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

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

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

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

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

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

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

Rewrite the left-hand-side of Eq. (*):
$$S_i = \text{sgn}\left(\sum_{j=1}^{N} w_{ij} \zeta_j^{(\nu)}\right) = \text{sgn}\left[\sum_{\substack{j=1 \\ j \neq i}}^{N} \left(\frac{1}{N} \sum_{\mu=1}^{P} \zeta_i^{(\mu)} \zeta_j^{(\mu)}\right) \zeta_j^{(\nu)}\right] =$$

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

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

Rewrite the right hand side of (#):

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

Stability condition:

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

Stability condition satisfied when:

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

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 $\times$ $(-\varphi_i^{(\nu)})$)

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

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

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

$$C_i^{(\nu)} \approx -\frac{1}{N}\sum_{j=1}^{N}\sum_{\mu=1}^{P} \varphi_i^{(\mu)}\varphi_j^{(\mu)}\varphi_j^{(\nu)}\varphi_i^{(\nu)} \text{, for } N \gg 1$$

**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)

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

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

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

$$c_i^{(\nu)} \underset{N \gg 1}{\approx} - \frac{1}{N} \overbrace{\sum_{\substack{k=1 \\ i \neq j}}^{(p-1)(N-1)}}^{\text{random variables } (x_k)}_{\text{with } \pm 1}$$

$$\left[ (p-1)(N-1) \text{ terms} \right]$$

Since we assume $p \gg 1$ and $N \gg 1$, we can use
the Central limit theorem (patterns are <u>random</u>!)
Variables $x_k$ have the mean $\underline{0}$, and variance $\sigma_x^2 = 1$.
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 $\sigma^2$ of $c_i^{(\nu)}$ is :

$$\sigma^2 = \frac{1}{N^2} \cdot (N-1)(p-1)\, \sigma_x^2 \nearrow \approx \frac{p}{N}$$

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

It follows that

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

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

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

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

② Hopfield model: recognition of one pattern.

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

weight matrix: $\underline{\underline{w}} = \frac{1}{N}\underline{y}^{(1)}\,\underline{y}^{(1)T} = \frac{1}{4}\begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix}\begin{pmatrix} 1 & -1 & -1 & -1 \end{pmatrix} = \frac{1}{4}\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}$
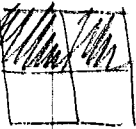
— Feeding in the $2^4$ possible patterns:

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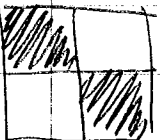
$\underline{S}_0 = \underline{y}^{(1)} = \begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix}$

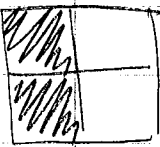2)  $\rightarrow \underline{S}_1 = sgn\left(\underline{\underline{w}}\,\underline{S}_0\right) = \frac{1}{4}\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ -1 \\ -1 \end{pmatrix}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

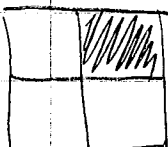$\underline{S}_0 = \begin{pmatrix} -1 \\ -1 \\ -1 \\ -1 \end{pmatrix}$

3) 

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

$\rightarrow S_1 = \text{sgn}(\underline{w}, \underline{S_o}) = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ -1 \\ -1 \end{pmatrix}\right] =$

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

4) 

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

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

5) 

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

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

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

6) 

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

$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ -1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

Orthogonal pattern to the stored pattern. The network doesnot restore the stored pattern, In fact, it retreives zero vector; failure of the network performance.

7)

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

$$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ 1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$
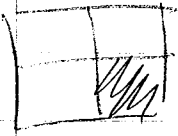
Same as Case 6 : orthogonal pattern.

8)

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

$$S_1 = \text{sgn}\left[\begin{pmatrix} -1 & 1 & 1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ -1 \\ 1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Same as Cases 6-7.

9)

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

$$S_1 = \text{sgn}\left[\begin{pmatrix} -1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ 1 \\ -1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$
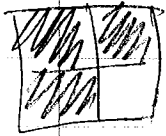
Same as Cases 6-8.

10)

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

$$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ -1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$
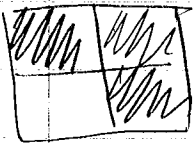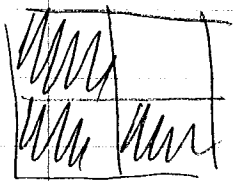
Same as Cases 6-9.

11)

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

$$S_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

Same as Cases 6-10.

**12)**



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

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

$$\boxed{\underline{S}_1 = -\underline{\xi}^{(1)}}$$

**13)**



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

$$\underline{S}_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} -1 \\ +1 \\ -1 \\ +1 \end{pmatrix}\right] = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix}$$

$$\boxed{\underline{S}_1 = -\underline{\xi}^{(1)}}$$

**14)**



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

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

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

**15)**



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

$$\underline{S}_1 = \text{sgn}\left[\begin{pmatrix} 1 & -1 & -1 & -1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \\ -1 & 1 & 1 & 1 \end{pmatrix}\begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}\right] = \begin{pmatrix} -1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \Rightarrow$$

$$\boxed{\underline{S}_1 = -\underline{\xi}^{(1)}}$$

**16)**



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

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

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)

cases 6-11

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).

[13] Back-propagation (two hidden layers)
- Two hidden layers.
- Input patterns $\underline{\xi}^{(\mu)} = (\xi_1, \xi_2, ..., \xi_N)^T$
- Target output $\zeta_1^{(\mu)}$
- Network output $O_1^{(\mu)}$

- 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: $O_1^{(\mu)} = g(b_1^{(\mu)})$, $b_1^{(\mu)} = \sum_k W_{1k} V_k^{(2,\mu)} - \Theta_1$

- 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:

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

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

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

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

$$= g'(b_1^{(\mu)}) \cdot V_k^{(2,\mu)} \qquad \left|\right| \text{Since } \frac{\partial W_{1\ell}}{\partial W_{1k}} = \delta_{\ell k}$$

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

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

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

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

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

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

- Second hidden layer

$$\delta w_{kj}^{(2)} = -\eta \frac{\partial H}{\partial w_{kj}^{(2)}} = -\eta \frac{\partial}{\partial w_{kj}^{(2)}} \left\{ \frac{1}{2} \sum_{\mu} (\mathcal{Y}_{1}^{(\mu)} - O_{1}^{(\mu)})^{2} \right\} =$$

$$= \eta \sum_{\mu} (\mathcal{Y}_{1}^{(\mu)} - O_{1}^{(\mu)}) \frac{\partial O_{1}^{(\mu)}}{\partial w_{kj}^{(2)}}$$

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

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

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

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

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

$$= g'(b_{1}^{(\mu)}) \cdot W_{1k} \, g'(b_{k}^{(2,\mu)}) \cdot V_{j}^{(1,\mu)}$$

$$\Rightarrow \delta w_{kj}^{(2)} = \eta \sum_{\mu} \underbrace{(\mathcal{Y}_{1}^{(\mu)} - O_{1}^{(\mu)}) \, g'(b_{1}^{(\mu)})}_{\delta_{1}^{(3,\mu)}} \cdot W_{1k} \, g'(b_{k}^{(2,\mu)}) V_{j}^{(1,\mu)}$$

$$\delta w_{kj}^{(2)} = \eta \sum_{\mu} \underbrace{\delta_{1}^{(3,\mu)} W_{1k} \, g'(b_{k}^{(2,\mu)})}_{\delta_{k}^{(2,\mu)}} V_{j}^{(1,\mu)}$$

$$\boxed{\delta w_{kj}^{(2)} = \eta \sum_{\mu} \delta_{k}^{(2,\mu)} V_{j}^{(1,\mu)}}$$

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

$$\delta\Theta_K^{(2)} = -\eta \frac{\partial H}{\partial\Theta_K^{(2)}} = \eta \sum_M \left(\varphi_1^{(M)} - O_1^{(M)}\right) \frac{\partial O_1^{(M)}}{\partial\Theta_K^{(2)}}$$

from previous page

$$\frac{\partial O_1^{(M)}}{\partial\Theta_K^{(2)}} \overset{\downarrow}{=} g'(b_1^{(M)}) \frac{\partial}{\partial\Theta_K^{(2)}} \left[\sum_\ell W_{1\ell} \, g\left(\sum_\lambda w_{\ell\lambda}^{(2)} V_\lambda^{(1,M)} - \Theta_\ell^{(2)}\right) - \Theta_1\right]$$

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

$$= - g'(b_1^{(M)}) \cdot W_{1K} \, g'(b_K^{(2,M)})$$

$$\Rightarrow \delta\Theta_K^{(2)} = -\eta \sum_M \underbrace{\left(\varphi_1^{(M)} - O_1^{(M)}\right) g'(b_1^{(M)})}_{\delta_1^{(3,M)}} W_{1K} \, g'(b_K^{(2,M)})$$

$$= -\eta \sum_M \underbrace{\delta_1^{(3,M)} W_{1K} \, g'(b_K^{(2,M)})}_{\delta_K^{(2,M)}}$$

$$\boxed{\delta\Theta_K^{(2)} = -\eta \sum_M \delta_K^{(2,M)}}$$

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:

$$\delta_K^{(2,M)} = \delta_1^{(3,M)} W_{1K} \, g'(b_K^{(2,M)})$$

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

layer:

$$\delta_j^{(1,M)} = \sum_K \delta_k^{(2,M)} w_{kj}^{(2)} g'\left(b_j^{(1,M)}\right)$$

The update formulae are, therefore, as follows:

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

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

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

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

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

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

Here we have the following:

$$\delta_1^{(3,M)} = \left(\zeta_1^{(M)} - O_1^{(M)}\right) g'\left(b_1^{(M)}\right), \quad b_1^{(M)} = \sum_K W_{1k} V_k^{(2,M)} - \Theta_1$$

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

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

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

(4) Backpropagation II — discussion of the implementation of the algorithm above. Explain how you program back-propagation.

(17) Oja's rule — Output $y = \sum_{j=1}^{N} w_j \xi_j \equiv \underline{w}^T \underline{\xi}$
$$\delta w_j = \eta y (\xi_j - y w_j)$$

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

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

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

For $\underline{w} = \underline{w}^*$, find $\langle y^2 \rangle = \langle \underbrace{\underline{w}^{*T} \mathbb{C} \underline{w}^*}_{\underline{w}^*} \rangle = \lambda_{max} \langle \underbrace{\underline{w}^{*T} \underline{w}^*}_{= 1} \rangle$

$\underbrace{\lambda_{max} \underline{w}^*}_{\text{(from ii)}}$  $\underbrace{= 1}_{\text{(from i)}}$

$$\Rightarrow \boxed{\langle y^2 \rangle_{\underline{w}^*} = \lambda_{max}}, \text{ where } \lambda_{max} \text{ is the maximum eigenvalue of } \mathbb{C}.$$

Since $\mathbb{C}$ is symmetric $(\langle \xi_i \xi_j \rangle = \langle \xi_j \xi_i \rangle)$ it has real eigenvalues $\lambda_\alpha$ and its eigenvectors $\underline{u}_\alpha$ are orthogonal:

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

Furthermore, all eigenvalues of $\mathbb{C}$ are positive, since

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

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

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

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

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

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

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

$$\Downarrow$$

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

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:

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

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

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

$$= \underbrace{\underline{\xi} \underline{\xi}^T \underline{w}^*}_{C}$$

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

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

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

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

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

$$\lambda = \underline{w}^{*T} \underbrace{C \underline{w}^*}_{\text{from } (\ast\#)} = \underline{w}^{*T} \lambda \underline{w}^* = \lambda \underline{w}^{*T} \underline{w}^* = \lambda |\underline{w}^*|^2$$

$$\Rightarrow \boxed{|\underline{w}^*|^2 = 1} \text{ Shown (i)}$$

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.

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

up to linear
order in $\varepsilon$      $= 0$ because $\underline{w}^*$ is steady
(previous page)

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

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

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

$$\left. - \langle \underline{w}^{*T} \underline{\xi} \, \underline{\varepsilon}^T \underline{\xi} \, \underline{w}^* \rangle \right]$$

Say that
$\underline{w}^* = \underline{u}_\alpha$
one of
the eigen-
vectors.
$\lambda_\alpha \underline{u}_\alpha$

$\mathbb{C}$

$$\Rightarrow \langle \delta(\underline{w}^* + \underline{\varepsilon}) \rangle \approx \eta \left[ \langle \underbrace{\underline{\xi} \, \underline{\xi}^T}_{} \underline{\varepsilon} \rangle - \langle \underline{\varepsilon}^T \underbrace{\underline{\xi} \, \underline{\xi}^T}_{=\lambda_\alpha \underline{u}_\alpha} \underline{w}^* \underbrace{\underline{w}^*}_{} \rangle \right.$$

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

$$\left. - \langle \underline{w}^{*T} \underbrace{(\underline{\xi} \, \underline{\xi}^T} \underline{w}^*) \underline{\varepsilon} \rangle - \langle \underline{w}^{*T} \underline{\xi} \underbrace{\underline{\xi}^T}_{} \underline{\varepsilon} \underline{w}^* \rangle \right]$$

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

$$= \eta \left[ \mathbb{C} \underline{\varepsilon} - \underline{\varepsilon}^T \underbrace{\lambda_\alpha \underline{u}_\alpha}_{=\lambda_\alpha} \underline{u}_\alpha \underbrace{}_{= \underline{\varepsilon}^T \underline{u}_\alpha} \right.$$

$$\left. - \underline{u}_\alpha^T \lambda_\alpha \underline{u}_\alpha \underline{\varepsilon} - \lambda_\alpha \underline{u}_\alpha^T \underline{\varepsilon} \underline{u}_\alpha \right]$$

$$= \eta \left[ \mathbb{C} \underline{\varepsilon} - 2\lambda_\alpha (\underline{\varepsilon}^T \underline{u}_\alpha) \underline{u}_\alpha - \lambda_\alpha \underline{\varepsilon} \right)$$

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

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

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

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

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:

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

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 $\phi$, and $|\underline{w}^*|^2 = 1$

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

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

where $\quad S_i = \sum_{j=1}^{N} w_{ij} \, \xi_j$ .

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.