**1. One-step error probability in deterministic Hopfield model.**

- Update rule: $S_i \leftarrow sgn\left(\sum_{j=1}^{N} w_{ij} S_j\right)$

- Weights: $\begin{cases} w_{ij} = \frac{1}{N} \sum_{\mu=1}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)}, & \text{for } i \neq j \\ w_{ii} = 0 \end{cases}$

- Input patterns: $\underline{\varphi}^{(\nu)}$; $\varphi_i^{(\nu)}$ — bit $i$ of input pattern $\underline{\varphi}^{(\nu)}$; $\varphi_i^{(\nu)} = +1$ or $-1$.

**a)** Condition for bit $\varphi_i^{(\nu)}$ to be stable after a single step of asynchronous update?

Apply $\underline{\varphi}^{(\nu)}$, obtain:
$$S_i = sgn\left[\sum_{j=1}^{N} w_{ij} \varphi_j^{(\nu)}\right]$$

For stability of $\varphi_i^{(\nu)}$ require: $\boxed{S_i \overset{!}{=} \varphi_i^{(\nu)}}$ (*)

Rewrite the left-hand-side of Eq. (*):

$$S_i = sgn\left(\sum_{j=1}^{N} w_{ij} \varphi_j^{(\nu)}\right) = sgn\left[\sum_{\substack{j=1 \\ j \neq i}}^{N} \left(\frac{1}{N} \sum_{\mu=1}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)}\right) \varphi_j^{(\nu)}\right]$$

$$= sgn\left[\frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \varphi_i^{(\nu)} \underbrace{\varphi_j^{(\nu)} \varphi_j^{(\nu)}}_{=1} + \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)} \varphi_j^{(\nu)}\right]$$

$$S_i = sgn\left[\frac{N-1}{N} \varphi_i^{(\nu)} + \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)} \varphi_j^{(\nu)}\right]$$ (#)

Rewrite the right-hand side of (#):

$$\text{RHS of (#)} = sgn\left[\varphi_i^{(\nu)} - \frac{1}{N}\varphi_i^{(\nu)} + \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)} \varphi_j^{(\nu)}\right]$$

$\underbrace{\qquad\qquad}_{\text{"cross-talk term"}}$

Stability condition:

(**) $$\boxed{\varphi_i^{(\nu)} \overset{!}{=} sgn\left[\left(1 - \frac{1}{N}\right)\varphi_i^{(\nu)} + \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)} \varphi_j^{(\nu)}\right]}$$

Stability condition satisfied when:

$$\boxed{\left|-\frac{1}{N}\varphi_i^{(\nu)} + \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)} \varphi_j^{(\nu)}\right| < 1}$$

Alternatively, one can define $C_i^{(\nu)}$ as follows:

$$C_i^{(\nu)} = \frac{1}{N} - \frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)} \varphi_j^{(\nu)} \varphi_i^{(\nu)}$$

$\uparrow$

(= cross-talk term × $(-\varphi_i^{(\nu)})$)

Multiply (**) by $(-\varphi_i^{(\nu)})$ and rewrite the stability condition (**) as follows:

$$\boxed{-1 \overset{!}{=} sgn\left(-1 + C_i^{(\nu)}\right)}$$

This condition is satisfied for $\boxed{C_i^{(\nu)} < 1}$.

Note: no limits were taken so far. In the limit of $N \gg 1$, $C_i^{(\nu)}$ is:
$$C_i^{(\nu)} \approx -\frac{1}{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} \sum_{\mu \neq \nu}^{P} \varphi_i^{(\mu)} \varphi_j^{(\mu)} \varphi_j^{(\nu)} \varphi_i^{(\nu)}, \quad \text{for } N \gg 1$$

**b)** Random patterns: $y_i^{(\mu)} = \begin{cases} +1, & \text{with prob. } \frac{1}{2}, \\ -1, & \text{with prob. } \frac{1}{2}. \end{cases}$

Bit $y_i^{(\nu)}$ is stable after a single step of asynchronous update if $c_i^{(\nu)} < 1$ (task a).

Therefore, the probability that $y_i^{(\nu)}$ is unstable is: (Perror)

$$\boxed{P_{error} = Prob\left( c_i^{(\nu)} > 1 \right)}$$

To evaluate Perror, consider $c_i^{(\nu)}$:

$$c_i^{(\nu)} = \frac{1}{N} - \frac{1}{N-1} \sum_{\substack{j=1 \\ j\neq i}}^{N} \sum_{\substack{\mu=1 \\ \mu\neq \nu}}^{p} y_i^{(\mu)} y_j^{(\mu)} y_i^{(\nu)} y_j^{(\nu)} \implies$$

$$\xrightarrow{N\gg 1} \quad c_i^{(\nu)} \approx -\frac{1}{N} \sum_{\substack{k=1 \\ k=1}}^{(p-1)(N-1)} \text{random variables } (x_k) \text{ with } \pm 1$$

$$\boxed{(p-1)(N-1) \text{ terms}}$$

- Since we assume $p\gg 1$ and $N\gg 1$, we can use the Central limit theorem (patterns are random!)
- Variables $x_k$ have the mean $\underline{0}$, and variance $\delta_k^2 = 1$.

It follows that $c_i^{(\nu)}$ has the following properties:
- $c_i^{(\nu)}$ is approximately Gaussian distributed,
- the mean of $c_i^{(\nu)}$ is equal to $\underline{0}$ (since the mean of the random variables $x_k$ is $\underline{0}$)
- the variance $\delta^2$ of $c_i^{(\nu)}$ is:

$$\delta^2 = \frac{1}{N^2} \cdot (N-1)(p-1)\, \delta_x^2 \approx \frac{p}{N}$$

$$\implies \delta^2 \approx \frac{p}{N} \quad (\text{since } p\gg 1, N\gg 1)$$

---

It follows that

$$\boxed{erf(z) = \frac{2}{\sqrt{\pi}} \int_0^{z} e^{-y^2} dy}$$

$$P_{error} = \int_1^{\infty} \frac{1}{\sqrt{2\pi \delta^2}} e^{-\frac{x^2}{2\delta^2}} dx = \frac{1}{2}\left[ 1 - erf\left(\frac{1}{\sqrt{2}\,\delta}\right)\right]$$

$$\underbrace{\phantom{\frac{1}{\sqrt{2\pi\delta^2}}}}_{\text{Gaussian distribution}}$$

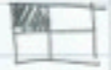$$\implies P_{error} = \frac{1}{2}\left[ 1 - erf\left(\frac{1}{\sqrt{2\frac{p}{N}}}\right)\right]$$

$$\boxed{P_{error} = \frac{1}{2}\left[ 1 - erf\left(\sqrt{\frac{N}{2p}}\right)\right]}$$
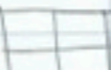
---

② Hopfield model: recognition of one pattern.

Stored pattern: $\underline{y}^{(1)} = \begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix}$

weight matrix: $\underline{\underline{w}} = \frac{1}{N} \underline{y}^{(1)} \underline{y}^{(1)T} = \frac{1}{4}\begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix}(1\ -1\ -1\ -1) = \frac{1}{4}\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}$

- Feeding in the $2^4$ possible patterns:

1)  $\to \underline{S}_1 = sgn\left(\underline{\underline{w}}\, \underline{y}^{(1)}\right) = \frac{1}{4} \underline{y}^{(1)} \underline{y}^{(1)T} \underline{y}^{(1)} = \frac{1}{4} \cdot 4\, \underline{y}^{(1)} = \underline{y}^{(1)}$

$\underline{S}_0 = \underline{y}^{(1)} = \begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix}$

2)  $\to \underline{S}_1 = sgn\left(\underline{\underline{w}}\, \underline{S}_0\right) = \frac{1}{4}\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix} = \underline{y}^{(1)}$

$\underline{S}_0 = \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix}$

3) $S_0 = \begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix}$ → $\underline{S}_1 = sgn(w\,\underline{S}_0) = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix}\right] =$

$$= \begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix} = \underline{y}^{(1)}$$

4) $\underline{S}_0 = \begin{pmatrix} -1 \\ 1 \\ -1 \\ 1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ -1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = \underline{y}^{(1)}$

5) $\underline{S}_0 = \begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ 1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix}$

$$= \underline{y}^{(1)}$$

6) $\underline{S}_0 = \begin{pmatrix} -1 \\ 1 \\ -1 \\ -1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ -1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Orthogonal pattern to the stored pattern. The network does not restore the stored pattern. In fact, it retrieves zero redor; failure of the network performance.

7) $\underline{S}_0 = \begin{pmatrix} -1 \\ -1 \\ 1 \\ -1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ -1 \\ 1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Same as case 6: orthogonal pattern.

8) $\underline{S}_0 = \begin{pmatrix} -1 \\ -1 \\ -1 \\ 1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ -1 \\ -1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Same as cases 6-7.

9) $\underline{S}_0 = \begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Same as cases 6-8.

10) $\underline{S}_0 = \begin{pmatrix} 1 \\ 1 \\ -1 \\ 1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ -1 \\ 1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Same as cases 6-9.

11) $\underline{S}_0 = \begin{pmatrix} 1 \\ -1 \\ 1 \\ 1 \end{pmatrix}$  $\underline{S}_1 = sgn\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Same as cases 6-10.

12)


$S_0 = \begin{pmatrix} -1 \\ 1 \\ 1 \\ -1 \end{pmatrix}$

$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ 1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$
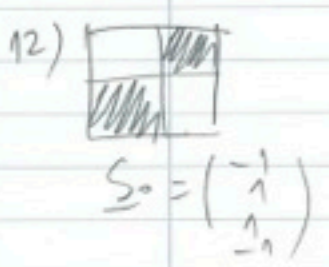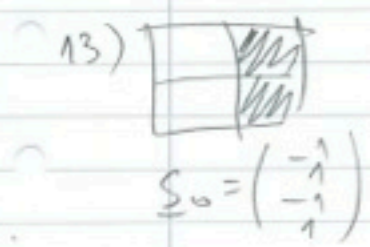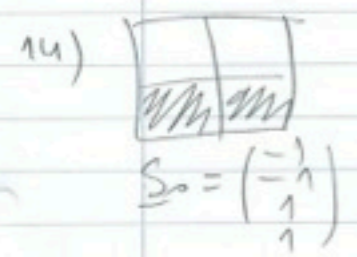
$$\boxed{S_1 = -\underline{\varphi}^{(1)}}$$

13)


$S_0 = \begin{pmatrix} -1 \\ -1 \\ 1 \\ 1 \end{pmatrix}$

$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} +1 \\ -1 \\ -1 \\ +1 \end{pmatrix}\right] = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$

$$\boxed{S_1 = -\underline{\varphi}^{(1)}}$$

14)


$S_0 = \begin{pmatrix} -1 \\ -1 \\ 1 \\ 1 \end{pmatrix}$

$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ -1 \\ 1 \\ 1 \end{pmatrix}\right] =$

$= \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix} = -\underline{\varphi}^{(1)}$

15)


$S_0 = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$

$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \Rightarrow$

$$\boxed{S_1 = -\underline{\varphi}^{(1)}}$$
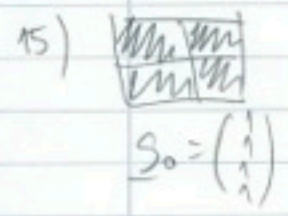
16)


$S_0 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$

$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$

---

In summary: ① In the first 5 cases, the network retrieves the stored pattern.
Note: in cases 2,3,4,5, the pattern that was fed had only one distorted bit in comparison to the stored pattern.
Case 1: fed pattern = stored pattern.

② In cases when more than 2 bits are distorted, the network retrieves the inverted version of the stored pattern (cases 12-16)

③ When exactly $\frac{N}{2} = 2$ bits are distorted, the network fails: unable to deal with patterns orthogonal to the stored pattern (due to Hebb's rule).

cases 6-11

3 Back-propagation I.
- Two hidden layers.
- Input patterns $\underline{\xi}^{(\mu)} = (\xi_1, \xi_2, ..., \xi_N)^T$
- Target output $\varphi_i^{(\mu)}$
- Network output $O_i^{(\mu)}$

- First hidden layer: $V_j^{(1,\mu)} = g(b_j^{(1,\mu)})$, $b_j^{(1,\mu)} = \sum_i w_{ji}^{(1)} \xi_i^{(\mu)} - \theta_j^{(1)}$

- Second hidden layer: $V_k^{(2,\mu)} = g(b_k^{(2,\mu)})$, $b_k^{(2,\mu)} = \sum_j w_{kj}^{(2)} V_j^{(1,\mu)} - \theta_k^{(2)}$

- Output layer: $O_i^{(\mu)} = g(b_i^{(\mu)})$, $b_i^{(\mu)} = \sum_k W_{ik} V_k^{(2,\mu)} - \Theta_i$

−Energy function: $H = \frac{1}{2} \sum_\mu \left( y_1^{(\mu)} - O_1^{(\mu)} \right)^2$

− Gradient-descent: find the parameters that minimise $H$.

− Start from the output layer:

$$\delta W_{1K} = -\eta \frac{\partial H}{\partial W_{1K}} = -\eta \frac{\partial}{\partial W_{1K}} \left\{ \frac{1}{2} \sum_\mu \left[ y_1^{(\mu)} - g(b_1^{(\mu)}) \right]^2 \right\} =$$

$$= -\eta \left[ \sum_\mu \left[ y_1^{(\mu)} - \underbrace{g(b_1^{(\mu)})}_{O_1^{(\mu)}} \right] \left( -\frac{\partial g(b_1^{(\mu)})}{\partial W_{1K}} \right) \right] =$$

$$= \eta \left[ \sum_\mu \left[ y_1^{(\mu)} - O_1^{(\mu)} \right] \cdot \frac{\partial g(b_1^{(\mu)})}{\partial W_{1K}} \right]$$

$$\frac{\partial g(b_1^{(\mu)})}{\partial W_{1K}} = \frac{\partial}{\partial W_{1K}} \left( g\left( \sum_\ell W_{1\ell} V_\ell^{(2,\mu)} - \Theta_1 \right) \right) =$$

$$= g'(b_1^{(\mu)}) \cdot V_K^{(2,\mu)} \qquad // \text{ Since } \frac{\partial \overline{W}_{1\ell}}{\partial W_{1K}} = \delta_{\ell K}$$

$$\boxed{\delta W_{1K} = \eta \sum_\mu \left[ y_1^{(\mu)} - O_1^{(\mu)} \right] \cdot g'(b_1^{(\mu)}) \cdot V_K^{(2,\mu)}} \boxed{= \eta \sum_\mu \delta_1^{(3,\mu)} V_K^{(2,\mu)}}$$

$$\delta \Theta_1 = -\eta \frac{\partial H}{\partial \Theta_1} = -\eta \frac{\partial}{\partial \Theta_1} \left\{ \frac{1}{2} \sum_\mu \left[ y_1^{(\mu)} - g(b_1^{(\mu)}) \right]^2 \right\} =$$

$$= -\eta \sum_\mu \left( y_1^{(\mu)} - O_1^{(\mu)} \right) \cdot \left( -\frac{\partial g(b_1^{(\mu)})}{\partial \Theta_1} \right) =$$

$$= \eta \sum_\mu \left( y_1^{(\mu)} - O_1^{(\mu)} \right) \cdot g'(b_1^{(\mu)}) \cdot (-1)$$

$$\Rightarrow \boxed{\delta \Theta_1 = -\eta \sum_\mu \left( y_1^{(\mu)} - O_1^{(\mu)} \right) \cdot g'(b_1^{(\mu)})} \boxed{= -\eta \sum_\mu \delta_1^{(3,\mu)}}$$

(left margin vertical) $\delta_1^{(3,\mu)} = \left( y_1^{(\mu)} - O_1^{(\mu)} \right) g'(b_1^{(\mu)})$

−Second hidden layer

$$\delta w_{kj}^{(2)} = -\eta \frac{\partial H}{\partial w_{kj}^{(2)}} = -\eta \frac{\partial}{\partial w_{kj}^{(2)}} \left\{ \frac{1}{2} \sum_\mu \left( y_1^{(\mu)} - O_1^{(\mu)} \right)^2 \right\} =$$

$$= \eta \sum_\mu \left( y_1^{(\mu)} - O_1^{(\mu)} \right) \frac{\partial O_1^{(\mu)}}{\partial w_{kj}^{(2)}}$$

$$O_1^{(\mu)} = g(b_1^{(\mu)}) = g\left[ \sum_\ell W_{1\ell} V_\ell^{(2,\mu)} - \Theta_1 \right] =$$

$$= g\left[ \sum_\ell W_{1\ell} g(b_\ell^{(2,\mu)}) - \Theta_1 \right] =$$

$$= g\left[ \sum_\ell W_{1\ell} g\left( \sum_s w_{\ell s}^{(2)} V_s^{(1,\mu)} - \Theta_\ell^{(2)} \right) - \Theta_1 \right]$$

$$\Rightarrow \frac{\partial O_1^{(\mu)}}{\partial w_{kj}^{(2)}} = g'(b_1^{(\mu)}) \cdot \frac{\partial}{\partial w_{kj}^{(2)}} \left( \sum_\ell W_{1\ell} g\left( \underbrace{\sum_s w_{\ell s}^{(2)} V_\wedge^{(1,\mu)} - \Theta_\ell^{(2)}}_{b_\ell^{(2,\mu)}} \right) - \Theta_1 \right)$$

$$= g'(b_1^{(\mu)}) \cdot \sum_\ell W_{1\ell} g'(b_\ell^{(2,\mu)}) \cdot \underbrace{\frac{\partial b_\ell^{(2,\mu)}}{\partial w_{kj}^{(2)}}}_{= \sum_s V_s^{(1,\mu)} \delta_{k\ell} \delta_{js}} =$$

$$= g'(b_1^{(\mu)}) \cdot W_{1K} g'(b_k^{(2,\mu)}) \cdot V_j^{(1,\mu)}$$

$$\Rightarrow \delta w_{kj}^{(2)} = \eta \sum_\mu \underbrace{\left( y_1^{(\mu)} - O_1^{(\mu)} \right) g'(b_1^{(\mu)})}_{\delta_1^{(3,\mu)}} \cdot W_{1K} g'(b_k^{(2,\mu)}) V_j^{(1,\mu)}$$

$$\delta w_{kj}^{(2)} = \eta \sum_\mu \underbrace{\delta_1^{(3,\mu)} W_{1k} g'(b_k^{(2,\mu)})}_{\delta_k^{(2,\mu)}} V_j^{(1,\mu)}$$

$$\boxed{\delta w_{kj}^{(2)} = \eta \sum_\mu \delta_k^{(2,\mu)} V_j^{(1,\mu)}}$$

Thresholds $\Theta_k^{(2)}$:

$$\delta \Theta_k^{(2)} = -\eta \frac{\partial H}{\partial \Theta_k^{(2)}} = \eta \sum_\mu \left( y_1^{(\mu)} - O_1^{(\mu)} \right) \frac{\partial O_1^{(\mu)}}{\partial \Theta_k^{(2)}}$$

from previous page

$$\frac{\partial O_1^{(\mu)}}{\partial \Theta_k^{(2)}} \stackrel{\downarrow}{=} g'(b_1^{(\mu)}) \frac{\partial}{\partial \Theta_k^{(2)}} \left[ \sum_\ell W_{1\ell} \, g\left( \sum_s w_{\ell s}^{(2)} V_s^{(1,\mu)} - \Theta_\ell^{(2)} \right) - \Theta_1 \right] =$$

$$= g'(b_1^{(\mu)}) \sum_\ell W_{1\ell} \, g'(b_\ell^{(2,\mu)}) (-1) \delta_{\ell k}$$

$$= -g'(b_1^{(\mu)}) \cdot W_{1k} \, g'(b_k^{(2,\mu)})$$

$$\Rightarrow \delta \Theta_k^{(2)} = -\eta \sum_\mu \underbrace{\left( y_1^{(\mu)} - O_1^{(\mu)} \right) g'(b_1^{(\mu)})}_{\delta_1^{(3,\mu)}} \, W_{1k} \, g'(b_k^{(2,\mu)})$$

$$= -\eta \sum_\mu \underbrace{\delta_1^{(3,\mu)} \, W_{1k} \, g'(b_k^{(2,\mu)})}_{\delta_k^{(2,\mu)}}$$

$$\boxed{\delta \Theta_k^{(2)} = -\eta \sum_\mu \delta_k^{(2,\mu)}}$$

For the first hidden layer we should proceed as above.
Alternatively, we note that $\delta$'s for the 3rd and
2nd layer obey the following relation:

$$\delta_k^{(2,\mu)} = \delta_1^{(3,\mu)} \, W_{1k} \, g'(b_k^{(2,\mu)})$$

We can use this to find the $\delta$'s for the first hidden

layer:

$$\delta_j^{(1,\mu)} = \sum_k \delta_k^{(2,\mu)} \, w_{kj}^{(2)} \, g'(b_j^{(1,\mu)})$$

The update formulae are, therefore, as follows:

Output layer: $\delta W_{1k} = \eta \sum_\mu \delta_1^{(3,\mu)} V_k^{(2,\mu)}$

$$\delta \Theta_1 = -\eta \sum_\mu \delta_1^{(3,\mu)}$$

Second hidden layer: $\delta w_{kj}^{(2)} = \eta \left( \sum_\mu \right) \delta_k^{(2,\mu)} V_j^{(1,\mu)}$

$$\delta \Theta_k^{(2)} = -\eta \sum_\mu \delta_k^{(2,\mu)}$$

First hidden layer: $\delta w_{ji}^{(1)} = \eta \sum_\mu \delta_j^{(1,\mu)} \xi_i^{(\mu)}$

$$\delta \Theta_j^{(1)} = -\eta \sum_\mu \delta_j^{(1,\mu)}$$

Summation over $\mu$ only for batch mode. Otherwise $\eta$: no Summation!

Here we have the following:

$$\delta_1^{(3,\mu)} = \left( y_1^{(\mu)} - O_1^{(\mu)} \right) g'(b_1^{(\mu)}), \quad b_1^{(\mu)} = \sum_k W_{1k} V_k - \Theta_1$$

$$\delta_k^{(2,\mu)} = \delta_1^{(3,\mu)} W_{1k} \, g'(b_k^{(2,\mu)}), \quad b_k^{(2,\mu)} = \sum_j w_{kj}^{(2)} V_j^{(1,\mu)} - \Theta_k^{(2)}$$

$$\delta_j^{(1,\mu)} = \sum_k \delta_k^{(2,\mu)} w_{kj}^{(2)} \, g'(b_j^{(1,\mu)}), \quad b_j^{(1,\mu)} = \sum_i w_{ji}^{(1)} \xi_i^{(\mu)} - \Theta_j^{(1)}$$

④ Backpropagation II — discussion of the implementation of the algorithm above. Explain how you program backpropagation.

Output $y = \sum_{j=1}^{N} w_j \xi_j = \underline{w}^T \underline{\xi}$

⑤ Oja's rule — $\delta w_j = \eta y (\xi_j - y w_j)$

a) Prove that $\underline{w}^*$ maximises $\langle y^2 \rangle$ using that

$|\underline{w}^*|^2 = 1$ and $\underline{w}^*$ is the leading eigenvector of $C$, with elements $C_{ij} = \langle \xi_i \xi_j \rangle$.

$\langle y^2 \rangle = \langle (\underline{w}^T \underline{\xi})(\underline{\xi}^T \underline{w}) \rangle = \langle \underline{w}^T C \underline{w} \rangle$

For $\underline{w} = \underline{w}^*$, find $\langle y^2 \rangle_{\underline{w}^*} = \langle \underline{w}^{*T} \underbrace{C \underline{w}^*}_{\substack{\lambda_{max} \underline{w}^* \\ (\text{from } ii)}} \rangle = \lambda_{max} \underbrace{\langle \underline{w}^{*T} \underline{w}^* \rangle}_{\substack{= 1 \\ (\text{from } i)}}$

$\Rightarrow \boxed{\langle y^2 \rangle = \lambda_{max}}$, where $\lambda_{max}$ is the maximum eigenvalue of $C$.

Since $C$ is symmetric $(\langle \xi_i \xi_j \rangle = \langle \xi_j \xi_i \rangle)$ it has real eigenvalues and its eigenvectors are orthogonal:

$\underline{u}_\alpha \underline{u}_\beta = \delta_{\alpha\beta}$, where $\delta_{\alpha\beta} = \begin{cases} 1, & \text{for } \alpha = \beta \\ 0, & \text{otherwise} \end{cases}$

Furthermore, all eigenvalues of $C$ are positive, since

$\lambda_\alpha = \underline{u}_\alpha^T C \underline{u}_\alpha = \underline{u}_\alpha^T \langle \underline{\xi} \underline{\xi}^T \rangle \underline{u}_\alpha = \langle \underline{u}_\alpha^T \underline{\xi} \underline{\xi}^T \underline{u}_\alpha \rangle =$

$= \langle |\underline{u}_\alpha^T \underline{\xi}|^2 \rangle \geq 0$

For any unit vector $\underline{w} = \sum_\alpha K_\alpha \underline{u}_\alpha$ that can be represented as a linear combination of the eigenvectors $\underline{u}_\alpha$ with coefficients $K_\alpha$ (assuming that $|\underline{w}|^2 = 1$) we find

$\langle y^2 \rangle_{\underline{w}} = \langle (\sum_\alpha K_\alpha \underline{u}_\alpha)^T C (\sum_\beta K_\beta \underline{u}_\beta) \rangle = \langle \sum_\alpha (K_\alpha \underline{u}_\alpha)^T (\sum_\beta K_\beta \lambda_\beta \underline{u}_\beta) \rangle =$

$= \langle \sum_{\alpha\beta} K_\alpha K_\beta \lambda_\beta \underbrace{\underline{u}_\alpha^T \underline{u}_\beta}_{\delta_{\alpha\beta}} \rangle = \langle \sum_\alpha (K_\alpha)^2 \lambda_\alpha \rangle \leq \lambda_{max} \langle \sum_\alpha |K_\alpha|^2 \rangle$

From $|\underline{w}|^2 = 1$, we find $\sum_\alpha (K_\alpha)^2 = 1$

Therefore: $\langle y^2 \rangle_{\underline{w}} \leq \lambda_{max} \langle \sum_\alpha |K_\alpha|^2 \rangle = \lambda_{max}$

$\Downarrow$

$\boxed{\langle y^2 \rangle_{\underline{w}} \leq \lambda_{max}}$ and $\underline{\langle y^2 \rangle_{\underline{w}^*} = \lambda_{max}}$

This shows that $\langle y^2 \rangle_{\underline{w}^*}$ is maximal in comparison to $\langle y^2 \rangle$ evaluated for any other $\underline{w}$ such that $|\underline{w}|^2 = 1$.

b) Assume that $\underline{w}^*$ is a steady state. In other words:

$\langle \delta \underline{w} \rangle_{\underline{w}^*} = 0$

$\Rightarrow \langle \eta y (\underline{\xi} - y \underline{w}) \rangle_{\underline{w}^*} = 0$

$\Rightarrow \langle \underline{w}^{*T} \underline{\xi} (\underline{\xi} - \underline{w}^{*T} \underline{\xi} \underline{w}^*) \rangle = 0 \quad | (\underline{w}^{*T} \underline{\xi}) \underline{\xi} = \underline{\xi} (\underline{w}^{*T} \underline{\xi})$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad = \underline{\xi} \underline{\xi}^T \underline{w}^*$

$\langle \underbrace{\underline{\xi} \underline{\xi}^T}_{C} \underline{w}^* - \underline{w}^{*T} \underbrace{\underline{\xi} \underline{\xi}^T}_{C} \underline{w}^* \underline{w}^* \rangle = 0$

$C \underline{w}^* - \underbrace{[\underline{w}^{*T} C \underline{w}^*]}_{\text{scalar; let's call it } \lambda} \underline{w}^* = 0$

$(**) \Rightarrow \boxed{C \underline{w}^* = \lambda \underline{w}^*} \Rightarrow$ Thus, $\underline{w}^*$ is an eigenvector of $C$, with eigenvalue

$\boxed{\lambda = \underline{w}^{*T} C \underline{w}^*}$

Norm of $\underline{w}^*$ (property $i$)

$\lambda = \underline{w}^{*T} \underbrace{C \underline{w}^*}_{\text{from } (**)} = \underline{w}^{*T} \lambda \underline{w}^* = \lambda \underline{w}^{*T} \underline{w}^* = \lambda |\underline{w}^*|^2$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \Rightarrow \boxed{|\underline{w}^*|^2 = 1} \; \text{shown} \; ①$

Now we must show that $\underline{w}^*$ has the maximum eigenvalue $\lambda_{max}$. Note: in order for the network to converge to a steady state, this steady state needs to be stable. Otherwise, the network would not converge to it.

Therefore, check the stability of $\underline{w}^*$.

Evaluate: $\langle \delta \underline{w} \rangle$ at $\underline{w} = \underline{w}^* + \underline{\varepsilon}$, where $|\underline{\varepsilon}|$ is small.

$$\langle \delta(\underline{w}^* + \underline{\varepsilon}) \rangle = \eta \langle (\underline{w}^* + \underline{\varepsilon})^T \underline{\xi} [\underline{\xi} - (\underline{w}^* + \underline{\varepsilon})^T \underline{\xi}(\underline{w}^* + \underline{\varepsilon})] \rangle$$

up to linear order in $\underline{\varepsilon}$

$= 0$ because $\underline{w}^*$ is steady (previous page)

$$\approx \eta [\langle \underline{w}^{*T} \underline{\xi}(\underline{\xi} - \underline{w}^{*T}\underline{\xi}\,\underline{w}^*) \rangle$$

$$+ \langle \underline{\varepsilon}^T \underline{\xi}\,\underline{\xi} \rangle - \langle \underline{\varepsilon}^T \underline{\xi}(\underline{w}^{*T}\underline{\xi}\,\underline{w}^*) \rangle$$

$= (\underline{\xi}\,\underline{\xi}^T \underline{\varepsilon})$

say $\underline{w}^* = \underline{u}_\alpha$ one of the eigenvectors.

$$- \langle \underline{w}^{*T}\underline{\xi}\,\underline{w}^{*T}\underline{\xi}\,\underline{\varepsilon} \rangle$$

$$- \langle \underline{w}^{*T}\underline{\xi}\,\underline{\varepsilon}^T\underline{\xi}\,\underline{w}^* \rangle$$

$\lambda_\alpha \underline{u}_\alpha$

$c$

$= \underline{u}_\alpha$

$$\Rightarrow \langle \delta(\underline{w}^* + \underline{\varepsilon}) \rangle \approx \eta [\langle \underline{\xi}\,\underline{\xi}^T \underline{\varepsilon} \rangle - \langle \underline{\varepsilon}^T \underline{\xi}\,\underline{\xi}^T \underline{w}^* (\underline{w}^*) \rangle$$

$= \lambda_\alpha \underline{u}_\alpha$

$$- \langle \underline{w}^{*T}\underline{\xi}\,\underline{\xi}^T\underline{w}^* \underline{\varepsilon} \rangle - \langle \underline{w}^{*T}\underline{\xi}\,\underline{\xi}^T \underline{\varepsilon}\,\underline{w}^* \rangle$$

$$= \eta [c\underline{\varepsilon} - \underline{\varepsilon}^T \lambda_\alpha \underline{u}_\alpha \underline{u}_\alpha$$

$= \lambda_\alpha$    $= \underline{\varepsilon}^T \underline{u}_\alpha$

$$- \underline{u}_\alpha^T \lambda_\alpha \underline{u}_\alpha \underline{\varepsilon} - \lambda_\alpha \underline{u}_\alpha^T \underline{\varepsilon}\,\underline{u}_\alpha ]$$

$$= \eta [c\underline{\varepsilon} - 2\lambda_\alpha(\underline{\varepsilon}^T \underline{u}_\alpha)\underline{u}_\alpha - \lambda_\alpha \underline{\varepsilon}]$$

Multiply both sides by $\underline{u}_\beta^T$. Find:

$= \lambda_\beta \underline{u}_\beta^T$

$$\underline{u}_\beta^T \langle \delta(\underline{w}^* + \underline{\varepsilon}) \rangle = \eta\left( \underline{u}_\beta^T c\,\underline{\varepsilon} - 2\lambda_\alpha(\underline{\varepsilon}^T \underline{u}_\alpha)\underline{u}_\beta^T \underline{u}_\alpha - \lambda_\alpha \underline{u}_\beta^T \underline{\varepsilon} \right)$$

$$= \eta(\lambda_\beta - 2\lambda_\alpha \delta_{\alpha\beta} - \lambda_\alpha)\underline{u}_\beta^T \underline{\varepsilon}$$

Recall: $\lambda_\alpha$ is the eigenvalue assigned to $\underline{w}^*$. Assume that this is not the maximal eigenvalue. In this case, thus, there will be at least one $\beta$ with $\lambda_\beta > \lambda_\alpha$. In this case, it follows that an initially small fluctuation around $\underline{w}^*$ (denoted by $\underline{\varepsilon}$ above) will grow! This is because the right-hand-side of the equation above is, in this case, positive:

$$\lambda_\beta > \lambda_\alpha \Rightarrow (\lambda_\beta - \underset{=0}{2\lambda_\alpha \delta_{\alpha\beta}} - \lambda_\alpha) = \lambda_\beta - \lambda_\alpha > 0$$

Therefore, in this case $\underline{w}^*$ is not the weight vector to which the network converges.

What happens if $\lambda_\alpha$ is the maximum eigenvalue? From the above argument, find that $\underline{\varepsilon}$ will shrink in size in all directions $\underline{u}_\beta$ ($\beta \neq \alpha$). What happens in the direction $\underline{u}_\alpha = \underline{w}^*$? In this direction $\underline{\varepsilon}$ also shrinks because the right-hand-side of the equation above is negative:

$$\lambda_\alpha - 2\lambda_\alpha - \lambda_\alpha = -2\lambda_\alpha < 0$$

Thus, we have shown that if the network converges to $\underline{w}^*$, then $\underline{w}^*$ is the leading eigenvector of $c$, and $|\underline{w}^*|^2 = 1$.

c) Generalisation of Oja's rule for learning M principal components for zero-mean data

$$\delta w_{ij} = \eta \, \varsigma_i \left( \xi_j - \sum_{k=1}^{M} \varsigma_k \, w_{kj} \right)$$

where $\varsigma_i = \sum_{j=1}^{N} w_{ij} \xi_j$ .

When $M=1$, this rule reduces to the rule (5) in the exam text.
Weight decay (second term in the rule) assures that the weight vectors remain normalised.