# Boltzmann machines (BM)

Consider

$$E(\vec{x}) = -\frac{1}{2}\sum_{i,j} x_i \omega_{ij} x_j = -\frac{1}{2}\vec{x}^T W \vec{x},$$

$$P(\vec{x}) = \frac{1}{Z} e^{-E(\vec{x})}$$

Stochastic Hopfield network (aka Boltzmann machine, BM) ← actually implements Boltzmann distr'n

## Activity rule:

$$a_i = \sum_j \omega_{ij} x_j,$$

Gibbs sampling

$$x_i = \begin{cases} +1, & \text{prob. } q_i = \dfrac{1}{1+e^{-2a_i}} = \dfrac{e^{a_i}}{e^{a_i}+e^{-a_i}} \\ -1, & \text{prob. } 1-q_i = \dfrac{e^{-a_i}}{e^{a_i}+e^{-a_i}} \end{cases}$$

→ stochastic update

cf. p.402 31.1

Consider

$$E(\vec{x}) = -\frac{1}{2}J \sum_{\substack{m,n \\ m \neq n}} x_m x_n - H\sum_n x_n$$

↗ spin glass

Then $\ell_n = J\sum_{\substack{m \\ m \neq n}} x_m + H$ is the local field for spin $n$.

Indeed, for 2 spins

$$E = -\frac{J}{2}(x_1 x_2 + x_2 x_1) - H(x_1 + x_2) =$$

$$= -x_2(\underbrace{Jx_1 + H}_{\ell_2}) - Hx_1 =$$

$$= -x_1(\underbrace{Jx_2 + H}_{\ell_1}) - Hx_2$$

In general,
$$E = -x_n \ell_n + \text{const}(x_n)$$

<u>Gibbs</u> <u>sampling</u>: select spin $n$ at random

$$\begin{cases} P(S_n = +1 \mid b_n) = \dfrac{e^{+\beta b_n}}{e^{\beta b_n} + e^{-\beta b_n}} = \dfrac{1}{1 + e^{-2\beta b_n}}, \\[4mm] \underset{\substack{\uparrow \\ \text{all other spins} \\ \text{fixed}}}{} \\[4mm] P(S_n = -1 \mid b_n) = 1 - P(S_n = +1 \mid b_n) \end{cases}$$

Use these probabilities to set the spin state: $\pm 1$.

This converges to Boltzmann equilibrium.

<u>Metropolis</u> <u>sampling</u>:

Compute $\quad \Delta E = \begin{cases} x_n = 1 \Rightarrow x_n = -1 : E = -b_n \times \text{const} \Rightarrow b_n + \text{const} : \Delta E = 2b_n \\ x_n = -1 \Rightarrow x_n = 1 : E = b_n + \text{const} \Rightarrow -b_n + \text{const} : \\ \qquad\qquad\qquad\qquad\qquad \Delta E = -2b_n \end{cases}$

So, $\quad \Delta E = 2 b_n x_n$.

$P(\text{accept spin flip}) = \begin{cases} 1 & \Delta E \le 0 \\ e^{-\beta \Delta E} & \Delta E > 0 \end{cases}$

This converges to Boltzmann eq'm as well.

Now, given $\underset{\text{a set of } N}{\vee}$ examples $\{\vec{x}^{(n)}\}_1^N$, we might adjust weights $W$ s.t. the likelihood of generating those examples from the Boltzmann distribution $P(\vec{x})$ is maximized:

$$\mathcal{L} = \prod_{n=1}^{N} P(\vec{x}^{(n)}) \quad , \quad \text{or}$$

$$\log \mathcal{L} = \sum_{n=1}^{N} \log P(\vec{x}^{(n)}) = \sum_{n} \left[ \frac{1}{2} \vec{x}^{(n)T} W \vec{x}^{(n)} - \right.$$

$$\left. - \log Z \right].$$

We need

$$\frac{\partial}{\partial w_{ij}} \log Z = \frac{1}{Z} \frac{\partial}{\partial w_{ij}} \left\{ \sum_{\vec{x}} e^{-E(\vec{x})} \right\} =$$

$$= -\sum_{\vec{x}} P(\vec{x}) \frac{\partial}{\partial w_{ij}} E(\vec{x}) = \sum_{\vec{x}} x_i x_j P(\vec{x}) =$$

$$= \langle x_i x_j \rangle_P .$$

Then

$$\frac{\partial}{\partial w_{ij}} \log \mathcal{L} = \underbrace{\sum_{n} x_i^{(n)} x_j^{(n)}}_{\equiv \, N \langle x_i x_j \rangle_D} - N \langle x_i x_j \rangle_P =$$

$$= N \left[ \underbrace{\langle x_i x_j \rangle_D}_{\substack{\text{empirical} \\ \text{2-point} \\ \text{correl'n}}} - \underbrace{\langle x_i x_j \rangle_P}_{\substack{\text{model} \\ \text{2-point} \\ \text{correl'n}}} \right].$$

If $\dfrac{\partial}{\partial w_{ij}} \log \mathcal{L} = 0 \implies \underset{\substack{\nearrow \\ \text{compute} \\ \text{directly}}}{\langle x_i x_j \rangle_D} = \underset{\substack{\uparrow \\ \text{estimate} \\ \text{by gibbs} \\ \text{sampling}}}{\langle x_i x_j \rangle_P}$

Otherwise,

$$\langle x_i x_j \rangle_D - \langle x_i x_j \rangle_P \text{ provides the}$$

<u>gradient</u> for optimization algorithms.

Note that if $W=0 \implies E(\vec{x})=0$, $\forall \vec{x}$.

like the $\beta=0$
limit
(infinite $T$)

Then

$$\langle x_i x_j \rangle_P = \langle x_i \rangle_P \langle x_j \rangle_P = 0,$$

since all spins are equally likely
to be up or down.

If the weights are adjusted by the
gradient descent,

$$w_{ij}^{(\tau+1)} = w_{ij}^{(\tau)} + \eta \frac{\partial}{\partial w_{ij}} \log \mathcal{I} \Big|_{w_{ij}^{(\tau)}}$$

learning rate, $>0$

guaranteed to increase $\log \mathcal{I}$ if
the step is small:

$$\log \mathcal{I}(w_{ij}^{(\tau+1)}) \simeq \log \mathcal{I}(w_{ij}^{(\tau)}) +$$

$$+ \eta \underbrace{\frac{\partial}{\partial w_{ij}} \log \mathcal{I}\Big|_{w_{ij}^{(\tau)}} \times \frac{\partial}{\partial w_{ij}} \log \mathcal{I}\Big|_{w_{ij}^{(\tau)}}}_{\geq 0} .$$

$$\geq 0 \text{ if } \eta > 0$$

Thus in the $W=0$ case,

$$w_{ij}^{(1)} = \underbrace{w_{ij}^{(0)}}_{=0, \text{ say}} + \eta \sum_n x_i^{(n)} x_j^{(n)} \quad \leftarrow$$

Hebbian
learning
rule is
recovered in
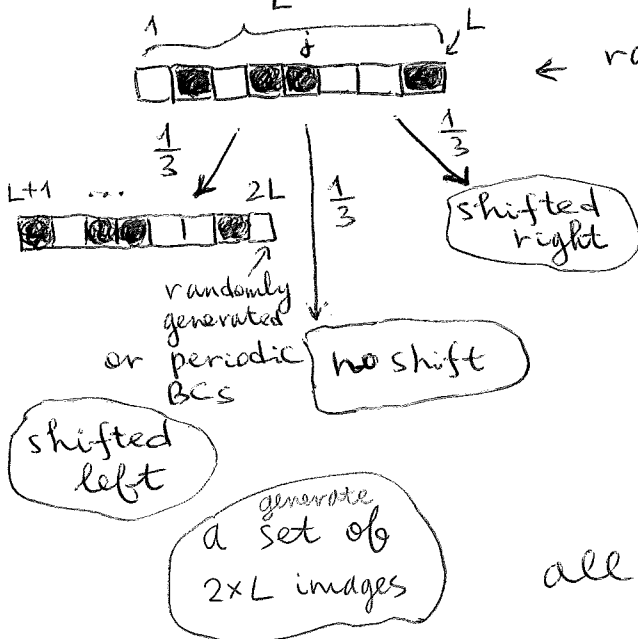1 iteration

# Poetic interpretation of BM learning:

When the BM is "awake", it measures

$$\text{i.e.}\uparrow \\ \text{gets input} \\ \text{from the world}$$

real-world correlations $\langle x_i x_j \rangle_D$ & uses them to adjust the weights. When it is "asleep", it does not adjust the weights — it "dreams" about the world & computes $\langle x_i x_j \rangle_P$ (i.e., its "idea" of the world). When $\langle x_i x_j \rangle_D = \langle x_i x_j \rangle_P$, the two views are balanced.

————o————

However, the "world" is represented by just two-point correlations $\langle x_i x_j \rangle_D$, seems to be too poor to really capture the richness of the world.

For example, consider a "shifter ensemble" of images:



← randomly generated pixels

randomly generated or periodic BCs

shifted left

shifted right

no shift

generate a set of $2 \times L$ images

Then, away from the boundaries:

$$\begin{cases} \langle x_j x_{j+L} \rangle = \frac{1}{3} & \text{unshifted} \\ \langle x_j x_{j+L-1} \rangle = \frac{1}{3} & \text{left} \\ \langle x_j x_{j+L+1} \rangle = \frac{1}{3} & \text{right} \end{cases}$$

all others are $= 0$

This seems too poor to describe the images $\Rightarrow$ need higher-order statistics:

$$P(\vec{x}) = \frac{1}{Z} e^{\frac{1}{2} \sum_{ij} w_{ij} x_i x_j + \frac{1}{6} \sum_{ijk} v_{ijk} x_i x_j x_k + \dots}$$

higher-order BM

Can get $\quad \dfrac{\partial}{\partial w_{ij}} \log Z \quad , \quad \dfrac{\partial}{\partial v_{ijk}} \log Z$, etc.

do gibbs sampling

[ But there are too many parameters ]
[ in general.

<u>Idea</u>: (due to Hinton & Sejnowski, 1986) introduce hidden variables to model higher-order correlations.

$\underline{\underline{BM}} \quad \underline{\underline{with}} \quad \underline{\underline{hidden}} \quad \underline{\underline{units}}$ [restricted BM]

$\vec{y} = $  $\left\{ \begin{array}{l} \vec{x} \\ \vec{h} \end{array} \right\}$ visible nodes state $(M_1)$ vector

hidden nodes state $(M_2)$ vector

$\uparrow$
node states, either
visible or hidden
$(M_1 + M_2)$
vector

In particular, when visible nodes are "clamped" at $\vec{x}^{(n)} \Rightarrow \vec{y}^{(n)} \equiv (\vec{x}^{(n)}, \vec{h})$.

Then $\quad P(\vec{x}^{(n)}) = \sum_{h} P(\vec{x}^{(n)}, \vec{h}) = \frac{1}{Z} \underbrace{\sum_{h} e^{\frac{1}{2} \vec{y}^{(n)T} W \vec{y}^{(n)}}}_{\equiv Z_{\vec{x}^{(n)}} \leftarrow \text{partial partition function}}$

$Z = \sum_{\vec{x}, \vec{h}} e^{\frac{1}{2} \vec{y}^T W \vec{y}}$

as before, consider

$$\frac{\partial}{\partial w_{ij}} \log \mathcal{L} = \sum_n \frac{\partial}{\partial w_{ij}} \left\{ \log Z_{\vec{x}^{(n)}} - \log Z \right\} \ominus$$

$$\mathcal{L} = \prod_{n=1}^{N} P(\vec{x}^{(n)})$$

$$\ominus \quad \sum_n \left\{ \frac{1}{Z_{\vec{x}^{(n)}}} \sum_h y_i^{(n)} y_j^{(n)} e^{\frac{1}{2} \vec{y}^{(n)T} W \vec{y}^{(n)}} - \right.$$

$$\left. - \underbrace{\langle x_i x_j \rangle_{P(\vec{x}, \vec{h})}}_{\text{as before}} \right\} \boxminus$$

$$\frac{\sum_h y_i^{(n)} y_j^{(n)} e^{\frac{1}{2} \vec{y}^{(n)T} W \vec{y}^{(n)}}}{\underbrace{\sum_h e^{\frac{1}{2} \vec{y}^{(n)T} W \vec{y}^{(n)}}}_{Z_{\vec{x}^{(n)}}}} \equiv \sum_h y_i^{(n)} y_j^{(n)} P(\vec{h} | \vec{x}^{(n)}) =$$
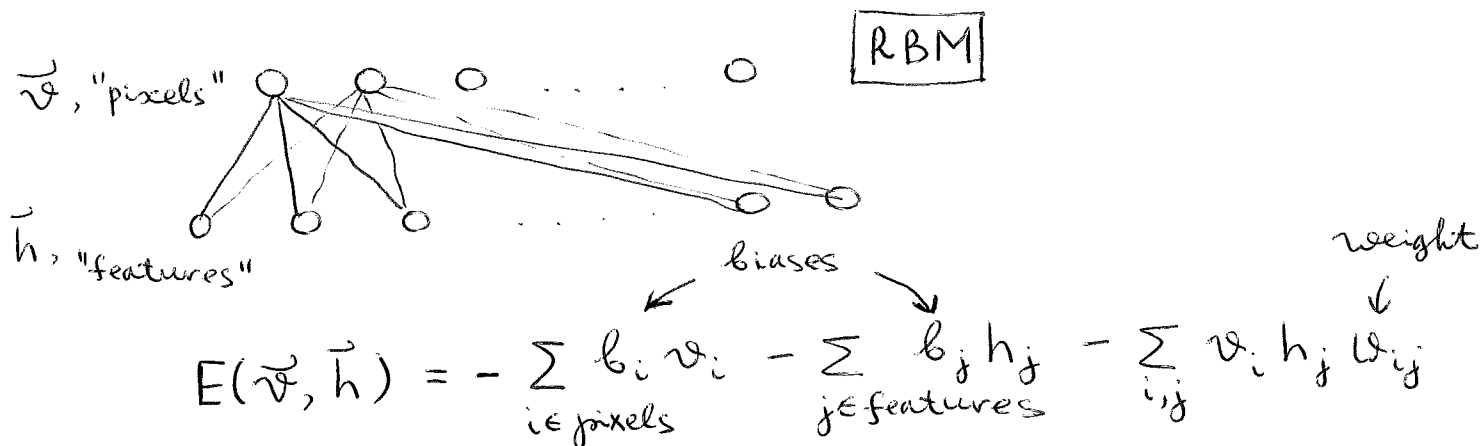
$$= \langle y_i y_j \rangle_{P(\vec{h} | \vec{x}^{(n)})}$$

$$\boxminus \quad \sum_n \left\{ \underbrace{\langle y_i y_j \rangle_{P(\vec{h} | \vec{x}^{(n)})}}_{\substack{\text{estimate by} \\ \text{gibbs sampling} \\ \text{with } \vec{x}^{(n)} \text{ fixed} \\ \text{(only hidden spins} \\ \text{flipped)}}} - \underbrace{\langle y_i y_j \rangle_{P(\vec{x}, \vec{h})}}_{\substack{\text{estimate by unrestricted} \\ \text{gibbs sampling} \\ \text{(both visible \&} \\ \text{hidden spins} \\ \text{flipped)}}} \right\}$$

# Application of BM in neural networks (NN)

Hinton & Salakhutdinov,
Science 2006

Idea: build a multi-layer NN, pre-train intermediate layers using BMs, then refine the weights by backpropagation.

Consider data that can be represented as binary vectors, e.g. images (0,1) (or vector of spins)

$\vec{v}$, "pixels"

$\vec{h}$, "features"

$\boxed{RBM}$

biases          weight

$$E(\vec{v}, \vec{h}) = -\sum_{i \in pixels} b_i v_i - \sum_{j \in features} b_j h_j - \sum_{i,j} v_i h_j w_{ij}$$

given pixel states,

(1) driven by data

$$h_j = \begin{cases} 1 & , \quad \sigma(b_j + \sum_i v_i w_{ij}) \\ 0 & , \quad \text{otherwise} \end{cases} \quad (*)$$

$\forall j$

record $v_i h_j$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

(2) $\forall i$

$$v_i = \begin{cases} 1 & , \quad \sigma(b_i + \sum_j h_j \overset{w_{ji}}{w_{ij}}) \\ 0 & , \quad \text{otherwise} \end{cases} \quad (**)$$

"confabulation"

(3) $\forall j$
$$h_j = \begin{cases} 1 & , \sigma(b_j + \sum_i v_i w_{ij}) \\ 0 & , \text{otherwise} \end{cases}$$

driven
by confabulation              record $v_i h_j$

Repeat many times, compute
$$\langle v_i h_j \rangle_{data} \quad \& \quad \langle v_i h_j \rangle_{recon}$$

Finally, adjust weights:
$$\Delta w_{ij} = \eta ( \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{recon} )$$

learning
rate

——— o ———

Iterate to convergence.

$$\begin{bmatrix} \text{Next, make the hidden units the} \\ \text{visible units of the next RBM.} \end{bmatrix}$$

Note: $E(\vec{v}, \vec{h}) = - \sum_i v_i \underbrace{[ b_i + \sum_j h_j w_{ij} ]}_{\text{local field for } v_i} + \text{const}(\vec{v})$

Then

$$P(v_i = +1) = \frac{e^{(b_i + \sum_j h_j w_{ij})}}{e^{(b_i + \sum_j h_j w_{ij})} + 1} =$$

all other
spins fixed

$v_i = 1$ state          $v_i = 0$ state

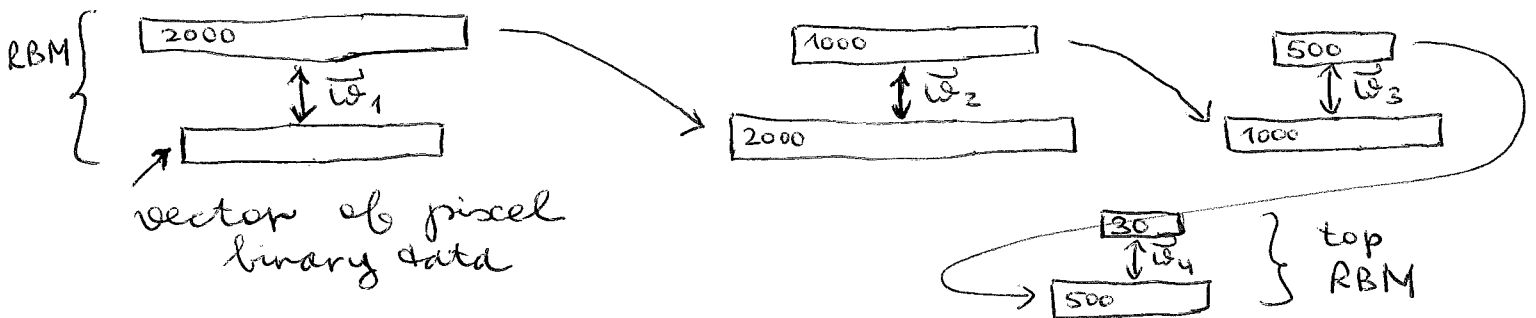$$= \sigma( b_i + \sum_j h_j w_{ij} )$$

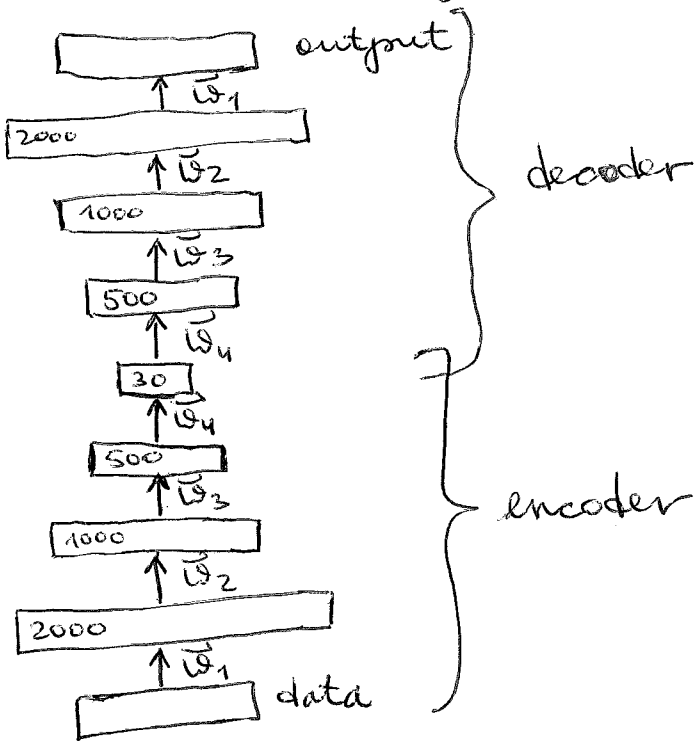$$P(v_i = 0) = 1 - P(v_i = +1)$$

Same as (**)

Likewise,

$$E(\vec{v}, \vec{h}) = -\sum_j h_j \left[ b_j + \underbrace{\sum_i v_i w_{ij}}_{\text{local field for } h_j} \right] + \text{const}(\vec{h})$$

leading to (*)

---

Finally, the whole architecture:



RBM { 
2000
$\updownarrow \vec{w}_1$

vector of pixel binary data

1000
$\updownarrow \vec{w}_2$
2000

500
$\updownarrow \vec{w}_3$
1000

30
$\updownarrow \vec{w}_4$
500
} top RBM

Unrolling:



output
$\uparrow \vec{w}_1$
2000
$\uparrow \vec{w}_2$
1000
$\uparrow \vec{w}_3$
500
$\uparrow \vec{w}_4$
30
$\uparrow \vec{w}_4$
500
$\uparrow \vec{w}_3$
1000
$\uparrow \vec{w}_2$
2000
$\uparrow \vec{w}_1$
data

} decoder

} encoder

For backpropagation, replace stochastic units with $\sigma$-units with local fields as activations

Minimize the error between output & data by backpropagation with conjugate gradients used on $10^3$ data vectors at a time.

-10-