

Problem 1 : Backpropagation.

$$\delta W_{mn} = -\eta \frac{\partial H}{\partial W_{mn}} = -\eta \sum_{p \neq m} \frac{\partial H}{\partial \sigma_p} \frac{\partial \sigma_p}{\partial W_{mn}}$$

$$= -\eta \sum_p \frac{\partial H}{\partial \sigma_p} \sigma'(B_p) \sum_2 \frac{\partial W_{p2}}{\partial W_{mn}} x_2$$

$$= +\eta \sum_{np} \left(\frac{t_p^{(n)}}{\sigma_p^{(n)}} - \frac{(1-t_p^{(n)})}{1-\sigma_p^{(n)}} \right) \sigma'(B_p) \sum_2 \delta_{pm} \delta_{qn} x_2$$

$$= \eta \sum_m \left(\frac{t_m^{(n)} - \sigma_m^{(n)}}{\sigma_m^{(n)} (1-\sigma_m^{(n)})} \right) \sigma'(B_m) (1-\sigma(B_m)) x_n$$

using $\tilde{\sigma}_m^{(n)} = \sigma(B_m)$

$$= \eta \sum_m (t_m^{(n)} - \sigma_m^{(n)}) x_n$$

Similarly by chain rule, obtain.

$$\delta w_{mn} = \eta \sum_{np} (t_p^{(n)} - \sigma_p^{(n)}) w_{pm} \sigma'(B_m) x_n$$

Problem 2: Restricted Boltzmann Machine.

Definitions: η : learning rate
 $\delta w_{mn}^{(u)}$: weight update for the weight
 w_{mn} , corresponding to pattern "u".

h_m : state of the m^{th} hidden neuron

$x_n^{(u)}$: input

v_n : state of the n^{th} visible neuron.

How to perform the averages:

$\langle h_m v_n \rangle_{\text{data}}$ is computed by averaging over all states
of the hidden neurons, given that pattern
"u" is applied to the visible neurons.

$\langle h_m v_n \rangle_{\text{model}}$ can be simplified in a similar
manner to the $\langle \cdot \rangle_{\text{data}}$ average, but is
computed in simulations by M-C sampling.

$$\begin{aligned}
 \langle h_m x_n^{(u)} \rangle_{\text{data}} &= \sum_{\substack{h_i=0,1 \\ h_m=0,1}} h_m x_n^{(u)} \prod_{i=1}^I p(h_i | V=x^{(u)}) \\
 &= \sum_{h_m=0,1} h_m p(h_m | V=x) \\
 &= 1 \cdot p(b_m^{(u)}) + 0 \cdot (1 - p(b_m^{(u)})) \\
 &\Rightarrow p(b_m^{(u)})
 \end{aligned}$$

where $p(b_m^{(u)}) = \frac{1}{1 + e^{-2b_m^{(u)}}}$

similarly,

$$\langle h_m v_n \rangle_{\text{model}} = \left\langle \frac{v_n}{1 + e^{-2b_m^{(u)}}} \right\rangle_{\text{model}}$$

Contrast: we find a sigmoid dependence on $b_m^{(u)}$
 instead of a tanh dependence, reflecting that
 the 0/1 neurons do not have a "0" mean.

Whereas 1/-1 neurons do.

choose the following filters:

$$1. \begin{bmatrix} +1 & +1 \\ -1 & -1 \end{bmatrix}$$

$$2. \begin{bmatrix} -1 & -1 \\ 1 & 1 \end{bmatrix}$$

$$3. \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}$$

$$4. \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}$$

choose thresholds = 0, stride = 1, padding = 0.

Use max pooling, 2×2 layers

Let bars be monochromatic columns, and stripes be monochromatic rows.

for bars, filters 1,2 output 0, whereas
(+ max pool)

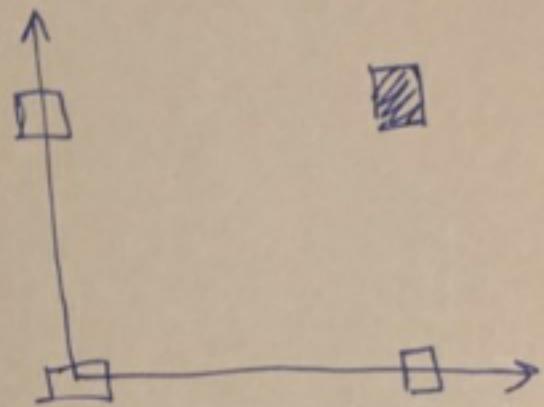
3,4 + max pool output 2, and vice versa for

stripes. Let x_i be the output of the i^{th} filter +
max pool.

Then a F.C. layer, $g(w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + \theta)$
with $w_1 = w_2 = +1, w_3 = w_4 = -1, \theta = 0$, and ~~act~~
step activation "g" outputs 0 for bars & 1 for stripes.

4. Boolean AND problem

Input		Output
0	0	-1
0	1	-1
1	0	-1
1	1	1



$$\begin{aligned}\square &= -1 \\ \blacksquare &= +1\end{aligned}$$

f) $H = \frac{1}{2} \left[(-w_1 - w_2 + \theta)^2 + (-1 - w_1 + \theta)^2 + (-1 - w_2 + \theta)^2 + (-1 + \theta)^2 \right]$

Let $\frac{\partial H}{\partial w_1}, \frac{\partial H}{\partial w_2}, \frac{\partial H}{\partial \theta} > 0$.

$$\frac{\partial H}{\partial w_1} = 2w_1 + w_2 - 2\theta > 0$$

$$\frac{\partial H}{\partial w_2} = w_1 + 2w_2 - 2\theta > 0$$

$$\frac{\partial H}{\partial \theta} = -2 - 2w_1 - 2w_2 + 4\theta = 0$$

Solution: $w_1 = 1, w_2 = 1, \theta = 3/2$

(c)

See Book section 5.3.

Reason: oneslap matrix not invertible for l.d. w/ infew.

(d) No, do not solve the Boolean AND.

The above procedure is equivalent to computing the pseudo inverse for the oneslap matrix. Thus essentially works as a best fit

⑤ Auto encoders

- write output of a layer as $b_i^{(n)} = w_{ij} x_j^{(n)} - \theta_i$
then, average over n and set-assume inputs have mean,

$$\langle b_i^{(n)} \rangle_n = w_{ij} \langle x_j^{(n)} \rangle_n - \theta_i$$

$$\langle b_i^{(n)} \rangle_n = -\theta_i$$

choosing $\theta_i = 0$, sets $\langle b_i^{(n)} \rangle_n > 0$

$$b) H = \frac{1}{2} \sum_{i,n} \left(t_i^{(n)} - \sum_{j,k} (w_d)_{ij} (w_e)_{jk} x_k^{(n)} \right)^2$$

$$f = \frac{1}{2} \| \hat{\mathbf{x}} - W_d W_c \mathbf{x} \|^2$$

- write $\hat{\mathbf{x}}_n = W_c \mathbf{x}$

- $\mathbf{x} = n \times N$

- assume that \mathbf{x} has rank n .

In the given case, $|W_d| = n \times 2$, $n > 2$,
 so can have rank at most 2.

the matrix $W_d \hat{\mathbf{x}}_n$ that minimises f
 is the best rank 2 approximation of \mathbf{x} .
 can be solved using SVD.

- $\mathbf{x} = U_n \Sigma_n V_n^T$

$$W_d \hat{\mathbf{x}}_n = U_p \Sigma_p V_p^T ; p=2.$$

thus the solution can be written

$$W_d = V_p T^{-1} ; \hat{\mathbf{x}}_n = T \Sigma_p V_p^T$$

for arbit, $p \times p$ T .

finally, need to solve

$$W_e \mathbf{x} = \pi \sum_p V_p^T$$

$\underset{P \times n}{\mathbf{x}}$ $\underset{n \times N}{\mathbf{x}}$ $\underset{P \times P}{\Sigma_p}$ $\underset{P \times P}{V_p}$ $\underset{P \times N}{\mathbf{x}}$

takes SVD pseudo inverse

$$\mathbf{A}^{\dagger} \mathbf{x}$$

$$W_c = \pi \underbrace{\sum_p V_p^T V_n \Sigma_n^{-1} U_n^T}_{\frac{1}{\sigma_p} \underset{P \times P}{\mathbf{x}} \underset{P \times N}{\mathbf{x}} \underset{N \times N}{\Sigma_n} \underset{N \times N}{\mathbf{x}} \underset{N \times N}{U_n^T}}$$

6. $P_{\text{error}}^{t=1} = P_{\text{obs}} (C_i^{(v)} > 1)$

using CLT, $C_i^{(v)}$ is normally distributed

with mean=0 & variance $\sigma_c^2 \approx \frac{P}{N}$

$$P_{\text{error}}^{t=1} = \int_1^{\infty} \frac{1}{\sqrt{2\pi \sigma_c^2}} e^{-\frac{c^2}{2\sigma_c^2}}$$

$$= [1 - \text{erf}\left(\frac{1}{\sqrt{2\sigma_c^2}}\right)]$$