

# 4. Adaline

- 4.1 Charakterisierung und Lernzyklus von Adaline
- 4.2 Adaline Proportional Lernregel
- 4.3 Adaline Gradientenabstieg Lernregel
- 4.4 Anwendungen von Adaline
- 4.5 Lineare Regression durch Batch-Lernen

# 4.1 Charakterisierung und Lernzyklus von Adaline

## Überblick:

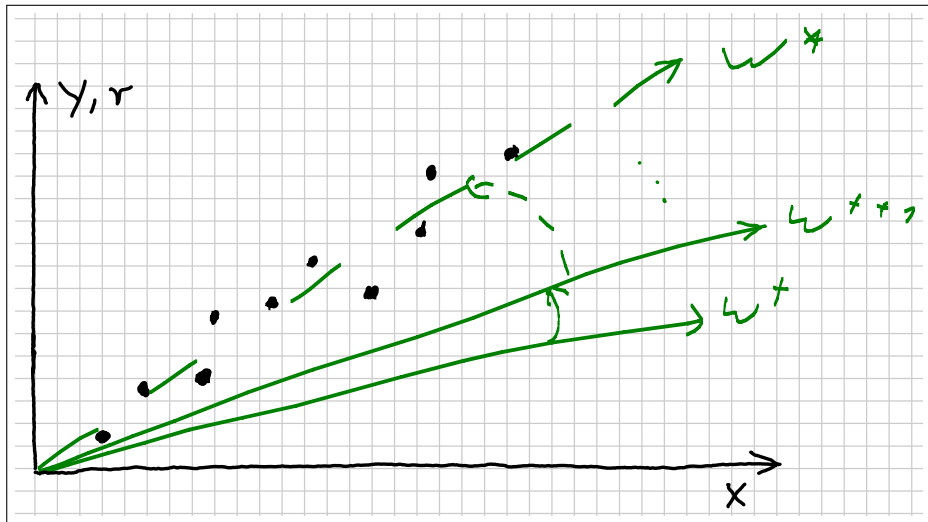
- Charakterisierung von Adaline
- Interpretation des Gewichtsvektors
- Graphische Illustration Lernzyklus

# Charakterisierung von Adaline

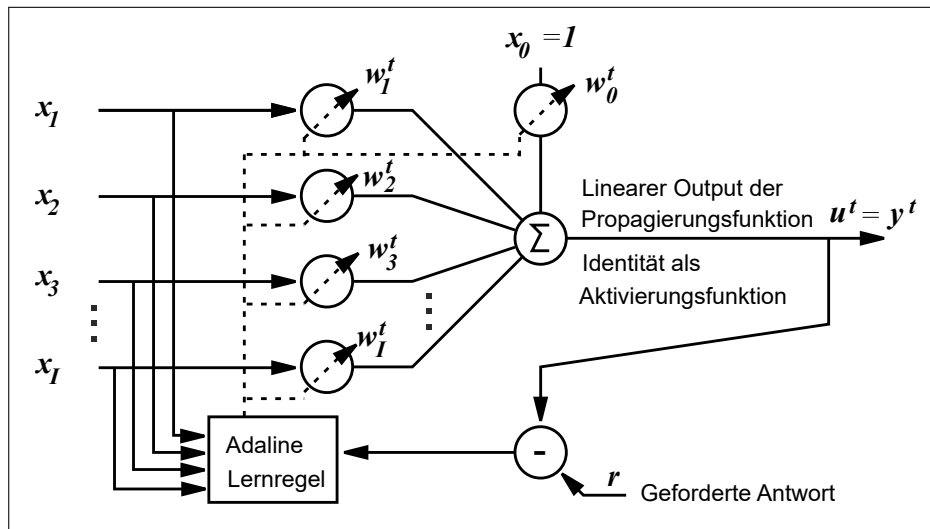
## **Adaptive linear element:** Widrow, Hoff (1960)

- Ziel des Lernens ist eine lineare Funktion, mit Minimierung der Abweichung zwischen gelernter Funktion und den Trainingsdaten.
- Input, Gewichtung, Propagierungsfunktion wie Perzeptron, Aktivierungsfunktion ist die Identität.
- Für das Lernen wird der lineare Output der Propagierungsfunktion verwendet.

# Interpretation des Gewichtsvektors



# Graphische Illustration Lernzyklus



## 4.2 Adaline Proportional Lernregel

### Überblick:

- Anpassung der Gewichte
- Änderung des Fehlers
- Graphische Illustration zur Anpassung der Gewichte
- Minimale Störung des bisherigen Lernerfolges

# Anpassung der Gewichte

Anpassung der Gewichte eines einzelnen Adaline:

$$\mathbf{w}^{t+1} := \mathbf{w}^t + \alpha \frac{d^t \mathbf{x}}{\|\mathbf{x}\|^2}$$

$$d^t := r - \underbrace{\mathbf{w}^{t,T} \mathbf{x}}_{u^t}$$

$r$  ist die geforderte reelle Ausgabe zum Eingabevektor  $\mathbf{x}$ .

Die Normierung geschieht durch Division mit der quadrierten Länge des Eingabevektors.

Hinweis: Im Unterkapitel wird das erweiterte Neuronenmodell mit dem Gewichtsvektor  $\mathbf{w}^e$  angenommen, aber zur kürzeren Schreibweise das Symbol  $e$  weggelassen.

# Änderung des Fehlers

Die Veränderung der Gewichte bewirkt folgende Änderung des Fehlers:

$$\begin{aligned}\Delta d^t &:= (r - w^{t+1,T} x) - (r - w^{t,T} x) \\ &= - \underbrace{(w^{t+1} - w^t)^T x}_{\Delta w^t}\end{aligned}$$

Speziell bei Adaline Proportional Lernregel gilt:

$$\Delta w^t := \alpha \frac{d^t x}{\|x\|^2}$$

$$\Delta d^t := -\alpha \frac{d^t x^T}{\|x\|^2} x = -\alpha \cdot d^t$$

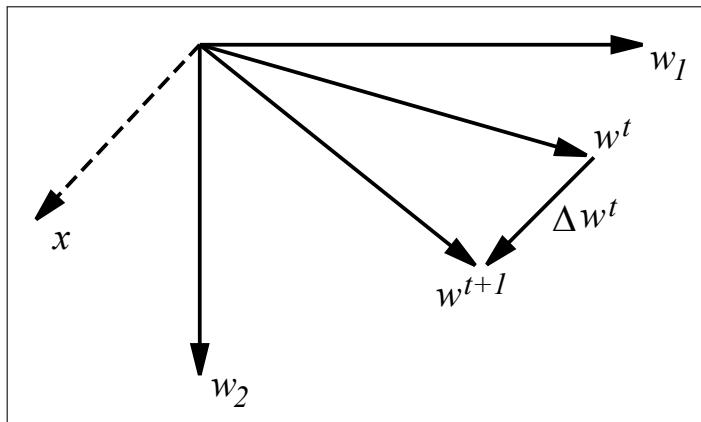


# Änderung des Fehlers

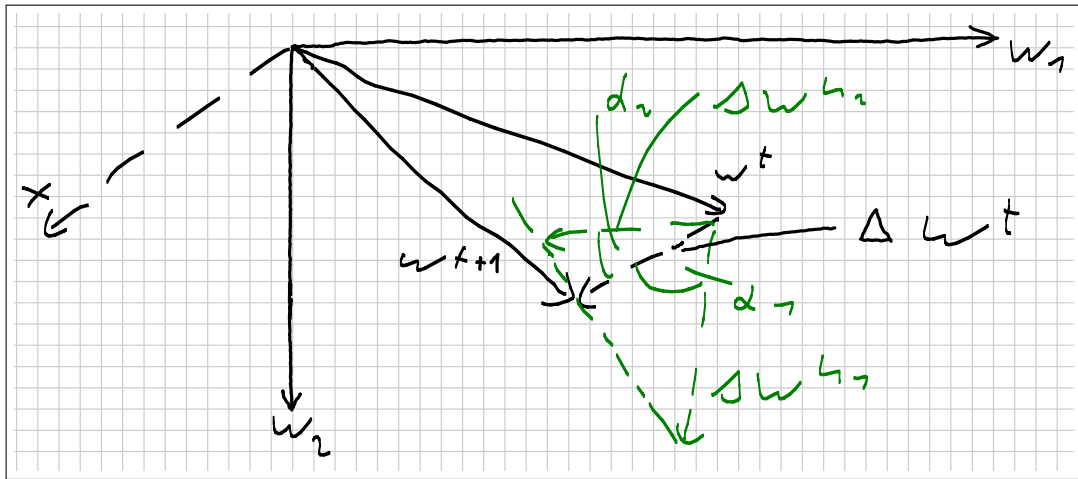
## Resumee:

- Lineare Fehlerreduzierung, da proportional zu Fehler, reduziert Anteil  $\alpha$  von  $d^t$ .
- Im Allgemeinen  $0 < \alpha < 1$ .  
Falls speziell  $\alpha = 1$ , dann Fehlerkorrektur.
- $\Delta w^t$  parallel zu  $x$

# Graphische Illustration zur Anpassung der Gewichte



# Graphische Illustration zur Anpassung der Gewichte



# Graphische Illustration zur Anpassung der Gewichte

$$\Delta w^t \cdot x = \|\Delta w^t\| \cdot \|x\| \cdot \underbrace{\cos 0}_1$$

$$\Delta w^{h_1} \cdot x = \underbrace{\|\Delta w^{h_1}\|} \cdot \|x\| \cdot \underbrace{\cos \alpha_1}$$

$$\Delta w^{h_2} \cdot x = \underbrace{\|\Delta w^{h_2}\|} \cdot \|x\| \cdot \underbrace{\cos \alpha_2}$$

$$\|\Delta w^{h_1}\| \cdot \cos \alpha_1 = \|\Delta w^t\| = \|\Delta w^{h_2}\| \cdot \cos \alpha_2$$

# Minimale Störung des bisherigen Lernerfolges

Es gibt verschiedene Varianten der Gewichtsänderung, welche zu gleicher Reduzierung des Fehlers  $d^t$  um  $\Delta d^t$  führen: z.B.

$$|\Delta d^t| = |\Delta w^{t,T} x| = |\Delta w^{h_1,T} x| = |\Delta w^{h_2,T} x|$$

Wahl von  $\Delta w^t$ , der parallel zu  $x$  ist, bewirkt minimale Gewichtsänderung, da kürzester Verschiebungsvektor.

⇒ Prinzip der minimalen Störung des bisherigen Lernerfolges.

## 4.3 Adaline Gradientenabstieg Lernregel

### Überblick:

- Ensemble von Input-/Output-Tupeln
- Quadratische Fehlerfunktion
- Gradientenabstieg an Fehlerfunktion
- Direkte Ermittlung von Minimum der Fehlerfunktion
- Gradientenabstieg Lernregel

# Ensemble von Input-/Output-Tupeln

Elementarfehler bei einem Input-/Output-Tupel:

$$\begin{aligned}(d^t)^2 &= (r - w^{t,T}x)^2 \\ &= r^2 - 2rx^Tw^t + w^{t,T}xx^Tw^t\end{aligned}$$

Eigentlich ist aber das Ensemble von Input-/Output-Tupeln, d.h. das Ensemble von zugehörigen Elementarfehlern, zu beachten.

Hinweis: Im Unterkapitel wird das erweiterte Neuronenmodell mit dem Gewichtsvektor  $w^e$  angenommen, aber zur kürzeren Schreibweise das Symbol  $e$  weggelassen.

# Ensemble von Input-/Output-Tupeln

Übergang auf Erwartungswert  $E$  des Fehlers für alle Elemente aus  $\Omega_T$ .

$$\underbrace{E((d^t)^2)}_{=: \text{MSE}(w^t)} = E((r)^2) - 2 \cdot \underbrace{E(rx^T)}_{=: V^T} \cdot w^t + w^{t,T} \cdot \underbrace{E(xx^T)}_{=: A} \cdot w^t$$

Mit  $V$  ein  $(I + 1)$ -dimensionaler Spaltenvektor,

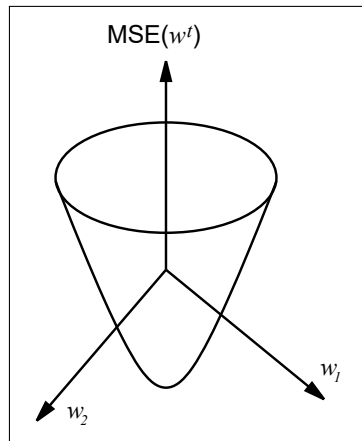
$A$  eine quadratische  $(I + 1) \times (I + 1)$  Matrix,

MSE die Mean Squared Error Funktion (im Gewichtsraum definiert).



# Quadratische Fehlerfunktion

Die Fehlerfunktion  $MSE(w^t)$  ist quadratisch in den Gewichten, d.h. konvex mit globalem Minimum. Die Position eines Punktes entspricht der Wahl der Gewichte mit zugehörigem MSE-Wert.



# Gradientenabstieg an Fehlerfunktion

- Ableiten der Fehlerfunktion (Fehlerfläche) liefert Gradientenvektor.
- Abstieg an Fehlerfläche in der negierten Gradientenrichtung, d.h. maximales Gefälle.
- Entsprechende Anpassung der Gewichte führt zum Minimum der Fehlerfunktion.

# Direkte Ermittlung von Minimum der Fehlerfunktion

Falls komplettes Ensemble von Input-/Output-Tupeln vorliegt, und die Fehlerfunktion quadratisch ist, dann kann das Minimum der Fehlerfunktion, d.h. der optimale Gewichtsvektor, direkt ermittelt werden.

$$\nabla \text{MSE}(\mathbf{w}^t) := \begin{pmatrix} \frac{\partial \text{MSE}(\mathbf{w}^t)}{\partial w_0} \\ \vdots \\ \frac{\partial \text{MSE}(\mathbf{w}^t)}{\partial w_I} \end{pmatrix} = -2 \cdot \mathbf{V} + 2 \cdot \mathbf{A} \cdot \mathbf{w}^t \stackrel{!}{=} \vec{0}$$
$$\implies \mathbf{w}^* := \mathbf{A}^{-1} \cdot \mathbf{V}$$

# Gradientenabstieg Lernregel

Falls aber wie beim Online-Lernen die Input-/Output-Tupel nur sequentiell auftreten, dann wird für das aktuelle Tupel der momentane Gradientenvektor berechnet und damit der Gewichtsvektor angepasst.

$$\begin{aligned} w^{t+1} &:= w^t - \mu \left( \frac{\partial (d^t)^2}{\partial w^t} \right) \\ &= w^t + 2\mu d^t x \end{aligned}$$

Hinweis: Beim Gradientenabstieg hat man eine garantierte aber langsame Konvergenz zum (nächsten lokalen) Minimum.

## 4.4 Anwendungen von Adaline

### Überblick:

- Adaline als Adaptiver Filter
- Adaline zur System-Modellierung
- Adaline zur Prädiktion

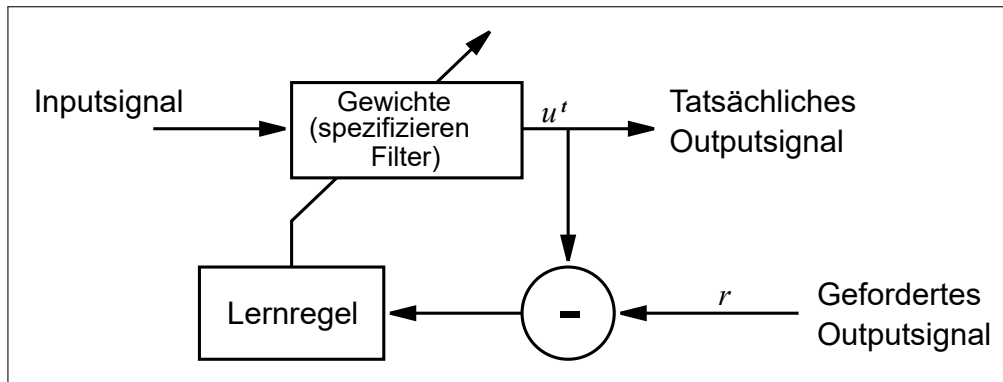
# Adaline als Adaptiver Filter

Aufgabe: Ermitteln eines Operators (Filters) für Verarbeitung eines Signals (z.B. eines Bildes).

Lösung:

- Gefordertes Output-Signal mit tatsächlichem Output-Signal vergleichen.
- Gewichtsadaption (Filteradaption) aufgrund berechneter Differenz.

# Adaline als Adaptiver Filter



# Adaline als Adaptiver Filter

Beispiel für einen gelernten Filter (Operatormuster, Faltungskern), um Grauwertkanten (in horizontaler Richtung) zu extrahieren. Siehe Grundlagen der Bildverarbeitung.

1	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1
1	1	1	2	2	2	2	2	2	1	1	1
1	1	1	2	2	2	2	2	2	1	1	1
1	1	1	2	2	2	2	2	2	1	1	1
1	1	1	2	2	2	2	2	2	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1
1	1	1	1	1	1	1	1	1	1	1	1

Faltung mit

1	0	-1
---	---	----

0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	-1	-1	0	0	0	0	1	1	0	0
0	0	-1	-1	0	0	0	0	1	1	0	0
0	0	-1	-1	0	0	0	0	1	1	0	0
0	0	-1	-1	0	0	0	0	1	1	0	0
0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

Relevant bei sog. Convolutional Neural Networks (Deep Learning).



# Adaline zur System-Modellierung

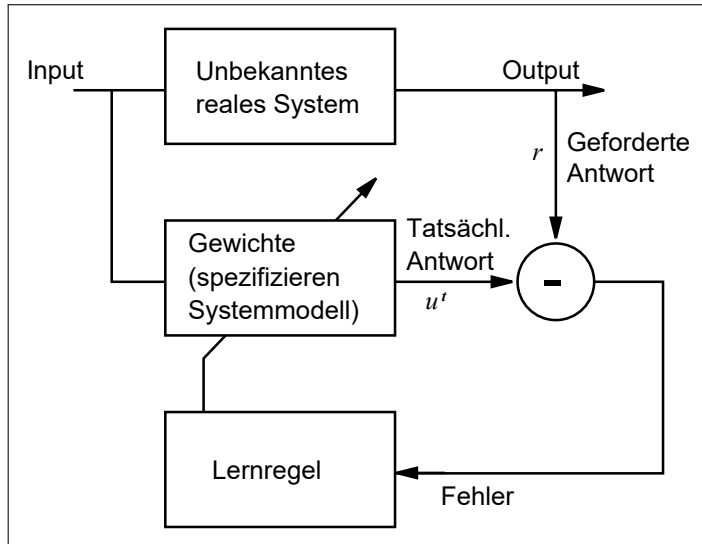
Synonym: System-Modellierung = System-Identifikation

Aufgabe: Aus beobachtbarer Input-Output-Relation soll ein System modelliert werden.

Lösung:

- Gemeinsamer Input für Adaline und unbekanntem System.
- Tatsächlicher Output des unbekannten Systems sollte der Geforderte Output des Systemmodells sein.

# Adaline zur System-Modellierung



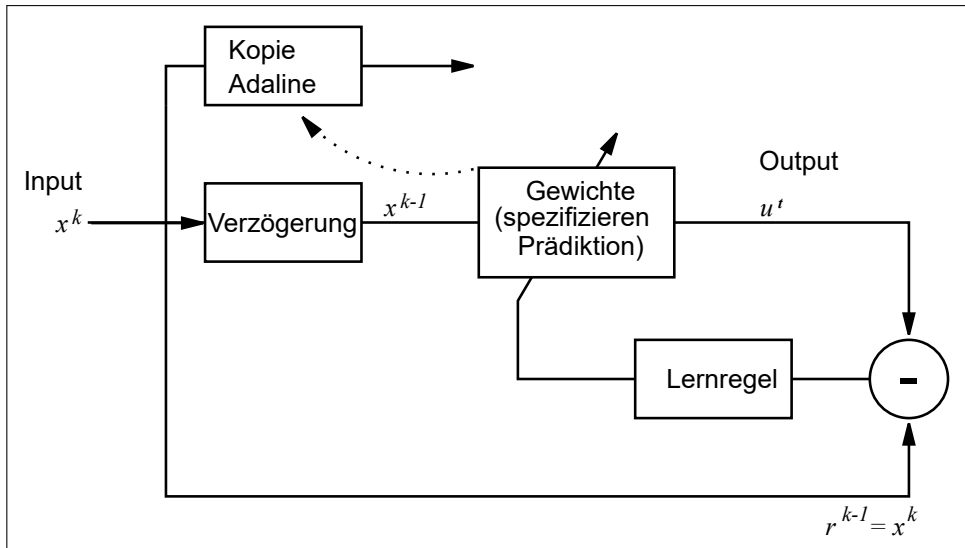
# Adaline zur Prädiktion

Aufgabe: Schätzung späterer Signale aus früheren Signalen.

Lösung:

- Früheres Signal als Input, späteres Signal als geforderter Output, Gewichtsadaption.
- Kopie der optimalen Gewichte als fertiges Adaline zur Prädiktion.

# Adaline zur Prädiktion



## 4.5 Lineare Regression durch Batch-Lernen

### Überblick:

- Aufgabenstellung der Linearen Regression
- Gleichungssystem, Parameteroptimierung
- Dimension Eingaberaum versus Anzahl Trainingselemente

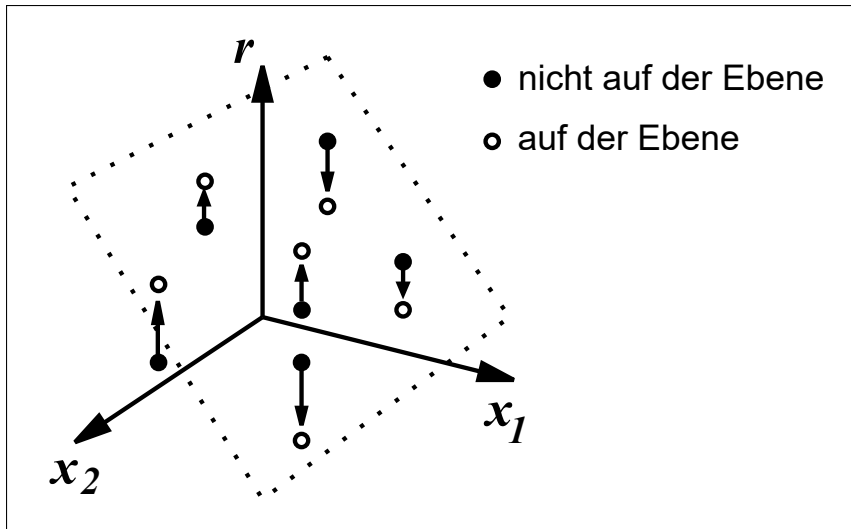
# Aufgabenstellung der Linearen Regression

Annahme: Adaline; und alle Trainingselemente vorhanden für Batch-Lernen.

Gegeben: Trainingselemente  $\mathbf{x}^1, \dots, \mathbf{x}^M$  ( $I$ -dimensional) mit zugehörigen skalaren Werten  $r^1, \dots, r^M$ .

Gesucht: Parameter der  $I$ -dimensionalen Ebene, so daß  $(\mathbf{x}^m, r^m)$  zur Ebene minimalen Abstand hat, für alle  $m \in \{1, \dots, M\}$ .

# Aufgabenstellung der Linearen Regression



# Gleichungssystem, Parameteroptimierung

Gesucht sind optimale Werte für die Parameter  $w_0, w_1, \dots, w_I$  in den Formeln  $r^m = w_0 x_0^m + w_1 x_1^m + w_2 x_2^m + \dots + w_I x_I^m + \delta^m$ , mit  $x_0^m = 1$ , und  $m \in \{1, \dots, M\}$ , sodaß die Fehler  $\delta^m$  minimal sind.

Matrixschreibweise:

Sei die Menge von  $M$  Trainingslementen  $(x^{m,T}, r^m)$  repräsentiert in Matrix  $\mathbf{H}$  und Vektor  $\mathbf{R}$ , sowie die Gewichte  $w_i$  in Vektor  $\mathbf{w}^e$  und die Fehler  $\delta^m$  in Vektor  $\Delta$ .

Dann gilt folgendes System von Gleichungen:  $\mathbf{R} = \mathbf{H} \cdot \mathbf{w}^e + \Delta$



# Gleichungssystem, Parameteroptimierung

Definitionen:

$$\mathbf{H} := \begin{pmatrix} 1 & \mathbf{x}_1^1 & \cdots & x_I^1 \\ \vdots & & \vdots & \\ 1 & \mathbf{x}_1^M & \cdots & x_I^M \end{pmatrix}; \quad (M \times (I + 1)) \text{ Eingabedaten-Matrix}$$

$$\mathbf{R} := (r^1, \dots, r^M)^T; \quad M - \text{dimensionaler Solldaten-Vektor}$$

$$\mathbf{w}^e := (w_0, w_1, \dots, w_I)^T; \quad (I + 1) - \text{dimens. Gewichts-Vektor}$$

$$\Delta := (\delta^1, \dots, \delta^M)^T; \quad M - \text{dimensionaler Fehler-Vektor}$$

# Dimension Eingaberaum, Anzahl Trainingselemente

**Fall a:**  $(I + 1) = M = \text{Rang}(H)$

Trainingselement  $(x^{m,T}, r^m)$  liefert Hyperebene im Gewichtsraum.

Betrifft Dualität Eingabe-/Gewichtsraum (siehe auch Unterkapitel 3.5).

Insgesamt erhält man  $M$  Hyperebenen der Form:

$$(1, x^{m,T}) \cdot w^e = r^m, m \in \{1, \dots, M\}.$$

Eindeutige Lösung  $w^{e,*}$ , als Schnittpunkt der  $M$  Hyperebenen im Gewichtsraum.

Gleichungssystem  $H \cdot w^e = R$ , Lösung  $w^{e,*} := H^{-1} \cdot R$

Fehlervektor  $\Delta$  ist der Nullvektor, d.h. Fehlerfunktion  $D(w^{e,*}) = 0$

# Dimension Eingaberaum, Anzahl Trainingselemente

**Fall b:**  $\text{Rang}(\mathbf{H}) = (I + 1) < M$

Überbestimmtes Gleichungssystem:  $\mathbf{R} = \mathbf{H} \cdot \mathbf{w}^e + \Delta$

Es existiert keine fehlerfreie, aber eine eindeutig bestimmte Lösung  $\mathbf{w}^{e,*}$ .  
Diese wird gefunden am globalen Minimum der Fehlerfunktion

$$D(\mathbf{w}^e) := \sum_{m=1}^M \left( r^m - \sum_{i=0}^I w_i x_i^m \right)^2$$

Bestimmung von Vektor  $\mathbf{w}^{e,*}$  durch Minimierung von  $D(\mathbf{w}^e)$ .  
Dabei gilt unter Verwendung von Matrizen:

$$D(\mathbf{w}^e) = \|\Delta\|^2 := (\mathbf{R} - \mathbf{H} \cdot \mathbf{w}^e)^T \cdot (\mathbf{R} - \mathbf{H} \cdot \mathbf{w}^e)$$

# Dimension Eingaberaum, Anzahl Trainingselemente

Lösung: Partielle Ableitungen der Fehlerfunktion  $D(w^e)$  nach den Parametern  $w_i$ , und diese jeweils Null-Setzen.

$$\begin{aligned}\nabla_{w^e}[(R - H \cdot w^e)^T \cdot (R - H \cdot w^e)] = \\ -2 \cdot H^T \cdot R + 2 \cdot H^T \cdot H \cdot w^e \stackrel{!}{=} \vec{0}\end{aligned}$$

$$\text{Daraus folgt } w^{e,*} := \underbrace{(H^T \cdot H)^{-1} \cdot H^T}_{\text{Pseudo-Inverse von } H} \cdot R$$

Bei linearer Regression und quadratischer Fehlerfunktion führt der Gradientenabstieg zur Pseudo-Inversen. Die Lösung  $w^{e,*}$  ergibt sich direkt aus obiger Formel.

# Dimension Eingaberaum, Anzahl Trainingselemente

**Fall c:**  $\text{Rang}(H) = M < (I + 1)$

Unterbestimmtes Gleichungssystem:  $R = H \cdot w^e + \Delta$

Es existiert keine eindeutige Lösung, sondern sogar mehrere Lösungen für  $w^e$ , und alle diese Lösungen sind fehlerfrei. Die Fehlerfunktion hat an diesen Stellen also den Wert Null.

Hinweis zur Bestimmung dieser Lösungen: Singulärwertzerlegung.