

Adaptive lineare Transformationen AS2-3

PCA-Transformation

PCA-Netze und Weissen

ICA-Transformation

Hebb'sches Lernen

$$\Delta w = w_i(t) - w_i(t-1) = \gamma_i(t) y_i x \quad \text{Iterative Hebb'sche Lernregel}$$

$$\Delta W = W(t) - W(t-1) = \gamma(t) y x^T$$

$$W = W(1) + W(2) + W(3) + \dots$$

Problem: ex. kein „Vergessen“, $w \rightarrow \infty$

Unendliches Wachstum ??

Hebb'sches Lernen - Ergänzungen

Lösung 1: lin. Term, Abklingen der Synapsen

$$\mathbf{w}(t) = \mathbf{w}(t-1) + \gamma(t) \mathbf{y}\mathbf{x} \quad \text{Iterative Hebb'sche Lernregel}$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) + \gamma(t) (\mathbf{y}\mathbf{x} - \mathbf{w}(t-1)) \quad \text{Abklingen des Gewichts}$$

$$\gamma^{-1} \frac{\partial}{\partial t} \mathbf{w} \approx \mathbf{w}(t) - \mathbf{w}(t-1) = \mathbf{y}\mathbf{x} - \mathbf{w}(t-1) \quad \text{Diff.gleichung mit } \tau = 1/\gamma$$

$$\text{Erwartetes Ziel bei lin. System } \mathbf{y} = \mathbf{w}^T \mathbf{x} \quad \mathbf{C}_{xx} := \langle \mathbf{x}\mathbf{x}^T \rangle$$

$$\frac{\partial}{\partial t} \mathbf{w} = \gamma_1 \langle \mathbf{x}\mathbf{y} \rangle - \gamma_2 \mathbf{w} = \gamma_1 \langle \mathbf{x}\mathbf{x}^T \mathbf{w} \rangle - \mathbf{w} = \gamma_1 \langle \mathbf{x}\mathbf{x}^T \rangle \mathbf{w} - \gamma_2 \mathbf{w} = \mathbf{0}$$

$$\Leftrightarrow \mathbf{C}_{xx} \mathbf{w} = \lambda \mathbf{w} \quad \text{bei Fixpunkt } \mathbf{w} = \mathbf{w}^* \text{ Eigenvektor von } \mathbf{C}_{xx}$$

\mathbf{w}^* stabil? 1-dim Beispiel: **Nein!**

$$\frac{\partial}{\partial t} \mathbf{w} = \gamma_1 \mathbf{x}\mathbf{y} - \gamma_2 \mathbf{w}^n \quad \text{nicht-lin. Abklingterm } m > 1$$

Hebb'sches Lernen - Ergänzungen

Lösung : Normierung der Gewichte

$$\mathbf{w}(t) = \mathbf{w}(t-1) + \gamma(t) \mathbf{y}\mathbf{x} \quad \text{mit } |\mathbf{w}(t)| = 1$$

Wie?

$$\hat{\mathbf{w}}(t) = \mathbf{w}(t-1) + \gamma(t) \mathbf{y}\mathbf{x}$$

$$\mathbf{w}(t) = \frac{\hat{\mathbf{w}}}{|\hat{\mathbf{w}}|} = \frac{\mathbf{w}(t-1) + \gamma(t) \mathbf{y}\mathbf{x}}{|\mathbf{w}(t-1) + \gamma(t) \mathbf{y}\mathbf{x}|}$$

Wohin konvergiert $\mathbf{w}(t)$?

Lernen: beschränkte Hebb'sche Regel

Konvergenzziel?

$$\mathbf{w}_i(t) = \mathbf{w}_i(t-1) + \gamma_i(t) \langle \mathbf{y}_i, \mathbf{x} \rangle \quad \text{Hebb'sche Lernregel mit NB } |\mathbf{w}| = \text{const.}$$

$$= \mathbf{w}_i(t-1) + \gamma_i(t) \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}} \quad \text{Gradientenaufstieg}$$

Also:

$$\frac{\partial R(\mathbf{w})}{\partial \mathbf{w}} = \langle \mathbf{x}\mathbf{y} \rangle = \langle \mathbf{x}\mathbf{x}^T \rangle \mathbf{w} = \mathbf{A}\mathbf{w} \quad \text{bei } \mathbf{y} = \mathbf{x}^T \mathbf{w} \text{ lin. Neuron}$$

$$\text{Zielfunktion } R(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{A}\mathbf{w} \quad \text{mit } |\mathbf{w}| = \text{const} = 1$$

Extremwert von $L(\mathbf{w}, \mu) = R(\mathbf{w}) + \mu \cdot (\mathbf{w}^2 - 1)$ Lagrangefunktion

$$\text{bei } \frac{\partial L}{\partial \mathbf{w}}(\mathbf{w}, \mu) = \frac{\partial R(\mathbf{w})}{\partial \mathbf{w}} + \mu 2\mathbf{w} = \mathbf{A}\mathbf{w} + \mu 2\mathbf{w} = \mathbf{0}$$

$$\mathbf{A}\mathbf{w} = \lambda \mathbf{w} \quad \text{Eigenwertgleichung} \quad \text{mit } \lambda := -2\mu$$

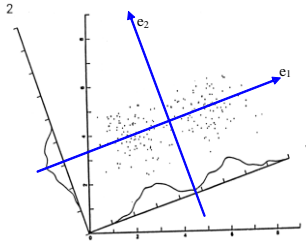
$$\mathbf{w} \rightarrow \text{EV}(\mathbf{A}) \text{ mit EW } \lambda = \max$$

Principal Component Analysis PCA

Zerlegung in orthogonale Eigenvektoren = Basisvektoren

„Hauptkomponentenanalyse“, „principal component analysis“, „Karhunen-Loève-Entwicklung“, „Hotelling-Transformation“, ...

Eigenvektoren – Wozu?



Merkmals-
transformation
auf
Hauptrichtungen

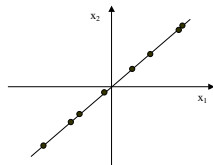
Principal Component Analysis PCA

Transformation auf Unkorreliertheit

$$\langle (x_1 - \langle x_1 \rangle) (x_2 - \langle x_2 \rangle) \rangle = 0$$

Unkorreliertheit von x_1, x_2

Beispiel



Rauschfrei korrelierte Daten

$$\mathbf{x} = (x_1, x_2) \text{ mit } x_2 = a x_1$$

Rechnung: EV, EW = ?

Dekorrelation und Unabhängigkeit

• **DEF dekorreliert:** $\langle (y_i - \langle y_i \rangle) (y_j - \langle y_j \rangle) \rangle = \langle y_i - \langle y_i \rangle \rangle \langle y_j - \langle y_j \rangle \rangle = 0 \quad i \neq j$

• **Satz: PCA dekorreliert Daten**

DEF PCA: $\mathbf{y} = \mathbf{w}^T \mathbf{x}$ mit $\mathbf{C}_{xx} \mathbf{w} = \lambda \mathbf{w}$ Eigenvektoren

⇒ Mit PCA gilt

$$y_i = y_i - \langle y_i \rangle, \quad \mathbf{x}' = \mathbf{x} - \bar{\mathbf{x}}$$

$$\langle y_i y_j \rangle = \langle \mathbf{w}_i^T \mathbf{x}' \mathbf{x}'^T \mathbf{w}_j \rangle = \mathbf{w}_i^T \langle \mathbf{x}' \mathbf{x}'^T \rangle \mathbf{w}_j = \mathbf{w}_i^T \mathbf{C}_{xx} \mathbf{w}_j = \mathbf{w}_i^T \lambda_j \mathbf{w}_j$$

$$= \begin{cases} \lambda_i & \text{bei } i = j \\ 0 & \text{bei } i \neq j \end{cases} \quad \text{da ja } \mathbf{w}_i^T \mathbf{w}_j = \begin{cases} 1 & \text{bei } i = j \\ 0 & \text{bei } i \neq j \end{cases} \text{ gilt.}$$

• Daten sind **unabhängig** ⇒ Daten sind **dekorreliert**

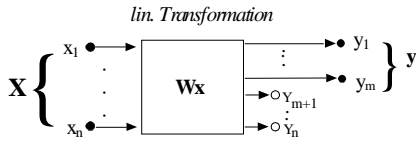
DEF **unabhängig:** $P(x_i, x_j) = P(x_i)P(x_j)$ x_i, x_j Zufallsvariable

$$\Rightarrow \langle y_i y_j \rangle = \sum_{y_i, y_j} P(y_i, y_j) y_i y_j = \sum_{y_i} P(y_i) P(y_j) y_i y_j = \sum_{y_i} P(y_i) y_i \sum_{y_j} P(y_j) y_j = \langle y_i \rangle \langle y_j \rangle = 0 \text{ bei } \langle y_i \rangle = 0 \quad i \neq j$$

Aber: umgekehrt gilt nicht: dekorreliert ist **nicht** unabhängig!

Transformation mit minimalem MSE

Allgemeine Situation



$$\min_{\mathbf{w}} R(\mathbf{W}) = \min \langle (\mathbf{x} - \hat{\mathbf{x}})^2 \rangle \quad \text{least mean squared error (LMSE)}$$

Wann minimal ?

Transformation mit minimalem MSE

Minimaler Rekonstruktionsfehler

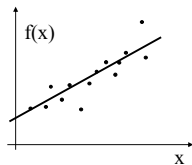
$$\min_{\mathbf{w}} R(\mathbf{W}) = \min \langle (\mathbf{x} - \hat{\mathbf{x}})^2 \rangle \quad \text{least mean squared error (LMSE)}$$

$$\mathbf{x} = \sum_{i=1}^m y_i \mathbf{w}_i + \sum_{i=m+1}^n y_i \mathbf{w}_i \quad \hat{\mathbf{x}} = \sum_{i=1}^m y_i \mathbf{w}_i + \sum_{i=m+1}^n c_i \mathbf{w}_i \quad y_i = \mathbf{x}^T \mathbf{w}_i$$

- Was ist die beste Schätzung für die Konstanten c_i ?
 $\min R(c_i) = ?$ **Rechnung!** Anhang B
- Bei welchen Basisvektoren \mathbf{w}_i ist der Fehler minimal?
 $\min R(\mathbf{w}_i) = ?$ **Rechnung!** Anhang B

Transformation mit minimalem MSE

m Messungen

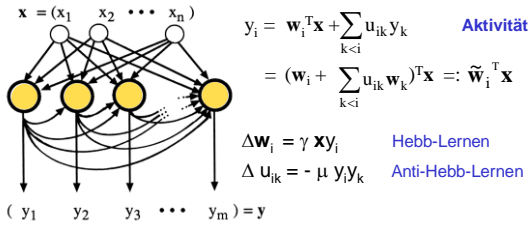


- Modellierung als Gerade
 $y = f(x) = y_0 + ax$
- Beispiel: Ökonomie
Konsum $y = f(\text{Einkommen } x)$
= Konsumsokkel + $a \cdot \text{Einkommen}$
- Problem: Ergebnis hängt vom Koord.system ab

PCA Netze durch laterale Inhibition

Asymmetrische Netze

Rubner, Tavan 1990



Anti-Hebb auch aus Prinzip „kleinste gemeins. Information“ $H(y_i, y_j)$

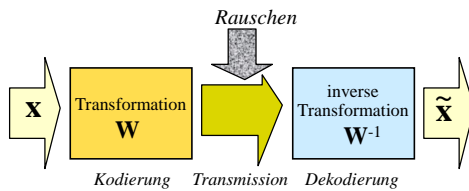
Whitening Filter

Shannon: Whitening für alle Frequenzen, d.h. alle diskreten Signalfrequenzen
Übertragung auf parallele Signale y_i : gleiche Varianz aller durch Transformation W .

Anhebung zu geringer Amplituden: Wähle W so,

$$\text{dass } \langle y_i y_j \rangle = 1 \text{ bei } i = j, \text{ sonst } = 0; \quad \text{also } \langle y y^T \rangle = I$$

Absenkung der Amplituden: durch inverse Matrix W^{-1}



Whitening Filter

Anhebung bei parallelen Signalen

Wenn für die Transformation W eine orthonormale Basis $M^{-1} = M^T$ gewählt wird, ist das Ziel des Lernens

mit $I = \langle y y^T \rangle = \langle W x x^T W^T \rangle = W \langle x x^T \rangle W^T = W A W^T$

auch $W^T I = W^T (W A W^T) = A W^T$

bzw. $w_k = A w_k$ **Eigenvektoren w_k von A mit $\lambda = 1$**

Also:

- Signal zentrieren und PCA durchführen. Wir erhalten orthonormale Eigenvektoren e_i mit Eigenwerten λ_i .
- e_i normieren: $w_i = e_i / \lambda_i^{1/2}$ so dass $\|w_i\|^2 = 1 / \lambda_i$.
Es ergibt sich orthogonale Transformationsmatrix W mit Zeilenvektoren w_i

Sie erfüllt $\langle y_i y_j \rangle = \langle w_i^T x x^T w_j \rangle = w_i^T A w_j = w_i^T w_j \lambda_i = 1$

Whitening Filter

Abseitung (Rücktransformation) $W^{-1} = ?$

Wenn $B = (\lambda_1 w_1, \dots, \lambda_n w_n)$ gewählt wird, ist mit $|w_i|^2 = \lambda_i^{-1}$

$$\text{mit } W \cdot B = \begin{pmatrix} w_1^T & \dots & w_n^T \\ \dots & \dots & \dots \\ w_n^T & \dots & w_n^T \end{pmatrix} \begin{pmatrix} \lambda_1 w_1 & \dots & \lambda_n w_n \\ \dots & \dots & \dots \\ \lambda_1 w_1 & \dots & \lambda_n w_n \end{pmatrix} = \begin{pmatrix} 1 & \dots & 0 \\ 0 & \ddots & 0 \\ 0 & \dots & 1 \end{pmatrix} = I,$$

Also ist $W^{-1} = B$ mit den Spalten aus den Zeilen von W .

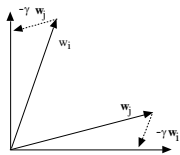
Also: Rücktransformation

- Aus der PCA haben wir e_i, λ_i mit $|e_i|^2 = 1$ und so die Matrix W mit $|w_i|^2 = \lambda_i^{-1}$
- Basis b_i bilden aus W : $b_i = (\lambda_1 w_{1i}, \lambda_2 w_{2i}, \dots, \lambda_n w_{ni})$

Orthonormalisierende Netze

Heuristische Methode

Silva, Almeida 1991



Ziel: Projektion $a_{ij} := w_j^T w_i$ eines Basisvektors w_i auf einen anderen w_j vermindern

$$\begin{aligned} w_i(t) &= w_i(t-1) - \gamma(t) a_{ij} w_j(t-1) \\ &= w_i(t-1) - \gamma(t) (w_j^T(t-1) w_i(t-1)) w_j(t-1) \\ &\quad + \gamma(1 - \langle y_i, y_j \rangle) w_i(t) \quad \text{Normierung} \end{aligned}$$

Ziel: orthogonal im Datenraum: $a_{ij} := \langle y_i, y_j \rangle = 0$

Alle Einflüsse

$$w_i(t) = w_i(t-1) - \gamma(t) \sum_{j \neq i} \langle y_i, y_j \rangle w_j(t-1) + \gamma(1 - \langle y_i, y_i \rangle) w_i(t)$$

$$W(t+1) = W(t) - \gamma C_{yy}(t) W(t) + \gamma W(t) \quad \text{Matrixversion}$$

Orthonormalisierende Netze

Konvergenz der Heuristischen Methode

Silva, Almeida 1991

Beh: $C_{yy} \rightarrow I$ mit $C_{yy} = \langle yy^T \rangle = W \langle xx^T \rangle W^T = W C_{xx} W^T$

Bew: Einsetzen der Matrixversion $W(t+1) = W(t) - \gamma C_{yy}(t) W(t) + \gamma W(t)$

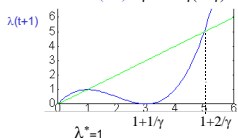
$$\Rightarrow C_{yy}(t+1) = \gamma^2 C_{yy}^3 - 2\gamma(1+\gamma)C_{yy}^2 + (1+\gamma)^2 C_{yy}$$

Darstellung im Eigenvektorraum

$$C_{yy} e_i = e_i \lambda_i \Rightarrow C_{yy} E = E \Lambda \Rightarrow C_{yy} = E \Lambda E^T$$

ändert die Eigenvektoren nicht, sondern nur die Eigenwerte λ zu

$$\lambda(t+1) = \gamma^2 \lambda^3 - 2\gamma(1+\gamma)\lambda^2 + (1+\gamma)^2 \lambda \quad \text{Fixpunktgleichung}$$



Konvergenz $\lambda \rightarrow 1 \Rightarrow C_{yy} \rightarrow I$

q.e.d.

Ex. keine Zielfunktion!

DEF Information

$$I \sim n = \text{ld}(2^n) = \text{ld} \text{ (Zahl der möglichen Daten)}$$

$$I \sim \text{ld}(1/P) \quad [\text{Bit}]$$

DEF $I(X) := \ln(1/P(x^k)) = -\ln(P(x^k))$ *Information*

DEF $H(X) := \sum_k P(x^k) I(x^k) = \langle I(x^k) \rangle_k$ *Entropie*

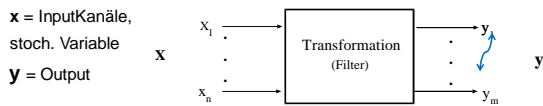
$$H(X) := \int_{-\infty}^{+\infty} p(x) \ln p(x)^{-1} dx \quad \textit{differenzielle Entropie}$$

Frage: Wieviel Information hat eine 32-bit floating-point Zahl?

DEF $I(X;Y) = H(X) + H(Y) - H(X,Y)$ *Transinformation mutual information*

ICA - Algorithmen 1a

Ziel: *minimale Transinformation zwischen den Ausgaben y_i*



Transinformation $I(y_1; y_2) = H(y_1) + H(y_2) - H(y_1, y_2)$

minimal bei $I(y_1; y_2) = 0$ bzw. maximaler Entropie $H(y_1, y_2) = H(y_1) + H(y_2)$
bzw. $p(y_1, y_2) = p(y_1) \cdot p(y_2)$ *stochastische Unabhängigkeit* der Variablen y_i

$$W(t+1) = W(t) - \frac{\partial}{\partial W} \gamma I(y_1; y_2; \dots; y_n) \quad \textit{Gradientenabstieg}$$

(Amari, Young, Cichocki 1996)

Entwicklung von $p(y_1, y_2, \dots, y_n)$ in $I(y_1; y_2; \dots; y_n)$ nach höheren Momenten

$$W(t+1) = W(t) - (I - f(y)) y^T W(t) \quad \text{mit } f_i(y_i) = \frac{3}{4} y^{11} + \frac{25}{4} y^9 - \frac{14}{3} y^7 - \frac{47}{4} y^5 + \frac{29}{4} y^3$$

ICA - Algorithmen 1b

Ziel: *maximale Transinformation* (Bell, Sejnowski 1995)
zwischen Eingabe und Ausgabe

Transinformation $I(X;Y) = H(Y) - H(Y|X)$ aus Def.
 $H(Y|X)$ ist konstant im determinist. System

Maximierung von $I(X;Y)$ durch Maximierung von $H(Y)$

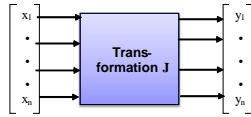
$$R(w) := H(Y) = H(X) + \int_{-\infty}^{+\infty} p(x) \ln |\det J| dx \quad \text{mit } J = \left(\frac{\partial y_i}{\partial x_j} \right)$$

Nicht-lin. Ausgabe y , z.B. $y = \tanh(z)$

Informationstransformation

$H(\mathbf{Y}) = ?$

$H(\mathbf{Y}) = H(\mathbf{Y}(\mathbf{X}))$



Transformation
kontinuierlicher
Zufallsvariabler

$$H(\mathbf{Y}) = - \int_{-\infty}^{+\infty} p(\mathbf{y}) \ln p(\mathbf{y}) d\mathbf{y}$$

$$d\mathbf{y} = |\det \mathbf{J}| d\mathbf{x} \\ \Rightarrow p(\mathbf{y}(\mathbf{x})) = p(\mathbf{x}) |\det \mathbf{J}|^{-1}$$

$$= - \int_{-\infty}^{+\infty} (p(\mathbf{x}) |\det \mathbf{J}|^{-1}) \ln (p(\mathbf{x}) \cdot |\det \mathbf{J}|^{-1}) |\det \mathbf{J}| d\mathbf{x}$$

$$= - \int_{-\infty}^{+\infty} p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} + \int_{-\infty}^{+\infty} p(\mathbf{x}) \ln |\det \mathbf{J}| d\mathbf{x}$$

$$\mathbf{y}(\mathbf{x}) = \mathbf{W}\mathbf{x} \\ \Rightarrow \mathbf{J} = \mathbf{W}$$

$$H(\mathbf{Y}) = H(\mathbf{X}) + \ln |\det(\mathbf{W})|$$

ICA - Algorithmen 1b

Ziel: *maximale Transinformation* (Bell, Sejnowski 1995)
zwischen Eingabe und Ausgabe

Gradientenregel

$$\mathbf{W}_{(t+1)} = \mathbf{W}_{(t)} + \gamma \left(\left[\mathbf{W}^T \right]^{-1} - 2\mathbf{y}\mathbf{x}^T \right) \quad \text{Rechnung: 1-dim Fall (Kap.3.4.1)}$$

„Natürl.“ Gradient

Amarı 1985

$$\frac{d\mathbf{W}}{dt} = \Delta\mathbf{W} = \gamma(t) \frac{\partial R(\mathbf{W})}{\partial \mathbf{W}} \mathbf{W}^T \mathbf{W} = \gamma(\mathbf{I} - 2\mathbf{y}\mathbf{z}^T) \mathbf{W}$$

Statist. Momente und Kurtosis

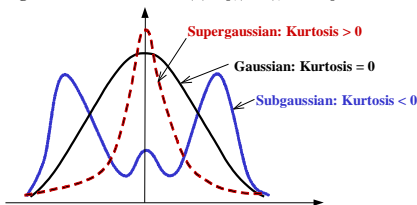
Momente einer Zufallsvariablen x :

$$\alpha_i = \langle x^i \rangle, \quad \text{z.B. } \alpha_1 = \langle x \rangle \quad \text{Mittelwert}$$

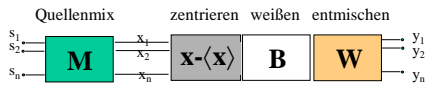
Zentrale Momente einer Zufallsvariablen x :

$$m_k = \langle (x - \alpha_1)^k \rangle, \quad \text{z.B. } m_2 = \langle (x - \alpha_1)^2 \rangle \quad \text{Varianz}$$

Wölbungsmaß Kurtosis: $\text{kurt}(x) = [\langle (x - \alpha_1)^4 \rangle - 3m_2^2] / m_2^2$



ICA-Algorithmen: Vorverarbeitungsfolge



Zentrieren

$$\langle x \rangle = 0 \quad \langle x - \langle x \rangle \rangle^2 = 1$$

- Mittelwertbildung, z.B. iterativ durch $w_0(t+1) = w_0(t) - \gamma (w_0 \cdot x)$, $\gamma = 1/t$

Weiß

- PCA durchführen: w_i Eigenvektoren von $C = \langle xx^T \rangle$ mit $|w_i| = 1$ und Eigenwerten λ_i
- Gewichtsvektoren w_i normieren zu $w_i / \lambda_i^{1/2}$. Dies führt zu $\langle y^2 \rangle = w_i^T \langle xx^T \rangle w_i = w_i^T \lambda_i w_i = 1$

Entmischen

- ICA Algorithmen, z.B. minimale Transinformation, maximale Kurtosis etc.
- Speziell:** dekorrelierte x benötigen nur eine orthogonale Matrix W (Vereinfachung)

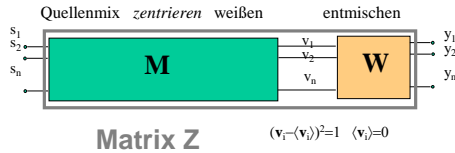
ICA – Algorithmen 2

Ziel: extreme Kurtosis

(Delfosse, Loubaton 1995)

Extrema bei $s_j = \text{unabh. Komp.}$ und $z_j = \pm 1$

$$\text{kurt}(y) = \text{kurt}(w^T v) = \text{kurt}(w^T M s) = \text{kurt}(z^T s) = \sum_{j=1}^n z_j^4 \text{kurt}(s_j)$$



ICA – Algorithmen 2

Ziel: extreme Kurtosis bei $y = w^T v$

$$R(w) = \langle (w^T v)^4 \rangle - 3 \langle (w^T v)^2 \rangle^2 = \min_w$$

$$w(t+1) = w(t) + \gamma \text{grad } R(w)$$

$$= w(t) + \gamma 4 \left(\langle (w^T v)^3 v \rangle - 3 |w|^2 w \right)$$

Bei $|w| = 1$ ist die Richtung gegeben durch

$$w(t+1) = \alpha \left(\langle (w^T v)^3 v \rangle - \beta w \right)$$

- Lernalgorithmus** für einzelnes Neuron (Hyvarinen, Oja 1996)

$$w(t+1) = \langle (w^T v)^3 v \rangle - 3 w \quad \text{Fixpunktalgorithmus}$$

mit $|w| = 1$

ICA – Algorithmen 2

- Sequentielle Extraktion aller Komponenten

Gegeben: Trainingsmenge $\{v(0)\}$

$$w_1(t+1) = \langle (w_1^T v)^3 v \rangle - 3 w_1 \quad \text{mit } |w_1| = 1$$

Konvergenz zum 1. ICA-Vektor.

Dann neue Trainingsmenge durch $v(1) = v(0) - w_1 y_1$

$$w_2(t+1) = \langle (w_2^T v)^3 v \rangle - 3 w_2 \quad \text{mit } |w_2| = 1$$

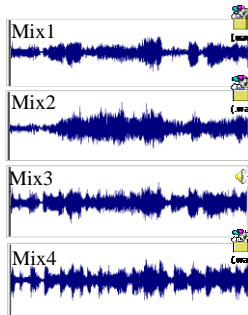
Konvergenz zum 2. ICA-Vektor, usw.

- Schnellere Konvergenz: Orthogonalisierung

$$w_i(t+1) = w_i(t) - \sum_{j=1}^{i-1} (w_i^T w_j) w_j \quad j < i$$

ICA-Anwendung: Audioanalyse

Mischung



entmischte Quellen

