

Machine Intelligence 2

1.2 Online-PCA / Hebbian Learning

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

SS 2018

Recap: PCA

assuming centered data ($\underline{\mathbf{m}}_a = 0$), we have complex features $\underline{\mathbf{u}}_a = \underline{\mathbf{X}}\mathbf{e}_a$

$$\sigma_\alpha^2 = \frac{1}{p} \underline{\mathbf{u}}_a^T \underline{\mathbf{u}}_a = \frac{1}{p} (\mathbf{e}_a^T \underline{\mathbf{X}}^T) \cdot (\underline{\mathbf{X}}\mathbf{e}_a) = \mathbf{e}_a^T \underline{\mathbf{C}}\mathbf{e}_a$$

Goal:

$$\mathbf{e}_a^* = \underset{\mathbf{e}_a}{\operatorname{argmax}} (\sigma_a^2) \quad \text{with} \quad \|\mathbf{e}_a\| = 1$$

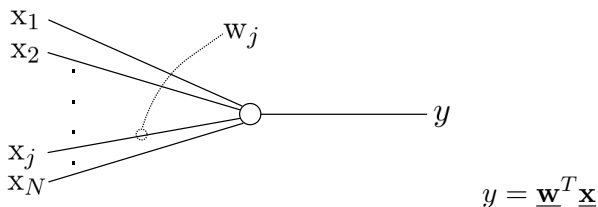
$$\underbrace{\mathbf{e}_a^T \underline{\mathbf{C}}\mathbf{e}_a}_{\text{objective}} - \lambda \underbrace{(\mathbf{e}_a^T \mathbf{e}_a - 1)}_{\text{constraints}} \stackrel{!}{=} \max$$

Constrained optimization \rightsquigarrow Eigenvalue problem

$$\underline{\mathbf{C}}\mathbf{e}_a = \lambda \mathbf{e}_a$$

\Rightarrow **Principal Components:** normalized eigenvectors \mathbf{e}_α of $\underline{\mathbf{C}}$

Linear connectionist neurons



observations: $\underline{\mathbf{x}}^{(\alpha)}, \quad \alpha = 1, \dots, p, \quad \underline{\mathbf{x}}^{(\alpha)} \in \mathbb{R}^N$

Hebbian learning (Donald Hebb, 1949): "fire together - wire together"

Hebbian learning

initialization of weights (e.g. to small numbers)

choose learning rate ε

begin loop

choose an observation $\underline{\mathbf{x}}^{(\alpha)}$

change weights according to:

$$\Delta \underline{\mathbf{w}} = \varepsilon y(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}}) \underline{\mathbf{x}}^{(\alpha)}$$

end

⇒ weights increase (decrease), if input and output are correlated (anticorrelated)

Proposition

Weight vector converges to the Principal Component with the largest eigenvalue.

Hebbian learning

assumption: centered data

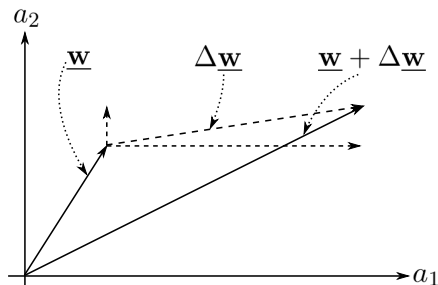
small learning steps \rightarrow average over patterns:

$$\begin{aligned}
 \Delta \underline{\mathbf{w}}_j &\approx \frac{\varepsilon}{p} \sum_{\alpha=1}^p y(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}}) \mathbf{x}_j^{(\alpha)} \\
 &= \frac{\varepsilon}{p} \sum_{\alpha=1}^p \sum_{k=1}^N w_k x_k^{(\alpha)} x_j^{(\alpha)} \\
 &= \varepsilon \sum_{k=1}^N w_k \left\{ \frac{1}{p} \sum_{\alpha=1}^p x_k^{(\alpha)} x_j^{(\alpha)} \right\} \\
 &= \varepsilon \sum_{k=1}^N w_k C_{kj} \\
 \Delta \underline{\mathbf{w}} &\approx \varepsilon \underline{\mathbf{C}} \underline{\mathbf{w}}
 \end{aligned}$$

Hebbian learning

eigenvectors of $\underline{\mathbf{C}}$: $\underline{\mathbf{e}}_1, \underline{\mathbf{e}}_2, \dots, \underline{\mathbf{e}}_N$ corresponding eigenvalues: $\lambda_1 > \lambda_2 > \dots > \lambda_N$

$$\underline{\mathbf{w}} = a_1 \underline{\mathbf{e}}_1 + a_2 \underline{\mathbf{e}}_2 + \dots + a_N \underline{\mathbf{e}}_N$$



$$\Delta \underline{\mathbf{w}} = \varepsilon \underline{\mathbf{C}} \underline{\mathbf{w}}$$

$$\Delta a_j = \varepsilon \lambda_j a_j$$

(see blackboard)

consequence

$$\blacksquare t \rightarrow \infty \quad \rightsquigarrow |\underline{\mathbf{w}}| \rightarrow \infty$$

$$\blacksquare \underline{\mathbf{e}}_{\underline{\mathbf{w}}} = \frac{\underline{\mathbf{w}}}{|\underline{\mathbf{w}}|} \text{ converges to } \underline{\mathbf{e}}_1 \text{ (eigenvector with the largest eigenvalue)}$$

Implicit normalization: Oja's rule

- adaptive tracking of the direction of largest variance: "on-line" PCA
- implicit normalization

Oja's rule

$$\Delta \mathbf{w}_j = \varepsilon y(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}}) \left\{ \underbrace{\mathbf{x}_j^{(\alpha)}}_{\text{Hebbian learning}} - \underbrace{y(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}}) \mathbf{w}_j}_{\text{decay term}} \right\}$$

Proposition

Oja's rule converges to the unit vector which points into the direction of the largest variance.

Derivation of Oja's rule

- Let $y^{(\alpha)} = y(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}}(t))$.
- Normalization of $|\underline{\mathbf{w}}(t)| = 1 \ \forall t$ is achieved by the learning rule

$$\underline{\mathbf{w}}(t+1) = \frac{\overbrace{\underline{\mathbf{w}}(t) + \varepsilon y^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)}}^{\text{Hebbian learning}}}{\underbrace{|\underline{\mathbf{w}}(t) + \varepsilon y^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)}|}_{\text{Euclidean weights normalization}}}$$

However

\leadsto Multiplicative constraint requires computation of the norm in each step

Derivation of Oja's rule

- Small learning step ε :

Taylor expansion around $\varepsilon = 0$ gives (\rightarrow calculation: exercise sheet)

$$\underline{\mathbf{w}}(t+1) = \underline{\mathbf{w}}(t) + \varepsilon \left\{ y^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{w}}(t) y^{(\alpha)} \left(\underline{\mathbf{w}}(t)^T \underline{\mathbf{x}}^{(\alpha)} \right) \right\} + \mathcal{O}(\varepsilon^2)$$

Oja's rule

$$\Delta \mathbf{w}_j = \varepsilon y(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}}) \left\{ \underbrace{x_j^{(\alpha)}}_{\text{Hebbian learning}} - \underbrace{y(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}}) \mathbf{w}_j}_{\text{decay term}} \right\}$$

Oja's rule = Hebbian Learning with weight normalization

Convergence properties of Oja's rule

a) The learning rule analyzes the covariance matrix $\underline{\mathbf{C}}$ of the data

Small learning steps \rightsquigarrow average over all patterns

$$\Delta \underline{\mathbf{w}} \approx \frac{1}{p} \sum_{\alpha=1}^p \overbrace{\varepsilon y^{(\alpha)} (\underline{\mathbf{x}}^{(\alpha)} - y^{(\alpha)} \underline{\mathbf{w}})}^{\text{Oja's rule}} = \varepsilon \left(\underbrace{\underline{\mathbf{C}} \underline{\mathbf{w}}}_{\text{Hebbian rule}} - \underbrace{(\underline{\mathbf{w}}^T \underline{\mathbf{C}} \underline{\mathbf{w}})}_{\text{decay term}} \underline{\mathbf{w}} \right)^{\geq 0}$$

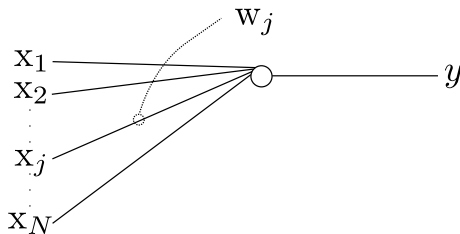
b) stationary states $\underline{\mathbf{w}}^*$ of Oja's rule $\hat{=}$ normalized eigenvector $\underline{\mathbf{e}}_j$ of $\underline{\mathbf{C}}$

Proof: See supplementary material

c) the stationary state $\underline{\mathbf{w}}^* = \underline{\mathbf{e}}_j$ is stable if and only if $\underline{\mathbf{e}}_j = \pm \underline{\mathbf{e}}_1$, i.e., if $\underline{\mathbf{e}}_j$ is the eigenvector with the largest eigenvalue λ_1

Proof: See supplementary material

Hebbian PCA (generalized Hebbian algorithm)

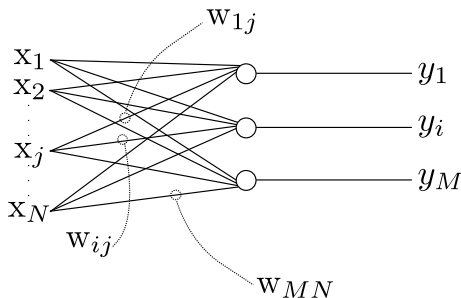


One linear neuron: $y = \underline{\mathbf{w}}^T \underline{\mathbf{x}}$

observations: $\underline{\mathbf{x}}^{(\alpha)}$, $\alpha = 1, \dots, p$, $\underline{\mathbf{x}}^{(\alpha)} \in \mathbb{R}^N$

Weight vector converges to the first principal component

Hebbian PCA (generalized Hebbian algorithm)



M linear neurons: $y_i = \underline{\mathbf{w}}_i^T \underline{\mathbf{x}} = \sum_{j=1}^N w_{ij} x_j, i = 1, \dots, M$

observations: $\underline{\mathbf{x}}^{(\alpha)}, \alpha = 1, \dots, p, \underline{\mathbf{x}}^{(\alpha)} \in \mathbb{R}^N$

The (feedforward) neural network extracts the M PCs with the largest eigenvalues

→ online-PCA for data with time-varying statistics

Hebbian PCA (generalized Hebbian algorithm)

Extended learning (Sanger's rule)

$$\Delta \mathbf{w}_{ij} = \varepsilon y_i \left\{ \underbrace{\mathbf{x}_j}_{\text{Hebbian rule}} - \underbrace{\sum_{k=1}^i \mathbf{w}_{kj} y_k}_{\text{is added to Oja's rule}} \right\}$$

→ weights converge to the M eigenvectors with the largest eigenvalues

$$\underline{\mathbf{w}}_1 \rightarrow \underline{\mathbf{e}}_1$$

$$\underline{\mathbf{w}}_2 \rightarrow \underline{\mathbf{e}}_2$$

$$\vdots$$

$$\underline{\mathbf{w}}_M \rightarrow \underline{\mathbf{e}}_M$$

→ $y_i = \underline{\mathbf{e}}_i^T \underline{\mathbf{x}} =: a_i$ after learning

Learning: Oja's rule & Gram-Schmidt orthonormalization

Sanger's rule: $\Delta w_{ij} = \varepsilon y_i \left\{ x_j - \sum_{k=1}^i w_{kj} y_k \right\}$

- Define $\hat{x}_j^{(i)} := x_j - \sum_{k=1}^{i-1} w_{kj} y_k$
- Then $\Delta w_{ij} = \varepsilon y_i \left\{ \hat{x}_j^{(i)} - y_j w_{ij} \right\} \longrightarrow$ Oja's rule with modified input

Case $i = 1$:

$$\begin{aligned} \hat{x}_j^{(1)} = x_j &\leadsto \text{original form of Oja's rule} \\ &\leadsto \underline{w}_1 \text{ converges to eigenvector } \pm \underline{e}_1 \end{aligned}$$

Learning: Oja's rule & Gram-Schmidt orthonormalization

Sanger's rule: $\Delta \mathbf{w}_{ij} = \varepsilon y_i \left\{ \mathbf{x}_j - \sum_{k=1}^i \mathbf{w}_{kj} y_k \right\}$

- Define $\hat{\mathbf{x}}_j^{(i)} := \mathbf{x}_j - \sum_{k=1}^{i-1} \mathbf{w}_{kj} y_k$
- Then $\Delta \mathbf{w}_{ij} = \varepsilon y_i \left\{ \hat{\mathbf{x}}_j^{(i)} - y_j \mathbf{w}_{ij} \right\} \rightarrow$ Oja's rule with modified input

Case $i = 2$:

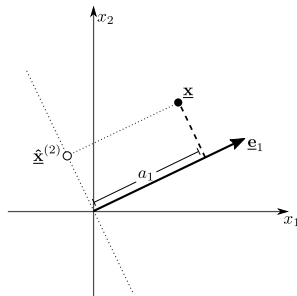
$$\hat{\mathbf{x}}_j^{(2)} = \mathbf{x}_j - \mathbf{w}_{1j} y_1$$

$$\underline{\mathbf{w}}_1 = \underline{\mathbf{e}}_1 \rightarrow y_1 = \underline{\mathbf{x}}^T \underline{\mathbf{e}}_1 =: a_1$$

$$\& \hat{\mathbf{x}}_j^{(2)} = \mathbf{x}_j - (\underline{\mathbf{e}}_1)_j a_1$$

$\leadsto \hat{\mathbf{x}}^{(2)}$ is the projection of $\underline{\mathbf{x}}$ onto subspace orthogonal to $\underline{\mathbf{e}}_1$:

$\leadsto \underline{\mathbf{w}}_2$ converges to $\pm \underline{\mathbf{e}}_2$ by Oja's rule since $\underline{\mathbf{e}}_2$ is the direction of largest variance in that subspace



Learning: Oja's rule & Gram-Schmidt orthonormalization

Sanger's rule: $\Delta \mathbf{w}_{ij} = \varepsilon y_i \left\{ \mathbf{x}_j - \sum_{k=1}^i \mathbf{w}_{kj} y_k \right\}$

- Define $\hat{\mathbf{x}}_j^{(i)} := \mathbf{x}_j - \sum_{k=1}^{i-1} \mathbf{w}_{kj} y_k$
- Then $\Delta \mathbf{w}_{ij} = \varepsilon y_i \left\{ \hat{\mathbf{x}}_j^{(i)} - y_j \mathbf{w}_{ij} \right\} \rightarrow$ Oja's rule with modified input

Case $i = 3$:

$$\hat{\mathbf{x}}_j^{(3)} = \mathbf{x}_j - \mathbf{w}_{1j} y_1 - \mathbf{w}_{2j} y_2$$

$$\underline{\mathbf{w}}_1 = \underline{\mathbf{e}}_1 \ \& \ \underline{\mathbf{w}}_2 = \underline{\mathbf{e}}_2 \rightarrow y_1 = a_1, y_2 = \underline{\mathbf{x}}^T \underline{\mathbf{e}}_2 =: a_2$$

$$\& \ \hat{\mathbf{x}}_j^{(3)} = \mathbf{x}_j - a_1 (\underline{\mathbf{e}}_1)_j - a_2 (\underline{\mathbf{e}}_2)_j$$

$\leadsto \hat{\underline{\mathbf{x}}}^{(3)}$ is the projection of $\underline{\mathbf{x}}$ onto subspace orthogonal to $\text{span}\{\underline{\mathbf{e}}_1, \underline{\mathbf{e}}_2\}$

$\leadsto \underline{\mathbf{w}}_3$ converges to $\pm \underline{\mathbf{e}}_3$ by Oja's rule

Learning: Oja's rule & Gram-Schmidt orthonormalization

Sanger's rule: $\Delta \mathbf{w}_{ij} = \varepsilon y_i \left\{ \mathbf{x}_j - \sum_{k=1}^i \mathbf{w}_{kj} y_k \right\}$

- Define $\hat{\mathbf{x}}_j^{(i)} := \mathbf{x}_j - \sum_{k=1}^{i-1} \mathbf{w}_{kj} y_k$

- Then $\Delta \mathbf{w}_{ij} = \varepsilon y_i \left\{ \hat{\mathbf{x}}_j^{(i)} - y_j \mathbf{w}_{ij} \right\} \rightarrow$ Oja's rule with modified input

⋮

Case $i = M$:

$$\hat{\mathbf{x}}_j^{(M)} = \mathbf{x}_j - \sum_{k=1}^{M-1} \mathbf{w}_{kj} y_k$$

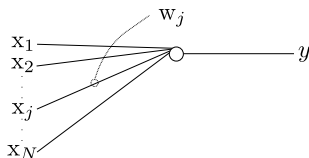
$$\underline{\mathbf{w}}_k = \underline{\mathbf{e}}_k \text{ for } k = 1, \dots, M-1 \rightarrow y_k = \underline{\mathbf{x}}^T \underline{\mathbf{e}}_k =: a_k$$

$$\& \hat{\mathbf{x}}_j^{(M)} = \mathbf{x}_j - \sum_{k=1}^{M-1} a_k (\underline{\mathbf{e}}_k)_j$$

$\leadsto \hat{\underline{\mathbf{x}}}^{(M)}$ is the proj. of $\underline{\mathbf{x}}$ onto subspace orthogonal to $\text{span} \{ \underline{\mathbf{e}}_1, \dots, \underline{\mathbf{e}}_{M-1} \}$

$\leadsto \underline{\mathbf{w}}_M$ converges to $\pm \underline{\mathbf{e}}_M$ by Oja's rule

Summary of Hebbian learning



$$y = \underline{\mathbf{w}}^T \underline{\mathbf{x}}$$

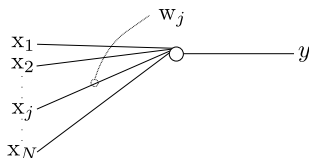
1. Hebbian learning without constraint

$$\underbrace{\Delta \underline{\mathbf{w}} = \varepsilon y \underline{\mathbf{x}}}_{\text{Hebb's rule}} \leadsto \lim_{t \rightarrow \infty} \underline{\mathbf{w}} \parallel \underline{\mathbf{e}}_1 \text{ (orthogonal)}$$

\leadsto weights converge to direction of largest variance in the data

\leadsto but: $|\underline{\mathbf{w}}| \rightarrow \infty$ for $t \rightarrow \infty$

Summary of Hebbian learning



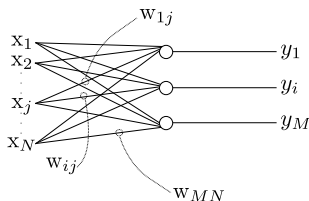
$$y = \underline{\mathbf{w}}^T \underline{\mathbf{x}}$$

II. Hebbian learning with normalization

$$\underbrace{\Delta \underline{\mathbf{w}} = \varepsilon y (\underline{\mathbf{x}} - y \underline{\mathbf{w}})}_{\text{Oja's rule}} \leadsto \lim_{t \rightarrow \infty} \underline{\mathbf{w}} \in \{+\underline{\mathbf{e}}_1, -\underline{\mathbf{e}}_1\}$$

\leadsto weights remain finite: $|\underline{\mathbf{w}}| = 1$

Summary of Hebbian learning



$$\mathbf{y} = \underline{\mathbf{w}}^T \underline{\mathbf{x}}$$

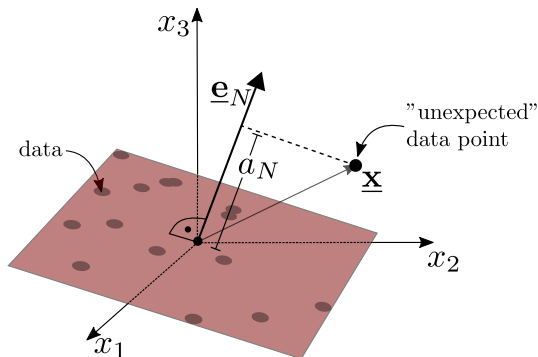
III. Hebbian PCA with M neurons and normalization

$$\underbrace{\Delta w_{ij} = \varepsilon y_i \left\{ x_j - \sum_{k=1}^i w_{kj} y_k \right\}}_{\text{Sanger's rule}} \rightsquigarrow \lim_{t \rightarrow \infty} \underline{\mathbf{w}}_i \in \{+\underline{\mathbf{e}}_i, -\underline{\mathbf{e}}_i\}, i = 1, \dots, M$$

\rightsquigarrow combination of Oja's rule & Gram-Schmidt-orthonormalization

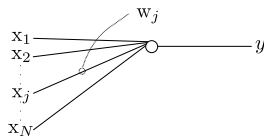
Novelty Filter

Reminder



$y = \underline{e}_N^T \underline{x} =: a_N$ after learning \leadsto projection onto smallest PC
 \leadsto large output for unexpected data \rightarrow Novelty Filter

Novelty Filter: On-line Learning



$$y = \underline{\mathbf{w}}^T \underline{\mathbf{x}}$$

Anti-Hebbian rule:

$$\Delta \mathbf{w}_j = \underbrace{\text{"Anti"-Hebbian}}_{-} \varepsilon y^{(\alpha)} \mathbf{x}_j^{(\alpha)}$$

Novelty Filter: On-line Learning

Conjecture:

\mathbf{w} converges to the direction of smallest eigenvector.

Proof:

Learning rule:

$$\Delta \mathbf{w}_j = -\varepsilon y^{(\alpha)} \mathbf{x}_j^{(\alpha)}$$

Assume small learning steps \rightarrow average over all patterns

$$\Delta \mathbf{w}_j = -\frac{\varepsilon}{p} \sum_{\alpha=1}^p y^{(\alpha)} \mathbf{x}_j^{(\alpha)} = -\frac{\varepsilon}{p} \sum_{\alpha=1}^p \mathbf{x}_j^{(\alpha)} \sum_{k=1}^N x_k^{(\alpha)} \mathbf{w}_k = -\varepsilon \sum_{k=1}^N C_{jk} \mathbf{w}_k$$

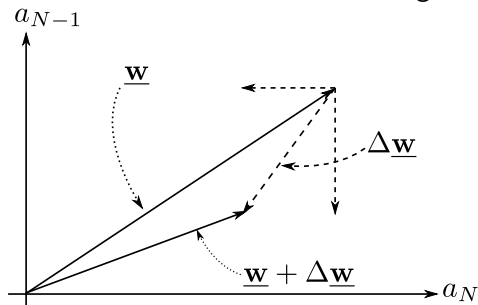
$$\Delta \underline{\mathbf{w}} = -\varepsilon \underline{\mathbf{C}} \underline{\mathbf{w}}$$

Novelty Filter: On-line Learning

Proof cont.:

$$\Delta \underline{\mathbf{w}} = -\varepsilon \underline{\mathbf{C}} \underline{\mathbf{w}}$$

Transformation into eigenbasis of covariance matrix:

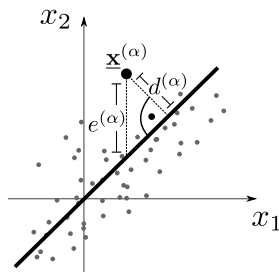


$$\underline{\mathbf{w}} = a_1 \underline{\mathbf{e}}_1 + a_2 \underline{\mathbf{e}}_2 + \cdots + a_N \underline{\mathbf{e}}_N$$

$$\Delta a_j = -\varepsilon \lambda_j a_j$$

- ~ for $\lambda_j > 0 : a_j \rightarrow 0$, constraints required
- ~ for $\lambda_j = 0 : a_j$ remains unchanged
- ~ weights converge to the eigenvector with the smallest eigenvalue

Novelty Filter and linear regression



ordinary least squares:

$$\frac{1}{p} \sum_{\alpha=1}^p \left(e^{(\alpha)} \right)^2 \stackrel{!}{=} \min.$$

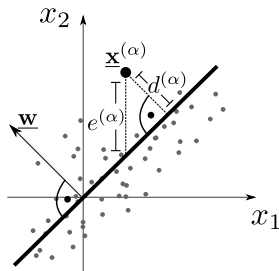
~> correct if data points are noisy along x_2 -component only

~> wrong if data points are also noisy along x_1 -component

Novelty Filter and linear regression

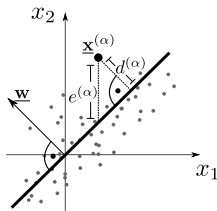
total least squares:

$$\frac{1}{p} \sum_{\alpha=1}^p \left(d^{(\alpha)} \right)^2 \stackrel{!}{=} \min.$$



tacit assumption: same variance noise
centered data $\rightarrow \underline{\mathbf{w}}^T \underline{\mathbf{x}} = 0$

Novelty Filter and linear regression



Cost function:

$$\mathbb{E}(\underline{w}) = \frac{1}{p} \sum_{\alpha=1}^p \left(d^{(\alpha)} \right)^2 \stackrel{!}{=} \min_{\underline{w}} \quad \text{s.t. } |\underline{w}| = 1$$

$$\underline{w}^T \underline{C} \underline{w} \stackrel{!}{=} \min_{\underline{w}} \quad \text{s.t. } |\underline{w}| = 1$$

solution:

\underline{w} is the normalized eigenvector to the smallest eigenvalue of the covariance matrix.

Novelty Filter with normalization

$$\Delta \underline{\mathbf{w}} = -\varepsilon \frac{y^{(\alpha)} \left\{ \underline{\mathbf{x}}^{(\alpha)} - y^{(\alpha)} \underline{\mathbf{w}} \right\}}{\left| \underline{\mathbf{w}} - \varepsilon y^{(\alpha)} \left\{ \underline{\mathbf{x}}^{(\alpha)} - y^{(\alpha)} \underline{\mathbf{w}} \right\} \right|}$$

Anti-Hebbian version of Oja's rule:

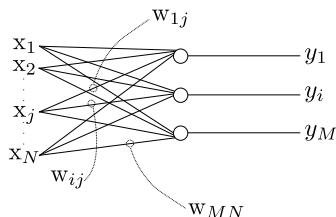
$$\Delta w_j = -\varepsilon y^{(\alpha)} \left\{ \overbrace{x_j^{(\alpha)}}^{\text{Anti-Hebbian learning}} - y^{(\alpha)} w_j \right\} + \varepsilon \underbrace{\left\{ 1 - \sum_{k=1}^N \overbrace{w_k^2}^{=|\underline{\mathbf{w}}|^2} \right\} w_j}_{\text{normalization}}$$

$\leadsto \underline{\mathbf{w}}$ converges to $\pm \underline{\mathbf{e}}_N$

Feedforward network as a Novelty Filter

Extension of the learning rule to N neurons:

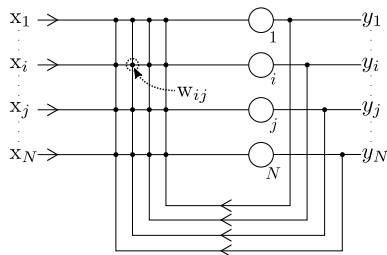
$$\Delta w_{ij} = -\varepsilon y_i^{(\alpha)} \left\{ x_j^{(\alpha)} - \underbrace{\sum_{k=1}^{i-1} w_{kj} y_k^{(\alpha)}}_{\text{is added}} \right\} + \varepsilon \left\{ 1 - \sum_{k=1}^N w_{ik}^2 \right\} w_{ij}$$



\leadsto result: $\underline{\mathbf{w}}_1 \rightarrow \underline{\mathbf{e}}_N$ (PC with smallest eigenvalue)
 $\underline{\mathbf{w}}_2 \rightarrow \underline{\mathbf{e}}_{N-1}$ \vdots
 \vdots \vdots
 $\underline{\mathbf{w}}_M \rightarrow \underline{\mathbf{e}}_{N-M+1}$ (PC with largest eigenvalue, if $N = M$)

Sequential calculation of eigenvectors (cf. Sanger's rule).

Recurrent network as a Novelty Filter



$$y_i^{(\alpha)}(t+1) = \sum_{j=1}^N w_{ij} y_j^{(\alpha)}(t) + x_i^{(\alpha)}$$

$$\underline{\mathbf{y}}^{(\alpha)}(t+1) = \underline{\mathbf{W}} \underline{\mathbf{y}}^{(\alpha)}(t) + \underline{\mathbf{x}}^{(\alpha)}$$

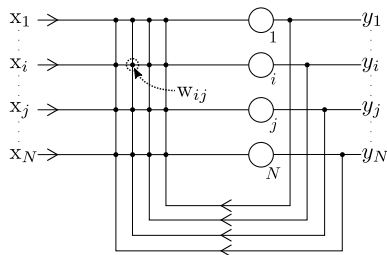
Stationary state:

$$\tilde{\underline{\mathbf{y}}}^{(\alpha)}(t+1) = \tilde{\underline{\mathbf{y}}}^{(\alpha)}(t) =: \tilde{\underline{\mathbf{y}}}^{(\alpha)}$$

$$\tilde{\underline{\mathbf{y}}}^{(\alpha)} = (\underline{\mathbf{I}} - \underline{\mathbf{W}})^{-1} \underline{\mathbf{x}}^{(\alpha)}$$

Convergence is guaranteed, if weight matrix is symmetric.

Recurrent network as a Novelty Filter



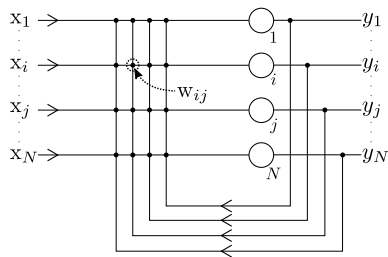
$$y_i^{(\alpha)}(t+1) = \sum_{j=1}^N w_{ij} y_j^{(\alpha)}(t) + x_i^{(\alpha)}$$

$$\underline{\mathbf{y}}^{(\alpha)}(t+1) = \underline{\mathbf{W}} \underline{\mathbf{y}}^{(\alpha)}(t) + \underline{\mathbf{x}}^{(\alpha)}$$

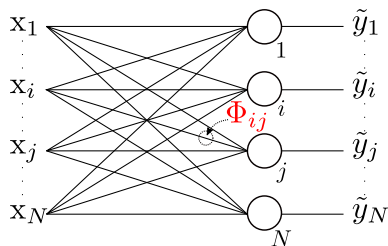
Learning rule:

$$\Delta w_{ij} = -\varepsilon \tilde{y}_i \tilde{y}_j$$

Recurrent network as a Novelty Filter

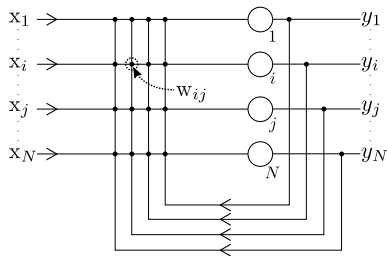


$$\underline{\tilde{\mathbf{y}}}^{(\alpha)} = (\underline{\mathbf{I}} - \underline{\mathbf{W}})^{-1} \underline{\mathbf{x}}^{(\alpha)}$$

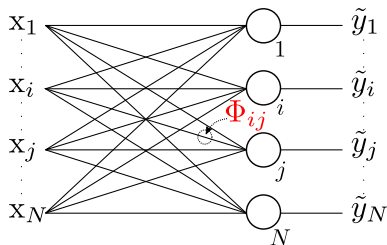


$$\underline{\tilde{\mathbf{y}}}^{(\alpha)} = \underline{\Phi} \underline{\mathbf{x}}^{(\alpha)}$$

Recurrent network as a Novelty Filter



$$\Delta \underline{\mathbf{w}} = -\varepsilon \tilde{\underline{\mathbf{y}}}^{(\alpha)} \left(\tilde{\underline{\mathbf{y}}}^{(\alpha)} \right)^T$$



$$\Delta \underline{\Phi} = -\varepsilon \underline{\Phi}^2 \underline{\mathbf{x}}^{(\alpha)} \left(\underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\Phi}^2$$

Recurrent network as a Novelty Filter

initialization of $\underline{\Phi}$ with identity matrix, $\underline{\Phi} = \underline{\mathbf{I}}$

repeat

choose an observation $\underline{\mathbf{x}}^{(\alpha)}$

change weight matrix according to:

$$\Delta \underline{\Phi} = -\varepsilon \underline{\Phi}^2 \underline{\mathbf{x}}^{(\alpha)} \left(\underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\Phi}^2$$

until *convergence*

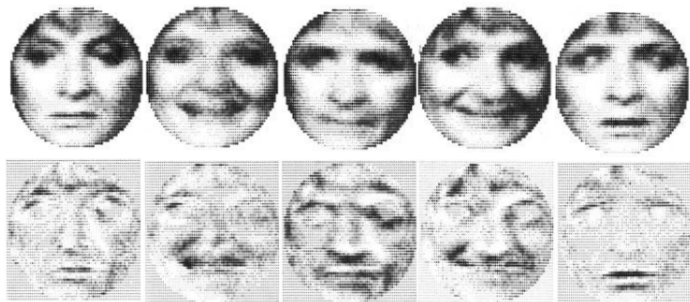
\rightsquigarrow $\underline{\Phi}$ converges to a matrix, which projects onto the subspace orthogonal to the training data

example:

training data $\left\{ \underline{\mathbf{x}}^{(1)}, \dots, \underline{\mathbf{x}}^{(p)} \right\} \subseteq \text{span} \left\{ \underline{\mathbf{e}}_1, \dots, \underline{\mathbf{e}}_{N-1} \right\} \quad \underline{\mathbf{x}}^{(\alpha)} \in \mathbb{R}^N$

$$\underline{\Phi} \xrightarrow[t \rightarrow \infty]{} \underline{\mathbf{I}} - \sum_{k=1}^{N-1} \underline{\mathbf{e}}_k \underline{\mathbf{e}}_k^T \rightsquigarrow \underline{\Phi} \underline{\mathbf{e}}_j = \underline{\mathbf{e}}_N \delta_{jN} \rightsquigarrow \underline{\Phi} \underline{\mathbf{x}} = \underbrace{\left(\underline{\mathbf{e}}_N^T \underline{\mathbf{x}} \right)}_{=: a_N} \underline{\mathbf{e}}_N$$

Recurrent network as a Novelty Filter

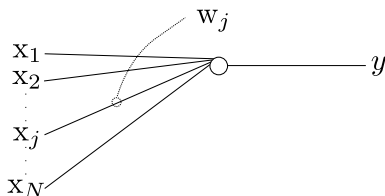


(taken from Kohonen 1989)

- training data: "neutral" facial expressions (not shown here)
 - ↪ Φ projects into space orthogonal to this data
- top row: faces $\underline{x}^{(\beta)}$ with different expressions
- bottom row: projection $\Phi \underline{x}^{(\beta)}$

Novelty Filter: Summary

One neuron:



$$y = \underline{\mathbf{w}}^T \underline{\mathbf{x}}$$

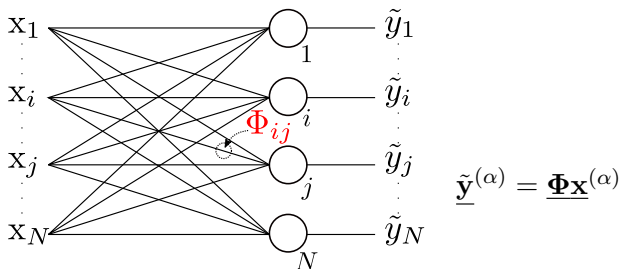
Anti-Hebbian learning rule

$$\Delta w_j = -\varepsilon y^{(\alpha)} \left\{ \underbrace{x_j^{(\alpha)}}_{\text{Anti-Hebbian learning}} - y^{(\alpha)} w_j \right\} + \varepsilon \underbrace{\left\{ 1 - \sum_{k=1}^N w_k^2 \right\}}_{\text{normalization}} w_j$$

$\sum_{k=1}^N w_k^2 = |\underline{\mathbf{w}}|^2$

Novelty Filter: Summary

N neurons:



Learning rule:

$$\Delta \underline{\Phi} = -\varepsilon \underline{\Phi}^2 \underline{\mathbf{x}}^{(\alpha)} \left(\underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\Phi}^2$$

\leadsto $\underline{\Phi}$ converges to projection matrix onto subspace orthogonal to training data