Neuronale Verfahren zur Funktionsapproximation

- Klassifikation mit dem Perceptron von Rosenblatt
- Vom Perceptron zum Multilagen-Perceptron
- Error-Backpropagation Lernregel
- Radiale Basisfunktionen-Netze

Das Rosenblatt-Perceptron (1962)

Vorher: Regression. Jetzt: Binäre Klassifikation

Neuronale Struktur:
 Das binäre Modell-Neuron

$$\hat{y}(\mathbf{x}) = \Theta(\mathbf{w}^T \mathbf{x} + b), \quad y = \hat{y} + n$$

• Wiederholung:

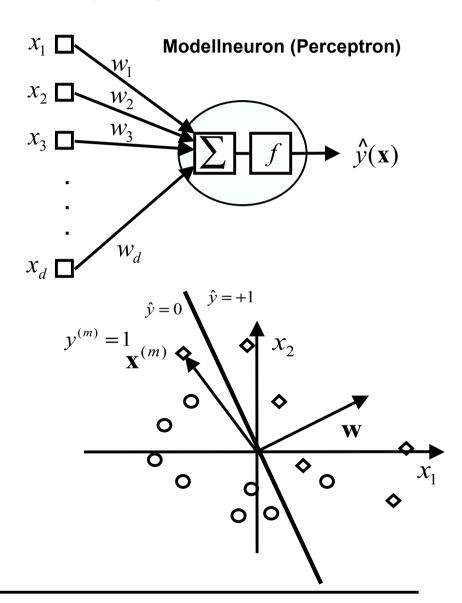
Das Perceptron kann als binärer Klassifikator dienen

- -- Geg: Daten in zwei Klassen, y = 1, 0
- -- Ges: Klassifikationsregel als Fkt. von x

-- Def:
$$x_{d+1} = 1$$
, $w_{d+1} = b \Rightarrow \hat{y} = \Theta(\mathbf{w}^T \mathbf{x})$

- Lernen: Minimieren einer Fehlerfunktion. Hier: Gewichteter Falsch-Klassifikationsfehler
- Korrekte Klassifikation: Kein Fehler
 Falsch-Klassifikation: positiver Fehler

$$F_P(w) = -\sum_{m=1}^{M} (y^{(m)} - \hat{y}(\mathbf{x}^{(m)})) \mathbf{w}^T \mathbf{x}^{(m)}$$

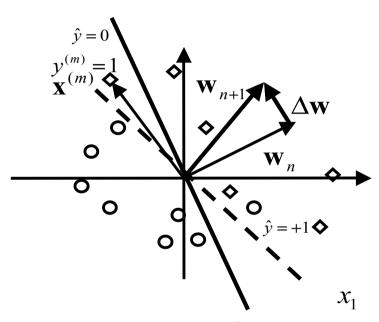


• Perceptron-Lernregel: Gradientenabstieg

$$\Delta \mathbf{w} = -\eta \nabla_{\mathbf{w}} F_P(\mathbf{w}) = \eta \sum_{m=1}^{M} (y^{(m)} - \hat{y}(\mathbf{x}^{(m)})) \mathbf{x}^{(m)}$$
 (Batch-Modus)

Online-Lernregel:

- 1. Wähle beliebiges \mathbf{w}_0
- 2. Wähle Muster m
- 3: $\mathbf{w}_{n+1} = \mathbf{w}_n + \eta(y^{(m)} \hat{y}(\mathbf{x}^{(m)}))\mathbf{x}^{(m)}$ (d.h. tue nichts bei korrekter Klassifikation, biege w bei Falsch-Klassifikation)
- 4: Bis beste Klassifikation gehe zu 2

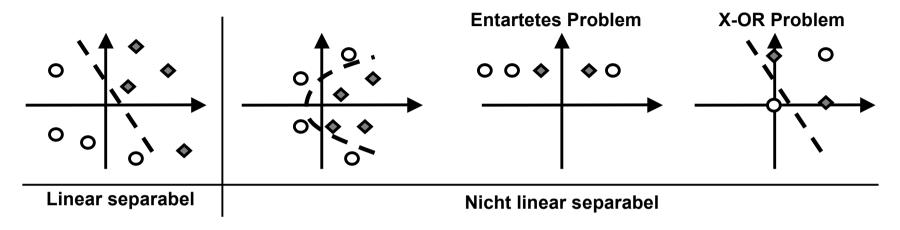


• Bem:

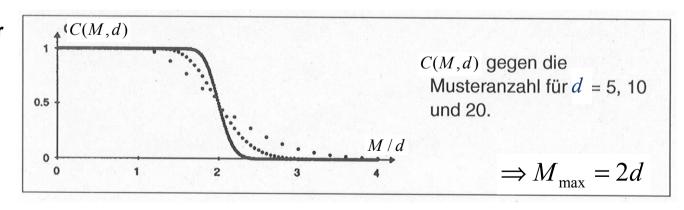
- -- Für linear separable Probleme Konvergenzgarantie in endl. vielen Schritten
- -- Funktioniert nur für linear separable Probleme
- -- Erweiterbar auf kontinuierliche Outputs

Lineare Separabilität:

Klassen können durch Hyperebene vollständig getrennt werden



- Wie wahrscheinlich sind *M* Muster in *d* Dimensionen linear separabel?
 - -- Wähle zufällig $M\,d$ dimensionale Muster
 - -- Verteile zufällig Klassenlabels
 - -- Bestimme Wa. für lin. Separabilität
 - -- Erg. f. hohe
 Dimensionen:



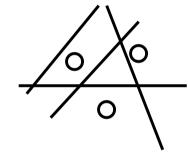
Behandlung nicht linear separabler Probleme 1: Dimensionalität

Lineare Klassifikation entspricht einer bestimmten Modellkomplexität

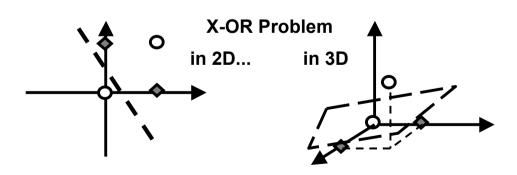
- -- Wie komplex sind die Funktionen, die ein solches Modell implementieren kann?
- --> "Zerschmettern" von M Datenpunkten: Fähigkeit, alle $\mathbf{2}^M$ möglichen Funktionen zu implementieren.

Def: <u>Vapnik-Chervonenkis-Dimension</u> (VC-Dimension) eines Modells: Größtes M, das das Modell zerschmettern kann

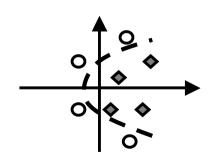
VC-Dimension e. linearen Klassifikators in d Dimensionen: d+1



Lineare Klassifikation in höherer Dimension ... oder ...



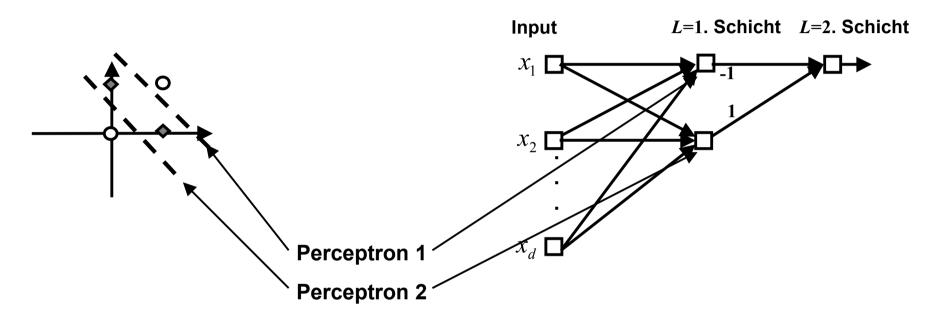
Nichtlineare Klassifikation



Vgl: Support Vector Machine (siehe später)

Vom Perceptron zum Multilagen-Perceptron

Behandlung nicht linear separabler Probleme 2: Mehrere Lagen



Perceptron:

- Regression "komponentenweise nichtlinearer Funktionen"
- Klassifikation linear sparabler
 Probleme

Multilagen Perceptron:

- Allgemeine
 Funktionsapproximation
- Regression und Klassifikation



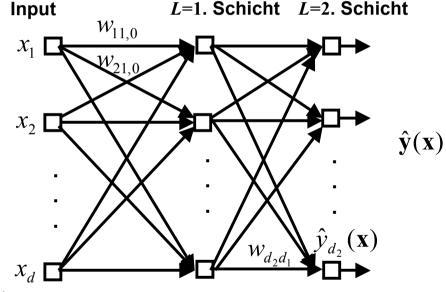
Das Multilagen-Perceptron (MLP) für Regression

- Motivation:
 - Nervenzellen im Gehirn sind hintereinandergeschaltet
 - Geschachtelte nichtlineare Transformationen sind m\u00e4chtiger als eine einzige
- Aufbau
 - Feed-forward
 - Jedes Neuron gibt nichtlineare Funktion g des summierten Inputs weiter

$$y_{i,0} = x_i, \quad d_0 = d$$

$$y_{i,L+1} = g \left(\sum_{j=1}^{d_L} w_{ij,L} y_{j,L} - \theta_i \right)$$

- mit $d_L \to d_L + 1$, $y_{d_L + 1, L} = -1$, $w_{id_L + 1, L} = \theta_i$ $y_{i, L + 1} = g\left(\sum_{i=1}^{d_L} w_{ii, L} y_{i, L}\right)$
- Bsp: zwei Schichten: $\hat{y}_p \equiv y_{p,L=2} = g \sum_{q=1}^{d_1} u_{pq} g \sum_{r=1}^{d} v_{qr} x_r$



g(x)

Transferfunktion

Sigmoide

Als Regressionsmodell geschrieben

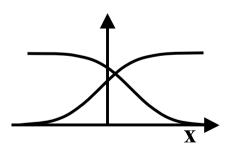
$$y_{i} = \hat{y}_{i}(\mathbf{x}) + n_{i} = f_{i}(\mathbf{x} \mid \mathbf{w}) + n_{i} \equiv g \sum_{j=1}^{d_{L}} w_{ij,L} g \sum_{k=1}^{d_{L-1}} w_{jk,L-1} \dots g \sum_{q=1}^{d} w_{pq,0} x_{q} + n_{i}$$

• Wichtigkeit der Nichtlinearität: Falls g linear, entspricht MLP einschichtigem Fall

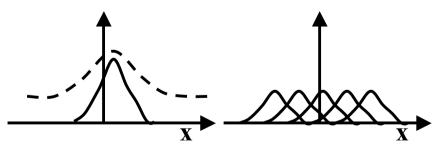
$$y_{i} = \sum_{j=1}^{d_{1}} w_{ij,1} \sum_{k=1}^{d} w_{jk,0} x_{k} = \sum_{k=1}^{d} \left(\sum_{j=1}^{d_{1}} w_{ij,1} w_{jk,0} \right) x_{k} = \sum_{k=1}^{d} u_{ik} x_{k}$$

Universale Approximationseigenschaft

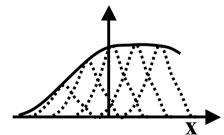
Ein dreischichtiges Netz kann jede beliebige kontinuierliche Funktion approximieren Anschaulich:



1. Schicht
Pool von kontinuierlichen
Perceptrons



2. Schicht
Lokalisierte Antwort durch Kombination von
Peceptron-Outputs plus sigmoide Transformation



3. Schicht
Kombination zur
gewünschten Funktion

• Unklar

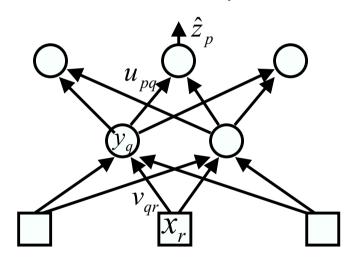
- Wieviele versteckte Neuronen werden benötigt
- Wie sichert man sich gegen Überfitten ab? (=> Kreuzvalidierung)
- Wie trainiert man die versteckten Neuronen?
 - => Error-Backpropagation-Algorithmus (Rumelhart, Werbos, `80er Jahre)
- Jetzt: Betrachte zweischichtiges Netzwerk

$$\hat{z}_p = g \left(\sum_q u_{pq} y_q \right) = g \left(\sum_q u_{pq} g \left(\sum_r v_{qr} x_r \right) \right)$$

- Datensatz $D = \{\mathbf{x}^{(m)}, \mathbf{z}^{(m)}\}, \ m = 1, \dots, M$
- Fehlerfunktion:

$$F(\mathbf{w}) = \sum_{m} \sum_{p} (z_{p}^{(m)} - \hat{z}_{p}^{(m)})^{2} = \sum_{m} \sum_{p} \left(z_{p}^{(m)} - g \left(\sum_{q} u_{pq} g \left(\sum_{r} v_{qr} x_{r}^{(m)} \right) \right) \right), \quad \mathbf{w} \equiv (\mathbf{u}, \mathbf{v})$$

• Modular:
$$F(\mathbf{w}) = \sum_{m} E^{(m)}(\mathbf{w})$$
 $E^{(m)}(\mathbf{w}) = \sum_{p} \left(z_{p}^{(m)} - g \left(\sum_{q} u_{pq} g \left(\sum_{r} v_{qr} x_{r}^{(m)} \right) \right) \right)$



- **Backpropagation Algorithmus: Herleitung**
 - --> Geschickte Anwendung des Gradientenabstiegs $\Delta \mathbf{w} = -\eta \nabla_{\mathbf{w}} E(\mathbf{w})$
 - Def: Inputs $a_p = \sum_q u_{pq} y_q$, $b_q = \sum_r v_{qr} x_r$
 - Def. Fehler: $E(\mathbf{w}) = \sum_{p} \left(z_{p} g \left(\sum_{q} u_{pq} g \left(\sum_{r} v_{qr} x_{r} \right) \right) \right) =$ (ohne index m) $= \sum_{p} \left(z_{p} - g \left(\sum_{q} u_{pq} g(b_{q}) \right) \right) = \sum_{p} \left(z_{p} - g(a_{p}) \right)^{2} = \sum_{p} \left(z_{p} - \hat{z}_{p} \right)^{2}$
 - Partielle Ableitungen für Versteckt-zu-Output Gewichte

$$-\frac{\partial E}{\partial u_{ij}} = 2(z_i - \hat{z}_i) \frac{\partial \hat{z}_i}{\partial u_{ij}} = 2(\underline{z_i - \hat{z}_i}) \underline{g'(a_i)} \frac{\partial a_i}{\partial u_{ij}} = \delta_i y_j$$

$$\delta_i = 2(z_i - \hat{z}_i) \underline{g'(a_i)}$$
"Error"

$$\delta_i = 2(z_i - \hat{z}_i)g'(a_i)$$
"Error"

Partielle Ableitungen für Input-zu-Versteckt Gewichte

$$-\frac{\partial E}{\partial v_{jk}} = \sum_{p} \underbrace{2(z_{p} - \hat{z}_{p})g'(a_{p})}_{\delta_{p}} \frac{\partial a_{p}}{\partial v_{jk}} = \sum_{p} \delta_{p} \frac{\partial}{\partial v_{jk}} \left(\sum_{q} u_{pq} g \left(\sum_{r} v_{qr} x_{r} \right) \right)$$

$$= \underbrace{\sum_{p} \delta_{p} u_{pj} g'(b_{j})}_{=:\delta_{j}} \underbrace{\frac{\partial b_{j}}{\partial v_{jk}}}_{\delta_{k}} = \delta_{j} x_{k},$$

$$\delta_{j} = \sum_{p} \delta_{p} u_{pj} g'(b_{j})$$

• Backprop-Algorithmus:

Nach Zufallsbelegung der Gewichte auf kleine Werte

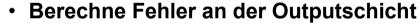
• Präsentiere Muster $\mathbf{X}^{(m)}$

Signal forward-propagation:

 Berechne und speichere von jeder Schicht die Inputs und Aktivitäten, also

$$b_{q} = \sum_{r} v_{qr} x_{r}^{(m)}, \ y_{q} = g(b_{q})$$

$$a_{p} = \sum_{r} u_{pq} y_{q}, \ \hat{z}_{p} = g(a_{p})$$



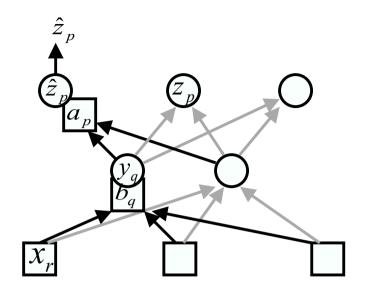
$$\delta_p = 2(z_p - \hat{z}_p)g'(a_p)$$

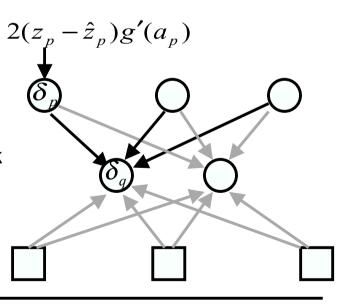
Error-Backpropagation:

propagiere Fehler durch versteckte Schichten zurück

$$\delta_q = \sum_p \delta_p u_{pq} g'(b_q)$$

• Lerne mit Regel:
$$\Delta u_{ij} = -\eta \frac{\partial E}{\partial u_{ij}} = \eta \delta_i y_j$$
, $\Delta v_{jk} = \eta \delta_j x_k$

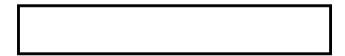




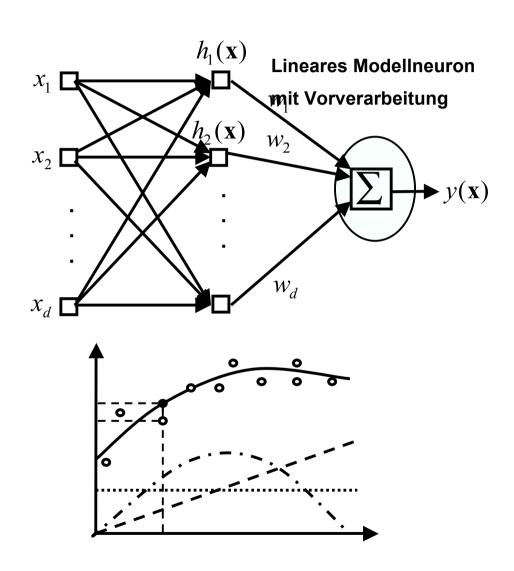
Wiederholung: Das lineare Modell

$$y(\mathbf{x}) = h_1(\mathbf{x})w_1 + h_2(\mathbf{x})w_2 + \dots + h_d(\mathbf{x})w_d$$
$$= \sum_{i=1}^d h_i(\mathbf{x})w_i$$

- Funktionen h(x) sind von Hand vorgegeben, werden nicht gelernt
- Aber: Einfache Lösung



analytisch berechenbar!



Radiale Basisfunktionen (RBF) Netzwerke

Motivation:

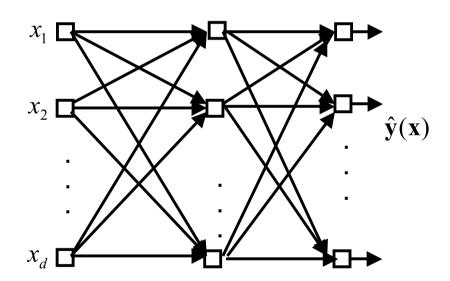
- Hybrid zwischen linearem Modell und universalem Funktionsapproximator
- Approximiere Output als lineare Summe nichtlinearer (Gauss-) Funktionen

Aufbau

- 1. Schicht Gauss-Funktionen (anstatt h(x)) approximieren Segmente des Inputs
- 2. Schicht: Linearkombination der Gaussfunktionen

Input

- 1. Schicht
- 2. Schicht



- Bem: -- 1. Schicht nichtlinear, aber Gaussisch => effizient zu optimieren
 - -- 1. Schicht entspricht Schicht 1,2 des universalen Approximators
 - -- 2. Schicht linear => einfache analytische ML Lösung

• Optimierung

- Optimiere μ , Σ der Gaussfunktionen nur basierend auf Input-Daten
 - -- Entspricht Mixture of Gaussian Dichteschätzer
 - Verwende optimale Gaussfunktionen als Modellfunktionen eines linearen Modells

