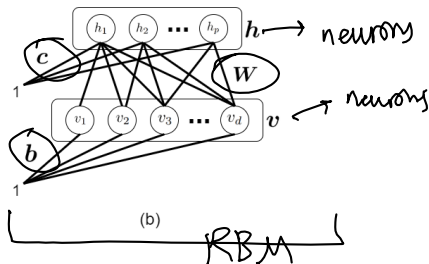
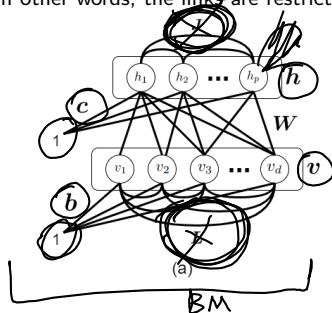
Structure of Restricted Boltzmann Machine

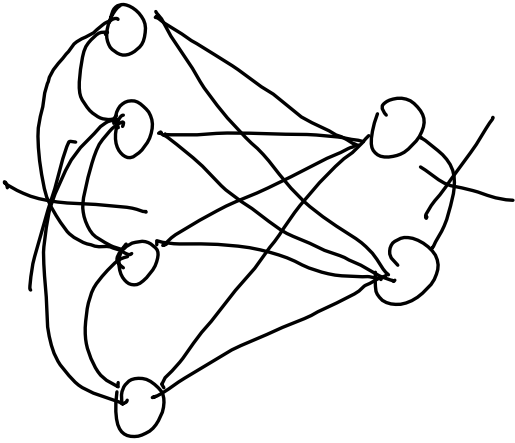
- **Boltzmann Machine (BM)** is a generative model and a Probabilistic Graphical Model (PGM) [34] which is a building block of many probabilistic models.
- Its name is because of the Boltzmann distribution [1, 2] used in this model.
- It was first introduced to be used in machine learning in (1983-1985) [8, 9] and then in (2002-2004) [11, 12].
- A BM consists of a visible (or observation) layer $\mathbf{v} = [v_1, \dots, v_d] \in \mathbb{R}^d$ and a hidden layer $\mathbf{h} = [h_1, \dots, h_p] \in \mathbb{R}^p$.
- The visible layer is the layer that we can see; for example, it can be the layer of data. The hidden layer is the layer of latent variables which represent meaningful features or embeddings for the visible data.
- In other words, there is a meaningful connection between the hidden and visible layers although their dimensionality might differ, i.e., $d \neq p$.



Structure of Restricted Boltzmann Machine

- Each of the elements of \mathbf{v} and \mathbf{h} also have a bias. There are also links between the elements of \mathbf{v} as well as between the elements of \mathbf{h} [35].
- Let w_{ij} denote the link between v_i and h_j , and l_{ij} be the link between v_i and v_j , and j_{ij} be the link between h_i and h_j , and b_i be the bias link for v_i , and c_i be the bias link for h_i .
- The dimensionality of these links are $\mathbf{W} = [w_{ij}] \in \mathbb{R}^{d \times p}$, $\mathbf{L} = [l_{ij}] \in \mathbb{R}^{d \times d}$, $\mathbf{J} = [j_{ij}] \in \mathbb{R}^{p \times p}$, $\mathbf{b} = [b_1, \dots, b_d] \in \mathbb{R}^d$, and $\mathbf{c} = [c_1, \dots, c_p] \in \mathbb{R}^p$. Note that \mathbf{W} is a symmetric matrix, i.e., $w_{ij} = w_{ji}$. Also, as there is no link from a node to itself, the diagonal elements of \mathbf{L} and \mathbf{J} are zero, i.e., $l_{ii} = j_{ii} = 0, \forall i$.
- Restricted Boltzmann Machine (RBM)** is BM which does not have links within a layer, i.e., there is no any link between the elements of \mathbf{v} and no any link between the elements of \mathbf{h} . In other words, the links are restricted in RBM to be $\mathbf{L} = \mathbf{J} = \mathbf{0}$.





Structure of Restricted Boltzmann Machine

- Recall that RBM is an Ising model. As we saw in Eq. (8):

$$E(x) := \mathcal{H}(x) = - \sum_{(i,j)} \underbrace{J_{ij} x_i x_j}_{\text{interaction}}$$

the energy of an Ising model can be modeled as [8, 9]:

$$\mathbb{R} \ni E(\mathbf{v}, \mathbf{h}) := \underbrace{-\mathbf{b}^\top \mathbf{v}}_{\text{bias}} - \underbrace{\mathbf{c}^\top \mathbf{h}}_{\text{bias}} - \underbrace{\mathbf{v}^\top \mathbf{W} \mathbf{h}}_{\text{interaction}} \quad (11)$$

which is based on interactions between linked units.

- As introduced in Eq. (1):

$$\mathbb{P}(x) = \frac{e^{-\beta E(x)}}{Z},$$

the visible and hidden variables make a joint Boltzmann distribution [36]:

$$\mathbb{P}(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} \exp(-\beta E(\mathbf{v}, \mathbf{h})) \stackrel{(11)}{=} \frac{1}{Z} \exp(\mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h}), \quad (12)$$

where Z is the partition function:

$$Z := \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h})). \quad (13)$$

- According to Lemma 1, the BM and RBM try to reduce the energy of model. Training the BM or RBM reduces its energy [8, 9].

Conditional Distributions

Conditional Distributions

$$h_i \perp h_j \mid v$$
$$v_i \perp v_j \mid h$$

Lemma (Conditional Independence of Variables)

In RBM, given the visible variables, the hidden variables are conditionally independent. Likewise, given the hidden variables, the visible variables are conditionally independent. This does not hold in BM because of the links within each layer.

Proof:

- According to the Bayes' rule, we have:

$$\begin{aligned} \underbrace{\mathbb{P}(h|v)} &= \frac{\mathbb{P}(h, v)}{\mathbb{P}(v)} = \frac{\mathbb{P}(v, h)}{\sum_{h \in \mathbb{R}^p} \mathbb{P}(v, h)} \stackrel{(12)}{=} \frac{\frac{1}{Z} \exp(\mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h})}{\sum_{h \in \mathbb{R}^p} \frac{1}{Z} \exp(\mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h})} \\ &= \frac{\frac{1}{Z} \exp(\mathbf{b}^\top \mathbf{v}) \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})}{\frac{1}{Z} \sum_{h \in \mathbb{R}^p} \exp(\mathbf{b}^\top \mathbf{v}) \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})} \\ &\stackrel{(a)}{=} \frac{\exp(\mathbf{b}^\top \mathbf{v}) \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})}{\cancel{\exp(\mathbf{b}^\top \mathbf{v})} \sum_{h \in \mathbb{R}^p} \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})} = \frac{\exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})}{\sum_{h \in \mathbb{R}^p} \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})}, \end{aligned}$$

where (a) is because the term $\exp(\mathbf{b}^\top \mathbf{v})$ does not have \mathbf{h} in it.

- Note that $\sum_{h \in \mathbb{R}^p}$ denotes summation over all possible p -dimensional hidden variables for the sake of marginalization.

Conditional Distributions

- We had:

$$\star \mathbb{P}(\mathbf{h}|\mathbf{v}) = \frac{\exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})}{\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})}$$

- Let $Z' = \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(\mathbf{c}^\top \mathbf{h}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})$. Hence:

$$\begin{aligned} \star \mathbb{P}(\mathbf{h}|\mathbf{v}) &= \frac{1}{Z'} \exp(\mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h}) = \frac{1}{Z'} \exp\left(\sum_{j=1}^p c_j h_j + \sum_{j=1}^p \mathbf{v}^\top \mathbf{W}_{:j} h_j\right) \\ &= \frac{1}{Z'} \prod_{j=1}^p \exp(c_j h_j + \mathbf{v}^\top \mathbf{W}_{:j} h_j), \quad \leftarrow e^{a+b} = e^a e^b \end{aligned} \quad (14)$$

where $\mathbf{W}_{:j} \in \mathbb{R}^d$ denotes the j -th column of matrix \mathbf{W} .

- The Eq. (14) shows that given the visible variables, the hidden variables are conditionally independent because their joint distribution is the product of every distribution.

Conditional Distributions

- We can write similar expressions for the probability $\mathbb{P}(\mathbf{v}|\mathbf{h})$:

$$\begin{aligned}\cancel{\star} \quad \mathbb{P}(\mathbf{v}|\mathbf{h}) &= \frac{1}{Z''} \exp(\mathbf{b}^\top \mathbf{v} + \mathbf{v}^\top \mathbf{W} \mathbf{h}) = \frac{1}{Z''} \exp\left(\sum_{i=1}^d b_i v_i + \sum_{i=1}^d v_i \mathbf{W}_{i:} \mathbf{h}\right) \\ &= \frac{1}{Z''} \prod_{i=1}^d \exp(b_i v_i + v_i \mathbf{W}_{i:} \mathbf{h}),\end{aligned}\tag{15}$$

where $\mathbf{W}_{i:} \in \mathbb{R}^P$ denotes the i -th row of matrix \mathbf{W} and $Z'' := \sum_{\mathbf{v} \in \mathbb{R}^d} \exp(\mathbf{b}^\top \mathbf{v}) \exp(\mathbf{v}^\top \mathbf{W} \mathbf{h})$.

- This equation shows that given the hidden variables, the visible variables are conditionally independent. Q.E.D.

Conditional Distributions

- According to Eq. (14):

$$\star \mathbb{P}(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{j=1}^p \exp(c_j h_j + \mathbf{v}^\top \mathbf{W}_{:j} h_j) \star$$

$\frac{1}{P(\mathbf{v})}$ (pointing to $\frac{1}{Z'}$)
 $P(\mathbf{h}, \mathbf{v}) = \prod_{j=1}^p P(h_j, \mathbf{v})$ (pointing to the product)
 $P(h_j, \mathbf{v})$ (circled)

and considering the rule $\mathbb{P}(\mathbf{h}|\mathbf{v}) \propto \mathbb{P}(\mathbf{h}, \mathbf{v}) P(\mathbf{v})$, we have:

$$\mathbb{P}(\mathbf{h}|\mathbf{v}) = \frac{1}{Z'} \prod_{j=1}^p \exp(c_j h_j + \mathbf{v}^\top \mathbf{W}_{:j} h_j) = \frac{1}{Z'} \prod_{j=1}^p \mathbb{P}(h_j, \mathbf{v})$$

jth hidden variable (pointing to the product)

$$\Rightarrow \mathbb{P}(h_j, \mathbf{v}) = \exp(c_j h_j + \mathbf{v}^\top \mathbf{W}_{:j} h_j) = \exp(c_j h_j + \sum_{i=1}^d v_i w_{ij} h_j). \quad (16)$$

v_1, \dots, v_d (pointing to \mathbf{v})

Conditional Distributions

- Similarly, according to Eq. (15):

$$\star \quad \mathbb{P}(\mathbf{v}|\mathbf{h}) = \frac{1}{Z''} \prod_{i=1}^d \exp(b_i v_i + v_i \mathbf{W}_{i:} \mathbf{h}), \quad (17)$$

and considering the rule $\underbrace{\mathbb{P}(\mathbf{v}|\mathbf{h}) = \mathbb{P}(\mathbf{h}, \mathbf{v})/\mathbb{P}(\mathbf{h})}$, we have:

$$\mathbb{P}(\mathbf{h}|\mathbf{v}) = \frac{1}{Z''} \prod_{i=1}^d \exp(b_i v_i + v_i \mathbf{W}_{i:} \mathbf{h}) = \frac{1}{Z''} \prod_{i=1}^d \mathbb{P}(\mathbf{h}, v_i)$$

$$\Rightarrow \mathbb{P}(\mathbf{h}, v_i) = \exp(b_i v_i + v_i \mathbf{W}_{i:} \mathbf{h}) = \exp(b_i v_i + \sum_{j=1}^p v_i w_{ij} h_j). \quad (18)$$

- We will use these equations later.

h_1, \dots, h_p

Sampling Hidden and Visible Variables

Gibbs Sampling

- We can use Gibbs sampling for sampling and generating the hidden and visible units.
- If ν denotes the iteration index of Gibbs sampling, we iteratively sample:

$$h^{(\nu)} \sim \mathbb{P}(h | v^{(\nu)}), \quad (19)$$

$$v^{(\nu+1)} \sim \mathbb{P}(v | h^{(\nu)}), \quad (20)$$

until the burn-in convergence.

- In Gibbs sampling, only several iterations of Gibbs sampling are usually sufficient.
- After the burn-in, the samples are approximate samples from the joint distribution $\mathbb{P}(v, h)$.

$$\mathbb{P}(v, h)$$

Gibbs Sampling

- As the variables are conditionally independent, this Gibbs sampling can be implemented as in this algorithm:

```
1 Input: visible dataset  $v$ , (initialization: optional)
2 Get initialization or do random initialization of  $v$ 
3 while until burn-in do
4   for  $j$  from 1 to  $p$  do
5      $h_j^{(\nu)} \sim \mathbb{P}(h_j | v^{(\nu)})$ 
6   for  $i$  from 1 to  $d$  do
7      $v_i^{(\nu+1)} \sim \mathbb{P}(v_i | h^{(\nu)})$ 
```

Handwritten annotations: $\rightarrow h \sim P(h|v)$ next to line 5, and $\rightarrow v \sim P(v|h)$ next to line 7.

Algorithm 1: Gibbs sampling in RBM

- In this algorithm, $h_j^{(\nu)} \sim \mathbb{P}(h_j | v^{(\nu)})$ can be implemented as drawing a sample from uniform distribution $u \sim U[0, 1]$ and comparing it to the value of Probability Density Function (PDF), $\mathbb{P}(h_j | v^{(\nu)})$. If u is less than or equal to this value, we have $h_j = 1$; otherwise, we have $h_j = 0$.
- Implementation of sampling v_i has a similar procedure.
- Alternatively, we can use inverse of cumulative distribution function of these distributions for drawing samples (see [37] for more details about sampling).

Generations and Evaluations by Gibbs Sampling

- Gibbs sampling for generating both observation and hidden units is used for both training and evaluation phases of RBM.
- Use of Gibbs sampling in training RBM will be explained later.
- After the RBM model is trained, we can generate any number of p -dimensional hidden variables as a meaningful representation of the d -dimensional observation using Gibbs sampling.
- Moreover, using Gibbs sampling, we can generate other d -dimensional observations in addition to the original dataset. These new generated observations are d -dimensional representations for the p -dimensional hidden variables.
- This shows that BM and RBM are generative models.

**Training Restricted
Boltzmann Machine by
Maximum Likelihood
Estimation**

Training Restricted Boltzmann Machine by Maximum Likelihood Estimation

- The weights of links which are \mathbf{W} , \mathbf{b} , and \mathbf{c} should be learned so that we can use them for sampling/generating the hidden and visible units. Consider a dataset of n visible vectors $\{\mathbf{v}_i \in \mathbb{R}^d\}_{i=1}^n$.
- Note that \mathbf{v}_i should not be confused with v_i where the former is the i -th visible data instance and the latter is the i -th visible unit. We denote the j -th dimension of \mathbf{v}_i by $v_{i,j}$; in other words, $\mathbf{v}_i = [v_{i,1}, \dots, v_{i,d}]^\top$.
- The log-likelihood of the visible data is:

$$\sum_y P(x, y) = P(x)$$

$$\log\left(\frac{a}{b}\right) = \log a - \log b$$

$$\begin{aligned} \ell(\mathbf{W}, \mathbf{b}, \mathbf{c}) &= \sum_{i=1}^n \log(\mathbb{P}(\mathbf{v}_i)) = \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \mathbb{P}(\mathbf{v}_i, \mathbf{h}) \right) \\ &\stackrel{(12)}{=} \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \left(\frac{1}{Z} \right) \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) = \sum_{i=1}^n \log \left(\left(\frac{1}{Z} \right) \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) \\ &= \sum_{i=1}^n \left[\log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) - \log Z \right] = \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) - n \log Z \\ &\stackrel{(13)}{=} \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) - n \log \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h})). \end{aligned} \quad (21)$$

Training Restricted Boltzmann Machine by Maximum Likelihood Estimation

- We found:

$$\cancel{\star} \ell(\mathbf{W}, \mathbf{b}, \mathbf{c}) = \overbrace{\sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right)} - \overbrace{\log \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h}))}.$$

- We use Maximum Likelihood Estimation (MLE) for finding the parameters $\theta := \{\mathbf{W}, \mathbf{b}, \mathbf{c}\}$. The derivative of log-likelihood with respect to parameter θ is:

$$\nabla_{\theta} \ell(\theta) = \overbrace{\cancel{\nabla_{\theta}} \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right)} - \overbrace{\cancel{\nabla_{\theta}} \log \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h}))}. \quad (22)$$

Training Restricted Boltzmann Machine by Maximum Likelihood Estimation

- We had:

$$\nabla_{\theta} \ell(\theta) = \nabla_{\theta} \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) - n \nabla_{\theta} \log \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h})).$$

- The first term of this derivative is:

$$\begin{aligned} \nabla_{\theta} \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) &= \sum_{i=1}^n \nabla_{\theta} \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) \\ &= \sum_{i=1}^n \frac{\nabla_{\theta} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h}))}{\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h}))} = \sum_{i=1}^n \frac{\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \nabla_{\theta}(-E(\mathbf{v}_i, \mathbf{h}))}{\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h}))} \\ &\stackrel{(a)}{=} \sum_{i=1}^n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\nabla_{\theta}(-E(\mathbf{v}_i, \mathbf{h}))], \end{aligned} \quad (23)$$

where (a) is because the definition of expectation is $\mathbb{E}_{\sim \mathbb{P}}[\mathbf{x}] := \sum_{i=1} \mathbb{P}(\mathbf{x}_i) \mathbf{x}_i$. However, if \mathbb{P} is not an actual distribution and does not sum to one, we should normalize it to behave like a distribution in the expectation: $\mathbb{E}_{\sim \mathbb{P}}[\mathbf{x}] := (\sum_{i=1} \mathbb{P}(\mathbf{x}_i) \mathbf{x}_i) / (\sum_{i=1} \mathbb{P}(\mathbf{x}_i))$.

Training Restricted Boltzmann Machine by Maximum Likelihood Estimation

- We had:

$$\nabla_{\theta} \ell(\theta) = \nabla_{\theta} \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) - \overbrace{n \nabla_{\theta} \log \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h}))}.$$

- The second term of the derivative of log-likelihood is:

$$\begin{aligned} -n \nabla_{\theta} \left(\log \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h})) \right) &= -n \frac{\nabla_{\theta} \sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h}))}{\sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h}))} \\ &= -n \frac{\sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \nabla_{\theta} \exp(-E(\mathbf{v}, \mathbf{h}))}{\sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h}))} \\ &= -n \frac{\sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h})) \nabla_{\theta} (-E(\mathbf{v}, \mathbf{h}))}{\sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h}))} \stackrel{(a)}{=} -n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\nabla_{\theta} (-E(\mathbf{v}, \mathbf{h}))], \quad (24) \end{aligned}$$

where (a) is for the definition of expectation which was already explained above.

- In summary, the derivative of log-likelihood is:

$$\longrightarrow \nabla_{\theta} \ell(\theta) = \sum_{i=1}^n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h} | \mathbf{v}_i)} [\nabla_{\theta} (-E(\mathbf{v}_i, \mathbf{h}))] - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\nabla_{\theta} (-E(\mathbf{v}, \mathbf{h}))]. \quad (25)$$

- Setting this derivative to zero does not give us a closed-form solution. Hence, we should learn the parameters iteratively using gradient ascent for MLE.

Training Restricted Boltzmann Machine by Maximum Likelihood Estimation

- Recall Eqs. (11) and (25):

$$\mathbb{R} \ni E(\mathbf{v}, \mathbf{h}) := -\mathbf{b}^\top \mathbf{v} - \mathbf{c}^\top \mathbf{h} - \mathbf{v}^\top \mathbf{W} \mathbf{h},$$

$$\nabla_{\theta} \ell(\theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{h} \sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\nabla_{\theta} (-E(\mathbf{v}_i, \mathbf{h}))] - n \mathbb{E}_{\mathbf{h} \sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\nabla_{\theta} (-E(\mathbf{v}, \mathbf{h}))].$$

- Now, consider each of the parameters $\theta = \{\mathbf{W}, \mathbf{b}, \mathbf{c}\}$. The derivative w.r.t. these parameters in Eq. (25) are:

$$\begin{aligned} \nabla_{\mathbf{W}} (-E(\mathbf{v}, \mathbf{h})) &\stackrel{(11)}{=} \frac{\partial}{\partial \mathbf{W}} (\mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h}) = \mathbf{v} \mathbf{h}^\top \\ \nabla_{\mathbf{b}} (-E(\mathbf{v}, \mathbf{h})) &\stackrel{(11)}{=} \frac{\partial}{\partial \mathbf{b}} (\mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h}) = \mathbf{v} \\ \nabla_{\mathbf{c}} (-E(\mathbf{v}, \mathbf{h})) &\stackrel{(11)}{=} \frac{\partial}{\partial \mathbf{c}} (\mathbf{b}^\top \mathbf{v} + \mathbf{c}^\top \mathbf{h} + \mathbf{v}^\top \mathbf{W} \mathbf{h}) = \mathbf{h} \end{aligned}$$

Training Restricted Boltzmann Machine by Maximum Likelihood Estimation

- Therefore, Eq. (25) for these parameters becomes:

$$\begin{aligned}\nabla_{\mathbf{W}} \ell(\theta) &= \sum_{i=1}^n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\mathbf{v} \mathbf{h}_i^\top] - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{v} \mathbf{h}^\top] \\ &= \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\mathbf{h}^\top] - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{v} \mathbf{h}^\top],\end{aligned}$$

$$\nabla_{\mathbf{b}} \ell(\theta) = \sum_{i=1}^n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\mathbf{v}_i] - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{v}] = \sum_{i=1}^n \mathbf{v}_i - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{v}],$$

$$\nabla_{\mathbf{c}} \ell(\theta) = \sum_{i=1}^n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\mathbf{h}] - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{h}].$$

Training Restricted Boltzmann Machine by Maximum Likelihood Estimation

- If we define:

$$\hat{\mathbf{h}}_i := \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)}[\mathbf{h}], \quad (26)$$

we can summarize these derivatives as:

$$\left\{ \begin{array}{l} \mathbb{R}^{d \times p} \ni \nabla_{\mathbf{W}} \ell(\theta) = \sum_{i=1}^n \mathbf{v}_i \hat{\mathbf{h}}_i^{\top} - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})}[\mathbf{v} \mathbf{h}^{\top}], \end{array} \right. \quad (27)$$

$$\left\{ \begin{array}{l} \mathbb{R}^d \ni \nabla_{\mathbf{b}} \ell(\theta) = \sum_{i=1}^n \mathbf{v}_i - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})}[\mathbf{v}], \end{array} \right. \quad (28)$$

$$\left\{ \begin{array}{l} \mathbb{R}^p \ni \nabla_{\mathbf{c}} \ell(\theta) = \sum_{i=1}^n \hat{\mathbf{h}}_i - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})}[\mathbf{h}]. \end{array} \right. \quad (29)$$

- Setting these derivatives to zero does not give a closed form solution. Hence, we need to find the solution iteratively using gradient ascent where the above gradients are used.
- In the derivatives of log-likelihood, we have two types of expectation. The conditional expectation $\mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)}[\cdot]$ is based on the observation or data which is \mathbf{v}_i . The joint expectation $\mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})}[\cdot]$, however, has nothing to do with the observation and is merely about the RBM model.

Contrastive Divergence

Contrastive Divergence

- According to Eq. (23):

$$\star \quad \nabla_{\theta} \sum_{i=1}^n \log \left(\sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}_i, \mathbf{h})) \right) = \sum_{i=1}^n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\nabla_{\theta} (-E(\mathbf{v}_i, \mathbf{h}))],$$

the conditional expectation used in Eq. (26), $\hat{\mathbf{h}}_i := \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)}[\mathbf{h}]$, includes one summation.

- Moreover, according to Eq. (24):

$$-n \nabla_{\theta} \log \left(\sum_{\mathbf{v} \in \mathbb{R}^d} \sum_{\mathbf{h} \in \mathbb{R}^p} \exp(-E(\mathbf{v}, \mathbf{h})) \right) = -n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\nabla_{\theta} (-E(\mathbf{v}, \mathbf{h}))],$$

the joint expectations used in Eqs. (27), (28), and (29):

$$\left\{ \begin{array}{l} \mathbb{R}^{d \times p} \ni \nabla_{\mathbf{w}} \ell(\theta) = \sum_{i=1}^n \mathbf{v}_i \hat{\mathbf{h}}_i^{\top} - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{v} \mathbf{h}^{\top}], \\ \mathbb{R}^d \ni \nabla_{\mathbf{b}} \ell(\theta) = \sum_{i=1}^n \mathbf{v}_i - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{v}], \\ \mathbb{R}^p \ni \nabla_{\mathbf{c}} \ell(\theta) = \sum_{i=1}^n \hat{\mathbf{h}}_i - n \mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\mathbf{h}]. \end{array} \right.$$

contain two summations.

Contrastive Divergence

- This double-summation makes computation of the joint expectation intractable because it sums over all possible values for both hidden and visible units. Therefore, exact computation of MLE is hard and we should approximate it. One way to approximate computation of joint expectations in MLE is contrastive divergence (2002) [11].
- Contrastive divergence improves the efficiency and reduces the variance of estimation in RBM [11, 12].
- The idea of contrastive divergence is as follows. First, we obtain a point $\tilde{\mathbf{v}}$ using Gibbs sampling starting from the observation \mathbf{v}_i . Then, we compute expectation by using only that one point $\tilde{\mathbf{v}}$.
- The intuitive reason for why contrastive divergence works is explained in the following. We need to minimize the gradients to find the solution of MLE. In the joint expectations in Eqs. (27), (28), and (29), rather than considering all possible values of observations, contrastive divergence considers only one of the data points (observations). If this observation is a wrong belief which we do not wish to see in generation of observations by RBM, contrastive divergence is performing a task which is called negative sampling [36]. In negative sampling, we say rather than training the model to not generate all wrong observations, we train it iteratively but less ambitiously in every iteration. Each iteration tries to teach the model to not generate only one of the wrong outputs. Gradually, the model learns to generate correct observations by avoiding to generate these negative samples.

Contrastive Divergence

- Let $\tilde{\mathbf{h}} = [\tilde{h}_1, \dots, \tilde{h}_m]^\top$ be the corresponding sampled \mathbf{h} to $\tilde{\mathbf{v}} = [\tilde{v}_1, \dots, \tilde{v}_m]^\top$ in Gibbs sampling.
- According to the above explanations, contrastive divergence approximates the joint expectation in the derivative of log-likelihood, Eq. (25), by Monte-Carlo approximation [37] evaluated at $\tilde{\mathbf{v}}_i$ and $\tilde{\mathbf{h}}_i$ for the i -th observation and hidden units where $\tilde{\mathbf{v}}_i$ and $\tilde{\mathbf{h}}_i$ are found by Gibbs sampling. Hence:

$$\mathbb{E}_{\sim \mathbb{P}(\mathbf{h}, \mathbf{v})} [\nabla_{\theta} (-E(\mathbf{v}, \mathbf{h}))] \approx \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} (-E(\mathbf{v}_i, \mathbf{h}_i)) \Big|_{\mathbf{v}_i = \tilde{\mathbf{v}}_i, \mathbf{h}_i = \tilde{\mathbf{h}}_i} \quad (30)$$

- Experiments have shown that a small number of iterations in Gibbs sampling suffice for contrastive divergence. Paper [11] even uses one iteration of Gibbs sampling for this task.

Contrastive Divergence

- By the approximation in Eq. (30), the Eqs. (27), (28), and (29) become:

$$\nabla_{\mathbf{w}} \ell(\theta) = \sum_{i=1}^n \underbrace{(\mathbf{v}_i \hat{\mathbf{h}}_i^\top)}_{\text{data}} - \sum_{i=1}^n \underbrace{(\tilde{\mathbf{v}}_i \tilde{\mathbf{h}}_i^\top)}_{\text{recon.}}, \quad \approx 0 \quad (31)$$

$$\nabla_{\mathbf{b}} \ell(\theta) = \sum_{i=1}^n \underbrace{\mathbf{v}_i}_{\text{data}} - \sum_{i=1}^n \underbrace{\tilde{\mathbf{v}}_i}_{\text{recon.}} = 0 \quad (32)$$

$$\nabla_{\mathbf{c}} \ell(\theta) = \sum_{i=1}^n \underbrace{\mathbf{h}_i}_{\text{data}} - \sum_{i=1}^n \underbrace{\tilde{\mathbf{h}}_i}_{\text{recon.}} \approx 0 \quad (33)$$

- These equations make sense because **when the observation variable and hidden variable given the observation variable become equal to the approximations by Gibbs sampling, the gradient should be zero and the training should stop.**
- Note that some works in the literature restate Eqs. (31), (32), and (33) as [11, 36, 38]:

$$\forall i, j: \nabla_{w_{ij}} \ell(\theta) = \langle \mathbf{v}_i \mathbf{h}_j \rangle_{\text{data}} - \langle \mathbf{v}_i \mathbf{h}_j \rangle_{\text{recon.}}, \quad (34)$$

$$\forall i: \nabla_{b_i} \ell(\theta) = \langle \mathbf{v}_i \rangle_{\text{data}} - \langle \mathbf{v}_i \rangle_{\text{recon.}}, \quad (35)$$

$$\forall j: \nabla_{c_j} \ell(\theta) = \langle \mathbf{h}_j \rangle_{\text{data}} - \langle \mathbf{h}_j \rangle_{\text{recon.}}, \quad (36)$$

where $\langle \cdot \rangle_{\text{data}}$ and $\langle \cdot \rangle_{\text{recon.}}$ denote expectation over data and reconstruction of data, respectively.

Contrastive Divergence

```

1 Input: training data  $\{x_i\}_{i=1}^n$ 
2 Randomly initialize  $W, b, c$ 
3 while not converged do
4   Sample a mini-batch  $\{v_1, \dots, v_m\}$  from
     training dataset  $\{x_i\}_{i=1}^n$  (n.b. we may set
      $m = n$ )
5   // Gibbs sampling for each data point:
6   Initialize  $\hat{v}_i^{(0)} \leftarrow v_i$  for all  $i \in \{1, \dots, m\}$ 
7   for  $i$  from 1 to  $m$  do
8     Algorithm 1  $\leftarrow \hat{v}_i^{(0)}$ 
9      $\{h_i\}_{i=1}^p, \{v_i\}_{i=1}^d \leftarrow$  Last iteration of
       Algorithm 1
10     $\tilde{h}_i \leftarrow [h_1, \dots, h_p]^\top$ 
11     $\tilde{v}_i \leftarrow [v_1, \dots, v_d]^\top$ 
12     $\tilde{h}_i \leftarrow \mathbb{E}_{\sim \mathbb{P}(h|v_i)}[h]$ 
13    // gradients:
14     $\nabla_W \ell(\theta) \leftarrow \sum_{i=1}^m v_i \hat{h}_i^\top - \sum_{i=1}^m \tilde{h}_i \tilde{v}_i^\top$ 
15     $\nabla_b \ell(\theta) \leftarrow \sum_{i=1}^m v_i - \sum_{i=1}^m \tilde{v}_i$ 
16     $\nabla_c \ell(\theta) \leftarrow \sum_{i=1}^m \hat{h}_i - \sum_{i=1}^m \tilde{h}_i$ 
17    // gradient descent for updating solution:
18     $W \leftarrow W + \eta \nabla_W \ell(\theta)$ 
19     $b \leftarrow b + \eta \nabla_b \ell(\theta)$ 
20     $c \leftarrow c + \eta \nabla_c \ell(\theta)$ 
21 Return  $W, b, c$ 

```

Algorithm 2: Training RBM using contrastive divergence

Boltzmann Machine

Boltzmann Machine

- So far, we introduced and explained RBM. BM has more links compared to RBM [35].
- Here, we briefly introduce training of BM. As was explained before, BM has additional links $\mathbf{L} = [l_{ij}] \in \mathbb{R}^{d \times d}$ and $\mathbf{J} = [j_{ij}] \in \mathbb{R}^{p \times p}$. The weights $\mathbf{W} \in \mathbb{R}^{d \times p}$ and biases $\mathbf{b} \in \mathbb{R}^d$ and $\mathbf{c} \in \mathbb{R}^p$ are trained by gradient descent using the gradients in Eqs. (31), (32), and (33). The **additional weights** \mathbf{L} and \mathbf{J} are updated similarly using the following gradients [35]:

$$\left\{ \begin{array}{l} \nabla_{\mathbf{L}} \ell(\theta) = \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^\top - \sum_{i=1}^n \tilde{\mathbf{v}}_i \tilde{\mathbf{v}}_i^\top, \end{array} \right. \quad (37)$$

$$\left\{ \begin{array}{l} \nabla_{\mathbf{J}} \ell(\theta) = \sum_{i=1}^n \mathbb{E}_{\mathbf{h} \sim \mathbb{P}(\mathbf{h}|\mathbf{v}_i)} [\mathbf{h} \mathbf{h}^\top] - \sum_{i=1}^n \tilde{\mathbf{h}}_i \tilde{\mathbf{h}}_i^\top. \end{array} \right. \quad (38)$$

- These equations can be restated as:

$$\forall i, j: \nabla_{l_{ij}} \ell(\theta) = \langle \mathbf{v}_i \mathbf{v}_j \rangle_{\text{data}} - \langle \mathbf{v}_i \mathbf{v}_j \rangle_{\text{recon.}}, \quad (39)$$

$$\forall i, j: \nabla_{j_{ij}} \ell(\theta) = \langle \mathbf{h}_i \mathbf{h}_j \rangle_{\text{data}} - \langle \mathbf{h}_i \mathbf{h}_j \rangle_{\text{recon.}}, \quad (40)$$

where $\langle \cdot \rangle_{\text{data}}$ and $\langle \cdot \rangle_{\text{recon.}}$ denote expectation over data and reconstruction of data, respectively.

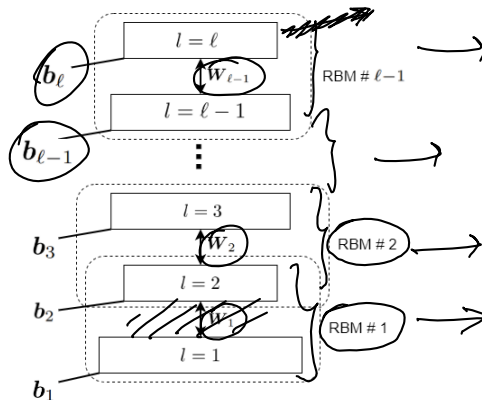
Deep Belief Network

Stacking RBM Models

- We can train a neural network using RBM training (2006) [14, 15]. Training a neural network using RBM training can result in very good initialization of weights for training network using backpropagation.
- Before the development of ReLU [18] and dropout [19], multilayer perceptron networks could not become deep for the problem of vanishing gradients. This was because random initial weights were not suitable enough for starting optimization in backpropagation, especially in deep networks. Therefore, a method was proposed for pre-training neural networks which initializes network to a suitable set of weights and then the pre-trained weights are fine-tuned using backpropagation [14, 15].
- A neural network consists of several layers. Let ℓ denote the **number of layers**, where the first layer gets the input data, and let p_ℓ be the **number of neurons** in the ℓ -th layer. By convention, we have $p_1 = d$.

Stacking RBM Models

- We can consider every two successive layers as one RBM.



Stacking RBM Models

- We start from the first pair of layers as an RBM and we introduce training dataset $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$ as the visible variable $\{\mathbf{v}_i\}_{i=1}^n$ of the first pair of layers. We train the weights and biases of this first layer as an RBM using the algorithm of training RBM.
- After training this RBM, we generate n p_2 -dimensional hidden variables using Gibbs sampling in the algorithm of sampling in RBM.
- Now, we **consider the hidden variables of the first RBM as the visible variables for the second RBM (the second pair of layers)**. Again, this RBM is trained by the algorithm of training RBM and, then, hidden variables are generated using Gibbs sampling in the algorithm of sampling in RBM.
- This procedure is **repeated until all pairs of layers are trained using RBM training**.
- This layer-wise training of neural network has a **greedy** approach [16]. This greedy training of layers prepares good **initialized weights and biases** for the whole neural network. After this initialization, we can fine-tune the weights and biases using **backpropagation** [39].
- The explained training algorithm was first proposed in [14, 15] and was used for dimensionality reduction.
- By increasing ℓ to any large number, the network becomes large and **deep**. As layers are trained one by one as RBM models, we can make the network as deep as we want without being worried for **vanishing gradients** because weights are initialized well for backpropagation.
- As this network can get deep and is pre-trained by belief propagation (RBM training), it is referred to as the **Deep Belief Network (DBN)** [15, 17].
- DBN can be seen as a **stack of RBM models**.

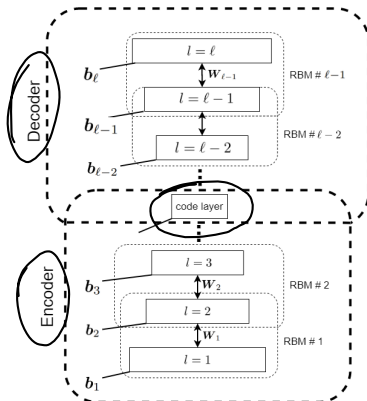
Stacking RBM Models

```
1 Input: training data  $\{x_i\}_{i=1}^n$ 
2 // pre-training:
3 for  $l$  from 1 to  $\ell - 1$  do
4     if  $l = 1$  then
5          $\{v_i\}_{i=1}^n \leftarrow \{x_i\}_{i=1}^n$ 
6     else
7         // generate  $n$  hidden variables of previous
7         RBM:
8          $\{h_i\}_{i=1}^n \leftarrow$  Algorithm 1 for  $(l - 1)$ -th
8         RBM  $\leftarrow \{v_i\}_{i=1}^n$ 
9          $\{v_i\}_{i=1}^n \leftarrow \{h_i\}_{i=1}^n$ 
10     $\underline{W_l, b_l, b_{l+1}} \leftarrow$  Algorithm 2 for  $l$ -th RBM
10     $\leftarrow \{v_i\}_{i=1}^n$ 
11 // fine-tuning using backpropagation:
12 Initialize network with weights  $\{W_l\}_{l=1}^{\ell-1}$  and
12 biases  $\{b_l\}_{l=2}^{\ell}$ .
13  $\{W_l\}_{l=1}^{\ell-1}, \{b_l\}_{l=1}^{\ell} \leftarrow$  Backpropagate the error
13 of loss fro several epochs.
```

Algorithm 3: Training a deep belief network

Stacking RBM Models

- Note that pre-training of DBN is an **unsupervised** task because RBM training is unsupervised. Fine-tuning of DBN can be either unsupervised or supervised depending on the loss function for backpropagation.
- If the DBN is an **autoencoder** with a low-dimensional middle layer in the network, both its pre-training and fine-tuning stages are unsupervised because the loss function of backpropagation is also a mean squared error. This DBN autoencoder can learn a low-dimensional embedding or representation of data and can be used for **dimensionality reduction** (2006) [14].



Acknowledgment

- This slide deck is based on our tutorial paper "Restricted boltzmann machine and deep belief network: Tutorial and survey" [40].
- For more information on RBM and DBN, refer to our tutorial paper [40].
- Some slides of this slide deck are inspired by teachings of Prof. Ali Ghodsi at University of Waterloo, Department of Statistics.
- Some slides of this slide deck (the Ising model part) are inspired by teachings of Prof. Mehdi Molkarai at University of Waterloo, Department of Statistics.

References

- [1] L. Boltzmann, "Studien über das Gleichgewicht der lebendigen Kraft," *Wissenschaftliche Abhandlungen*, vol. 1, pp. 49–96, 1868.
- [2] J. W. Gibbs, *Elementary principles in statistical mechanics*. Courier Corporation, 1902.
- [3] K. Huang, *Statistical Mechanics*. John Wiley & Sons, 1987.
- [4] W. Lenz, "Beiträge zum Verständnis der magnetischen Eigenschaften in festen Körpern," *Physikalische Z.*, vol. 21, pp. 613–615, 1920.
- [5] E. Ising, "Beitrag zur Theorie des Ferromagnetismus," *Zeitschrift für Physik*, vol. 31, no. 1, pp. 253–258, 1925.
- [6] W. A. Little, "The existence of persistent states in the brain," *Mathematical biosciences*, vol. 19, no. 1-2, pp. 101–120, 1974.
- [7] J. J. Hopfield, "Neural networks and physical systems with emergent collective computational abilities," *Proceedings of the national academy of sciences*, vol. 79, no. 8, pp. 2554–2558, 1982.
- [8] G. E. Hinton and T. J. Sejnowski, "Optimal perceptual inference," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, vol. 448, IEEE, 1983.

References (cont.)

- [9] D. H. Ackley, G. E. Hinton, and T. J. Sejnowski, "A learning algorithm for Boltzmann machines," *Cognitive science*, vol. 9, no. 1, pp. 147–169, 1985.
- [10] Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang, "A tutorial on energy-based learning," *Predicting structured data*, vol. 1, 2006.
- [11] G. E. Hinton, "Training products of experts by minimizing contrastive divergence," *Neural computation*, vol. 14, no. 8, pp. 1771–1800, 2002.
- [12] M. Welling, M. Rosen-Zvi, and G. E. Hinton, "Exponential family harmoniums with an application to information retrieval.," in *Advances in neural information processing systems*, vol. 4, pp. 1481–1488, 2004.
- [13] S. Geman and D. Geman, "Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images," *IEEE Transactions on pattern analysis and machine intelligence*, vol. PAMI-6, no. 6, pp. 721–741, 1984.
- [14] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, no. 5786, pp. 504–507, 2006.
- [15] G. E. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [16] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, "Greedy layer-wise training of deep networks," in *Advances in neural information processing systems*, pp. 153–160, 2007.

References (cont.)

- [17] G. E. Hinton, “Deep belief networks,” *Scholarpedia*, vol. 4, no. 5, p. 5947, 2009.
- [18] X. Glorot, A. Bordes, and Y. Bengio, “Deep sparse rectifier neural networks,” in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pp. 315–323, JMLR Workshop and Conference Proceedings, 2011.
- [19] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: a simple way to prevent neural networks from overfitting,” *The journal of machine learning research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [20] B. Ghoggh and M. Crowley, “The theory behind overfitting, cross validation, regularization, bagging, and boosting: tutorial,” *arXiv preprint arXiv:1905.12787*, 2019.
- [21] S. Carroll, *From eternity to here: the quest for the ultimate theory of time*. Penguin, 2010.
- [22] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi, “Optimization by simulated annealing,” *science*, vol. 220, no. 4598, pp. 671–680, 1983.
- [23] L. A. Passos and J. P. Papa, “Temperature-based deep Boltzmann machines,” *Neural Processing Letters*, vol. 48, no. 1, pp. 95–107, 2018.
- [24] D. Alberici, A. Barra, P. Contucci, and E. Mingione, “Annealing and replica-symmetry in deep Boltzmann machines,” *Journal of Statistical Physics*, vol. 180, no. 1, pp. 665–677, 2020.

References (cont.)

- [25] D. Alberici, P. Contucci, and E. Mingione, “Deep Boltzmann machines: rigorous results at arbitrary depth,” in *Annales Henri Poincaré*, pp. 1–24, Springer, 2021.
- [26] S. G. Brush, “History of the Lenz-Ising model,” *Reviews of modern physics*, vol. 39, no. 4, p. 883, 1967.
- [27] G. E. Hinton, “Boltzmann machine,” *Scholarpedia*, vol. 2, no. 5, p. 1668, 2007.
- [28] D. Hebb, *The Organization of Behavior*. Wiley & Sons, New York, 1949.
- [29] J. J. Hopfield, “Neurons with graded response have collective computational properties like those of two-state neurons,” *Proceedings of the national academy of sciences*, vol. 81, no. 10, pp. 3088–3092, 1984.
- [30] H. Ramsauer, B. Schäfl, J. Lehner, P. Seidl, M. Widrich, T. Adler, L. Gruber, M. Holzleitner, M. Pavlović, G. K. Sandve, et al., “Hopfield networks is all you need,” *arXiv preprint arXiv:2008.02217*, 2020.
- [31] D. Krotov and J. J. Hopfield, “Dense associative memory for pattern recognition,” *Advances in neural information processing systems*, vol. 29, pp. 1172–1180, 2016.
- [32] D. Krotov and J. Hopfield, “Large associative memory problem in neurobiology and machine learning,” in *International Conference on Learning Representations (ICLR)*, 2021.
- [33] D. Krotov, “Hierarchical associative memory,” *arXiv preprint arXiv:2107.06446*, 2021.

References (cont.)

- [34] C. M. Bishop, "Pattern recognition," *Machine learning*, vol. 128, no. 9, 2006.
- [35] R. Salakhutdinov and G. Hinton, "Deep Boltzmann machines," in *Artificial intelligence and statistics*, pp. 448–455, PMLR, 2009.
- [36] G. E. Hinton, "A practical guide to training restricted Boltzmann machines," in *Neural networks: Tricks of the trade*, pp. 599–619, Springer, 2012.
- [37] B. Ghojogh, H. Nekoei, A. Ghojogh, F. Karray, and M. Crowley, "Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review," *arXiv preprint arXiv:2011.00901*, 2020.
- [38] G. W. Taylor, G. E. Hinton, and S. T. Roweis, "Modeling human motion using binary latent variables," in *Advances in neural information processing systems*, pp. 1345–1352, 2007.
- [39] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533–536, 1986.
- [40] B. Ghojogh, A. Ghodsi, F. Karray, and M. Crowley, "Restricted boltzmann machine and deep belief network: Tutorial and survey," *arXiv preprint arXiv:2107.12521*, 2021.