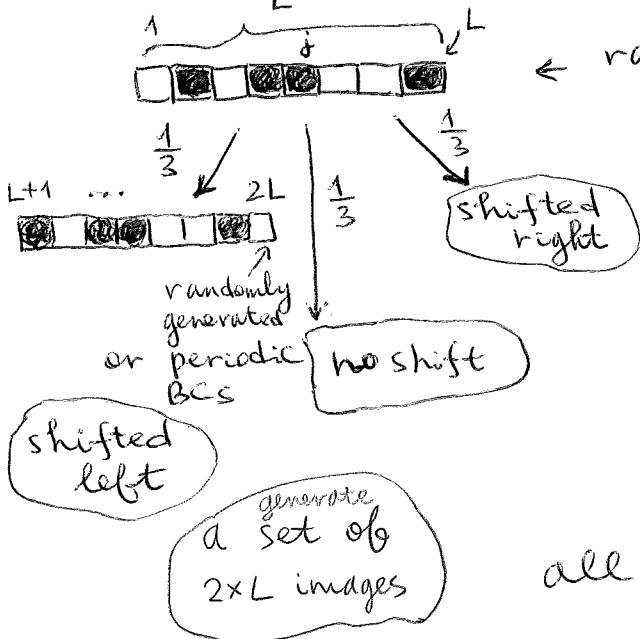# Poetic interpretation of BM learning:

When the BM is "awake", it measures

> i.e.↑
> gets input
> 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

shifted right

no shift

shifted left

randomly generated or periodic BCs

generate a set of 2×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:

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

$\uparrow$ higher-order BM

Can get $\frac{\partial}{\partial w_{ij}} \log Z$ , $\frac{\partial}{\partial v_{ijk}} \log Z$, 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.

<u>BM</u> <u>with</u> <u>hidden</u> <u>units</u> [restricted BM]

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

hidden nodes state ($M_2$) vector

$\vec{y} = $ (scribbled)

$\uparrow$
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 $P(\vec{x}^{(n)}) = \sum_{\vec{h}} P(\vec{x}^{(n)}, \vec{h}) = \frac{1}{Z} \sum_{\vec{h}} e^{\frac{1}{2}\vec{y}^{(n)T} W \vec{y}^{(n)}}$

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

$\underbrace{\qquad\qquad}$
$\equiv Z_{\vec{x}^{(n)}}$ partial
partition
function

As before, consider

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

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

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

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

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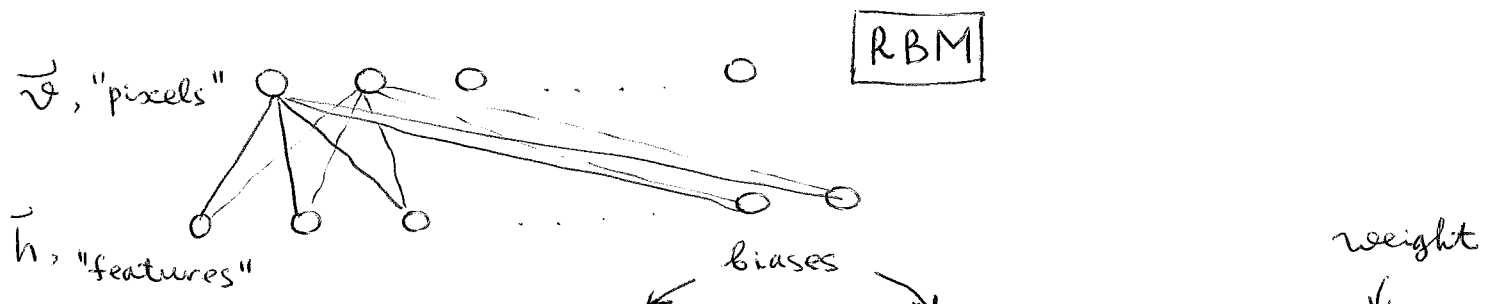
$$= \langle y_i y_j \rangle_{P(\vec{h}|\vec{x}^{(n)})}$$

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

# 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.

———o———

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



$\vec{v}$, "pixels"

$\vec{h}$, "features"

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

record $v_i h_j$

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

(2)
"confabulation"
$$\forall i \quad v_i = \begin{cases} 1, & \sigma(b_i + \sum_j h_j \overset{w_{ji}}{w_{ij}}) \\ 0, & \text{otherwise} \end{cases} \qquad (**)$$

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

driven
by confabulation $\qquad\qquad$ record $v_i h_j$

Repeat many times, compute

$\langle v_i h_j \rangle_{data}$ & $\langle v_i h_j \rangle_{recon}$

Finally, adjust weights:

$$\Delta \omega_{ij} = \eta \left( \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{recon} \right)$$

$\qquad\qquad \underset{\text{learning rate}}{\nearrow}$

——o——

Iterate to convergence.

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

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

Then
$\begin{cases} P(v_i = +1) = \dfrac{e^{\overset{\text{all other spins fixed}}{\curvearrowleft}(b_i + \sum_j h_j \omega_{ij})}}{e^{(b_i + \sum_j h_j \omega_{ij})} + 1} = \\[4ex] \qquad\qquad\qquad\quad \underset{v_i = 1 \text{ state}}{\nearrow} \qquad \underset{v_i = 0 \text{ state}}{\nearrow} \\[3ex] \qquad = \sigma(b_i + \sum_j h_j \omega_{ij}) \\[2ex] P(v_i = 0) = 1 - P(v_i = +1) \end{cases}$
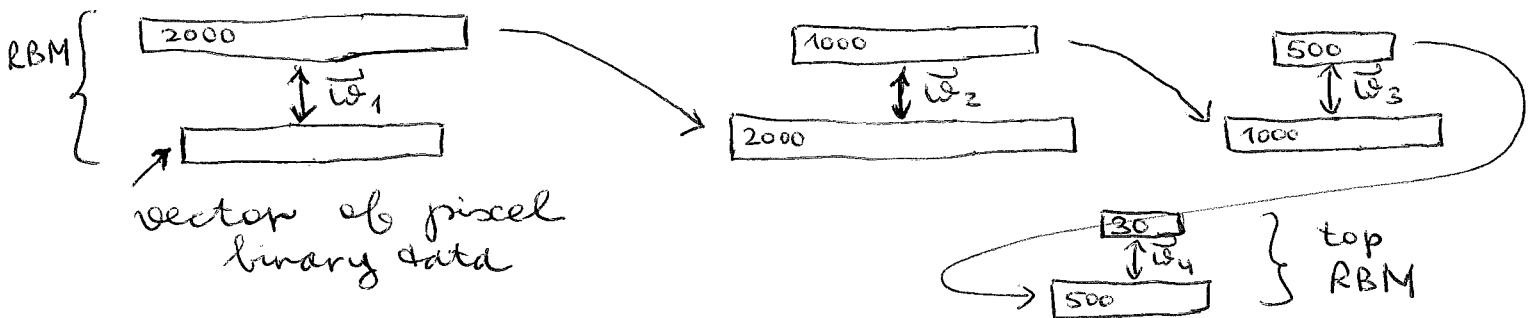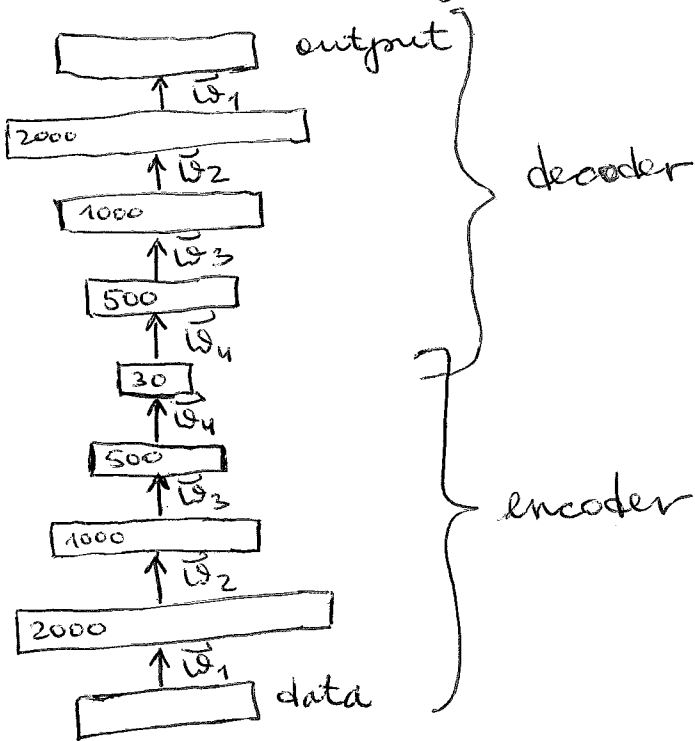
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ same as $(**)$

Likewise,

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

local field
for $h_j$

leading to (*)

---

Finally, the whole architecture:

RBM $\Big\{$

| 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 |

$\Big\}$ 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

$\Big\}$ decoder

$\Big\}$ 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.