

## 受限波尔兹曼机\*

张春霞<sup>1</sup>, 姬楠楠<sup>2</sup>, 王冠伟<sup>3</sup>

(1- 西安交通大学数学与统计学院, 西安 710049; 2- 长安大学理学院, 西安 710064;

3- 西安工业大学机电工程学院, 西安 710021)

**摘 要:** 受限波尔兹曼机(restricted Boltzmann machines, RBM)是一类具有两层结构、对称连接且无自反馈的随机神经网络模型, 层间全连接, 层内无连接. 近年来, 随着RBM的快速学习算法—对比散度的出现, 机器学习界掀起了研究RBM理论及应用的热潮. 实践表明, RBM是一种有效的特征提取方法, 用于初始化前馈神经网络可明显提高泛化能力, 堆叠多个RBM组成的深度信念网络能提取更抽象的特征. 鉴于RBM的优点及其在深度学习中的广泛应用, 本文对RBM的基本模型、学习算法、参数设置、评估方法、变形算法等进行了详细介绍, 最后探讨了RBM在未来值得研究的方向.

**关键词:** 机器学习; 深度学习; 受限波尔兹曼机; 对比散度; Gibbs 采样

**分类号:** AMS(2000) 92B20; 68T05

**中图分类号:** TP181; O235

**文献标识码:** A

### 1 引言

机器学习研究的主要任务是设计和开发计算机根据实际数据进行“智能学习”的算法, 使其自动发现隐藏在数据中的模式和规律. 目前, 各种机器学习算法在人工智能、工程应用、医学等诸多领域都扮演着非常重要的角色. 人工神经网络(artificial neural network, ANN)<sup>[1,2]</sup>作为一种通过模仿生物神经网络的结构和功能而建立起来的计算模型, 因其自学习、自组织、较好的容错性和优良的非线性逼近能力等优点, 而受到众多领域学者的广泛关注.

在诸多人工神经网络模型中, 波尔兹曼机(Boltzmann machine, BM)<sup>[3]</sup>是Hinton和Sejnowski于1986年提出的一种根植于统计力学的随机神经网络. 这种网络中的神经元是随机神经元, 其输出只有两种状态(未激活、激活), 一般用二进制的0和1表示, 状态的取值根据概率法则决定. 从功能上讲, BM是由随机神经元全连接组成的反馈神经网络, 且对称连接, 无自反馈, 包含一个可见层和一个隐层的BM模型, 如图1(a)所示.

BM具有强大的无监督学习能力, 能学习数据中复杂的规则. 但是, 拥有这种学习能力的代价是其训练(学习)过程耗时. 此外, BM所表示的分布不仅无法确切计算, 得到该分布的随机样本也很困难. 于是, Smolensky<sup>[4]</sup>引入了一种限制的波尔兹曼机(restricted

收稿日期: 2013-08-22. 作者简介: 张春霞(1980年6月生), 女, 博士, 副教授. 研究方向: 模式识别与集成学习.

\*基金项目: 国家重点基础研究发展计划973项目(2013CB329406); 国家自然科学基金重大研究计划(91230101); 国家自然科学基金(11201367); 中央高校基本科研业务费专项基金(xjj2011048).

Boltzmann machine, RBM). RBM 具有一个可见层, 一个隐层, 层内无连接, 其结构如图 1(b) 所示. RBM 具有很好的性质<sup>[5]</sup>: 在给定可见层单元状态时, 各隐单元的激活条件独立; 反之, 在给定隐单元状态时, 可见层单元的激活亦条件独立. 这样一来, 尽管 RBM 所表示的分布仍无法有效计算, 但通过 Gibbs 采样 (Gibbs sampling) 可以得到 RBM 所表示分布的随机样本. 此外, Roux 和 Bengio<sup>[6]</sup> 从理论上证明, 只要隐单元足够多, RBM 能够拟合任意离散分布. 自 Hinton<sup>[7]</sup> 于 2002 年提出了 RBM 的快速学习算法—对比散度 (contrastive divergence, CD) 之后, 机器学习界掀起了一轮研究 RBM、CD 算法的理论及应用的热潮. 理论方面, RBM 的 CD 快速学习算法促进了研究者们对随机近似理论、基于能量的模型、未归一化的统计模型的研究<sup>[8]</sup>. 应用方面, RBM 目前已被成功地应用于不同的机器学习问题<sup>[9-14]</sup>, 如分类、回归、降维、高维时间序列建模、图像特征提取、协同过滤等等.

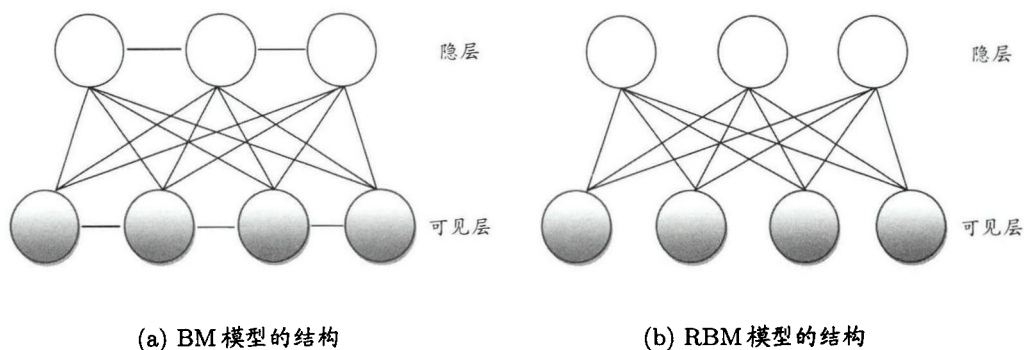


图 1: BM 和 RBM 模型的结构比较

2006 年, Hinton 等人<sup>[15]</sup> 提出了一种深度信念网络 (deep belief nets, DBN), 并给出了该模型的一个高效学习算法, 该算法至今仍是深度学习方法的主要框架. 在该算法中, 一个 DBN 模型被视为由若干个 RBM 堆叠在一起形成, 训练时可通过由低到高逐层训练 RBM 来实现:

- 1) 底部 RBM 以原始输入数据训练;
- 2) 将底部 RBM 抽取的特征作为顶部 RBM 的输入训练;
- 3) 过程 1) 和 2) 可以重复来训练所需要的尽可能多的层数.

由于 RBM 可以通过 CD 算法快速训练, 这一框架绕过了直接从整体上训练 DBN 的高复杂度, 从而将其化简为对多个 RBM 的训练问题. Hinton 建议, 经过这种方式训练后, 可以再通过传统的全局学习算法 (如反向传播算法) 对网络进行微调, 使模型收敛到局部最优点. 这种学习算法, 本质上等同于通过逐层 RBM 训练为模型寻找一个较好的参数初始值. 这样一来, 不仅解决了模型训练速度慢的问题, 大量试验结果也表明, 这种方式能够产生非常好的参数初始值, 使得模型的建模能力也大大提升. 自此, 机器学习领域又产生了一个新的研究方向—深度学习 (deep learning)<sup>[16-18]</sup>, 明确提出了设计面向人工智能的机器学习算法的目标.

当前, 以RBM为基本构成模块的DBN模型被认为是最有效的深度学习算法之一。鉴于RBM在深度学习领域中占据的核心位置及其本身的良好性质, 为了给RBM的初学者提供入门指导, 同时为设计与之相关的新算法提供参考, 本文将对RBM进行较为系统的介绍, 详细阐述其基本模型、具有代表性的快速学习算法、参数设置、评估方法及其变形算法, 最后对RBM在未来值得研究的方向进行探讨。

本文后续内容安排如下: 第2节介绍受限波尔兹曼机RBM的基本模型, 第3节详细阐述当前训练RBM的快速学习算法, 第4节讨论RBM的参数设置, 第5节给出评价RBM优劣的方法, 第6节简单介绍几种具有代表性的RBM变形算法, 第7节是总结与展望, 主要探讨RBM在未来值得研究的方向。

## 2 受限波尔兹曼机RBM的基本模型

RBM也可以被视为一个无向图(undirected graph)模型, 如图2所示。 $\mathbf{v}$ 为可见层, 表示观测数据,  $\mathbf{h}$ 为隐层, 可视为一些特征提取器(feature detectors),  $W$ 为两层之间的连接权重。Welling<sup>[19]</sup>指出, RBM中的隐单元和可见单元可以为任意的指数族单元(即给定隐单元(可见单元), 可见单元(隐单元)的分布可以为任意的指数族分布), 如softmax单元、高斯单元、泊松单元等等。这里, 为了讨论方便起见, 我们假设所有的可见单元和隐单元均为二值变量, 即对任意的 $i, j, v_i \in \{0, 1\}, h_j \in \{0, 1\}$ 。

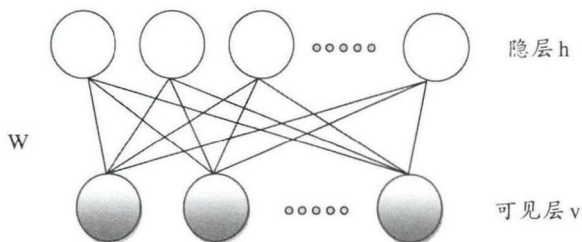


图2: RBM的图模型表示, 层内单元之间无连接

如果一个RBM有 $n$ 个可见单元和 $m$ 个隐单元, 用向量 $\mathbf{v}$ 和 $\mathbf{h}$ 分别表示可见单元和隐单元的状态。其中 $v_i$ 表示第 $i$ 个可见单元的状态,  $h_j$ 表示第 $j$ 个隐单元的状态。那么, 对于一组给定的状态 $(\mathbf{v}, \mathbf{h})$ , RBM作为一个系统所具备的能量定义为

$$E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta}) = - \sum_{i=1}^n a_i v_i - \sum_{j=1}^m b_j h_j - \sum_{i=1}^n \sum_{j=1}^m v_i W_{ij} h_j, \quad (1)$$

上式中 $\boldsymbol{\theta} = \{W_{ij}, a_i, b_j\}$ 是RBM的参数, 它们均为实数。其中 $W_{ij}$ 表示可见单元 $i$ 与隐单元 $j$ 之间的连接权重,  $a_i$ 表示可见单元 $i$ 的偏置(bias),  $b_j$ 表示隐单元 $j$ 的偏置。当参数确定时, 基于该能量函数, 我们可以得到 $(\mathbf{v}, \mathbf{h})$ 的联合概率分布

$$P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta}) = \frac{e^{-E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})}}{Z(\boldsymbol{\theta})}, \quad Z(\boldsymbol{\theta}) = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})}, \quad (2)$$

其中  $Z(\theta)$  为归一化因子 (也称为配分函数, partition function).

对于一个实际问题, 我们最关心的是由 RBM 所定义的关于观测数据  $\mathbf{v}$  的分布  $P(\mathbf{v} | \theta)$ , 即联合概率分布  $P(\mathbf{v}, \mathbf{h} | \theta)$  的边际分布, 也称为似然函数 (likelihood function)

$$P(\mathbf{v} | \theta) = \frac{1}{Z(\theta)} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h} | \theta)}. \quad (3)$$

为了确定该分布, 需要计算归一化因子  $Z(\theta)$ , 这需要  $2^{n+m}$  次计算. 因此, 即使通过训练可以得到模型的参数  $W_{ij}$ ,  $a_i$  和  $b_j$ , 我们仍旧无法有效地计算由这些参数所确定的分布.

但是, 由 RBM 的特殊结构 (即层间有连接, 层内无连接) 可知: 当给定可见单元的状态时, 各隐单元的激活状态之间是条件独立的. 此时, 第  $j$  个隐单元的激活概率为

$$P(h_j = 1 | \mathbf{v}, \theta) = \sigma\left(b_j + \sum_i v_i W_{ij}\right), \quad (4)$$

其中  $\sigma(x) = \frac{1}{1+\exp(-x)}$  为 sigmoid 激活函数.

由于 RBM 的结构是对称的, 当给定隐单元的状态时, 各可见单元的激活状态之间也是条件独立的, 即第  $i$  个可见单元的激活概率为

$$P(v_i = 1 | \mathbf{h}, \theta) = \sigma\left(a_i + \sum_j W_{ij} h_j\right). \quad (5)$$

### 3 基于对比散度的 RBM 快速学习算法

学习 RBM 的任务是求出参数  $\theta$  的值, 以拟合给定的训练数据. 参数  $\theta$  可以通过最大化 RBM 在训练集 (假设包含  $T$  个样本) 上的对数似然函数学习得到, 即

$$\theta^* = \arg \max_{\theta} \mathcal{L}(\theta) = \arg \max_{\theta} \sum_{t=1}^T \log P(\mathbf{v}^{(t)} | \theta). \quad (6)$$

为了获得最优参数  $\theta^*$ , 我们可以使用随机梯度上升法 (stochastic gradient ascent) 求

$$\mathcal{L}(\theta) = \sum_{t=1}^T \log P(\mathbf{v}^{(t)} | \theta)$$

的最大值. 其中, 关键步骤是计算  $\log P(\mathbf{v}^{(t)} | \theta)$  关于各个模型参数的偏导数.

由于

$$\begin{aligned} \mathcal{L}(\theta) &= \sum_{t=1}^T \log P(\mathbf{v}^{(t)} | \theta) = \sum_{t=1}^T \log \sum_{\mathbf{h}} P(\mathbf{v}^{(t)}, \mathbf{h} | \theta) \\ &= \sum_{t=1}^T \log \frac{\sum_{\mathbf{h}} \exp[-E(\mathbf{v}^{(t)}, \mathbf{h} | \theta)]}{\sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp[-E(\mathbf{v}, \mathbf{h} | \theta)]} \\ &= \sum_{t=1}^T \left( \log \sum_{\mathbf{h}} \exp[-E(\mathbf{v}^{(t)}, \mathbf{h} | \theta)] - \log \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp[-E(\mathbf{v}, \mathbf{h} | \theta)] \right). \end{aligned} \quad (7)$$

令  $\theta$  表示  $\boldsymbol{\theta}$  中的某一个参数, 则对数似然函数关于  $\theta$  的梯度为

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta} &= \sum_{t=1}^T \frac{\partial}{\partial \theta} \left( \log \sum_{\mathbf{h}} \exp[-E(\mathbf{v}^{(t)}, \mathbf{h} | \boldsymbol{\theta})] - \log \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp[-E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})] \right) \\ &= \sum_{t=1}^T \left( \sum_{\mathbf{h}} \frac{\exp[-E(\mathbf{v}^{(t)}, \mathbf{h} | \boldsymbol{\theta})]}{\sum_{\mathbf{h}} \exp[-E(\mathbf{v}^{(t)}, \mathbf{h} | \boldsymbol{\theta})]} \times \frac{\partial(-E(\mathbf{v}^{(t)}, \mathbf{h} | \boldsymbol{\theta}))}{\partial \theta} \right. \\ &\quad \left. - \sum_{\mathbf{v}} \sum_{\mathbf{h}} \frac{\exp[-E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})]}{\sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp[-E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})]} \times \frac{\partial(-E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta}))}{\partial \theta} \right) \\ &= \sum_{t=1}^T \left( \left\langle \frac{\partial(-E(\mathbf{v}^{(t)}, \mathbf{h} | \boldsymbol{\theta}))}{\partial \theta} \right\rangle_{P(\mathbf{h} | \mathbf{v}^{(t)}, \boldsymbol{\theta})} - \left\langle \frac{\partial(-E(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta}))}{\partial \theta} \right\rangle_{P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})} \right), \quad (8)\end{aligned}$$

其中  $\langle \cdot \rangle_P$  表示求关于分布  $P$  的数学期望.  $P(\mathbf{h} | \mathbf{v}^{(t)}, \boldsymbol{\theta})$  表示在可见单元限定为已知的训练样本  $\mathbf{v}^{(t)}$  时, 隐层的概率分布, 故式 (8) 中的前一项比较容易计算.  $P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})$  表示可见单元与隐单元的联合分布, 由于归一化因子  $Z(\boldsymbol{\theta})$  的存在, 该分布很难获取, 导致我们无法直接计算式 (8) 中的第二项, 只能通过一些采样方法 (如 Gibbs 采样) 获取其近似值. 值得指出的是, 在最大化似然函数的过程中, 为提高计算效率, 上述偏导数在每一迭代步中的计算一般只基于部分而非所有的训练样本进行, 关于这部分内容我们将在后面讨论 RBM 的参数设置时详细阐述.

假设只有一个训练样本, 我们分别用 “data” 和 “model” 简记概率分布  $P(\mathbf{h} | \mathbf{v}^{(t)}, \boldsymbol{\theta})$  和  $P(\mathbf{v}, \mathbf{h} | \boldsymbol{\theta})$ , 则对数似然函数关于连接权重  $W_{ij}$ 、可见层单元的偏置  $a_i$  和隐层单元的偏置  $b_j$  的偏导数分别为

$$\begin{aligned}\frac{\partial \log P(\mathbf{v} | \boldsymbol{\theta})}{\partial W_{ij}} &= \langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{model}}, \\ \frac{\partial \log P(\mathbf{v} | \boldsymbol{\theta})}{\partial a_i} &= \langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{model}}, \\ \frac{\partial \log P(\mathbf{v} | \boldsymbol{\theta})}{\partial b_j} &= \langle h_j \rangle_{\text{data}} - \langle h_j \rangle_{\text{model}}.\end{aligned}$$

### 3.1 RBM 中的 Gibbs 采样

Gibbs 采样 (Gibbs sampling)<sup>[20]</sup> 是一种基于 MCMC 策略的采样方法. 对于一个  $K$  维随机向量  $\mathbf{X} = (X_1, X_2, \dots, X_K)$ , 假设我们无法求得  $\mathbf{X}$  的联合分布  $P(\mathbf{X})$ , 但知道给定  $\mathbf{X}$  的其他分量时, 其第  $k$  个分量  $X_k$  的条件分布, 即

$$P(X_k | X_{k-}), \quad X_{k-} = (X_1, X_2, \dots, X_{k-1}, X_{k+1}, \dots, X_K).$$

那么, 我们可以从  $\mathbf{X}$  的一个任意状态 (比如  $[x_1(0), x_2(0), \dots, x_K(0)]$ ) 开始, 利用上述条件分布, 迭代地对其分量依次采样, 随着采样次数的增加, 随机向量  $[x_1(n), x_2(n), \dots, x_K(n)]$  的概率分布将以  $n$  的几何级数的速度收敛于  $P(\mathbf{X})$ .

基于RBM模型的对称结构, 以及其中神经元状态的条件独立性, 我们可以使用Gibbs采样得到RBM所定义分布的随机样本. 在RBM中进行 $k$ 步Gibbs采样的具体算法为: 用一个训练样本(或可见层的任何随机化状态)初始化可见层的状态 $\mathbf{v}_0$ , 交替进行如下采样:

$$\begin{aligned} \mathbf{h}_0 &\sim P(\mathbf{h} | \mathbf{v}_0), & \mathbf{v}_1 &\sim P(\mathbf{v} | \mathbf{h}_0), \\ \mathbf{h}_1 &\sim P(\mathbf{h} | \mathbf{v}_1), & \mathbf{v}_2 &\sim P(\mathbf{v} | \mathbf{h}_1), \\ &\dots\dots, & \mathbf{v}_{k+1} &\sim P(\mathbf{v} | \mathbf{h}_k). \end{aligned}$$

在采样步数 $k$ 足够大的情况下, 可以得到服从RBM所定义分布的样本. 此外, 使用Gibbs采样也可以得到式(8)中第二项的一个近似.

### 3.2 基于对比散度的快速学习算法

尽管利用Gibbs采样我们可以得到对数似然函数关于未知参数梯度的近似, 但通常需要使用较大的采样步数, 这使得RBM的训练效率仍不高, 尤其是当观测数据的特征维数较高时.

2002年, Hinton<sup>[7]</sup>提出了RBM的一个快速学习算法, 即对比散度(contrastive divergence, CD). 与Gibbs采样不同, Hinton指出当使用训练数据初始化 $\mathbf{v}_0$ 时, 我们仅需要使用 $k$ (通常 $k=1$ )步Gibbs采样便可以得到足够好的近似. 在CD算法一开始, 可见单元的状态被设置成一个训练样本, 并利用式(4)计算所有隐层单元的二值状态. 在所有隐层单元的状态确定之后, 根据式(5)来确定第 $i$ 个可见单元 $v_i$ 取值为1的概率, 进而产生可见层的一个重构(reconstruction). 因此, 使用随机梯度上升法最大化对数似然函数在训练数据上的值时, 各参数的更新准则为

$$\begin{aligned} \Delta W_{ij} &= \epsilon (\langle v_i h_j \rangle_{\text{data}} - \langle v_i h_j \rangle_{\text{recon}}), \\ \Delta a_i &= \epsilon (\langle v_i \rangle_{\text{data}} - \langle v_i \rangle_{\text{recon}}), \\ \Delta b_j &= \epsilon (\langle h_j \rangle_{\text{data}} - \langle h_j \rangle_{\text{recon}}), \end{aligned}$$

这里 $\epsilon$ 是学习率(learning rate),  $\langle \cdot \rangle_{\text{recon}}$ 表示一步重构后模型定义分布.

在RBM中, 可见单元数一般等于训练数据的特征维数, 而隐单元数需要事先给定. 为了与前文记号一致, 假设可见单元数和隐单元数分别为 $n$ 和 $m$ . 令 $W$ 表示可见层与隐层间的连接权重矩阵( $m \times n$ 阶),  $\mathbf{a}$ ( $n$ 维列向量)和 $\mathbf{b}$ ( $m$ 维列向量)分别表示可见层与隐层的偏置向量.

#### 算法1 RBM的基于CD的快速学习算法主要步骤

- 
- 输入: 一个训练样本 $\mathbf{x}_0$ ; 隐层单元个数 $m$ ; 学习率 $\epsilon$ ; 最大训练周期 $T$ .
  - 输出: 连接权重矩阵 $W$ 、可见层的偏置向量 $\mathbf{a}$ 、隐层的偏置向量 $\mathbf{b}$ .
  - 训练阶段:
    - 初始化: 令可见层单元的初始状态 $\mathbf{v}_1 = \mathbf{x}_0$ ;  $W$ 、 $\mathbf{a}$ 和 $\mathbf{b}$ 为随机选取的较小数值.
    - For  $t = 1, 2, \dots, T$

For  $j = 1, 2, \dots, m$  (对所有隐单元)

计算  $P(h_{1j} = 1 | \mathbf{v}_1)$ , 即  $P(h_{1j} = 1 | \mathbf{v}_1) = \sigma(b_j + \sum_i v_{1i} W_{ij})$ ;  
从条件分布  $P(h_{1j} | \mathbf{v}_1)$  中抽取  $h_{1j} \in \{0, 1\}$ .

EndFor

For  $i = 1, 2, \dots, n$  (对所有可见单元)

计算  $P(v_{2i} = 1 | \mathbf{h}_1)$ , 即  $P(v_{2i} = 1 | \mathbf{h}_1) = \sigma(a_i + \sum_j W_{ij} h_{1j})$ ;  
从条件分布  $P(v_{2i} | \mathbf{h}_1)$  中抽取  $v_{2i} \in \{0, 1\}$ .

EndFor

For  $j = 1, 2, \dots, m$  (对所有隐单元)

计算  $P(h_{2j} = 1 | \mathbf{v}_2)$ , 即  $P(h_{2j} = 1 | \mathbf{v}_2) = \sigma(b_j + \sum_i v_{2i} W_{ij})$ ;

EndFor

按下式更新各个参数

- $W \leftarrow W + \epsilon(P(\mathbf{h}_{1.} = 1 | \mathbf{v}_1) \mathbf{v}_1^T - P(\mathbf{h}_{2.} = 1 | \mathbf{v}_2) \mathbf{v}_2^T)$ ;
- $\mathbf{a} \leftarrow \mathbf{a} + \epsilon(\mathbf{v}_1 - \mathbf{v}_2)$ ;
- $\mathbf{b} \leftarrow \mathbf{b} + \epsilon(P(\mathbf{h}_{1.} = 1 | \mathbf{v}_1) - P(\mathbf{h}_{2.} = 1 | \mathbf{v}_2))$ ;

EndFor

在上述算法中, 记号  $P(\mathbf{h}_{k.} = 1 | \mathbf{v}_k)$  ( $k = 1, 2$ ) 是  $m$  维列向量, 其第  $j$  个元素为  $P(h_{kj} = 1 | \mathbf{v}_k)$ .

尽管上述基于 CD 的学习算法是针对 RBM 的可见单元和隐层单元均为二值变量的情形提出的, 但很容易推广到可见层单元为高斯变量、可见层和隐层单元均为高斯变量等其他情形, 关于这方面的具体研究可参见文献 [21–25].

此外, 还有一些研究者在 CD 算法的基础上, 对其作了进一步改进. 例如 Tieleman<sup>[26]</sup> 提出了持续对比散度 (persistent contrastive divergence, PCD) 算法, 它与 CD 算法的区别在于: 首先, PCD 不再使用训练数据初始化 CD 算法中 Gibbs 采样的马氏链; 其次, PCD 算法中的学习率较小且不断衰减. 根据随机近似理论, 尽管每次更新参数后模型都发生了改变 (每次对于  $W$ ,  $\mathbf{a}$  和  $\mathbf{b}$  的更新, RBM 定义的分布都会发生改变), 但由于学习率较小且不断衰减, 则可认为那条马氏链产生的负样本是由当前 RBM 定义的分布的一个近似分布采样而来. Tieleman 和 Hinton<sup>[27]</sup> 进一步改进了 PCD 算法, 他们通过引入一组辅助参数以加快 PCD 中的马氏链的混合率, 提出了快速持续对比散度 (fast persistent contrastive divergence, FPCD) 算法. Desjardins 等<sup>[28]</sup> 提出了平行回火算法 (parallel tempering, PT), 该算法的主要思想是模型以不同温度平行运行多个 Markov 链, 假设每个链处于一个有序序列温度  $t_i$  上,  $t_0 < t_1 < \dots < t_i < \dots < t_{T-1} < t_T$ , 其中  $t_0 = 1$  是目标分布采样温度,  $t_T = \tau$  是高温. 在每个运行步骤中, PT 算法以一定的概率交换两个临近链之间的样本, 以增加 Markov 链的混合率, 进而快速逼近目标分布.

关于 RBM 的学习算法, 除了上述提到的基于 CD 的一些方法之外, 还有最大化拟似然函数 (maximum pseudo-likelihood)、比率匹配方法 (ratio matching)、状态翻转变

换算子 (flip-the-state transition operator) 等, 有兴趣的读者可参阅文献 [29–33] 查找关于 RBM 学习算法比较详细的阐述。

## 4 RBM 的参数设置

RBM 的训练通常是基于 CD 算法进行的, 但如何设置其中的一些参数 (如隐单元个数、学习率、参数的初始值等), 是需要有一定经验的。近来, 已有部分研究结果<sup>[34,35]</sup>表明: 对于特定的数据集和 RBM 结构, 如果参数设置不合适, RBM 将很难对真正的数据分布正确建模。因此, 对实际使用者 (尤其是初学者) 来说, 了解 RBM 中参数设置的一般规则非常重要。根据 Hinton<sup>[23]</sup> 的建议及我们通过数值试验所获经验, RBM 中的参数设置可参考以下规则。

**小批量数据及其容量** 对于连接权重、可见层和隐层偏置的更新, 虽然可以基于一个训练样本进行 (类似于在线学习的方式), 但计算量很大。为充分利用图形处理器 GPU (graphic processing unit) 或 Matlab 中矩阵之间相乘运算的优势, 可将训练集事先分成包含几十或几百个样本的小批量数据 (mini-batches) 以提高计算效率。同时, 为了避免在小批量数据的样本容量发生改变时, 学习率也必须做相应的修改, 通常的做法是在参数的更新过程中, 使用参数的平均梯度 (即总梯度除以数据容量), 即

$$\theta^{(t+1)} = \theta^{(t)} + \epsilon \left( \frac{1}{B} \sum_{t'=Bt+1}^{B(t+1)} \frac{\partial \log P(\mathbf{v}^{(t')} | \boldsymbol{\theta})}{\partial \theta} \right),$$

这里  $B$  表示小批量数据的容量, 其值不应设得太大。 $B=1$  表示参数更新以在线学习的方式进行, 而  $B=T$  则对应批处理方式。一般而言, 若训练集是包含来自不同类 (具有同等概率) 的样本, 理想的  $B$  应为总类数, 使得每批数据中都包含来自每个类的一个样本, 以减小梯度估计的抽样误差。对于其他数据集, 则可先随机化训练样本的次序, 再将其分为容量为 10 的倍数的小批量数据。

**学习率** 学习率若过大, 将导致重构误差急剧增加, 权重也会变得异常大。设置学习率的一般做法是先做权重更新和权重的直方图, 令权重更新量为权重的  $10^{-3}$  倍左右。如果有一个单元的输入值很大, 则权重更新应再小一些, 因为同一方向上较多小的波动很容易改变梯度的符号。相反地, 对于偏置, 其权重更新可以大一些。

**权重和偏置的初始值** 一般地, 连接权重  $W_{ij}$  可初始化为来自正态分布  $N(0, 0.01)$  的随机数, 隐单元的偏置  $b_j$  初始化为 0。对于第  $i$  个可见单元, 其偏置  $a_i$  通常初始化为  $\log[p_i/(1-p_i)]$ , 其中  $p_i$  表示训练样本中第  $i$  个特征处于激活状态所占的比率。

**动量学习率** 学习率  $\epsilon$  的选择至关重要,  $\epsilon$  大收敛速度快, 但是过大可能引起算法的不稳定;  $\epsilon$  小可避免不稳定情况的出现, 但收敛速度较慢。为克服这一矛盾, 一种具有代表性的思想是在参数更新式中增加动量项 (momentum), 使本次参数值修改的方向不完全由当前样本下的似然函数梯度方向决定, 而采用上一次参数值修改方向与本次梯度方向的组合。在某些情况下, 这可以避免算法过早地收敛到局部最优。以连接权重参数  $W_{ij}$  为



例, 其更新公式为

$$W_{ij}^{(t+1)} = kW_{ij}^{(t)} + \epsilon \frac{\partial \mathcal{L}}{\partial W_{ij}^{(t)}},$$

其中  $k$  为动量项学习率. 开始时,  $k$  可设为 0.5, 在重构误差处于平稳增加状态时,  $k$  可取为 0.9.

**权衰减** 使用权衰减 (weight-decay) 策略的主要目的是避免学习过程出现过拟合 (overfitting) 现象, 一般做法是在正常的梯度项后额外增加一项, 对较大的参数值作出惩罚. 最简单的罚函数是  $L_2$  函数  $(\lambda/2) \sum_i \sum_j W_{ij}^2$ , 即所有权重参数的平方和的  $1/2$  再乘上一个正则化系数  $\lambda$ ,  $\lambda$  在 RBM 中又称为权损失 (weight-cost) 系数. 重要的是, 惩罚项关于权重参数的梯度必须乘上学习率, 否则, 学习率的改变将导致优化的目标函数也发生改变. 在 RBM 中, 若使用  $L_2$  罚函数, 则  $\lambda$  可以取介于 0.01 与 0.0001 之间的任意值. 值得指出的是, 权衰减策略不需应用于可见层和隐层偏置, 因为它们不大可能导致过拟合, 且在某些情况下, 偏置的值还必须较大才行.

**隐单元个数** 如果我们关心的主要目标是避免过拟合而不是计算复杂度, 则可以先估算一下用一个好的模型描述一个数据所需的比特数, 用其乘上训练集容量. 基于所得的数, 选择比其低一个数量级的值作为隐元个数. 如果训练数据是高度冗余的 (比如数据集容量非常大), 则可以使用更少一些的隐元.

以上讨论的是 RBM 中一些常用的参数设置, 针对一个实际问题, 应使用什么类型的可见单元和隐单元, 在其中如何加入稀疏性使得隐单元只在少数情况下处于激活状态等问题的讨论, 可参见文献 [23,30,36].

## 5 RBM 的评估算法

对于一个已经学习得到或正在学习中的 RBM, 应通过何种指标评价其优劣呢? 显然, 最简单的指标就是该 RBM 在训练数据上的似然度

$$\mathcal{L}(\theta) = \sum_{t=1}^T \log P(\mathbf{v}^{(t)} | \theta).$$

但是,  $\mathcal{L}(\theta)$  的计算涉及到归一化常数  $Z(\theta)$ , 而这个值是无法通过数学方法直接解析得到的, 但我们又不可能枚举 RBM 的所有状态. 因此只能采用近似方法对 RBM 进行评估.

### 5.1 重构误差

所谓“重构误差” (reconstruction error), 就是以训练数据作为初始状态, 根据 RBM 的分布进行一次 Gibbs 采样后所获样本与原数据的差异 (一般用一范数或二范数来评估).

#### 算法 2 重构误差的计算

---

```
Error = 0                                % 初始化误差
for all  $\mathbf{v}^{(t)}$ ,  $t \in \{1, 2, \dots, T\}$  do % 对每个训练样本  $\mathbf{v}^{(t)}$  进行以下计算
```

```

h ~  $P(\cdot|\mathbf{v}^{(t)})$  % 对隐层采样
v ~  $P(\cdot|\mathbf{h})$  % 对可见层采样
Error = Error +  $\|\mathbf{v} - \mathbf{v}^{(t)}\|$  % 累计当前误差

end for
return Error % 返回总误差

```

重构误差能够在一定程度上反映RBM对训练数据的似然度, 不过并不完全可靠<sup>[23]</sup>. 但总的来说, 重构误差的计算十分简单, 因此, 在实践中非常有用.

## 5.2 退火式重要性采样

“退火式重要性采样”(annealed importance sampling, AIS)<sup>[37]</sup>是目前用于评估RBM的主流方法. 它的思想是采用“重要性采样”(importance sampling)<sup>[20]</sup>算法来估计RBM对数据的似然度. 这种算法的优点在于: 当目标分布十分陡峭时, 不直接对其进行采样, 而是引入另一个简单的分布, 在这个简单的分布上采样. 然后, 利用采样所获样本和两个分布之间的关系对原分布上的均值进行估算.

“重要性抽样”的基本思想如下: 假设我们要计算某个分布 $P_A(x)$ 的归一化常数 $Z_A$ , 那么, 我们可以引入另一个状态空间相同, 但更容易采样的分布 $P_B(x)$ , 并且事先知道它的归一化常数 $Z_B$ . 这时, 只要能计算出 $Z_A/Z_B$ 的值, 我们就可以算出原分布的归一化常数 $Z_A$ . 假设 $Z_A = \sum_x f(x)$ ,  $Z_B = \sum_x g(x)$ , 考虑它们的比例

$$\frac{Z_A}{Z_B} = \frac{\sum_x f(x)}{\sum_x g(x)} = \sum_x \frac{g(x)}{\sum_x g(x)} \frac{f(x)}{g(x)} = \left\langle \frac{f(x)}{g(x)} \right\rangle_{P_B},$$

上式表明 $Z_A/Z_B$ 最终等同于函数 $f(x)/g(x)$ 在引入的辅助分布 $P_B(x)$ 上的均值. 由于辅助分布上采样较容易, 这就绕过了传统MCMC可能面临的多模式问题. 不过, 如果两个分布的差别很大, 这个估计值的偏差就会很高, 导致估算的结果很不准确. AIS的想法是, 在两个分布中间进一步引入大量的中间分布, 使得相邻的两个分布十分相似, 这就克服了分布差别过大时造成的高偏差问题.

在评估RBM时, 我们可以引入一个非常简单的RBM, 使其归一化常数可以直接计算出来. 然后, 利用AIS, 估算两个RBM的归一化常数之比, 最后将这个比例乘上简单RBM的归一化常数, 即得到被评估RBM的归一化常数. 进而, RBM对训练数据的似然度即可顺利算出.

## 6 基本RBM模型的变形算法

自RBM的基本模型提出以来, 尤其是Hinton提出基于CD的快速学习算法之后, 研究者们针对RBM已发展了很多变形算法(如稀疏RBM, 稀疏组RBM, 分类RBM, 条件RBM等)<sup>[8,12,13,17,33,38-41]</sup>. 由于文章篇幅限制, 本节将对几种具有代表性的算法作一简单介绍.

### 稀疏受限波尔兹曼机 (sparse restricted Boltzmann machine, SRBM)

由于稀疏表示 (sparse representation) 模型符合生物视觉系统特性, 且能够提取图像的高级特征, 近年来在机器学习、图像处理、压缩感知等研究领域都得到了广泛关注. 一般而言, RBM 学习到的特征表示是分布式、非稀疏的. 在实际应用中, 隐单元只在少数情况下处于激活状态更容易解释 (相应隐单元仅被用来表示很小一部分训练数据), 且判别性能在某些情况下还会得到改进. Lee<sup>[38]</sup> 在对数似然函数的基础上, 引入了一个稀疏惩罚项, 以惩罚隐单元的平均激活概率偏离给定水平  $p$  所引起的损失, 提出了一种稀疏受限波尔兹曼机 (sparse restricted Boltzmann machines, SRBM). 给定训练数据  $\mathbf{v}^{(1)}, \mathbf{v}^{(2)}, \dots, \mathbf{v}^{(T)}$ , SRBM 的目标函数为

$$\underset{\{\omega_{ij}, a_i, b_j\}}{\text{minimize}} - \sum_{t=1}^T \log \sum_{\mathbf{h}} P(\mathbf{v}^{(t)}, \mathbf{h}^{(t)}) + \lambda \sum_{j=1}^m \left| p - \frac{1}{T} \sum_{t=1}^T \mathbb{E}[h_j^{(t)} | \mathbf{v}^{(t)}] \right|^2,$$

这里  $\mathbb{E}[\cdot]$  表示数据已知时的条件期望,  $\lambda$  是正则化系数,  $p$  是控制隐单元稀疏度的常数 (需事先指定). 在学习过程中, 可先基于对比散度的学习算法给出对数似然函数的梯度近似, 再利用正则化项的梯度进行梯度下降, 直至算法收敛. 在 MNIST 手写体数据集和自然图像上的试验结果表明, 稀疏 RBM 可以提取手写体的笔划特征及自然图像中类似于 Gabor 滤波的特征, 这与人脑 V1 区简单细胞感受野 (receptive fields) 十分相似. 更重要地, 堆叠两个稀疏 RBM 可以提取更抽象的特征. 文献 [38] 中试验结果表明, 对于自然图像, 两个稀疏 RBM 可以提取轮廓 (contours)、拐角 (angles) 以及边缘合并 (junctions of edges) 等特征, 这些特征与人脑 V2 区细胞的感受野十分相似. 较之前的稀疏表示方法, 堆叠稀疏 RBM 不但可以提取类似于 V1 区简单细胞的感受野, 而且能够提取类似于 V2 区细胞的感受野, 这促使机器学习研究向人工智能的目标迈进了一大步.

### 稀疏组受限波尔兹曼机 (sparse group restricted Boltzmann machine, SGRBM)

在实际问题中, 特征之间往往显示很强的统计相关性, 并且很多特征经常成组地出现 (或成组地不出现). 由于直接学习所有隐单元之间的统计相关性是困难的, 尤其是在高维问题建模时. 为简化该问题, 罗恒<sup>[13]</sup> 将组稀疏 (sparse group) 方法应用到 RBM 中, 提出了稀疏组受限波尔兹曼机. 其主要思想如下: 首先, 将隐单元划分到不重叠的组中, 只考虑组内隐单元状态的相关性; 其次, 不是去“学习”组内隐单元状态的相关性, 而是通过正则化方法惩罚组内隐单元的总激活程度, 从而使组内隐单元在学习过程中不再是条件独立的. 具体地, 给定训练数据, SGRBM 在似然函数中引入了一个惩罚项, 该项是关于隐单元激活概率的混合范数 ( $L_1/L_2$  范数). 假设一个 RBM 中有  $m$  个隐单元,  $\mathcal{H}$  表示隐单元下标组成的指标集, 即  $\mathcal{H} = \{1, 2, \dots, m\}$ . 将这些隐单元分成  $K$  组, 假设所有的组大小相同, 且互不重叠, 记第  $k$  组的指标集为  $\mathcal{G}_k$ ,  $\mathcal{G}_k \subset \mathcal{H}$ ,  $k = 1, 2, \dots, K$ . 给定分组  $\mathcal{G} = \{\mathcal{G}_1, \mathcal{G}_2, \dots, \mathcal{G}_K\}$  和观测数据  $\mathbf{v}^{(l)}$ , 第  $k$  组隐单元激活概率的  $L_2$  范数定义为

$$N_k(\mathbf{v}^{(l)}) = \left( \sum_{m \in \mathcal{G}_k} P(h_m = 1 | \mathbf{v}^{(l)})^2 \right)^{\frac{1}{2}},$$

该范数可解释为第  $k$  组隐单元在给定观测数据  $\mathbf{v}^{(l)}$  时的总体激活程度. 给定所有隐单元组的范数, 相应的  $L_1/L_2$  范数定义为

$$\sum_{k=1}^K |N_k(\mathbf{v}^{(l)})| = \sum_{k=1}^K \left( \sum_{m \in \mathcal{G}_k} P(h_m = 1 | \mathbf{v}^{(l)})^2 \right)^{\frac{1}{2}}.$$

于是, SGRBM 的目标函数

$$\underset{\{\omega_{ij}, a_i, b_j\}}{\text{minimize}} - \sum_{t=1}^T \log \sum_{\mathbf{h}} P(\mathbf{v}^{(t)}) + \lambda \sum_{k=1}^K |N_k(\mathbf{v}^{(l)})|,$$

其中  $\lambda$  是正则化常数. 为了学习得到模型参数, 罗恒建议采用如下迭代算法: 基于给定训练数据, 首先应用对比散度更新模型参数一次, 再使用正则化项的梯度更新参数一次, 直至算法收敛.

$L_1/L_2$  正则化的作用可以从组间和组内两个层面来解释. 在组间 ( $L_1$  范数), 为了最小化混合范数, 给定训练数据, 正则化项会鼓励隐单元的激活概率形成一种组稀疏的表示, 也就是大量隐单元组的  $L_2$  范数为 0. 由于隐单元的激活概率是非负的, 隐单元组的  $L_2$  范数为 0 便意味着组内所有隐单元的激活概率均为 0. 在组内 ( $L_2$  范数), 通常认为各分量均会受到同等程度的惩罚, 故而不会产生组内的稀疏表示. 但是由于 RBM 隐单元的激活概率的函数形式, 仍旧会产生一种组内的稀疏表示. 与 SRBM 相比, SGRBM 可以学习到更局部化的特征, 且可以达到更高的识别率.

#### 分类受限波尔兹曼机 (classification restricted Boltzmann machine, ClassRBM)

当使用 RBM 解决分类任务时, 最常见的做法是将 RBM 视为一个特征提取器 (feature detector): 使用观测数据 (忽略类标签) 训练 RBM, 然后以原训练数据在训练好的 RBM 的隐单元激活概率以及原有的类标签组成新的训练集, 进而使用其他常用的分类算法训练分类器. 由于 RBM 是采用无监督学习的方式训练的, 学习到的特征并不完全适合分类任务. Larochelle 等人<sup>[39,40]</sup> 指出 RBM 可直接用于解决有监督学习任务, 并提出了分类受限波尔兹曼机, 其主要思想是利用包含二值随机变量的隐单元来拟合输入特征与类标签的联合分布 (即把输入特征和类标签均作为可见单元进行学习). Larochelle 等人提出, ClassRBM 可以使用三种不同的训练目标进行学习, 其中的参数仍然可以基于对比散度进行训练. ClassRBM 使得分类过程得以简化 (不需要再训练另外的分类器), 保证了学习到的特征的判别能力, 并且可以以在线学习的方式进行训练, 可实时监测其学习到的特征表示的判别性能. 传统 RBM 的另一个用处是初始化深层神经网络, 此种分类方法需要分两个阶段进行训练, ClassRBM 与其相比, 省去了第二个训练阶段.

## 7 总结与展望

对比散度较好地解决了 RBM 的学习效率问题, 使得近些年 RBM 在许多领域都得到了广泛研究和应用. 本文对 RBM 的基本模型、基于对比散度的快速学习算法、参数设置、评估方法及其具有代表性的几种变形算法作了较为详细的介绍. RBM 为人们解决智能问题提供了一种强有力的工具, 并为其他领域的研究提供了新技术和新思路, 研究前景

十分广阔。尤其是随着深度神经网络的兴起, 借助RBM来学习深层网络逐渐成为深度学习研究中的主流, 也使得RBM在深度学习领域中逐渐占据核心地位。然而, 在RBM及其相关理论和学习算法的研究中, 仍有许多问题值得我们作进一步探讨。譬如, 如何提高RBM在无监督学习场景下所提取特征的辨别能力? 在不增加隐单元个数的情况下, 只利用RBM能量函数的非参数化形式能否提高其逼近性能? RBM能否用于图像分割、高维数据的聚类、缺失数据的重构等更广泛的实际应用? 这些问题的研究和探讨都将具有十分重要的理论和实际意义。

## 参考文献:

- [1] 叶世伟, 史忠植. 神经网络原理[M]. 北京: 机械工业出版社, 2006  
Ye S W, Shi Z Z. Neural Networks[M]. Beijing: China Machine Press, 2006
- [2] Haykin S. Neural Networks and Learning Machines (3rd Edition)[M]. New Jersey: Pearson Education, 2009
- [3] Hinton G E, Sejnowski T J. Learning and relearning in Boltzmann machines[C]// Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Cambridge, USA, 1986
- [4] Smolensky P. Information processing in dynamical systems: foundations of harmony theory[C]// Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Cambridge, USA, 1986
- [5] Freund Y, Haussler D. Unsupervised learning of distributions on binary vectors using two layer networks[R]. Santa Cruz: University of California, UCSC-CRL-94-25, 1994
- [6] Roux N L, Bengio Y. Representational power of restricted Boltzmann machines and deep belief networks[J]. Neural Computation, 2008, 20(6): 1631-1649
- [7] Hinton G E. Training products of experts by minimizing contrastive divergence[J]. Neural Computation, 2002, 14(8): 1771-1800
- [8] Cho K Y. Improved learning algorithms for restricted Boltzmann machines[D]. Espoo: Aalto University, 2011
- [9] Teh Y W, Hinton G E. Rate-coded restricted Boltzmann machines for face recognition[C]// Advances in Neural Information Processing Systems 13, MIT Press, 2001: 908-914
- [10] Salakhutdinov R, Mnih A, Hinton G E. Restricted Boltzmann machines for collaborative filtering[C]// Proceedings of the 24th International Conference on Machine Learning, Corvallis, OR, 2007: 791-798
- [11] 吴证, 周越, 杜春华, 等. 组合主成分分析的受限波尔兹曼机神经网络的降维方法[J]. 上海交通大学学报, 2008, 42(4): 559-563  
Wu Z, Zhou Y, Du C H, et al. Dimensionality reduction method based on restricted Boltzmann machine neural network with principal component analysis[J]. Journal of Shanghai Jiaotong University, 2008, 42(4): 559-563
- [12] 吴金龙. Netflix Prize 中的协同过滤算法[D]. 北京: 北京大学, 2010  
Wu J L. Collaborative filtering on the Netflix Prize dataset[D]. Beijing: Peking University, 2010
- [13] 罗恒. 基于协同过滤视角的受限波尔兹曼机研究[D]. 上海: 上海交通大学, 2011  
Luo H. Restricted boltzmann machines: a collaborative filtering perspective[D]. Shanghai: Shanghai Jiaotong University, 2011
- [14] 潘闻特, 申丽萍. 基于BM神经网络编码的生理信号情感识别[J]. 计算机工程与设计, 2012, 33(3): 1101-1106  
Pan W T, Shen L P. Emotion detection with BM network encoded physiological signals[J]. Computer Engineering and Design, 2012, 33(3): 1101-1106

- [15] Hinton G E, Osindero S, Teh Y W. A fast learning algorithm for deep belief nets[J]. *Neural Computation*, 2006, 18(7): 1527-1554
- [16] Hinton G E, Salakhutdinov R R. Reducing the dimensionality of data with neural networks[J]. *Science*, 2006, 313(5786): 504-507
- [17] Bengio Y. Learning deep architectures for AI[J]. *Foundations and Trends in Machine Learning*, 2009, 2(1): 1-127
- [18] Arel I, Rose D C, Karnowski T P. Deep machine learning—a new frontier in artificial intelligence research[J]. *IEEE Computational Intelligence Magazine*, 2010, 5(4): 13-18
- [19] Welling M, Rosen-Zvi M, Hinton G E. Exponential family harmoniums with an application to information retrieval[C]// *Advances in Neural Information Processing Systems 17*, Cambridge, USA, 2005: 1481-1488
- [20] Liu J S. *Monto Carlo Strategies in Scientific Computing*[M]. New York: Springer-Verlag, 2001
- [21] Chen H, Murray A F. Continuous restricted Boltzmann machine with an implementable training algorithm[J]. *IEE Proceedings Vision, Image and Signal Processing*, 2003, 150(3): 153-158
- [22] Nair V, Hinton G E. Rectified linear units improve restricted Boltzmann Machines[C]// *Proceedings of the 27th International Conference on Machine Learning*, Haifa, Israel, 2010: 807-814
- [23] Hinton G E. A practical guide to training restricted Boltzmann machines[R]. Montreal: Department of Computer Science, University of Toronto, 2010
- [24] Courville A, Bergstra J, Bengio Y. A spike and slab restricted Boltzmann machine[J]. *Journal of Machine Learning Research-Proceedings Track*, 2011, 15: 233-241
- [25] Tran T, Phung D Q, Venkatesh S. Mixed-variate restricted Boltzmann machines[J]. *Journal of Machine Learning Research-Proceedings Track*, 2011, 20: 213-229
- [26] Tieleman T. Training restricted boltzmann machines using approximations to the likelihood gradient[C]// *Proceedings of the 25th International Conference on Machine Learning*, Helsinki, Finland, 2008: 1064-1071
- [27] Tieleman T, Hinton G E. Using fast weights to improve persistent contrastive divergence[C]// *Proceedings of the 26th International Conference on Machine Learning*, Montreal, Canada, 2009: 1033-1040
- [28] Desjardins G, Courville A, Bengio Y, et al. Parallel tempering for training of restricted Boltzmann machines[C]// *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, New York, 2010: 145-152
- [29] Brügge K, Fischer A, Igel C. The flip-the-state transition operator for restricted Boltzmann machines[J]. *Machine Learning*, 2013, 93(1): 53-69
- [30] Fischer A, Igel C. Training restricted Boltzmann machines: an introduction[J]. *Pattern Recognition*, 2013, 47(1): 25-39
- [31] 刘建伟, 刘媛, 罗雄麟. 玻尔兹曼机研究进展[J]. *计算机研究与发展*, 2014, 51(1): 1-16  
Liu J W, Liu Y, Luo X L. Research and development of Boltzmann machine[J]. *Journal of Computer Research and Development*, 2014, 51(1): 1-16
- [32] Cho K H, Raiko T, Ilin A. Enhanced gradient for training restricted Boltzmann machines[J]. *Neural Computation*, 2013, 25(3): 805-831
- [33] Bengio Y, Courville A, Bincent P. Unsupervised feature learning and deep learning: a review and new perspectives[R]. Montreal: Department of Computer Science and Operations Research, University of Montreal, 2012
- [34] Schulz H, Müller A, Behnke S. Investigating convergence of restricted Boltzmann machine learning[C]// *NIPS 2010 Workshop on Deep Learning and Unsupervised Feature Learning*, Whistler, Canada, 2010: 1-9
- [35] Fischer A, Igel C. Empirical analysis of the divergence of Gibbs sampling based learning algorithms for restricted Boltzmann machines[C]// *Proceedings of the 20th International Conference on Artificial Neural*

Networks, Part III, Berlin, Springer-Verlag, 2010: 208-217

- [36] Gengio Y. Practical recommendations for gradient-based training of deep architectures[R]. Technical Report, Department of Computer Science and Operations Research, University of Montreal, 2012
- [37] Neal R M. Annealed importance sampling[J]. Statistics and Computing, 2001, 11(2): 125-139
- [38] Lee H, Ekanadham C, Ng A Y. Sparse deep belief net model for visual area V2[C]// Advances in Neural Information Processing Systems 20, Vancouver, Canada, 2008: 873-880
- [39] Larochelle H, Bengio Y. Classification using discriminative restricted Boltzmann machines[C]// Proceedings of the 25th International Conference on Machine Learning, Helsinki, Finland, 2008: 536-543
- [40] Larochelle H, et al. Learning algorithms for the classification restricted Boltzmann machine[J]. Journal of Machine Learning Research, 2012, 13: 643-669
- [41] Ji N N, et al. Discriminative restricted Boltzmann machine for invariant pattern recognition with linear transformation[J]. Pattern Recognition Letters, 2014, 45: 172-180

## Restricted Boltzmann Machines

ZHANG Chun-xia<sup>1</sup>, JI Nan-nan<sup>2</sup>, WANG Guan-wei<sup>3</sup>

(1- School of Mathematics and Statistics, Xi'an Jiaotong University, Xi'an 710049;

2- School of Science, Chang'an University, Xi'an 710064;

3- School of Mechatronic Engineering, Xi'an Technological University, Xi'an 710021)

**Abstract:** A restricted Boltzmann machine (RBM) is a particular type of random neural network model which has two-layer architecture, symmetric connections and no self-feedback. The two layers in an RBM are fully connected but there are no connections within the same layer. Recently, with the advent of a fast learning algorithm for RBMs (i.e., contrastive divergence), the machine learning community set off a surge to study the theory and applications of RBMs since it has many advantages. For example, a RBM provides us an effective tool to detect features. When a feed-forward neural network is initialized with an RBM, its generalization capability can be significantly improved. A deep belief network composed of several RBMs can detect more abstract features. Due to the advantages and wide applications of RBMs in deep learning, this paper attempts to provide a introductory guide for novice. It presents a detailed introduction of basic RBM model, its representative learning algorithm, parametric settings, evaluation methods, its variants and etc. Finally, some research directions of RBMs that are deserved to be further studied are discussed.

**Keywords:** machine learning; deep learning; restricted Boltzmann machine; contrastive divergence; Gibbs sampling

---

**Received:** 22 Aug 2013. **Accepted:** 19 May 2014.

**Foundation item:** The National Basic Research Program of China, 973 Program (2013CB329406); the Major Research Project of the National Natural Science Foundation of China (91230101); the National Natural Science Foundation of China (11201367); the Fundamental Research Funds for the Central Universities of China (xjj2011048).