

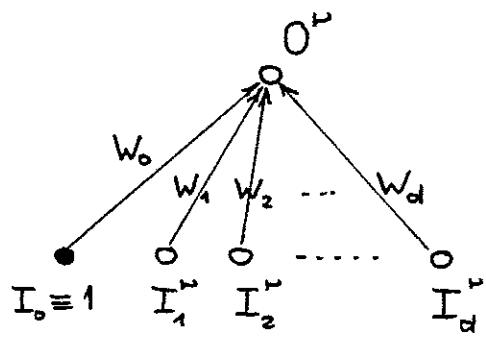
## V. PERCEPTRONS

- Literatur:
- Rosenblatt (1958, 1960), Minsky & Papert (1969)
  - N.J. Nilsson, "Learning Machines", McGraw-Hill (1965)

### V.1 Definition, Geometrische Interpretation

Ein Perceptron ist ein Feedforward-Netzwerk, das keine Hidden Units enthält. Es besteht also nur aus einer Schicht Input-Neuronen und einer Schicht Output-Neuronen.

Da die Outputs unabhängig voneinander sind, genügt es, sich auf Perceptrons mit nur einem Output zu beschränken:



$$w_i \text{ reell}, \leq 0$$

$$w_0 = \text{Schwelle} \quad [I_0 = 1 \quad \forall n]$$

$$I_i'' \text{ reell (ev. nur } \pm 1\text{)}$$

- Die Inputs sind Punkte  $\{I_1'', \dots, I_d''\}$  in  $\mathbb{R}^d$ , d.h. im d-dim. Euklidischen Raum,  $n=1, \dots, N$ .
- Vorgegebene Klassifizierung der N Inputpunkte:

$$D'' = \pm 1, \quad n = 1, \dots, N$$

[Klassifizierung in 2 Klassen:  $\rightarrow$  DICHOTOMIE]

- Problem: Bestimme  $w_1, \dots, w_d$  und  $w_0$  (Schwelle) so, dass  $O^r = D^r$  für  $r = 1, \dots, N$ .

$$O^r = \text{sign} \left( \sum_{i=1}^d w_i I_i^r + w_0 \right)$$

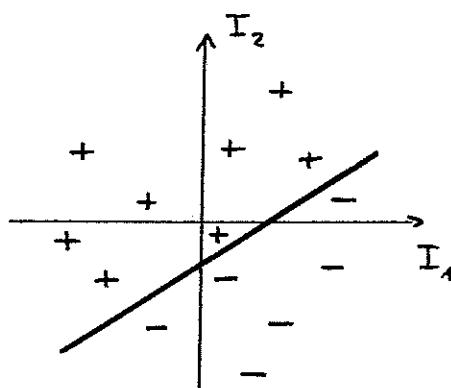
$$= \begin{cases} +1 & \text{falls } \sum_{i=1}^d w_i I_i^r + w_0 \geq 0 \\ -1 & \quad \quad \quad < 0 \end{cases}$$

- Die Gleichung  $\sum_{i=1}^d w_i I_i + w_0 = 0$  definiert eine  $(d-1)$ -dimensionale Hyperebene im  $d$ -dimensionalen Inputraum.

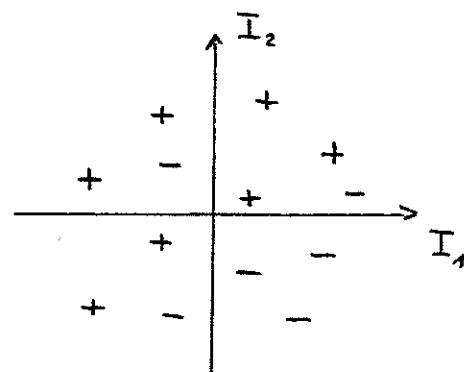
→ Klassifizierung mit Perceptron nur möglich, wenn die Klassen "+" und "-" durch eine Hyperebene getrennt werden können.

[LINEARE SEPARABILITÄT]

Beispiel ( $d=2$ ):



linear separabel



nicht linear separabel

## V.2 Kapazität eines Perceptrons

- N Punkte  $\underline{I}^1, \dots, \underline{I}^N$  im d-dim Euklid'schen Raum
  - $2^N$  mögliche Dichotomien
  - Wieviele davon sind linear separabel?
- $L(N, d) =$  Zahl der linear separablen Dichotomien von N Punkten in allgemeiner Lage in d Dimensionen.

Allgemeine Lage:  $N > d$ : Keine Untermenge von  $d+1$  Punkten liegt auf einer  $(d-1)$ -dim. Hyperebene

$N \leq d$ : Keine  $(N-2)$ -dim. Hyperebene enthält die N Punkte

- Berechnung von  $L(N, d)$ :

[Cover, IEEE Trans. on Electron. Computers, Vol. EC-14, pp. 326-334 (June 1965)]

Betrachte die  $L(N-1, d)$  linear separablen Dichotomien der Punkte  $\underline{I}^1, \dots, \underline{I}^{N-1}$ . Davon seien  $L_{\underline{I}^N}(N-1, d)$  mit einer Hyperebene realisierbar, die durch den Punkt  $\underline{I}^N$  geht.

$$\rightarrow L(N, d) = \underbrace{[L(N-1, d) - L_{\underline{I}^N}(N-1, d)]}_{\substack{\underline{I}^N \text{ ist entweder in} \\ \text{Klasse } "+" \text{ oder } "-"}} + \underbrace{2 L_{\underline{I}^N}(N-1, d)}_{\substack{\underline{I}^N \text{ kann in bei-} \\ \text{den Klassen sein} \\ (\text{beliebig kleine} \\ \text{Verschiebung der} \\ \text{Trennebene!})}}$$

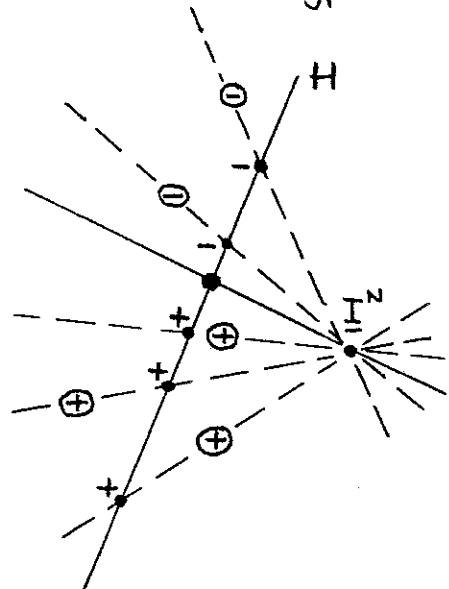
$\underline{I}^N$  ist entweder in  
Klasse "+" oder "-"

$\underline{I}^N$  kann in bei-  
den Klassen sein  
(beliebig kleine  
Verschiebung der  
Trennebene!)

$$\rightarrow L(N, d) = L(N-1, d) + L_{\underline{I}^N}(N-1, d)$$

$$\text{Behauptung: } L_{\underline{\mathbb{I}}^N}(N-1, d) = L(N-1, d-1)$$

Beweis: Schneide die  $N-1$  Geraden, die durch  $\underline{\mathbb{I}}^N$  und  $\underline{\mathbb{I}}^i$  gehen ( $i=1, \dots, N-1$ ), mit einer  $(d-1)$ -dim. Hyperebene  $H$ :



→  $N-1$  Punkte im  $d$ -dim. Raum genau dann durch Hyperebene durch  $\underline{\mathbb{I}}^N$  separierbar, wenn die  $N-1$  Projektionen auf  $H$  durch  $(d-2)$ -dim. Hyperebene trennbar sind.

q.e.d.

$$\rightarrow \text{Rekursion: } L(N, d) = L(N-1, d) + L(N-1, d-1)$$

Offensichtliche Randbed.:

$$L(1, d) = 2, \quad L(N, 1) = 2N$$

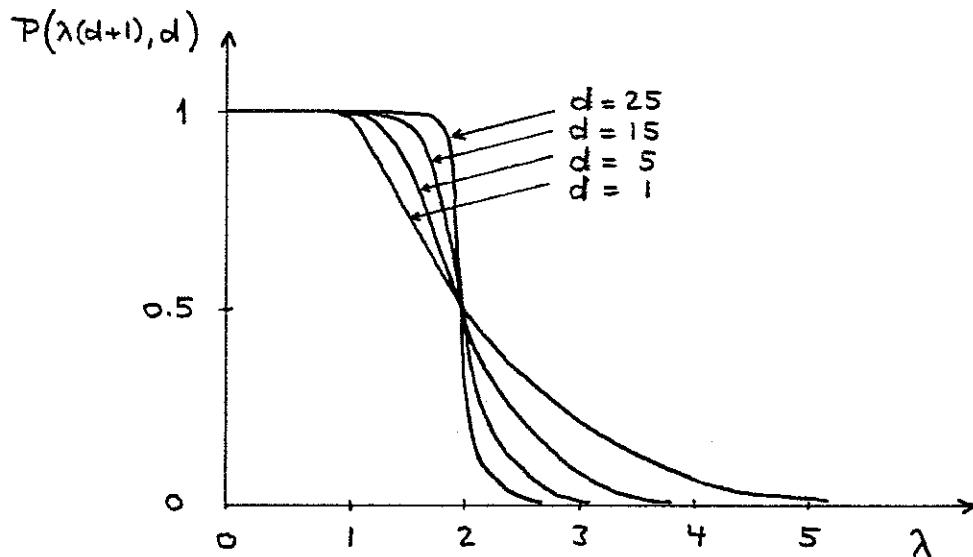
$$\rightarrow L(N, d) = \begin{cases} 2^N & \text{falls } N \leq d \\ 2 \sum_{i=0}^d \binom{N-1}{i} & \text{falls } N > d \end{cases}$$

Definition:

$P(N, d) =$  Wahrsch.keit, dass eine zufällig gewählte Dichotomie von  $N$  Punkten im  $d$ -dim. Raum linear separabel ist  
(d.h. durch ein Perceptron mit  $d$  Inputs und einer Schwelle realisiert werden kann).

$$\rightarrow P(N, d) = \frac{L(N, d)}{2^N} = \begin{cases} 1 & \text{falls } N \leq d \\ \frac{1}{2^{N-1}} \sum_{i=0}^d \binom{N-1}{i} & \text{falls } N > d \end{cases}$$

Setze  $N = \lambda(d+1)$  und betrachte  $P(\lambda(d+1), d)$ :



$$\rightarrow P(\lambda(d+1), d) = \frac{1}{2} \quad \text{für } \lambda = 2$$

$$\lim_{d \rightarrow \infty} P(\lambda(d+1), d) = \begin{cases} 0 & \text{falls } \lambda = 2 + \varepsilon \\ 1 & \text{falls } \lambda = 2 - \varepsilon \end{cases}$$

$\rightarrow$  KAPAZITÄT EINES PERCEPTRONS:

$$N_c = 2(d+1)$$

d.h. für grosse  $d$  gilt:

- Eine vorgegebene Dichotomie kann fast sicher durch ein Perceptron realisiert werden, falls  $N < 2(d+1)$
- und fast sicher nicht, falls  $N > 2(d+1)$

[  $N$  = Anzahl Lernbeispiele

$d$  = Dimension des Inputraums  
(Anzahl Inputs ohne Schwelle) ]

### V.3 Das Perceptron - Lernverfahren

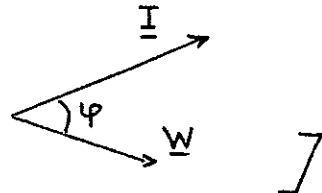
Notation:  $\underline{I}^n = (I_0=1, I_1^n, I_2^n, \dots, I_d^n)$ ,  $n=1, \dots, N$

$$\underline{W} = (W_0, W_1, W_2, \dots, W_d)$$

$$\rightarrow O^n = \text{sign}(\underline{W} \cdot \underline{I}^n) = \begin{cases} +1 & \underline{W} \cdot \underline{I}^n \geq 0 \\ -1 & \underline{W} \cdot \underline{I}^n < 0 \end{cases}$$

$$[\underline{W} \cdot \underline{I} = \sum_{i=0}^d W_i I_i = |\underline{W}| \cdot |\underline{I}| \cdot \cos \varphi$$

$\rightarrow$  Skalarprodukt von  
(d+1)-dimensionalen Vektoren



- Wie können die Gewichte  $W_0, W_1, \dots, W_d$  "gelernt" werden, sodass  $O^n = D^n$  ( $n=1, \dots, N$ ), für vorgegebene Klassifizierung  $D^n$  der Inputmuster?

$\rightarrow$  PERCEPTRON - LERNVERFAHREN:

- Die Inputbeispiele  $\underline{I}^n$  ( $n=1, \dots, N$ ) werden in zufälliger Reihenfolge immer wieder präsentiert.
- Bei jedem Schritt k wird  $O^n = \text{sign}(\underline{W} \cdot \underline{I}^n)$  berechnet und das aktuelle  $\underline{W} = \underline{W}(k)$  wie folgt geändert:

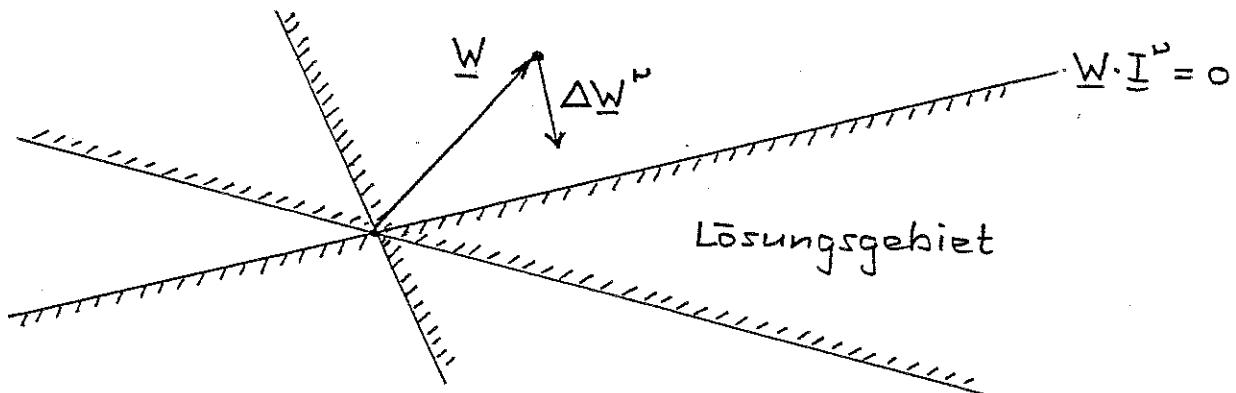
$$\underline{W}(k) \longrightarrow \underline{W}(k+1) = \underline{W}(k) + \Delta \underline{W}^n$$

$$\text{wobei } \Delta \underline{W}^n = \eta (D^n - O^n) \underline{I}^n, \quad \eta = \text{const}$$

$$[\Delta W_i^n = \eta (D^n - O^n) I_i^n, \quad i=0, 1, \dots, d]$$

## Konvergenzbeweis (geometrisch):

Voraussetzung: Es existiere ein Lösungsgebiet im  $\underline{W}$ -Raum. (Dieses ist durch Hyperebenen  $\underline{W} \cdot \underline{I}^r = 0$  begrenzt.)



Falls  $O^r(\underline{W}) \neq D^r$ , so wird  $\underline{W}$  bei Präsentation des Lernbeispiels  $r$  um  $\Delta\underline{W}^r = \eta(D^r - O^r)\underline{I}^r$  geändert.

Es ist leicht einzusehen, dass  $\Delta\underline{W}^r$  senkrecht zur Hyperebene  $\underline{W} \cdot \underline{I}^r = 0$  steht und in die Richtung des Lösungshalbraums für das  $r$ -te Inputbeispiel zeigt.

→ Da jedes  $r$  immer wieder auftritt, und da das Lösungsgebiet ein offener "Keil" ist, konvergiert das Verfahren immer in einer endlichen Anzahl Schritten gegen eine Lösung (mit Wahrscheinlichkeit 1).

Aber: Wenn keine Lösung für  $\underline{W}$  existiert, d.h. wenn das Problem nicht linear separabel ist, konvergiert das Verfahren nicht!  
(Auch nicht gegen eine optimale Lösung mit möglichst wenig Fehklassifizierungen!)

→ "POCKET-ALGORITHMUS"

(Während Lernprozess immer die Gewichte der bis anhin besten Lösung abspeichern.)

ANHANG: VOLLSTÄNDIGER KONVERGENZBEWEIS  
FÜR DAS PERCEPTRON-LERNVERFAHREN

Wahl der Anfangsgewichte:  $\underline{W} = \underline{W}^0 = \underline{0}$

(Der Beweis kann leicht auf beliebige Anfangsgewichte verallgemeinert werden!)

Gewichtsänderung beim Perceptron-Lernverfahren:

$$\Delta \underline{W} = \eta (\underline{D}^\nu - \underline{O}^\nu) \underline{I}^\nu = \eta (1 - \underline{D}^\nu \cdot \underline{O}^\nu) \underline{X}^\nu$$

$$\text{wobei: } \underline{D}^\nu = \pm 1, \quad \underline{O}^\nu = \text{sign}(\underline{W} \cdot \underline{I}^\nu) = \pm 1$$

$$\underline{X}^\nu = \underline{D}^\nu \cdot \underline{I}^\nu, \quad \nu = 1, \dots, N$$

$$\rightarrow \Delta \underline{W} = \begin{cases} 2\eta \underline{X}^\nu & \text{falls } \underline{O}^\nu \neq \underline{D}^\nu, \text{ d.h. falls } \underline{W} \cdot \underline{X}^\nu < 0 \\ 0 & \text{--- } \underline{O}^\nu = \underline{D}^\nu, \text{ --- } \underline{W} \cdot \underline{X}^\nu > 0 \end{cases}$$

$$\rightarrow \underline{W} = 2\eta \sum_{\nu=1}^N M^\nu \underline{X}^\nu,$$

wobei  $M^\nu = \text{Anzahl Präsentationen des Lernbeispiels } \nu, \text{ bei denen } \underline{O}^\nu \neq \underline{D}^\nu \text{ war.}$

Sei  $\underline{W}^*$  eine Lösung  $[\underline{W}^* \cdot \underline{X}^\nu > 0, \nu = 1, \dots, N]$ :

$$\rightarrow \underline{W} \cdot \underline{W}^* = 2\eta \sum_{\nu=1}^N M^\nu \underline{X}^\nu \cdot \underline{W}^* \geq 2\eta M \cdot D(\underline{W}^*) \cdot |\underline{W}^*| \quad (a)$$

$$\text{wobei: } M \equiv \sum_{\nu=1}^N M^\nu, \quad D(\underline{W}^*) = \frac{1}{|\underline{W}^*|} \min_\nu \underline{W}^* \cdot \underline{X}^\nu > 0$$

Für eine Gewichtsänderung aufgrund des Lernbeispiels  $\nu$  gilt andererseits:

$$\Delta |\underline{w}|^2 = (\underline{w} + 2\eta \underline{x}^\nu)^2 - \underline{w}^2 = 4\eta^2 (\underline{x}^\nu)^2 + 4\eta \underline{w} \cdot \underline{x}^\nu$$

$$\rightarrow \Delta |\underline{w}|^2 \leq 4\eta^2 (\underline{x}^\nu)^2 , \text{ da } \underline{w} \cdot \underline{x}^\nu < 0$$

$$\rightarrow |\underline{w}|^2 \leq \sum_{n=1}^N M^\nu \cdot 4\eta^2 (\underline{x}^\nu)^2 \leq 4\eta^2 M \cdot X_{\max}^2 \quad (b)$$

$$\text{wo bei: } X_{\max}^2 = \max_n (\underline{x}^\nu)^2$$

Aber:

$$\underline{w} \cdot \underline{w}^* = |\underline{w}| \cdot |\underline{w}^*| \cdot \cos \phi$$

wo bei  $\phi$  = Winkel zwischen  $\underline{w}$  und  $\underline{w}^*$

$$\rightarrow 1 \geq \cos \phi = \frac{\underline{w} \cdot \underline{w}^*}{|\underline{w}| \cdot |\underline{w}^*|} \geq \frac{2\eta M \cdot D(\underline{w}^*) \cdot |\underline{w}^*|}{2\eta \sqrt{M} \cdot \sqrt{X_{\max}^2} \cdot |\underline{w}^*|}$$

$\swarrow$  siehe (a) und (b)

$$\rightarrow 1 \geq \frac{\sqrt{M} \cdot D(\underline{w}^*)}{\sqrt{X_{\max}^2}}$$

$$\rightarrow M \leq \frac{X_{\max}^2}{D(\underline{w}^*)^2} , \text{ d.h. die Anzahl } M \text{ der Gewichtsänderungen während des Lernprozesses ist endlich!}$$


---

## PERCEPTRONS (ERGÄNZUNGEN):

Perceptron:  $O^v = \text{sign} \left( \sum_{i=0}^d w_i I_i^v \right)$

$$\rightarrow O^v = D^v, v=1, \dots, N \quad [O^v, D^v = \pm 1]$$



$$E^v = D^v \sum_{i=0}^d w_i I_i^v > 0, v=1, \dots, N$$

## Perceptron - Lernverfahren:

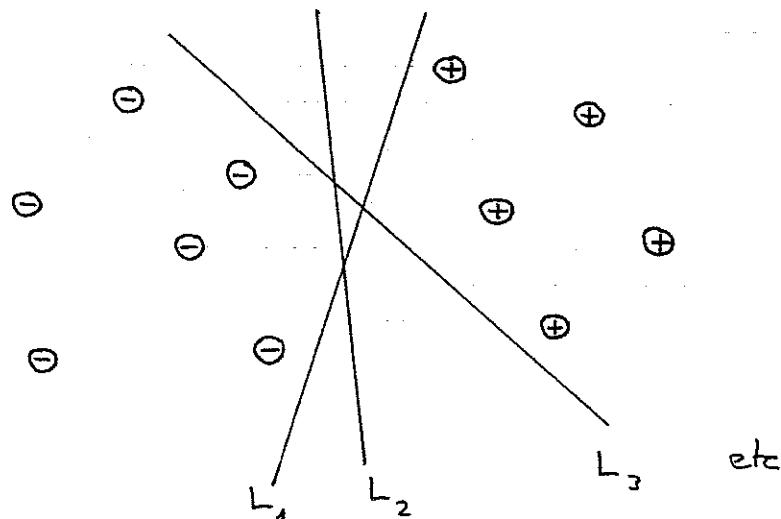
- $v$  zufällig oder der Reihe nach ausgewählt
- $\Delta w_i = \eta (D^v - O^v) I_i^v$

$$\rightarrow \Delta E^v = D^v \sum_{i=0}^d \Delta w_i I_i^v = \eta (1 - D^v O^v) \sum_{i=0}^d (I_i^v)^2$$

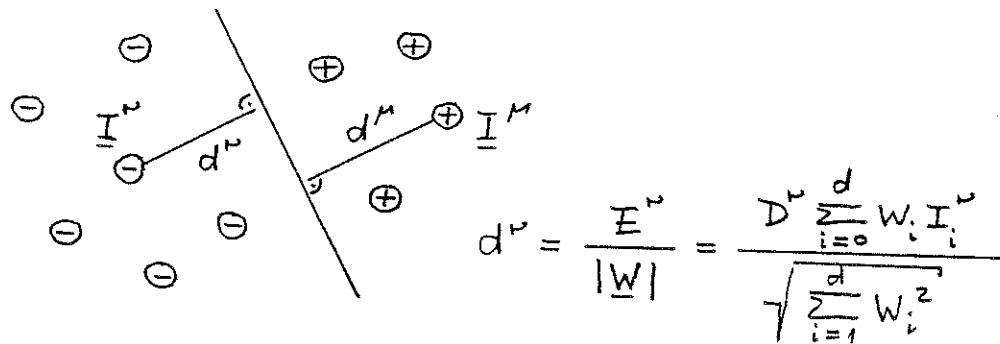
$$\rightarrow \Delta E^v > 0 \text{ falls } O^v \neq D^v \quad [=0 \text{ falls } O^v = D^v]$$

d.h.  $E^v$  wird immer wieder vergrößert, bis schliesslich alle  $E^v > 0$  [falls Lösung existiert!]

Es gibt aber viele verschiedene Lösungen:



## PERCEPTRON MIT MAXIMALER STABILITÄT



$$d^R = \frac{D^R \sum_{i=0}^D w_i I_i^R}{\sqrt{\sum_{i=1}^D w_i^2}}$$

→ Maximale Stabilität:

$$\min_p \{d^p\} = \max$$

### Lösungsmethoden:

i) Minover-Lernalgorithmus (für  $I_i^p = \pm 1$ )

[Krauth & Mézard, J. Phys. A 20, L745-752 (1987)]

ii) Formulierung als Optimierungsproblem

$$\cdot \sqrt{\sum_{i=1}^D w_i^2} = \min$$

$$\cdot \text{Nebenbedingungen: } D^p \sum_{i=0}^D w_i I_i^p \geq 1, p=1, \dots, N$$

→ "SUPPORT VECTOR CLASSIFIERS"

[siehe z.B.: T. Hastie et al

The Elements of Statistical Learning  
Springer Series in Statistics, 2001  
Chapters 4 and 12]

→ Verallgemeinerung auf nicht-lineare Trennflächen, nicht perfekt separierbare Klassifizierungsprobleme, und auf Regressionsprobleme:

"SUPPORT VECTOR MACHINES"

## V.4 Erweiterung des Inputraums

• Ursprünglicher Inputraum:  $I_i, i=1, \dots, d$

• Erweiterter Inputraum:

z.B.:  $I_i, I_i I_j, I_i I_j I_k, \dots$

oder:  $I_i, \sin(\pi I_i), \cos(\pi I_i), \sin(2\pi I_i), \dots$

→ Perceptrons mit erweitertem Inputraum können nichtlineare Trennflächen realisieren!

→ Probleme, die im ursprünglichen Inputraum nicht linear separabel sind, können im erweiterten Inputraum linear separabel werden!

→ Einfaches Perceptron-Lernverfahren anwendbar!  
(d.h. schnelleres Lernen als bei Verwendung von komplexeren Netzwerk-Architekturen)

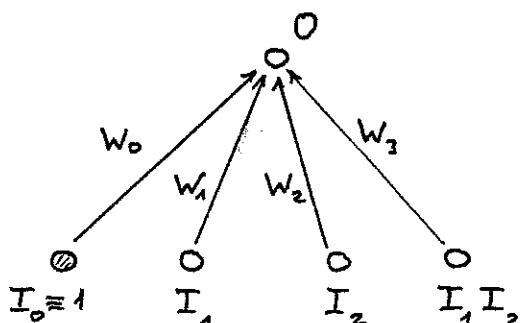
Aber: Geeignete Form der nichtlinearen Erweiterung nicht immer offensichtlich!

Beispiel: XOR-Problem

$I_1$	$I_2$	O
0	0	-1
1	0	+1
0	1	+1
1	1	-1

→ nicht linear separabel!

Mit erweitertem Inputraum aber linear separabel:



z.B.  $w_0 = -1$   
 $w_1 = 2 \quad w_2 = 2$   
 $w_3 = -4$

$[O = \text{sign} \left( \sum_{i=0}^3 w_i I_i \right)]$

## V.5 Minimierung des mittleren quadratischen Fehlers

$$F = \frac{1}{N} \sum_{n=1}^N (D^n - O^n)^2 = \min, \quad O^n = f\left(\sum_i w_i I_i^n\right)$$

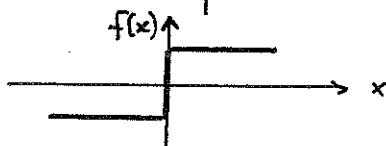
- EINFACHES GRADIENTENVERFAHREN:  
(zur Minimierung von F)

$$w_i \rightarrow w_i - \eta \frac{\partial F}{\partial w_i}$$

$$\frac{\partial F}{\partial w_i} = -\frac{1}{N} \sum_{n=1}^N 2(D^n - O^n) \cdot f'(\sum_j w_j I_j^n) \cdot I_i^n$$

Im Gegensatz zum Perceptron - Lernverfahren konvergiert dieses Verfahren immer (gegen ein lokales Minimum von F), falls  $\eta$  nicht zu gross gewählt wird.

Beim ursprünglichen Perceptron ist aber  $f = \text{sign}$ :



→  $f'$  existiert nicht, d.h. das Gradientenverfahren ist nicht anwendbar!

→ Änderung der Wahl der Aktivierungsfunktion f für das Output-Neuron:



oder Wahl eines anderen Optimierungsverfahrens zur Minimierung von F

[z.B. stochastisches Optimierungsverfahren]  
(siehe Kapitel VIII)]

## V.6 Analyse von Gradientenverfahren

### A) Lineares Output-Neuron [ADALINE]

$$\rightarrow O^u = \sum_i w_i I_i^u$$

$$\rightarrow F = \frac{1}{N} \sum_{r=1}^N [D^u - \sum_i w_i I_i^u]^2 = \frac{1}{N} \sum_{n=1}^N F_n$$

ist eine quadratische Funktion der  $w_i$

$\rightarrow$  Es existiert genau 1 Minimum!

#### Eindimensionales Beispiel:

$$F = A - 2Bw + Cw^2$$

$$\left[ \begin{array}{l} \rightarrow F_n = a_n - 2b_n w + c_n w^2 \\ A = \frac{1}{N} \sum_{n=1}^N a_n = \langle a_n \rangle, \text{ etc.} \end{array} \right]$$

$$\rightarrow F_{\min} = F(w^*) = A - \frac{B^2}{C}, \quad w^* = \frac{B}{C}$$

#### • Gradientenmethode ("Batch"):

$\rightarrow$  Änderung von  $w$  ist proportional zum negativen Gradienten von  $F$

( $F$  = Summe der quadratischen Outputfehler über alle Lernbeispiele)

$$w_{k+1} = w_k - \eta \frac{\partial F}{\partial w} \Big|_{w=w_k} = w_k + 2\eta (B - CW_k)$$

( $k$  = Iterationsindex)

Setze  $v_k = w_k - w^*$  ( $w^* = \frac{B}{C}$ ):

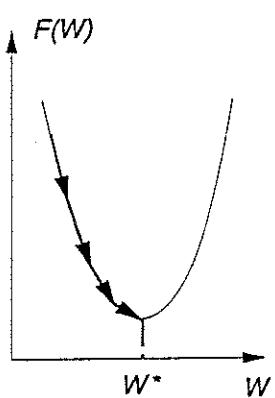
$$\rightarrow v_{k+1} = (1 - 2\eta C) v_k$$

→ Konvergenz falls  $\lim_{k \rightarrow \infty} (1 - 2\eta C)^k = 0$

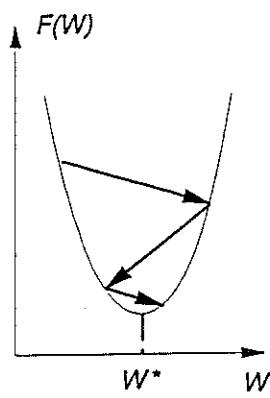
d.h. falls  $|1 - 2\eta C| < 1$

$$\rightarrow \boxed{0 < \eta < \frac{1}{2C}} \quad \text{oder} \quad \boxed{\frac{1}{2C} \leq \eta < \frac{1}{C}}$$

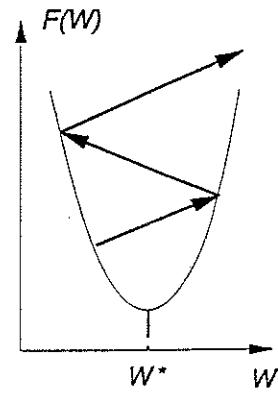
Im ersten Fall strebt  $W_k$  stetig gegen die Lösung  $W^*$ , im 2. Fall wird  $W^*$  in einer "Zickzackbewegung" angenähert ( $V_k = W_k - W^*$  wechselt bei jedem Iterationsschritt das Vorzeichen). Für  $\eta > \frac{1}{C}$  divergiert das Verfahren.



$$0 < \eta < 1/2C$$



$$1/2C < \eta < 1/C$$



$$\eta > 1/C$$

- LMS-Algorithmus ("Example-by-Example"):

[Widrow & Hoff 1960] (LMS = "least mean squares")

→ Lernbeispiele  $r$  werden in zufälliger Reihenfolge immer wieder präsentiert und  $W$  jedesmal um einen Betrag geändert, der proportional zum negativen Gradienten von  $F_r$  ist:

$$W_{k+1} = W_k - \eta \frac{\partial F_r}{\partial W} = W_k + 2\eta (b_r - c_r W_k)$$

( $r$  zufällig)

→ Wenn  $\eta$  genügend klein ist, dann strebt  $\{w_k\}$  gegen einen stationären Zufallsprozess mit:

$$\langle F \rangle = \langle a \rangle - 2\langle b \rangle \langle w \rangle + \langle c \rangle \langle w^2 \rangle$$

$$\langle w \rangle = \frac{\langle b \rangle}{\langle c \rangle} = w^*$$

$$\langle w^2 \rangle = \frac{\langle b \rangle^2 + \eta [\langle b^2 \rangle \langle c \rangle - 2\langle b \rangle \langle bc \rangle]}{\langle c \rangle [\langle c \rangle - \eta \langle c^2 \rangle]}$$

$$\rightarrow \langle F \rangle - F_{\min} \propto \eta \quad \text{falls } \eta \ll \frac{\langle c \rangle}{\langle c^2 \rangle}$$

### Bemerkung:

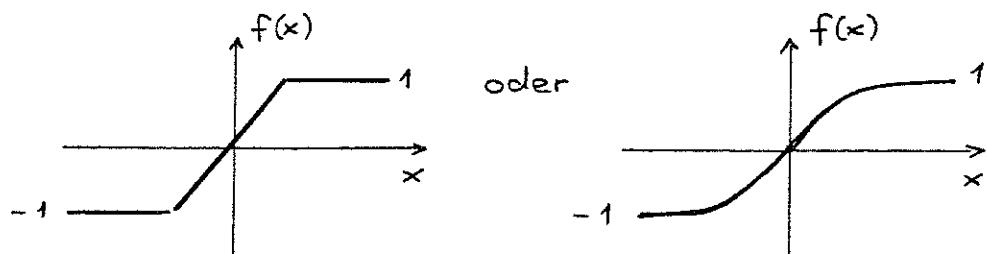
Die Konvergenzbeweise für die Gradientenmethode und für den LMS-Algorithmus lassen sich analog führen, wenn  $F$  von mehreren  $w_i$  abhängt.

[Vektor/Matrix-Notation!]

### B) Nichtlineares Output-Neuron

$$F = \frac{1}{2N} \sum_{k=1}^N (D^k - O^k)^2 = \min, \quad O^k = f(\sum_i w_i I_i^k)$$

wo bei z.B.:



$$\rightarrow \Delta W_i = -\eta \frac{\partial F}{\partial W_i} = \eta \frac{1}{N} \sum_r (D^r - O^r) \cdot f'(\sum_j W_j I_j^r) \cdot I_i^r$$

$$\text{bzw. } \Delta W_i = -\eta \frac{\partial F_r}{\partial W_i} = \eta (D^r - O^r) \cdot f'(\sum_j W_j I_j^r) \cdot I_i^r$$

- Auch im nichtlinearen Fall konvergieren die Verfahren, falls  $\eta$  genügend klein gewählt wird.
- Aber wenn  $f$  nichtlinear ist, können lokale Minima auftreten!  
→ Nur Konvergenz gegen nächstliegendes lokales Minimum garantiert!

### Bemerkungen:

- i) Wenn die "gewünschten" Outputs  $D^r = \pm 1$  sind, und die nichtlineare Aktivierungsfunktion  $f$  zwischen  $-1$  und  $+1$  variiert, dann

$$O^r = f(\sum_i W_i I_i^r) = \pm 1$$

nur möglich, wenn gewisse  $|W_i| \rightarrow \infty$ .

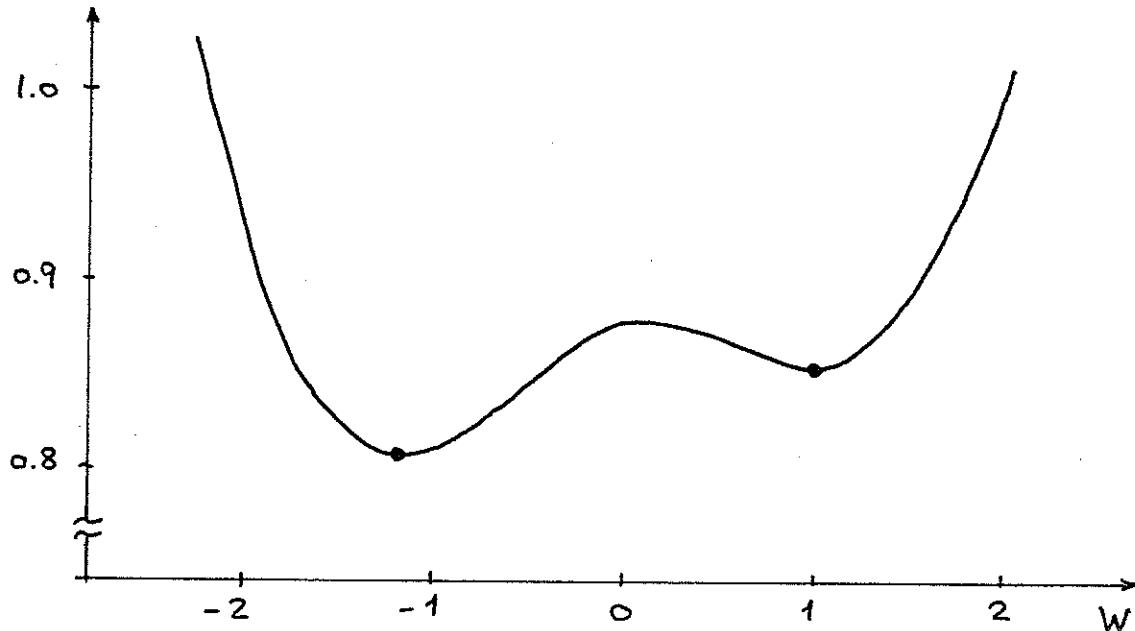
- Abhilfe:
- $D^r$  von  $\pm 1$  auf  $\pm 0.9$  reduzieren
  - $f$  in  $y$ -Richtung strecken, sodass ihre Maximalwerte z.B.  $\pm 1.1$  sind
  - Cutoff:  $F_r = \begin{cases} \frac{1}{2}(D^r - O^r)^2 & \text{if } |D^r - O^r| > \varepsilon \\ 0 & \text{if } |D^r - O^r| \leq \varepsilon \end{cases}$

- ii) Das Perceptron-Lernverfahren entspricht dem LMS-Algorithmus mit  $f' = 1$  (obwohl  $f = \text{sign}$ !)  
[→ Keine Konvergenz wenn  $F_{min} > 0$ !]

BEISPIEL MIT 2 LOKALEN MINIMA:

$$F = F_1 + F_2 = f(W-1.52)^2 + [f(W+1.52) + 0.04]^2 = \min$$

$$f(x) = \frac{1-e^{-x}}{1+e^{-x}}$$



a) "Batch": → Immer Konvergenz gegen dasjenige Minimum, das dem Anfangswert für  $W$  am nächsten liegt.  
(wenn  $\eta$  genügend klein)

b) "Example-by-Example":

→ Mit grosser Wahrscheinlichkeit Konvergenz gegen das tiefere Minimum (selbst wenn bei  $W=1$  gestartet wird!).

→ Example-by-Example Verfahren hat Vorteile bezüglich Steckenbleiben in schlechten lokalen Minima!