

Restricted Boltzmann Machines (RBM's)

Viktor Kuhn, Maximilian Lutz

August 8, 2019

Why do we need this

We want to understand high-dimensional data with rich structure.

"Understanding" → Density estimation (true $p(\mathbf{x})$ under data)

- generative models
- Missing value imputation (given partial \mathbf{x})
- Denoising (return original \mathbf{x} given damaged $\tilde{\mathbf{x}}$)
- ...

→ Density estimation (true $p(\mathbf{x})$ under data)

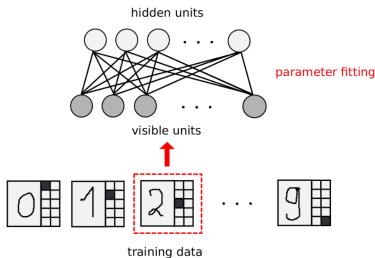
- generative models



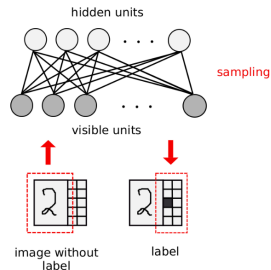
→ Density estimation (true $p(\mathbf{x})$ under data)

- Missing value imputation (given partial \mathbf{x})
→ classification is actually special case!

learning with labels



classification



→ Density estimation (true $p(\mathbf{x})$ under data)

- Denoising (return original \mathbf{x} given damaged $\tilde{\mathbf{x}}$)



Why this is hard

This requires modelling the underlying probability distribution.

→ fundamentally problematic

Naive approach to represent probability distribution has space complexity $\mathcal{O}(n^k)$ (where n is the cardinality of the state space of the probability distribution, k is the number of random variables)

Idea: use conditional structure of probability distribution to reduce complexity

→ motivation for so called **graphical models**

Graphical models

$G = (V, E)$, V nodes, E edges, defines an **undirected Graph**.

clique is a subset of V in which all nodes are pairwise connected. Call a clique **maximal**, if no node can be added without violating the definition of a clique.

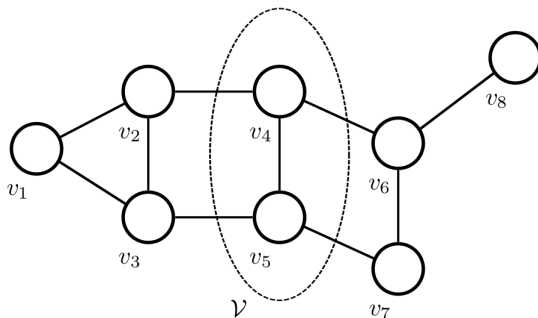


Figure 1: Undirected graph, $\{v_4, v_5\}$ clique, also maximal [6]

Graphical Models

- each node $v \in V$ corresponds to random variable X_v
- edges encode direct probabilistic interaction

Markov random fields

The set $\mathbf{X} = (X_v)_{v \in V}$ is called a **Markov random field (MRF)** iff

$$\forall \mathbf{x} : \forall v \in V : p(x_v | (x_w)_{w \in V \setminus \{v\}}) = p(x_v | (x_w)_{w \in \mathcal{N}_v}) \quad (1)$$

$$\mathcal{N}_v = \{w \in V : \{w, v\} \in E\} \text{ neighborhood of } v$$

In other words, in a Markov random field all X_v are conditionally independent of all other variables given their neighborhood.

Markov random fields

Theorem (Hammersley-Clifford) [15]

For a Markov random field the following equivalence holds

- 1) the distribution p is strictly positive
- 2) p factorizes over the Graph G , i.e.

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}) \quad (2)$$

where the product runs over all maximal cliques.

The $\psi_C(\mathbf{x})$ are often referred to as **clique potentials**.

Energy based models and Gibbs distribution

Rewrite (2):

$$p(\mathbf{x}) = \frac{1}{Z} \prod_{C \in \mathcal{C}} \psi_C(\mathbf{x}) = \frac{1}{Z} e^{\sum_{C \in \mathcal{C}} \ln \psi_C(\mathbf{x})} = \frac{1}{Z} e^{-E(\mathbf{x})} \quad (3)$$

→ **Energy based model (EBM)** (alt.: harmony based model)

underlying distribution: **Boltzmann distribution / Gibbs distribution**

→ + latent variables = **Boltzmann machine**

Problem

general MRFs virtually **untrainable**

Problem

general MRFs virtually **untrainable**

→ in the most general case (fully connected) no benefit for graphical model!

→ hard to do something clever without relying on some structure

Restricted Boltzmann machines (Harmonium)

Constraining structure → **bipartite graph**

one layer of visible & one layer of hidden/latent variables

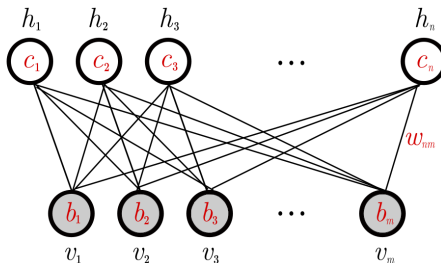
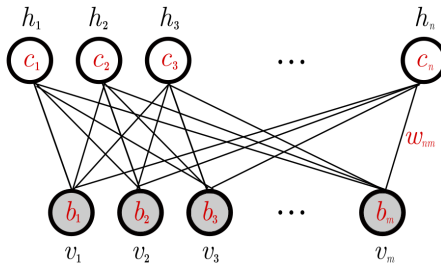
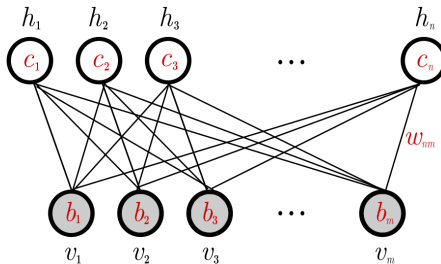


Figure 2: Restricted Boltzmann machine



binary random variables $(\mathbf{V}, \mathbf{H}) \rightarrow (\mathbf{v}, \mathbf{h}) \in \{0, 1\}^{m+n}$

$$E(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{j=1}^m b_j v_j - \sum_{i=1}^n c_i h_i \quad (4)$$



takeaway: RBMs are Boltzmann machines with structure & energy constraint

Properties of RBMs

no connections inside each layer

$$\begin{aligned}
 p(\mathbf{h}|\mathbf{v}) &= \prod_{i=1}^n p(h_i|\mathbf{v}) \\
 p(\mathbf{v}|\mathbf{h}) &= \prod_{j=1}^m p(v_j|\mathbf{h})
 \end{aligned}
 \tag{5}$$

Properties of RBMs

universal approximators for distributions on $\{0, 1\}^m$

→ additionally quite a few other results about universal approximator properties [11] [13]

Sampling from RBMs

How to access what we learn

Calculating partition function of RBM's intractable

→ **Markov chain Monte Carlo** (MCMC) methods provide ways to sample from $p(\mathbf{v}, \mathbf{h})$ without knowing Z

Markov Chains

A **Markov chain** is a family of random variables $X = \{X^k | k \in \mathbb{N}_0\}$ that, $\forall k \geq 0$ and $\forall j, i, i_0, \dots, i_{k-1} \in \Omega$ satisfies

$$\Pr(X^{(k+1)} = j | X^{(k)} = i, \dots, X^{(0)} = i_0) = \Pr(X^{(k+1)} = j | X^{(k)} = i) \quad (6)$$

Markov Chains

Transition matrix $\mathbf{P} = (\Pr(X^{(k+1)} = j | X^{(k)} = i))_{i,j \in \Omega} = (p_{ij})_{i,j \in \Omega}$

stationary distribution π for which $\pi^T = \pi^T \mathbf{P}$

Detailed balance condition $\pi(i)p_{ij} = \pi(j)p_{ji}$ sufficient for stationarity of π

Markov Chains

irreducible (can get from any state in Ω to any other in finite number of transitions)

$$\forall s_i, s_j \in \Omega : \exists n \in \mathbb{N} : P(X_n = s_i | X_0 = s_j) > 0$$

aperiodic (every state can reoccur at irregular times)

$$\gcd \{n > 0 : \Pr(X_n = i | X_0 = i) > 0\} = 1$$

\Rightarrow Markov chain over finite Ω converges to its stationary distribution

Gibbs sampling

Idea

Construct Markov chain by updating each variable based on its conditional distribution given state of the others

Gibbs sampling - Metropolis-Hastings

MRF $\mathbf{X} = (X_1, \dots, X_N)$ w.r.t. $G = (V, E)$ and $p(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})}$, then produce new state as follows:

1. Pick $X_i, i \in V$ at random with strictly positive $q(i)$
2. Transition between \mathbf{x} and \mathbf{y} , $\mathbf{x} \neq \mathbf{y}$ with $p_{\mathbf{x}\mathbf{y}} =$

$$\begin{cases} q(i)p(y_i | (x_v)_{v \in \mathcal{N}_i}), & \text{if } \exists i \in V \text{ so that } \forall v \in V \text{ with } v \neq i : x_v = y_v \\ 0, & \text{else} \end{cases}$$
 and stay in \mathbf{x} with $p_{\mathbf{x}\mathbf{x}} = \sum_{i \in V} q(i)p(x_i | (x_v)_{v \in \mathcal{N}_i})$

Gibbs sampling

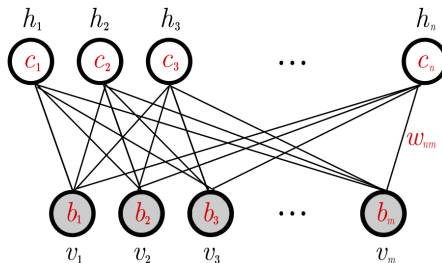
$p(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})}$ and conditional probabilities strictly positive

→ Chain irreducible and aperiodic

→ Chain converges to stationary distribution

→ $p(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})}$ stationary distribution (detailed balance)

k-step Block Gibbs sampling in RBM's



no intralayer connections in RBM's

→ sample states of all variables in one layer jointly

Algorithm 1 Block Gibbs sampling in RBM's

```

1: Input RBM ( $V_1, \dots, V_m, H_1, \dots, H_n$ ), Gibbs steps  $k$ 
2: Output sample ( $\mathbf{v}^{(k)}, \mathbf{h}^{(k)}$ ) after  $k$  Gibbs steps
3:  $\mathbf{v}^{(0)} \leftarrow \text{random\_vector}$ 
4: for  $t = 0, \dots, k - 1$  do
5:   for  $i = 1..n$  do
6:     Sample  $h_i^{(t)} \sim p(h_i | \mathbf{v}^{(t)})$ 
7:   end for
8:   for  $j = 1..m$  do
9:     Sample  $v_j^{(t+1)} \sim p(v_j | \mathbf{h}^{(t)})$ 
10:  end for
11: end for

```

Training RBMs

What to minimize

Gibbs sampling helps once we have parameters

But how to learn them?

→ minimize distance between (unknown) data distribution q and
RBM-distribution p in terms of **Kullback-Leibler divergence**

KL and log-likelihood

Minimizing

$$\text{KL}(q\|p) = \sum_{\mathbf{x} \in \Omega} q(\mathbf{x}) \ln \frac{q(\mathbf{x})}{p(\mathbf{x})} = \sum_{\mathbf{x} \in \Omega} q(\mathbf{x}) \ln q(\mathbf{x}) - \sum_{\mathbf{x} \in \Omega} q(\mathbf{x}) \ln p(\mathbf{x}) \quad (7)$$

\iff

Maximizing **log-likelihood**

$$\ln \mathcal{L}(\theta|S) = \ln \prod_{i=1}^{\ell} p(\mathbf{x}_i|\theta) = \sum_{i=1}^{\ell} \ln p(\mathbf{x}_i|\theta) \quad (8)$$

over training data $S = \{\mathbf{x}_1, \dots, \mathbf{x}_{\ell}\}$ w.r.t. parameters θ

Unsupervised RBM learning

→ Learning = gradient ascent on log-likelihood with learning rate

$\eta \in \mathbb{R}^+$,

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \eta \frac{\partial}{\partial \boldsymbol{\theta}^{(t)}} \left(\ln \mathcal{L} \left(\boldsymbol{\theta}^{(t)} | \mathcal{S} \right) \right) \quad (9)$$

Unsupervised RBM learning

For single training example \mathbf{v} ,

$$\ln \mathcal{L}(\theta|\mathbf{v}) = \ln \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} - \ln \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

$$\Rightarrow \frac{\partial \ln \mathcal{L}(\theta|\mathbf{v})}{\partial \theta} = \frac{\partial}{\partial \theta} \left(\ln \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \right) - \frac{\partial}{\partial \theta} \left(\ln \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \right)$$

$$= -\frac{1}{\sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}} \sum_{\mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} + \frac{1}{\sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}} \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta}$$

$$\frac{\partial \ln \mathcal{L}(\theta|\mathbf{v})}{\partial \theta} = - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta} \quad (10)$$

The first term in (10) is referred to as **positive phase** the second as **negative phase**.

Positive and negative phase

Positive phase

$$- \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta}$$

→ finds hidden configurations that work well with \mathbf{v} and lowers their energies

Negative Phase

$$\sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta}$$

→ finds joint configuration that are the "best competitors" and raises their energies

Positive phase in RBM's

→ generally very easy for RBM's exploiting independence property (5):

$$p(\mathbf{h}|\mathbf{v}) = \prod_{i=1}^n p(h_i|\mathbf{v})$$

$$p(\mathbf{v}|\mathbf{h}) = \prod_{j=1}^m p(v_j|\mathbf{h})$$

For example, positive Phase w.r.t. w_{ij} (and similarly for b_i , c_j):

Let \mathbf{h}_{-i} denote state of all hidden variables but the i -th one, then

$$\begin{aligned}
 - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} &= \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) h_i v_j = \sum_{\mathbf{h}} \prod_{k=1}^n p(h_k|\mathbf{v}) h_i v_j \\
 &= \sum_{h_i} \sum_{\mathbf{h}_{-i}} p(h_i|\mathbf{v}) p(\mathbf{h}_{-i}|\mathbf{v}) h_i v_j = \sum_{h_i} p(h_i|\mathbf{v}) h_i v_j \sum_{\mathbf{h}_{-i}} p(\mathbf{h}_{-i}|\mathbf{v}) \\
 &= p(H_i = 1|\mathbf{v}) v_j
 \end{aligned} \tag{11}$$

→ saves us from exponential sum in number of hidden variables

Negative phase in RBM's

We're not so lucky here!

→ same trick gets rid only of either $\sum_{\mathbf{v}}$ or $\sum_{\mathbf{h}}$ in $\sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}}$

For example, getting rid of $\sum_{\mathbf{h}}$ leaves:

$$\begin{aligned} \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} &= - \sum_{\mathbf{v}} p(\mathbf{v}) \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) h_i v_j \\ &= - \sum_{\mathbf{v}} p(\mathbf{v}) p(H_i = 1|\mathbf{v}) v_j \end{aligned} \quad (12)$$

→ sum remains exponential either in hidden or visible variables

$$\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta}|\mathbf{v})}{\partial w_{ij}} = - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}}$$

$$\begin{aligned}
 \frac{\partial \ln \mathcal{L}(\boldsymbol{\theta}|\mathbf{v})}{\partial w_{ij}} &= - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\
 &= p(H_i = 1|\mathbf{v}) v_j + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial \ln \mathcal{L}(\boldsymbol{\theta}|\mathbf{v})}{\partial w_{ij}} &= - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\
 &= p(H_i = 1|\mathbf{v}) v_j + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\
 &= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) h_i v_j
 \end{aligned}$$

$$\begin{aligned}
\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta}|\mathbf{v})}{\partial w_{ij}} &= - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\
&= p(H_i = 1|\mathbf{v}) v_j + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\
&= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) h_i v_j \\
&= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}} p(\mathbf{v}) p(H_i = 1|\mathbf{v}) v_j
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta}|\mathbf{v})}{\partial w_{ij}} &= - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\
&= p(H_i = 1|\mathbf{v}) v_j + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\
&= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) h_i v_j \\
&= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}} p(\mathbf{v}) p(H_i = 1|\mathbf{v}) v_j \\
&= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) h_i v_j
\end{aligned}$$

Training algorithms

For gradient ascent, need a way to compute second term:

$$\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta}|\mathbf{v})}{\partial w_{ij}} = p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) h_i v_j$$

→ approximate this sum by MCMC methods

$$\sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) h_i v_j$$

→ think of sum as expectation value

→ Gibbs sampling → approximate as average over a few samples

Algorithm 2 Naive MCMC gradient estimation

```

1: Input  $\text{RBM}(V_1, \dots, V_m, H_1, \dots, H_n)$ , training batch  $S$ 
2: Output  $\Delta w_{ij}$  for  $i = 1, \dots, n, j = 1, \dots, m$ 
3: init:  $\Delta w_{ij} = 0$  for  $i = 1, \dots, n, j = 1, \dots, m$ 
4: for all  $\mathbf{v} \in S$  do
5:   for  $l=1..\xi$  do
6:      $\mathbf{v}^{(0)} \leftarrow \text{random\_vector}$ 
7:     for  $t = 0, \dots, k-1$  do
8:       for  $i = 1..n$  do
9:         Sample  $h_i^{(t)} \sim p(h_i | \mathbf{v}^{(t)})$ 
10:      end for
11:      for  $j = 1..m$  do
12:        Sample  $v_j^{(t+1)} \sim p(v_j | \mathbf{h}^{(t)})$ 
13:      end for
14:    end for
15:  end for ▷ here  $\xi$   $(\mathbf{v}, \mathbf{h})$  pairs (approximately from  $p(\mathbf{v}, \mathbf{h})$ ) → call these  $(\mathbf{v}^l, \mathbf{h}^l)$ 
16:  for  $i = 1, \dots, n, j = 1, \dots, m$  do
17:     $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1 | \mathbf{v}) \cdot v_j - \frac{1}{\xi} \sum_{l=1..\xi} h_i^l v_j^l$ 
18:  end for
19: end for

```

Problem

requires running a lot of Gibbs chains (computationally expensive)

→ take a step back, tackle the exponential sum again

$$\begin{aligned}\frac{\partial \ln \mathcal{L}(\boldsymbol{\theta}|\mathbf{v})}{\partial w_{ij}} &= - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} \\ &= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}} p(\mathbf{v}) p(H_i = 1|\mathbf{v}) v_j\end{aligned}$$

Contrastive divergence $CD - k$

- initialize Gibbs chain by data distribution \rightarrow less time to burn in
- approximate the expectation value using a single sampled $\mathbf{v}^{(k)}$ (that is after k Gibbs steps)

Contrastive divergence $CD - k$

- initialize Gibbs chain by data distribution \rightarrow less time to burn in
- approximate the expectation value using a single sampled $\mathbf{v}^{(k)}$ (that is after k Gibbs steps)

$$\begin{aligned} \frac{\partial \ln \mathcal{L}(\theta|\mathbf{v})}{\partial w_{ij}} &= p(H_i = 1|\mathbf{v}) v_j - \sum_{\mathbf{v}} p(\mathbf{v}) p(H_i = 1|\mathbf{v}) v_j \\ &\approx p(H_i = 1|\mathbf{v}) v_j - p(H_i = 1|\mathbf{v}^{(k)}) v_j^{(k)} \end{aligned}$$

Algorithm 3 k-step contrastive divergence [1]

```

1: Input  $\text{RBM}(V_1, \dots, V_m, H_1, \dots, H_n)$ , training batch  $S$ 
2: Output  $\Delta w_{ij}, \Delta b_j$  and  $\Delta c_i$  for  $i = 1, \dots, n, j = 1, \dots, m$ 
3: init:  $\Delta w_{ij} = \Delta b_j = \Delta c_i = 0$  for  $i = 1, \dots, n, j = 1, \dots, m$ 
4: for all  $\mathbf{v} \in S$  do
5:    $\mathbf{v}^{(0)} \leftarrow \mathbf{v}$ 
6:   for  $t = 0, \dots, k - 1$  do
7:     for  $i = 1..n$  do
8:       Sample  $h_i^{(t)} \sim p(h_i | \mathbf{v}^{(t)})$ 
9:     end for
10:    for  $j = 1..m$  do
11:      Sample  $v_j^{(t+1)} \sim p(v_j | \mathbf{h}^{(t)})$ 
12:    end for
13:  end for
14:  for  $i = 1, \dots, n, j = 1, \dots, m$  do
15:     $\Delta w_{ij} \leftarrow \Delta w_{ij} + p(H_i = 1 | \mathbf{v}^{(0)}) \cdot v_j^{(0)} - p(H_i = 1 | \mathbf{v}^{(k)}) \cdot v_j^{(k)}$ 
16:  end for
17:  for  $j = 1..m$  do
18:     $\Delta b_j \leftarrow \Delta b_j + v_j^{(0)} - v_j^{(k)}$ 
19:  end for
20:  for  $j = 1..n$  do
21:     $\Delta c_i \leftarrow \Delta c_i + p(H_i = 1 | \mathbf{v}^{(0)}) - p(H_i = 1 | \mathbf{v}^{(k)})$ 
22:  end for
23: end for

```

Why does this even work?

- non-trivial and surprising when first introduced
- does not follow gradient of log-likelihood
 - important results give conditions for convergence and bounds on the expectation error, empirical observations support this further ([3, 4, 5, 17])

Why does this even work?

- non-trivial and surprising when first introduced
- does not follow gradient of log-likelihood
 - important results give conditions for convergence and bounds on the expectation error, empirical observations support this further ([3, 4, 5, 17])
- even more surprising: in practice, very often $k = 1$ [9]

Persistent contrastive divergence (PCD)

do not reinitialize Gibbs chains, but keep them permanently

→ update step is small enough so distribution keeps close to previous step

Where to go from here?

- better training
 - fast persistent contrastive divergence (FPCD)
 - parallel tempering
- other models
 - Deep Belief Networks
 - Conditional Boltzmann machines
 - ...

References i



D. H. Ackley, G. E. Hinton, and T. J. Sejnowski.

A learning algorithm for boltzmann machines.

Cognitive science, 9(1):147–169, 1985.



Y. Bengio.

Learning deep architectures for ai.

Foundations and Trends in Machine Learning, 2(1):1127, 2009.



Y. Bengio and O. Delalleau.

Justifying and generalizing contrastive divergence.

Neural computation, 21(6):1601–1621, 2009.



M. A. Carreira-Perpinan and G. E. Hinton.

On contrastive divergence learning.

In *Aistats*, volume 10, pages 33–40. Citeseer, 2005.

References ii



A. Fischer and C. Igel.

Bounding the bias of contrastive divergence learning.

Neural computation, 23(3):664–673, 2011.



A. Fischer and C. Igel.

Training restricted boltzmann machines: An introduction.

Pattern Recognition, 47(1):25–39, Jan. 2014.



I. Goodfellow, Y. Bengio, and A. Courville.

Deep learning.

MIT press, 2016.



G. E. Hinton.

Training products of experts by minimizing contrastive divergence.

Neural computation, 14(8):1771–1800, 2002.

References iii



G. E. Hinton.

A practical guide to training restricted boltzmann machines.

In *Neural networks: Tricks of the trade*, pages 599–619. Springer, 2012.



G. E. Hinton, S. Osindero, and Y.-W. Teh.

A fast learning algorithm for deep belief nets.

Neural computation, 18(7):1527–1554, 2006.



D. MacKay.

Failures of the one-step learning algorithm.

In *Available electronically at <http://www.inference.phy.cam.ac.uk/mackay/abstracts/gbm.html>*, 2001.

References iv



P. Mehta, M. Bukov, C.-H. Wang, A. G. Day, C. Richardson, C. K. Fisher, and D. J. Schwab.

A high-bias, low-variance introduction to machine learning for physicists.

Physics Reports, 2019.



G. Montufar and N. Ay.

Refinements of universal approximation results for deep belief networks and restricted boltzmann machines.

Neural Computation, 23(5):1306–1319, 2011.



V. Parmar and D. M. Suri.

Design exploration of hybrid cmos-oxram deep generative architectures.

01 2018.

References v



B. Pierre.

Markov chains: Gibbs fields, Monte Carlo simulation, and queues.

Springer, 2001.



P. Smolensky.

Information processing in dynamical systems: Foundations of harmony theory.

Technical report, Colorado Univ at Boulder Dept of Computer Science, 1986.



A. L. Yuille.

The convergence of contrastive divergences.

In *Advances in neural information processing systems*, pages 1593–1600, 2005.

Error Bound CD-k

Let p denote the marginal distribution of the visible units of an RBM and let q be the empirical distribution defined by a set of samples v_1, \dots, v_l . Then an upper bound on the expectation of the error of the CD-k approximation of the log-likelihood derivative w.r.t some RBM parameter θ_a is given by

$$\left| E_{q(v^{(0)})} \left[E_{p(v^{(k)} | v^{(0)})} \left[\frac{\partial \ln p(v^{(k)})}{\partial \theta_a} \right] \right] \right| \leq \frac{1}{2} |q - p| \left(1 - e^{-(m+n)\Delta} \right)^k$$

with

$$\begin{aligned} \Delta &= \max \left\{ \max_{l \in \{1, \dots, m\}} \vartheta_l, \max_{l \in \{1, \dots, n\}} \xi_l \right\} \\ \vartheta_l &= \max \left\{ \left| \sum_{i=1}^n I_{\{w_{il} > 0\}} w_{il} + b_l \right|, \left| \sum_{i=1}^n I_{\{w_{il} < 0\}} w_{il} + b_l \right| \right\} \\ \xi_l &= \max \left\{ \left| \sum_{j=1}^m I_{\{w_{lj} > 0\}} w_{lj} + c_l \right|, \left| \sum_{j=1}^m I_{\{w_{lj} < 0\}} w_{lj} + c_l \right| \right\} \end{aligned}$$

RBM's and neural networks

Why is $p(H_i = 1|\mathbf{v})$ in positive phase computationally feasible?

→ Z cancels out; explicit calculation shows:

$$p(H_i = 1|\mathbf{v}) = \text{sig}\left(\sum_{j=1}^m w_{ij}v_j + c_i\right)$$

$$p(V_j = 1|\mathbf{h}) = \text{sig}\left(\sum_{i=1}^n w_{ij}h_i + b_j\right)$$

Identify:

single variable → neuron

conditional probability → sigmoid activation function

RBM → **stochastic neural network**