

1 Begriffliche Grundlagen

Deduktion (logisches Schließen): Vom Allgemeinen aufs Spezifische schließen (*von generellen Axiomen auf spezielle Theoreme schließen*).

Induktion (Lernen durch Generalisierung von Beispielen): Vom Spezifischen aufs Allgemeine schließen (*von speziellen Beispielen auf generelle Regeln, Theorien schließen*).

↔ Schlussfolgerung über die Validität einer Theorie basiert auf Erfahrung

↔ Theorien sind nie "wahr" oder "existent"

↔ Theorien können falsifiziert werden, aber nicht verifiziert

↔ Theorien können nur durch Beobachtungen gestützt werden

Abduktion: Schließen ausgehend von Beispielen anhand von Regeln. Von einem speziellen Beispiel wird durch Rückwärtsanwendung von bekannten Regeln versucht daraus zu schließen (kann logisch richtig zu falschen Schlüssen führen).

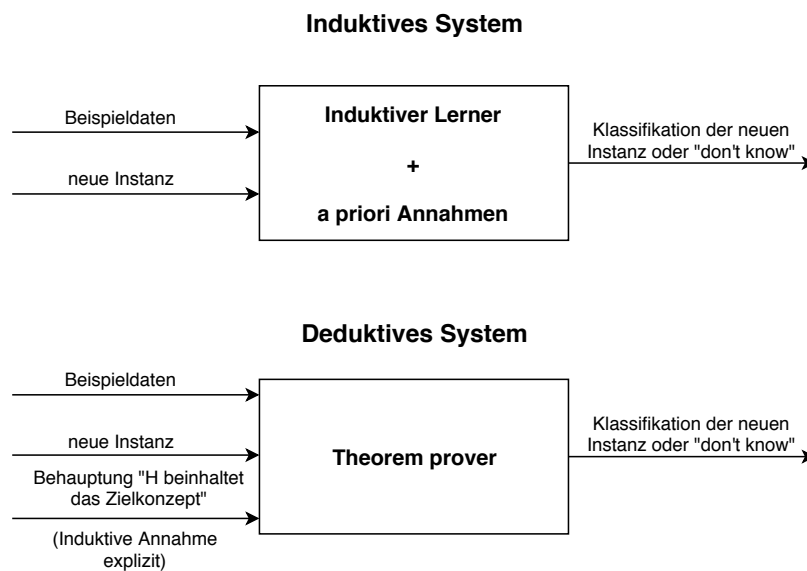


Abbildung 1: Induktives Lernen und äquivalentes deduktives System

Generalisierung (generalization): Gelerntes Modell auf neue Eingaben (die nicht zum Lernen benutzt wurden) anwenden.

Feature extraction: Eingabedaten werden Vorverarbeitet (pre-processing) damit das Problem einfacher wird (auch um Zeit-/Platzkomplexität zu reduzieren).

↔ Dimensionsreduktion

↔ Informationsverlust

Klassifikation (classification): Jede Eingabe wird einer endlichen Anzahl von diskreten Kategorien zugewiesen.

Regression (regression): Jede Eingabe wird einer oder mehreren (Vektor) stetigen Variablen zugewiesen.

Eingaben (inputs): Menge von Eingaben $X = \{x_1, \dots, x_n\}$ zu einer Funktion.

Ausgaben (targets): Menge von Sollausgaben $T = \{t_1, \dots, t_n\}$ zu einer Funktion.

Trainingsdaten (training set): Große Menge von Eingabedaten $X = \{x_1, \dots, x_n\}$, manchmal mit den zugehörigen Sollausgaben $T = \{t_1, \dots, t_n\}$.

Lineare Modelle (linear models): Funktionen, die linear in den unbekannten Parametern sind.

Inductive learning hypothesis: Jede Hypothese, die die Zielfunktion gut über eine ausreichend große Anzahl an Trainingsdaten approximiert, approximiert die Zielfunktion auch gut über neue Daten.

↔ Daten sind charakteristisch für generelle Regel/Funktion

↔ ist immer eine implizite Annahme

Learning bias: Bezeichnet die impliziten Annahmen über die Art der Generalisierung.

↔ ist immer notwendig!

Modellselektion (model selection): Auswahl einer parametrisierte Form des (Daten-) Modells.

Parameteroptimierung (parameter optimisation, learning): Optimierte (lerne) die Parameter eines ausgewählten Modells anhand der Trainingsdaten.

2 Überwachtes Lernen (supervised learning)

- gegeben eine Anzahl an Trainingsdaten (Eingaben, Sollausgaben)
- **Ziel**: finde Modell für die Daten
- **Ziel**: Generalisierung (Exploitation)
- verschiedenen Möglichkeiten für **Exploitation**
 - **deterministisch**: berechne neue Ausgabe für neue Eingabe
 - **probabilistisch**: erzeuge zufällige Ausgabe für neue Eingabe mit gelernter Wahrscheinlichkeitsverteilung (predictive distribution)

2.1 Regression (Deterministisch)

Datenmodell

- Polynom vom Grad M

- **Annahmen:**
 - ↪ Form des Modells: Polynom
 - ↪ Grad des Polynoms
 - ↪ bestimmt Komplexität des Modells
 - ↪ bestimmt die Anzahl der Parameter

Parameteroptimierung

- Summe der quadratischen Fehler:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2$$

- **Annahmen:**
 - ↪ quadratische Fehlerfunktion
 - ↪ "gute" Parameter für das Modell finden
- Betrachte $E(\mathbf{w})$ als Funktion der Parameter \mathbf{w}
- **Minimiere** quadratischen Fehler bezüglich \mathbf{w}
- wir erhalten durch Einsetzen der Trainingsdaten den **Trainingsfehler**
- Danach: teste mit Daten, die nicht zum Lernen verwendet wurden, wir erhalten den **Testfehler**
 - ↪ **Ziel:** Minimiere den Testfehler

Problem: Finden der richtigen Modellkomplexität.

- Fehlerminimierung allein funktioniert nicht
 - ↪ $E(\mathbf{w}) = 0$ ist bei verrauschten Daten unerwünscht
 - ↪ ist das Modell zu komplex (zu viele freie Parameter), dann kann es das Rauschen lernen $E(\mathbf{w}) \ll 1$
 - ↪ wenn Rauschen gelernt ist, dann ist die Generalisierung schlecht \Rightarrow **Overfitting**
- E_{RMS} zeigt Overfitting an

$$E_{RMS} = \sqrt{2E(\mathbf{w}^*)/N}$$

Maßnahmen gegen Overfitting:

- ↪ Modellselektion
- ↪ verändere #Daten relativ zu #Parameter

↪ stelle zusätzliche Bedingungen an die Parameter ⇒ **Regularisierung**

Regularisierung (hier):

- verbiete große Koeffizientenwerte
 - ↪ beugt Oszillation vor (Bias!)
 - ↪ macht die gelernte Funktion glatter (Bias!)

Regularisierung durch Bestrafung großer Koeffizienten:

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \underbrace{\frac{\lambda}{2} \|\mathbf{w}\|^2}_{\text{Regularisierungsterm}}$$

$\lambda = \text{wie wichtig ist Regularisierung}$

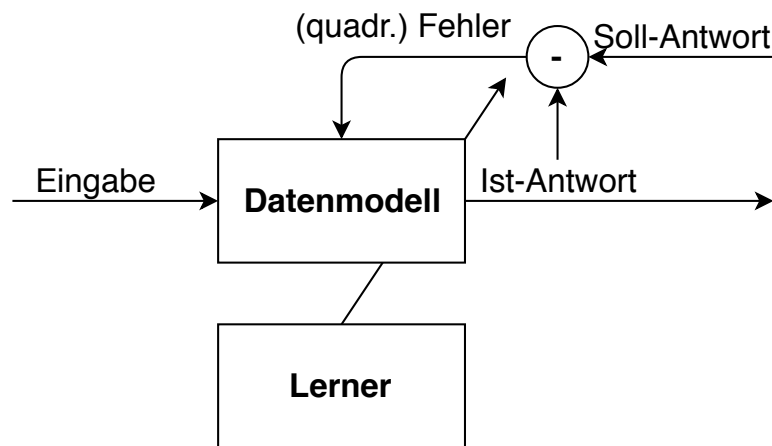


Abbildung 2: Vorgehensweise beim überwachten Lernen

2.2 Regression (Probabilistisch)

Modellierung durch Likelihood

- für gegebenes parametrisiertes Datenmodell $y(x, \mathbf{w})$
- **Likelihood** = bedingte Wahrscheinlichkeit eine Ausgabe t bei gegebenen Parametern \mathbf{w}, β und Eingabe x_0 zu erhalten: $p(t|x_0, \mathbf{w}, \beta)$.
 - ↪ ist nach Annahme gaußverteilt mit Mittelwert $y(x_0, \mathbf{w})$ und Varianz $\beta^{-1} = \sigma^2$ (1-dim. Fall):

$$p(t|x_0, \mathbf{w}, \beta) = N(t|y(x_0, \mathbf{w}), \sigma^2)$$

- $\hookrightarrow \beta$ heißt Präzision
- \hookrightarrow multidimensional:

$$p(\mathbf{t}|\mathbf{x}_0, \mathbf{w}, \Sigma^{-1}) = N(\mathbf{t}|\mathbf{y}(\mathbf{x}_0, \mathbf{w}), \Sigma^{-1})$$

Zugrundeliegende Annahmen

- es gibt nur einen Datensatz
- es gibt einen "wahren, idealen" Satz Modellparameter
- \mathbf{w}_{ML} , β_{ML} sind gute Approximationen
- Modellierung der Unsicherheit durch Störung resultiert in Verteilung
- Daten sind endlich und zufällig
- damit sind \mathbf{w}_{ML} , β_{ML} ebenfalls datenabhängig zufällig
- Unsicherheit in der Wahl des Datensatzes ist (noch) nicht modelliert

Von der Messung zur Wahrscheinlichkeit

- für gegebenes parametrisiertes Datenmodell $y(x, \mathbf{w})$ und gegebene Sollausgabe t

$$t = y(x, \mathbf{w}) + \underbrace{N(0, \beta^{-1})}_{\text{Rauschen}} \Leftrightarrow t - y(x, \mathbf{w}) \sim N(0, \beta^{-1}) \Leftrightarrow t \sim N(y(x, \mathbf{w}), \beta^{-1})$$

- **Datenmodell:** Normalverteilung um das Modell $y(x, \mathbf{w})$ mit Eingabe x und Parametern \mathbf{w}

Likelihood für alle Daten (Data-Likelihood)

Annahme: Daten unabhängig von einander erzeugt.

- \hookrightarrow gemeinsame Verteilung = Produkt der einzelnen Likelihoods:

$$\begin{aligned} L(\mathbf{w}) &= P(T|X, \mathbf{w}) \\ &= \prod_{n=1}^N N(t_n|y(x_n, \mathbf{w}), \beta^{-1}) \\ &= \prod \frac{1}{\mathcal{N}} e^{-\frac{(t_n - y(x_n, \mathbf{w}))^2}{2\sigma^2}} \end{aligned}$$

Parameteroptimierung

- Data-Likelihood ist ein stochastisches Datenmodell
- $L(\mathbf{w})$ ist Funktion aller Parameter des Datenmodells und von β
- Parameteroptimierung = Maximierung des Likelihood
 - ↪ maximiert Wahrscheinlichkeit, die gemessenen Ausgaben zu beobachten, gegeben die Modellparameter und Eingaben

Vorgehen:

- bilde $-\log L(\mathbf{w})$
 - ↪ Produkt wird zur Summe
- dann finde $\argmin -\log L(\mathbf{w})$
 - ↪ führt wieder zur Minimierung des quadratischen Fehlers!
 - ↪ wir erhalten die optimalen Parameter \mathbf{w}_{ML} , β_{ML}

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = - \underbrace{\frac{\beta}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln(\beta) - \frac{N}{2} \ln(2\pi)$$

↪ Berechne \mathbf{w}_{ML} indem $E(\mathbf{w})$ minimiert wird

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N (y(x_n, \mathbf{w}_{ML}) - t_n)^2$$

Generalisierung

- verwende die \mathbf{w}_{ML} , β_{ML} Parameter
- dann ist die optimale Output Verteilung

$$p(t|x, \mathbf{w}_{ML}, \beta_{ML}) = N(t|y(x, \mathbf{w}_{ML}), \beta_{ML}^{-1})$$

- die wahrscheinlichste Ausgabe für ein neues x ist der Mittelwert $y(x, \mathbf{w}_{ML})$
- aber: auch zufälliges Generieren von Ausgaben mit maximum Likelihood möglich (sampling, generatives Modell)
- β_{ML} gibt Konfidenz an

Bayes'scher Ansatz

Interpretiere Wahrscheinlichkeiten als Wissen/Unsicherheit über Parameter.

Annahme: Es gibt nur einen Datensatz, der unvollständig bekannt ist.

- verwende wieder stochastisches Datenmodell
- wenn neue/weitere Daten gemessen werden, dann ändert sich das Wissen/die Unsicherheit über die Parameter
- modelliere initiale Unsicherheit über die Parameter als $P(w)$
- $P(w)$ ist die a priori Wahrscheinlichkeit bevor die Daten beobachtet werden
- modelliere den Likelihood wie vorher
- dann berechne die a posteriori Wahrscheinlichkeit mit der Bayes Formel

$$P(w|D) = \frac{P(D|w)P(w)}{P(D)},$$

$$\underbrace{p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta)}_{\text{posterior}} \propto \underbrace{p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)}_{\text{likelihood}} \underbrace{p(\mathbf{w}|\alpha)}_{\text{prior}},$$

$$p(\mathbf{w}|\alpha) = N(\mathbf{w}|0, \alpha^{-1}I) = \frac{\alpha^{(M+1)/2}}{2\pi} e^{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}}.$$

↪ initiale Unsicherheit abhängig vom Hyperparameter α

Der **posterior** drückt das Wissen/die Unsicherheit über die Modellparameter aus nachdem die Daten beobachtet wurden.

↪ Daten werden selbst nur unter Unsicherheit beobachtet und damit durch stochastisches Datenmodell approximiert

Parameteroptimierung

- bilde $\mathbf{w}_{\mathbf{MAP}} = \operatorname{argmax} P(w|D) \Leftrightarrow \operatorname{argmin} -\log P(w|D)$

$$\hookrightarrow \beta \tilde{E}(\mathbf{w}) = \underbrace{\frac{\beta}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2}_{\text{likelihood}} + \underbrace{\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}}_{\text{prior}}$$

↪ Maximum a posteriori für gaußverteiltes stochastisches Datenmodell ist äquivalent zur Fehlerminimierung + Regularisierung

↪ **Overfitting** ist "automatisch" verhindert

Generalisierung

- Generalisiere durch $t_{new} = \underbrace{y(x_{new}, \mathbf{w}_{MAP})}_{\text{Mittelwert der posterior-Verteilung}}$
- möglich wäre auch "ziehen" eines Parameters \mathbf{w}' aus der posterior-Verteilung: $\mathbf{w}' \sim p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta)$ und Generalisierung durch $y(x_{new}, \mathbf{w}')$

Bayesian Predictive Distribution

Vollständige Berücksichtigung von Unsicherheit.

- bekannt: Data-Likelihood und Parameter posterior
- integriere über alle möglichen Parameter gewichtet mit ihrer Wahrscheinlichkeit:

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t})d\mathbf{w} = N(t|m(x), s^2(x))$$

- diese Wahrscheinlichkeit kann für lineare Modelle ebenfalls explizit berechnet werden und ist gaußverteilt
- dann generalisiere durch ziehen aus $\underbrace{p(t|x, \mathbf{x}, \mathbf{t})}_{\text{predictive distribution}}$

Bayesian Predictive Distribution vs. Maximum Likelihood

Der volle Bayes'sche Ansatz zeigt mehr Unsicherheit, da die Unsicherheit in den Parametern auch modelliert ist!

Inkrementelles Bayes'sches Lernen

Wenn Daten sequentiell beobachtet werden, wiederhole den Inferenzschritt durch Anwendung der Bayesregel für jeden Datenpunkt oder Datensatz D_1, D_2, D_3, \dots :

$$\begin{aligned} P(w|D_1) &= \frac{P(D_1|w)P(w)}{P(D_1)} \\ P(w|D_2) &= \frac{P(D_2|w)P(w|D_1)}{P(D_2)} \\ P(w|D_{k+1}) &= \frac{P(D_{k+1}|w)P(w|D_k)}{P(D_{k+1})} \end{aligned}$$

Hier wird der letzte Posterior zum neuen Prior!

2.3 Lineare vs. nichtlineare Modelle

Parameteroptimierung: Lineare Datenmodelle

Lineare Modelle haben die generelle Form

$$y(x_n, \mathbf{w}) = \sum_{m=0}^M w_m \Phi_m(x_n) = \mathbf{w}^T \Phi(x_n)$$

wobei $\phi(x_n) = (\Phi_1(x_n), \dots, \Phi_m(x_n))^T$

$\hookrightarrow \Phi_m(x)$ kann eine beliebige (nichtlineare) Funktion der Eingabe sein

Beispiele:

Polynom: $\Phi_j(x) = x^j$

j -th component: $\Phi(\mathbf{x}) = x_j$

RBF-net: $\Phi_j(x) = e^{-\frac{\|\mu_j - x\|^2}{2\sigma^2}}$

neural net: $\Phi_j(x) = \sigma(\frac{x - \mu_j}{s})$, $\sigma(a) = \frac{1}{1 + e^a}$

wenn μ_j, s_j zufällig gewählt werden: Extreme Learning Machine (ELM).

Fehlerminimierung

Maximum Likelihood

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \Phi(x))^2$$

Minimierung:

- bilde den Gradienten
- setze = 0
- löse auf

Wir erhalten das Resultat:

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = \Phi^\# \mathbf{t}$$

wobei $\Phi^\# = (\Phi^T \Phi)^{-1} \Phi^T$ die Moore-Penrose Pseudoinverse ist und

$$\Phi(X) = \begin{pmatrix} \Phi_1(x_1) & \dots & \Phi_M(x_1) \\ \vdots & \ddots & \vdots \\ \Phi_1(x_n) & \dots & \Phi_M(x_n) \end{pmatrix} \in \mathbb{R}^{N \times M} \text{ ist die sogenannte Designmatrix.}$$

Maximum A-Posteriori

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \Phi(x))^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- minimiere Fehler + Regularisierung
- berechne Gradienten, setze gleich null, löse

Resultat:

$$\mathbf{w}_{MAP} = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

(Erwünschter) Seiteneffekt: die Matrix wird numerisch stabiler.

Full Bayesian linear regression

Wir erhalten folgendes analytisches Resultat für die predictive distribution:

$$\begin{aligned} p(t|x, \mathbf{x}, \mathbf{t}) &= \int p(t|x, \mathbf{w}) p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w} = N(t|m(x), s^2(x)) \\ m(x) &= \beta(\Phi(x)^T S \sum_{n=1}^N \Phi(x_n) t_n) \\ s^2(x) &= \beta^{-1} + \Phi(x)^T S \Phi(x) \\ S^{-1} &= \alpha I + \beta \sum_{n=1}^N \Phi(x_n) \Phi(x_n)^T \end{aligned}$$

($\Phi(x_n) = (x_n^0, \dots, x_n^M)^T$ für Polynommodell).

Parameteroptimierung: nichtlineare Modelle

Problem: Minimiere quadratischen Fehler.

- ↔ keine analytische Lösung
- ↔ suche Minimum iterativ
- ↔ meist durch Gradientenabstieg
- ↔ notwendig für
 - direkte Minimierung des Fehlers als Funktion der Parameter

$$\mathbf{w}^* = \operatorname{argmin} E(\mathbf{w})$$

- Maximum Likelihood

$$\begin{aligned} \mathbf{w}_{ML} &= \operatorname{argmax} \log P(D|\mathbf{w}) \\ &\sim \operatorname{argmin} E(\mathbf{w}) \end{aligned}$$

- direkte Minimierung des Fehlers als Funktion der Parameter + Regularisierung

$$\mathbf{w} = \operatorname{argmin} [E(\mathbf{w}) + \lambda \|\mathbf{w}\|^2]$$

- Maximum A-Posteriori

$$\begin{aligned} \mathbf{w}_{\mathbf{MAP}} &= \operatorname{argmax} \log P(\mathbf{w}|D) \\ &\sim \operatorname{argmin} [E(\mathbf{w}) + \lambda \|\mathbf{w}\|^2] \end{aligned}$$

Gradientenabstieg = iterative Optimierung, die lokales Minimum findet:

- definiere Startpunkt \mathbf{w}_0 (guter Tipp oder zufällig)
- iteriere: $\mathbf{w}_{k+1} = \mathbf{w}_k - \mu \nabla E(\mathbf{w}_k)$
- bis $\nabla E(\mathbf{w}_k) \approx 0$

Gradientenabstieg führt zu einer allgemeinen Lernregel mit Lernrate μ .

Wobei $E(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2$. Das Modell $y(x, \mathbf{w})$ kann eine beliebige (nichtlineare) Funktion der Parameter und des Inputs sein. **Problem:**

- Minimum hängt von Startwert ab
- zu große Schrittweite kann Minimum verfehlen/oszillieren
- zu kleine Schrittweite konvergiert langsam

⇒ adaptive Schrittweitenbestimmung.

Problem: Berechnung der Gradienten mittels Kettenregel.

↔ wenn Modell komplex, muss Kettenregel rekursiv angewandt werden

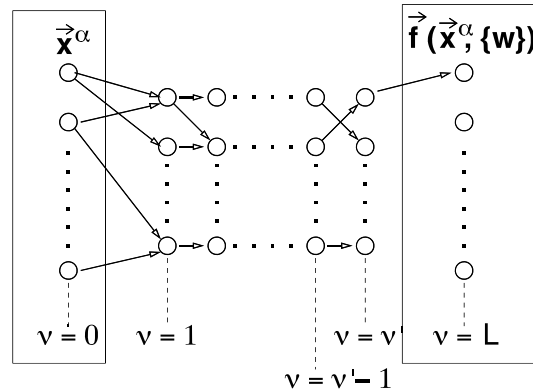
Multilayer Perceptron (MLP)

Abbildung 3: Multilayer Perceptron (MLP)

L = Anzahl der Schichten, $L-1$ innere Schichten ("hidden layers")

Für jeden Knoten:

$$s_i^{v+1} = \sigma\left(\sum w_{ij}s_j^v\right)$$

$$\sigma(a) = \tanh(a) \text{ oder } \sigma(a) = \frac{1}{1 + e^{-a}}$$

Modellselektion**Kapazität**

- Jede kontinuierliche Funktion kann theoretisch mit jeglicher Genauigkeit angenähert werden
- Die Anzahl notwendiger Neuronen (Knoten) ist aber unbekannt

Konkrete Modellparameter

- Anzahl der Schichten
 - ↔ eine innere Schicht ist ausreichend
 - ↔ häufig funktionieren 2 Schichten aber besser
 - ↔ modern: deep learning (viele, spezielle Schichten)
- Anzahl der Neuronen bestimmt Modellkomplexität
- Overfitting ist häufig ein Problem (Regularisierung notwendig)
- Feature extraction wichtig

Berechnung des Gradienten mittels Backpropagation

- Rekursive Formel zur Berechnung der Gradienten bzgl. der inneren Parameter w_{ij}
- Analytischer Ausdruck vorhanden (direkte Lösung möglich)
- Geringer Rechenaufwand
- **Backpropagation** = Methode zur Berechnung von Gradienten \Rightarrow funktioniert auch für andere Modelle

2.4 Lokale Modelle

Radial Basis Funktionen-Netz (RBF)

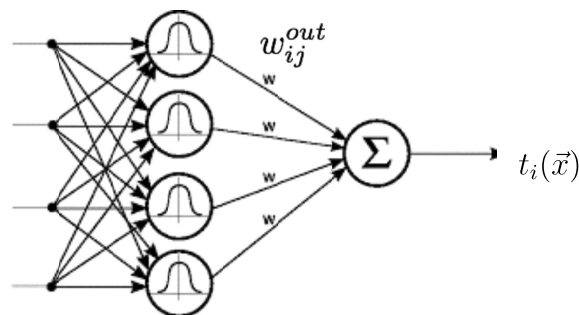


Abbildung 4: Radiales Basis Funktionen-Netz

- innere Schicht aus $m = 1 \dots M$ radialen Basisfunktionen (RBF, Gaussfunktionen)
- lineare Kombination an der Ausgabe mit Parametern (Gewichten w_m^{out})

$$t_i(\mathbf{x}) = \sum_{j=1}^J w_{ij}^{out} \cdot \Phi_j(\mathbf{x}, \mathbf{w}_j^{in}), \text{ mit } i \in 1, \dots, L, j \in 1, \dots, J$$

- RBF hängen nur vom Abstand des inputs zum Zentrum ab:

$$\Phi_j(\mathbf{x}, \mathbf{w}_j^{in}) = R(\|\mathbf{x} - \mathbf{w}_j^{in}\|), j \in 1, \dots, J$$

$$\hookrightarrow R() \text{ ist eine Gaussverteilung } R(s) = e^{-\frac{1}{2\sigma_j^2} s^2}$$

$$\hookrightarrow \text{RBF: } \Phi_m(\mathbf{x}, \mathbf{w}_j^{in}) = e^{-\frac{\|\mathbf{x} - \mathbf{w}_j^{in}\|^2}{2\sigma_j^2}}$$

RBF sind:

- universelle Approximatoren
- lineare Modelle in der Ausgabeschicht
- gut verstanden, da viele Methoden um die Parameter der Gaussfunktion zu wählen
- die RBF-Mittelwerte können auf (einige) Daten gesetzt werden
- Wahl von Sigma ("Breite") wichtig um gute Abdeckung der Daten zu erreichen

Berechnung:

- für L-dimensionale Ausgabe

$$t_i(\mathbf{x}) = \sum_{j=1}^M w_{ij} \cdot \Phi_j(\mathbf{x}, \mathbf{w}_j^i), \quad i \in 1, \dots, L, \quad j \in 1, \dots, M$$

- Zusammenfassung der Parameter in Vektor Θ
- und für 1-D Ausgabe: Linearkombination nichtlinearer RBF Φ

$$t(x) = \sum_{j=1}^M w_j \cdot \Phi_j(x, \Theta)$$

Lösung:

- für RBF wird der quadratische Fehler des linearen Modells zu

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^T \Phi(x_n))^2$$

- und die Minimierung hat, wie bekannt eine analytische Lösung
 \hookrightarrow Ableitung bilden, gleich Null setzen, umstellen fertig!
- wie vorher

$$\begin{aligned} \hat{\mathbf{w}} &= \operatorname{argmin} E(\mathbf{w}) = \|\mathbf{T} - \underline{\Phi} \cdot \mathbf{w}\| \\ \Leftrightarrow & -\underline{\Phi}^T (\mathbf{T} - \underline{\Phi} \cdot \mathbf{w}) = 0 \\ \Rightarrow & \hat{\mathbf{w}} = (\underline{\Phi}^T \cdot \underline{\Phi})^{-1} \underline{\Phi}^T \cdot \mathbf{T} \end{aligned}$$

- wobei \mathbf{T} alle Daten enthält (multidimensional, dann ist \mathbf{T} Matrix)
- $\underline{\Phi}$ ist die Design Matrix $\Phi_{ij} = \Phi_j(\mathbf{x}_j)$

Lineare gewichtete Summe von (nichtlinearen) Gaussfunktionen ist äquivalent zu Gaussgewichteter Summe von (konstanten) linearen Funktionen:

$$t(x) = \sum_{j=1}^M w_j \cdot \Phi_j(x, \Theta) \Leftrightarrow t(x) = \sum_{j=1}^M \Phi_j(x, \Theta)(w_j)$$

Gewichtete lineare Regression

Wir betrachten Regression nun mit zusätzlichen Gewichtungen für die einzelnen Datenpunkte.

Methode gewichteter kleinster Quadrate

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N d_n (t_n - \mathbf{a}^T x_n)^2$$

mit \mathbf{a} als Parametervektor für eine Gerade

Hinweis: Im einfachen linearen Modell werden hier die Datenpunkte direkt verwendet und nicht vorverarbeitet $\Phi_m(\mathbf{x}) = x_m$.

Die Lösung der Normalgleichung wird dann zu

$$\mathbf{a}^* = (\Phi^T D \Phi)^{-1} \Phi^T D \mathbf{t} = (X^T D X)^{-1} X^T D \mathbf{t}$$

mit $D_{nn} = d_n$, wobei D eine diagonale Skalierungs-Matrix darstellt $D = \begin{pmatrix} d_1 & \dots & \\ \vdots & \ddots & \vdots \\ & \dots & d_n \end{pmatrix}$.

Als Gewichtsfunktion wird die Dichtefunktion der Normalverteilung mit Mittelwert μ_1 und Kovarianz Σ genutzt:

$$d_n = \Phi(x_n, \Theta) = N(x_n | \mu_1, \Sigma) = e^{-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)}$$

damit tragen nur Datenpunkte in der Nähe des Erwartungswertes effektiv zum Fehler bei.

Lokale gewichtete Regression (LWR)

- bestimme mehrere Geraden-Segmente für verschiedene Gauß-Kurven, die jeweils mit μ_e, Σ_e parametrisiert sind

$$\mathbf{a}_e = (X^T D_e X)^{-1} X^T D_e \mathbf{t}, \text{ für } e = 1 \dots E$$

mit Gaußverteilter Gewichtsfunktion mit $\Theta = (\mu_e; \Sigma_e)$

- finale Modell wird dann als Summe überlagerter, linear skaliert Gauß-Funktionen bzw. aus überlagerten Geradensegmenten gebildet:

$$t(x) = \sum_{e=1}^E \Phi(x, \Theta_e) (\mathbf{a}_e^T x)$$

- es ist auch möglich die Geraden zu verschieben

$$t(x) = \sum_{e=1}^E \Phi(x, \Theta_e) (\mathbf{a}_e^T x + b_e)$$

indem eine erweiterte Regressions-Matrix verwendet wird

$$X = \begin{pmatrix} x_{11} & \dots & x_{1p} & 1 \\ \vdots & \ddots & \vdots & 1 \\ x_{n1} & \dots & x_{np} & 1 \end{pmatrix}, x_n \in R^p$$

Hinweis: hier ist X die (erweiterte) Designmatrix für die "innere Regression", da ja $\Phi_m(\mathbf{x}) = x_m$.

RBF ist Spezialfall von LWR ($E = M, b_e = w_e, a \equiv 0$).

↔ Das RBF kann als Summe sehr einfacher Modelle (einer Konstanten w_j) verstanden werden, die mit Gauss-Funktionen gewichtet sind.

Damit:

LWR: Überlagerung von Geraden, die lokal gewichtet sind.

↔ benötigt E Regressionsschritte und ist nichtlinear in den Parametern

↔ ist ein allgemeineres Modell

RBF: Überlagerung von konstanten Funktionen mit Gaußgewichten.

↔ benötigt nur einen Schritt, da linear in den Parametern

Vereinheitlichtes Modell (unified modell)

Viele Algorithmen verwenden die gleiche Modellform. Daher kann die gegebene Form

$$t(x) = \sum_{e=1}^E \Phi(x, \Theta_e)(\mathbf{a}_e^T x + b_e)$$

als vereinheitlichtes (unified) Modell verstanden werden.

Spezialfälle sind zum Beispiel:

- einfache lineare Regression
- RBF
- Gauß-gewichtete lineare Regression
- Gauß-gewichtete lokale lineare Regression
- auch probabilistische Modelle
- noch weitere Modelle

2.5 Modellierung von Unsicherheit (statistischer Ansatz)

Annahmen:

- \mathbf{w}^* ist echter Parameter eines parametrisierten Modells
- Daten sind verrauscht $t_n = f(x_n, \mathbf{w}) + \nu$
- durch einen der besprochenen Algorithmen (z.B. Maximum Likelihood) erhalten wir eine Schätzung für den echten Parameter \mathbf{w}^* :

$$\hat{\mathbf{w}}(Z), \quad Z = (X, T)$$

- ↔ $\hat{\mathbf{w}}(Z)$ wird **Schätzer** genannt und ist selbst eine Zufallsvariable, weil die Daten Z verrauscht sind.
- ↔ Versuch mit neuen Daten wiederholen um einen neuen Schätzer zu berechnen $\hat{\mathbf{w}}'(Z)$

Gewünschte Eigenschaften:

- ein Schätzer sollte **erwartungstreu** (unbiased) sein:
 - liefert im Mittel den "wahren" Parameter
 - d.h. der Schätzer macht keinen systematischen Fehler
 - Erinnerung: braucht die Annahme, dass es einen "wahren" Parameter \mathbf{w}^* gibt
 - falls es einen Bias (= systematische Abweichung) gibt, kann man diesen manchmal ausrechnen und diesen dann wiederum in den Schätzer integrieren
- ein Schätzer sollte minimale Varianz haben
 - geringe Varianz erhöht die Wahrscheinlichkeit bei einer/wenigen Durchführungen des Experimentes dem wahren Wert nahe zu kommen

Aber: Es gibt nicht immer einen erwartungstreuen Schätzer.

⇒ Schätzer soll erwartungstreu und eine minimale Varianz besitzen

$$\langle (\hat{\mathbf{w}} - \langle \hat{\mathbf{w}} \rangle)^2 \rangle \rightarrow \min.$$

↔ diese Schätzer werden als **UMV** (unbiased minimal variance) bezeichnet

Cramer-Rao Grenze

Es gibt eine untere Grenze für die Varianz von erwartungstreuen Schätzern, die mit der log-likelihood $l(\mathbf{w})$ verbunden ist:

$$\text{Var}(\hat{\mathbf{w}}) \geq \frac{1}{-\langle \frac{\delta^2 l(\mathbf{w})}{\delta \mathbf{w}^2} \rangle} = \frac{1}{\langle \left(\frac{\delta l}{\delta \mathbf{w}} \right)^2 \rangle}$$

Interpretation: Je flacher (gemessen durch die mittlere Krümmung, bzw. Ableitung) die log-likelihood Funktion, desto größer die minimal erreichbare Varianz.

⇒ Schätzer, die diese Grenze erreichen werden **effizient** genannt.

$$- < \frac{\delta^2 l(\mathbf{w})}{\delta \mathbf{w}^2} >_{\mathbf{z}} \begin{cases} \text{large} \\ \text{small} \end{cases} \Leftrightarrow \text{Var}(\hat{\mathbf{w}}) \begin{cases} \text{small} \\ \text{large} \end{cases}$$

Bestimmung von effizienten Schätzern:

- Die Cramer-Rao Grenze ist erreicht, wenn folgendes gilt

$$\text{Var}(\hat{\mathbf{w}}) < \left(\frac{\delta l}{\delta \mathbf{w}} \right)^2 > = 1$$

- Es gibt eine Zerlegung der folgenden Form

$$\frac{\delta \ln P(\mathbf{z}|\mathbf{w})}{\delta \mathbf{w}} = \lambda(\mathbf{w}) \cdot (\hat{\mathbf{w}}(\mathbf{z}) - \mathbf{w})$$

Mit den Funktionen $\lambda(\mathbf{w})$, $\hat{\mathbf{w}}(\mathbf{z})$ und dem wahren Parameter \mathbf{w}

- $\lambda(\mathbf{w}) = \text{Var}(\hat{\mathbf{w}}(\mathbf{z}))^{-1}$ und $\hat{\mathbf{w}}(\mathbf{z})$ sind effiziente Schätzer

Lineare Regression zur Schätzung

Betrachte das generelle, lineare, stochastische Modell mit nichtlinearer Eingangsfunktion $\Phi(\mathbf{x})$

$$y = w_0 + \sum_{j=1}^M w_j \Phi_j(\mathbf{x}) + \nu$$

Wir definieren \mathbf{x} als Datenvektor nach der Vorverarbeitung neu $\mathbf{x}_k = (\Phi_1(x_k), \dots, \Phi_M(x_k))^T \in \mathbb{R}^M$. Mit der folgenden Notation lässt sich das Modell sehr kompakt schreiben:

$$\begin{aligned} X &= \begin{pmatrix} 1 & 1 & \dots & 1 \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_k \end{pmatrix}^T \in \mathbb{R}^k \times \mathbb{R}^{M+1} \\ \mathbf{w} &= (w_0, \dots, w_k)^T \in \mathbb{R}^k \\ \boldsymbol{\nu} &= (\nu_1, \dots, \nu_k)^T \in \mathbb{R}^k \\ \mathbf{y} &= (y_1, \dots, y_k)^T \in \mathbb{R}^k \\ \mathbf{y} &= X\mathbf{w} + \boldsymbol{\nu} \end{aligned}$$

Dafür erhalten wir mit der Zerlegung die folgenden Schätzer

$$\begin{aligned} \Lambda(\mathbf{w}) &= \frac{1}{\sigma^2} X^T X \\ \mathbf{w}(z) &= (X^T X)^{-1} X^T \mathbf{y} \\ \Rightarrow \hat{\mathbf{w}} &= (X^T X)^{-1} X^T \mathbf{y} \end{aligned}$$

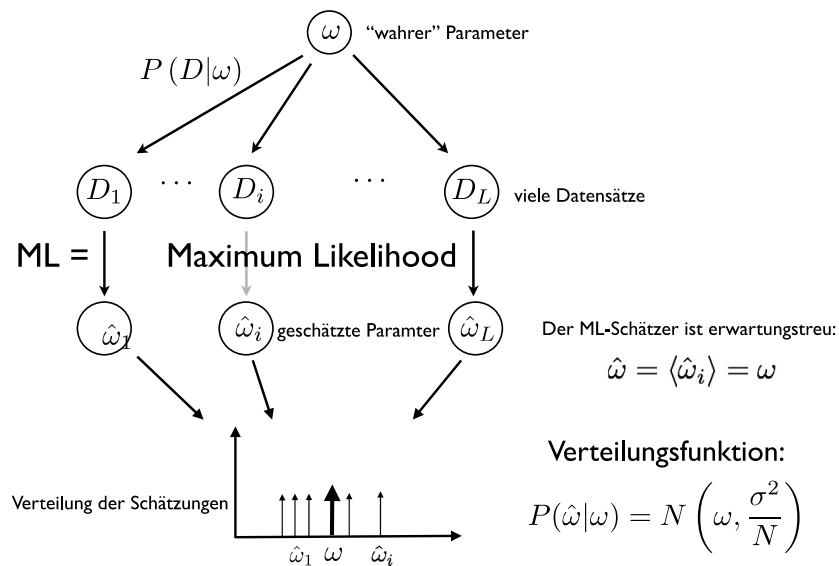


Abbildung 5: Die statistische Interpretation des Likelihood

Maximum Likelihood als Schätzer

In der Theorie liefert die Maximum Likelihood Methode UMV Schätzer, wenn ein Parameter gefunden werden kann, der die Likelihood maximiert

$$\begin{aligned}\hat{\mathbf{w}}(D) &= \operatorname{argmax} P(D|\mathbf{w}) \\ P(\hat{\mathbf{w}}|\mathbf{w}) &\sim N(\mathbf{w} | \frac{\sigma^2}{N})\end{aligned}$$

- der Schätzer ist dann selbst normalverteilt um den Mittelwert
- die Varianz nimmt mit zunehmender Datenmenge ab

In der Praxis wissen wir, dass wir den Fehler minimieren müssen aber häufig nur lokale Minima finden können.

2.6 Klassifikation

Aufgabe: Ordne Datum x einer von K diskreten, disjunkten Klassen zu.

- ↔ sofern die Klassen linear trennbar sind, können die Grenzen durch Hyperebenen dargestellt werden

Linear separierbarer Fall

Definition einer Hyperebene:

$$\mathbf{w}^T \mathbf{x} + w_0 = 0 \text{ (Gerade für } \mathbf{x} \in \mathbb{R}^2 \text{)}$$

Zuordnung (hier für zwei Klassen):

Klasse 0 : $\mathbf{w}^T \mathbf{x} + w_0 < 0$,

Klasse 1: $\mathbf{w}^T \mathbf{x} + w_0 > 0$,

Grenze : $\mathbf{w}^T \mathbf{x} + w_0 = 0$.

$K = 2$ Klassen:

Hier gilt für die Label/Sollausgaben t :

$$t \in \{-1, 1\} \text{ oder } t \in \{0, 1\}$$

\hookrightarrow Die Label können im letzteren Fall als Wahrscheinlichkeit interpretiert werden

$K > 2$ Klassen:

Hier wird gewöhnlich die 1-of- K Kodierung verwendet

$$\begin{aligned} \mathbf{t} &= (0, \dots, 0, 1, 0, \dots, 0) \in \mathbb{R}^K \\ t_k &= 1 \Leftrightarrow x \in C_k \end{aligned}$$

wobei C_k die k -te Klasse bezeichnet.

Drei Ansätze:

1. Finde eine direkte Abbildung f von den Eingängen auf die Klassen.
 - f heißt Diskriminatenfunktion (**Diskriminate**)
 - Wahrscheinlichkeiten werden nicht betrachtet

$$f(x) : x \rightarrow 0, 1$$

2. Bayes'scher Ansatz:

- a) Modellieren der Likelihood
- b) Modellierung der a-priori Wahrscheinlichkeiten für jede Klasse
- c) Berechne im nächsten Schritt den Posterior $p(C_k|x)$

3. Direkte Modellierung des Posteriors

Fisher-Diskriminante

Ansatz zur linearen Trennung (d.h. Diskriminante f ist Hyperebene)

- berechne Projektion der Daten auf eine Linie
- optimiere die Wahl der Linie um beste Trennung zu erreichen
- wähle Wert auf der Linie, zur Diskriminierung

Berechnung (zwei Klassen C_1, C_2 mit N_1, N_2 Punkten):

- Mittelwert $m_{1/2}$
- Projektion von 2D auf 1D mit linearem Modell

$$y(\mathbf{x} = \mathbf{w}^T \mathbf{x})$$

- \hookrightarrow Projizierte Punkte liegen auf Geraden in Richtung \mathbf{w}
 - \hookrightarrow Projektion verläuft senkrecht
- einfacher Ansatz: Abstand der projizierten Mittelwerte maximieren (maximize inter-class variance)

$$\max |m'_1 - m'_2| \Leftrightarrow \max |\mathbf{w}^T \cdot m_1 - \mathbf{w}^T \cdot m_2|$$

$\Rightarrow \mathbf{w} \propto (m_2 - m_1)$ (\mathbf{w} zeigt in Richtung $m_2 - m_1$) Beachte:

- Richtung, Länge von \mathbf{w} unwichtig
 - Projektionen können stark überlappen
- Besser: Zusätzlich die Varianz innerhalb der Klasse minimieren (minimize intra-class variance)

$$s_k^2 = \sum_{x_n \in C_k} (w_n x_n - m'_k)^2$$

Dann maximiere die folgende Kostenfunktion:

$$J(\mathbf{w}) = \frac{(m'_1 - m'_2)^2}{s_1^2 + s_2^2}$$

- $J(\mathbf{w})$ wird maximal, wenn gilt:

$$\begin{aligned} \mathbf{w} &\propto S_{\mathbf{w}}^{-1}(m_2 - m_1) \\ S_{\mathbf{w}} &= \sum_{x_n \in C_1} (x_n - m_1)(x_n - m_1)^T + \sum_{x_n \in C_2} (x_n - m_2)(x_n - m_2)^T \end{aligned}$$

Bemerkung: Abstand zwischen Mittelwerten nicht mehr maximal. Es lässt sich zeigen, dass das selbe Ergebnis durch Minimierung des quadratischen Fehlers erzielt werden kann

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T x_n + w_0 - t_n)^2 \text{ mit } t_n = \frac{N}{N_n}$$

Achtung: die traditionelle Bezeichnung "Fisher Diskriminante" ist etwas irreführend!

- Vektor \mathbf{w} gibt nur Projektionsgerade an, nicht die Klassengrenze
- Diskriminierungsfunktion notwendig

Optimale Klassifikation für neuen Datenpunkt $x \in C_1$ wenn $\mathbf{w}^T(x - m) > 0$ und sonst in der anderen Klasse.

Bayes'scher Ansatz

- modelliere die Likelihood $p(x|C_k)$ und Prior $p(C_k)$ (hier Likelihood = class conditional)
 - ↪ Indirekter Ansatz um die Wahrscheinlichkeit für eine Klasse (posterior) $p(C_k|x)$ zu finden
- Bayes Theorem für 2 Klassen:

$$p(C_1|x) = \frac{p(x|C_1)p(C_1)}{p(x|C_1)p(C_1) + p(x|C_2)p(C_2)} = \frac{1}{1 + e^{-a}}, \quad a = \ln \frac{p(x|C_1)p(C_1)}{p(x|C_2)p(C_2)}$$

- Generalisierung für $k > 2$ Klassen:

$$p(C_k|x) = \frac{p(x|C_k)p(C_k)}{\sum_j p(x|C_j)p(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)} \text{ (softmax Funktion), } a_j = \ln p(x|C_j)p(C_j)$$

- die Softmax Funktion kann als kontinuierliche differenzierbare Version des Maximums verstanden werden. Wenn $a_k \gg a_j$ dann gilt: $p(C_k|x) \approx 1, p(C_j|x) \approx 0, j \neq k$

Der Bayes'sche Ansatz mit sigmoiden Funktionen führt zu einer einfachen Lösung wenn die Klassen als Normalverteilungen modelliert werden.

- Annahme: gleiche Kovarianz

$$p(x|C_k) = \frac{1}{2\pi^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1}(x - \mu_k)\right)$$

- für den Posterior erhält man dann lineare Klassengrenzen:

$$p(C_1|x) = \sigma\left(\ln \frac{p(x|C_1)p(C_1)}{p(x|C_2)p(C_2)}\right)$$

- dies wird auch **generalized linear model** genannt, wobei die Ausgabe durch die sigmoid Funktion nichtlinear wird und als Wahrscheinlichkeit interpretierbar ist

$$\begin{aligned} p(C_1|x) &= \sigma(\mathbf{w}^T x + w_0) \\ \mathbf{w} &= \Sigma^{-1}(\mu_1 - \mu_2) \\ w_0 &= -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(C_1)}{p(C_2)} \end{aligned}$$

Anmerkung: Wegen der Annahme gleiche Kovarianz kürzen sich bei der Berechnung die Normalisierungsterme und quadratischen x -Terme heraus

Direkter Maximum Likelihood für Labels:

- direkter Ansatz um Daten auf Labels abzubilden
- Voraussetzung: gelabelte Daten (x_n, t_n)
- Annahmen:
 - wie vorher: $p(x|C_k)$ normalverteilt
 - Prior gegeben als $\pi = p(C_1) = 1 - p(C_2)$
- Damit wird Label-Data-Likelihood zu

$$p(\mathbf{t}|\pi, \mu_1, \mu_2, \Sigma) = L(\pi, \mu_1, \mu_2, \Sigma) = \prod_n [\pi N(\mathbf{x}|\mu_1, \Sigma)]^{t_n} [(1 - \pi) N(\mathbf{x}|\mu_2, \Sigma)]^{1-t_n}$$

- $\arg\max \log L()$ liefert optimale Parameter
- Likelihood zur Zuordnung von Labels auf neue Daten

Berechnung:

- zuerst nur für Terme die von π abhängen

$$\arg\max_{\pi} \sum_n (t_n \ln(\pi) + (1 - t_n) \ln(1 - \pi))$$

$$\Rightarrow \pi = \frac{1}{N} \sum_n t_n = \frac{N_1}{N} \Leftrightarrow 1 - \pi = \frac{N_2}{N}$$

- intuitiv richtig: A-priori \sim Anzahl der Punkte in der jeweiligen Klasse

- Mittelwert und Kovarianz

$$\begin{aligned}\mu_1 &= \frac{1}{N_1} \sum_{n=1}^N t_n x_n \\ \mu_2 &= \frac{1}{N_2} \sum_{n=1}^N (1 - t_n) x_n \\ \Sigma &= \frac{N_1}{N} s_1 + \frac{N_2}{N} s_2 \\ s_i &= \frac{1}{N_i} \sum_{x_n \in C_i} (x - \mu_i)(x - \mu_i)^T\end{aligned}$$

- intuitiv richtig: Summe der Intra-Klassen-Varianzen

Probabilistische Diskriminierung

- Posterior als Funktion von x modellieren

$$p(C_k|x) = y(x_n) = \sigma(\mathbf{w}^T x_n), x_n \in C_k$$

- diese Modelltypen nennen sich **logistische Regression**
- andere Herangehensweise: Vorgabe des Modells in Sigmoid Funktion
- optimiere Parameter der log-Likelihood

$$\begin{aligned}p(\mathbf{t}|\mathbf{w}) &= \prod_n y_n^{t_n} (1 - y_n)^{(1-t_n)} \\ \arg\max_{\mathbf{w}} \quad &- \log L(\mathbf{w}) \\ &= - \sum_{n=1}^N \ln(y_n^{t_n}) + \ln((1 - y_n)^{(1-t_n)})\end{aligned}$$

- Dies entspricht der Minimierung der sogenannten **cross-entropy**

$$E(\mathbf{w}) = - \sum_{n=1}^N \underbrace{t_n \ln(y_n)}_{0 \text{ für } t_n = 0} + \underbrace{(1 - t_n) \ln(1 - y_n)}_{0 \text{ für } t_n = 1}$$

- $y(x)$ ist dann direkt als Wahrscheinlichkeit $p(C_1|x)$ (posterior) zu interpretieren
- Optimierung mittels **gradient descent**
 - ↪ Analytische Lösung
 - ↪ Erweiterung auf mehrere Klassen

2.7 Konzeptlernen

- behandelt den allgemeinen Fall nicht numerischer Daten
- Abstraktes Framework zur Betrachtung von Mengen/Teilmengen
- $X :=$ Menge von Instanzen x
- $C :=$ **Konzept** = Untermenge von X = Funktion $c : X \rightarrow \{0, 1\}$ (**Indikator**) beschreibt Untermenge
- Anzahl der möglichen Konzepte $2^{|X|}$
- Beschreibung des Konzepts = Klassifikation aller Elemente von X

Konzeptlernen: Gewinnung einer booleschen Funktion aus den Instanzen und deren Label in den Trainingsdaten.

Äquivalent: Identifizierung einer bestimmten Untermenge aus den Trainingsdaten.

Hypothesenraum

- suche im gesamten Konzeptraum zu komplex
- schränke die Suche auf weniger Funktionen/Teilmengen ein
- $H :=$ Hypothesenraum von Funktionen h

$$h : X \rightarrow \{0, 1\}$$

- äquivalent: reduziere Anzahl möglicher Untermengen
- Hypothesenraum drückt Annahme über das Problem aus (Bias)
 - ermöglicht Generalisierung
 - induktives Hypothesenlernen
- Grundannahme (**inductive learning bias**): Jede Hypothese, die die Zielfunktion auf einem genügend großen Datensatz annähert ist auch eine gute Hypothese für den gesamten Raum X

Spezialfall: Beschreibung durch Attribute

- Datenbeschreibung als Tupel diskret-wertiger Attribute
 ("Wert(Attribut 1)", "Wert(Attribut 2)", ...)
 $N = \# \text{Kombinationen} = \# \text{Wert(Attribut 1)} \cdot \# \text{Wert(Attribut 2)} \dots \rightarrow 2^N$ Konzepte

Hypothesen

- Konjunktionen von Bedingungen an die Attribute
- "constraint(Attribut 1)" \wedge "constraint(Attribut 2)" \wedge ...
- mögliche constraints:
 - ? = alle Werte sind erlaubt
 - \emptyset = kein Wert ist erlaubt
 - v = nur der Wert v ist erlaubt

Bemerkung: Es gibt Hypothesen die semantisch äquivalent aber syntaktisch verschieden sind. Es gibt mehr Konzepte als semantisch verschiedene Hypothesen.

Allgemeinere oder genauso allgemeine (**generelle**) Hypothesen:

$$h_i \geq_g h_j$$

h_i ist genereller als h_j genau dann wenn $\forall x \in X : (h_j(x) = 1 \Rightarrow h_i(x) = 1)$.

Strikt mehr generelle Hypothesen:

"ist genereller als" = "schließt die gleiche oder Obermenge ein"

$$h_i \geq_g h_j$$

$$(h_i \geq_g h_j) \wedge (h_i \neq h_j)$$

$$\{x \in D \mid h_i(x) = 1\} \subseteq \{x \in D \mid h_j(x) = 1\}$$

Einfachster Algorithmus für attributbezogene Daten:

Algorithm 1 Generalisierung (nur positive Daten)

```

1: INIT  $h$  with  $\langle \emptyset, \emptyset, \dots \rangle$  // Initialisierung mit spezifischer Hypothese
2: for all positive examples  $x$  do // Iteriere über alle Trainingsdaten
3:   for all constraints  $a_i$  do // Iteriere über alle attribute Bedingungen
4:     if  $a_i$  is fulfilled by  $x$  then // Wenn Bedingung für Datenpunkt erfüllt ist
5:       OK
6:     else
7:       generalize  $a_i$  // Ansonsten wähle generellere Bedingung. Ersetze die
        Bedingung  $a_i$  mit der nächst generelleren Bedingung  $a_j$ , für die  $x$  die Bedingung
        erfüllt

```

Probleme:

- nur positive Beispiele werden verwendet
- Konvergenz kann nicht geprüft werden
- nicht robust gegen fehlerhafte Daten

Versionsraum

Eine Hypothese ist **konsistent** auf Beispielmenge $D \Leftrightarrow h(x) = c(x), \forall x \in D$. Das heißt die Hypothese klassifiziert korrekt für alle Datenpunkte in D . Die Menge aller konsistenten Hypothesen heißt **Versionsraum** $V_{H,D}$.

Der Versionsraum lässt sich durch zwei Hypothesensätze vollständig beschreiben:

$$\begin{aligned} G &:= \{g \in H \mid g \text{ ist konsistent und } \neg \exists g' : g' >_g g \wedge g' \text{ konsistent}\} \\ G &= \text{Satz generellster konsistenter Hypothesen} \\ S &:= \{s \in H \mid s \text{ ist konsistent und } \neg \exists s' : s >_g s' \wedge s' \text{ konsistent}\} \\ S &= \text{Satz spezifischster konsistenter Hypothesen} \end{aligned}$$

Damit wird der Versionsraum durch G und S begrenzt definiert:

$$V_{H,D} = \{h \in H \mid (\exists s \in S)(\exists g \in G) : g \geq_g h \geq_g s\}$$

Kandidateneliminierung

Eliminiere Hypothesen um Versionsraum zu erhalten. Kombiniere dazu die Ansätze Generalisierung und Spezialisierung.

Algorithm 2 Kandidateneliminierung

OUTPUT: Versionsraum

```

1: for all examples  $x$  do
2:   if example positive then
3:     generalize  $S$ 
4:   else if example negative then
5:     specialize  $G$ 
```

Konvergiert zu eine richtigen Konzept C , wenn $h = c \in H$ und alle Beispiele richtig klassifiziert werden.

Problem: Wenn ein Datenpunkt falsch gelabelt (inkonsistent) ist, wird die korrekte Hypothese fälschlicherweise gelöscht.

- Falls Versionsraum nur eine Hypothese enthält, verwende diese
- Falls Versionsraum noch viele Hypothesen zulässt
 - Generell mögliche Klassifikationsregel: Mehrheit + Konfidenz

Induktiver Bias

- L sei ein Konzeptlerner (Algorithmus)

- Für X gelte:

$$\begin{aligned}c &= \text{Konzept} : c : X \rightarrow \{0, 1\} \\ D &= \{x, c(x)\} = \text{Trainingsdatensatz (korrekt gelabelt)}\end{aligned}$$

- Bezeichne mit $L(x), c$ die Klassifizierung von $x \in X$ nach Training D
- Dann heißt ein willkürlicher Satz von Annahmen B induktiver Bias von L wenn für c und D gilt: $\forall x' \in X$ mit $(B \wedge D \wedge x')$ kann $L(x', D)$ durch Deduktion bestimmt werden

Entscheidungsbäume

- Instanzen durch diskrete Werte und Attribute verschrieben
- Klassifikation in mehreren Schritten
 - betrachte jeweils ein Attribut einzeln
- Handhabung korrupter Daten
- Gesamter Datensatz im Wurzelknoten (root)
- Knoten zerteilen Datensatz auf Basis der Attribute
- Blatt entspricht Klasse
 - Annahme: Instanzen gehören zu Klasse (bzw. Großteil der Instanzen)
- Änderung in der Auswertungsreihenfolge resultiert in unterschiedlichen Bäumen

Fragen:

- Was ist die optimale Reihenfolge?
- Welches Attribut sollte als nächstes Ausgewertet werden?
- Was heißt optimal?

ID3

1. Berechne den Informationsgehalt für jedes übrige Attribut des Unterbaums
2. Wähle Attribut, das größten Informationsgewinn (information gain) liefert

$$S = - \sum_{i=0,1} p_i \cdot \log_2 p_i$$

Zu 1.

Shannon Information Interpretation: Wie viel Information ist noch in den Instanzen enthalten?

Frage: Welches Attribut trennt die Instanzen bestmöglich indem es die meisten richtigen ja und nein Antworten erzeugt?

Zu 2.

Informationsgewinn entspricht der Abnahme an Entropy. Hohe Entropy bedeutet, dass die Instanzen nicht gut getrennt sind. Maximal für $p = 0,5$

Shannon Entropy

Frage: Wie bestimmt man Entropy in Entscheidungsbäumen?

Für gegebenen Wert des Attributs x ist die **Shannon Entropy**

$$S_{x=k} = - \sum_{i=0,1} P_{i|x=k} \cdot \log_2 P_{i|x=k}$$

Sie beschreibt die restliche Unsicherheit für gegebenen Wert k des Attributs in dem Unterbaum.

Merke: S ist klein für gute Trennung!

Summierung über alle möglichen Attribut-Werte beschreibt die gesamte Unsicherheit wenn Attribut ausgewertet wird

$$S_x = \sum_k P_{x=k} S_{x=k} = \sum_k \frac{|S_k|}{|S|} S_{x=k}$$

Wähle das Attribut, das den Informationsgewinn maximiert

$$\begin{aligned} \Delta_x S &= S - S_x \\ \hat{x} &= \operatorname{argmax}_x = \Delta_x S \end{aligned}$$

Anmerkungen:

- Ansatz auf > 2 Klassen übertragbar
- Erweiterte Ansätze erlauben kontinuierliche Attribute
- verbesserte Suche
- C 4.5, C 5.0
- decision forest

- Hypothesenraum des ID3 ist Menge möglicher Entscheidungsbäume
- Hypothesenraum ist vollständig:
 - enthält alle endlichen diskreten Klassifikationen
 - entspricht dem Konzeptraum
 - kein Bias durch a-priori Beschränkung der Hypothesen

Anmerkungen:

- Beginn mit "leerer" Hypothese
- nur eine einzelne Hypothese wird jeweils verfolgt
- Anzahl konsistenter Hypothesen kann nicht bestimmt werden
- robust gegen korrupte Daten

Induktiver Bias

ID3:

- unvollständige Suche
- Algorithmus schaut nicht voraus, ist "greedy"
- Bias durch Suchstrategie: **preference bias**

Candidate Elimination:

- unvollständiger Hypothesenraum
- vollständige Suche
- Bias durch begrenzten Hypothesenraum: **restriction bias**

Regression:

- Architekturdesign (**restriction bias**)
- Lern-Algorithmus (**preference bias**) (z.B. Gradientenverfahren findet nur lokale Minima)

Bayes'scher Ansatz

Verbinde Konzeptlernen und Bayes'sche Ansätze

- Nutzen von Prior, Likelihood und Posterior für Hypothesen

- Spezialfall: Diskrete Werte

$$P(D) = \sum_h P(D|h)P(h)$$

$$P(h|D) = \frac{P(D|h)P(h)}{\sum_h P(D|h)P(h)}$$

- **Annahme:** Hypothesen schließen sich gegenseitig aus (sind disjunkt)

Wie vorher: Betrachtung von Maximum Likelihood oder a-posteriori

$$h_{MAP} = \operatorname{argmax}_h P(h|D)$$

$$= \operatorname{argmax}_h \frac{P(D|h)P(h)}{P(D)}$$

$$= \operatorname{argmax}_h P(D|h)P(h)$$

Einfachster und naivster Ansatz: Berechne alles mit der Bayes-Regel

- Für alle Hypothesen $h \in H$, die Posteriors berechnen

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

$$\bar{h} = \operatorname{argmax}_h P(h|D)$$

Beispiel:

- D sei rauschfrei (korrekte Label)
- c stamme aus Hypothesenraum H
- kein a-priori Vorwissen

$$P(h) = \frac{1}{|H|} \Rightarrow \sum_{h \in H} P(h) = 1$$

- Likelihood

$$P(D|h) = \begin{cases} 1 & \text{wenn } y_i = h(x_i) \ \forall (x_i, y_i) \in D; \\ 0 & \text{sonst.} \end{cases}$$

- Zwei Fälle:

1. ist h nicht konsistent, dann ist $P(D|h) = 0$ und der Posterior ist $P(h|D) = \frac{0 \cdot P(h)}{P(D)} = 0$
2. ist h konsistent, dann ist der Posterior $P(h|D) = \frac{1 \cdot \frac{1}{|H|}}{\frac{|V_{H,D}|}{|H|}} = \frac{1}{|V_{H,D}|}$

Gegeben ein Hypothesenraum H , was ist die optimale Hypothese?

Problem: Maximum Likelihood gibt nicht das beste Ergebnis \Rightarrow Predictive Distribution nutzen

$$P(\bar{y}|D) = \sum_{h_i \in H} P(\bar{y}|h_i)P(h_i|D)$$

Optimale Klassifikation ist dann

$$\operatorname{argmax}_{\bar{y}} \sum_{h_i \in H} P(\bar{y}|h_i)P(h_i|D)$$

Keine andere Klassifikation kann im Mittel besser sein unter Annahme von gleichem H und a-priori Wissen.

Aber:

- unrealistisch alle Posteriors zu berechnen
- optimaler Klassifikator nicht in H enthalten

Naive Bayes

Verwende nun anstelle der Hypothesen direkt die Attribute. Naive Bayes bestimmt den Klassifikator direkt

- Annahme von N diskreten Attributen a_i
- Annahme von k Labeln y_i

$$\begin{aligned} y_{MAP} &= \operatorname{argmax}_{y_j} P(y_j|a_1, \dots, a_N) \\ &= \operatorname{argmax}_{y_j} \frac{P(a_1, \dots, a_N|y_j)P(y_j)}{P(a_1, \dots, a_N)} \\ &= \operatorname{argmax}_{y_j} P(a_1, \dots, a_N|y_j)P(y_j) \\ &= \operatorname{argmax}_{y_j} (\text{likelihood} \cdot \text{prior}) \end{aligned}$$

Bestimme nun den Likelihood und Prior aus den Daten

- Prior $P(y_j) \approx$ Frequenz der Klasse y_j in den Daten
- Likelihood
 - zählen nicht anwendbar, da nicht alle a_i für jedes Label vorhanden
 - nutze starke Annahme statistische Unabhängigkeit zwischen Attributen (in Klasse)

$$P(a_1, \dots, a_N|y_j) = \prod_{a_i} P(a_i|y_j)$$

- bestimme $P(a_i|y_j)$ dann durch zählen

\rightarrow Naive Bayes Klassifikator

$$\tilde{y} = \operatorname{argmax}_{y_j} \prod_{a_i} P(a_i|y_j)P(y_j)$$

2.8 Lazy learning

Bisher:

- Großer Aufwand für das Training
- Auswertung des gelernten Modells einfach und effizient

Alternativ:

- Speichere die Trainingsdaten oder Prototypen
- Großer Aufwand für die Auswertung
 \hookrightarrow rekombiniere gespeicherte Daten für die Ausgabe

K-nächste Nachbarn (k-nearest neighbour)

Naiv:

input: x_q

output: $f(\hat{x})$ mit $\hat{x} = \operatorname{argmin}_x \|x - x_q\|$ (Die Ausgabe von dem Beispiel, das der neuen Eingabe am nächsten ist)

Problem: Finde (berechne) den nächsten Nachbarn.

Besser:

- Bilde Mittelwert über k -nächsten Nachbarn

$$\hat{f}(x_q) = \frac{1}{k} \sum_{i=1}^k f(x_i)$$

- k - nächsten Nachbarn müssen erst gefunden werden, die ist Aufwendig

Problem: Finde k -nächsten Nachbarn.

Weighted k-nearest neighbours

- Bilde Mittelwert über k -nächsten Nachbarn und gewichte sie zusätzlich gemäß ihrer Entfernung zum Input

$$\hat{f}(x_q) = \frac{\sum_{i=1}^k w_i f(x_i)}{\sum_{i=1}^k w_i}, \text{ mit } w_i = d(x_i, x_q)^{-1} = \frac{1}{\text{euklidische Distanz}}$$

Problem:

- Finde k -nächste Nachbarn
- Definiere eine gute Metrik (kann schwierig sein für hochdimensionale Daten)

Bemerkung:

- die vorgestellten Methoden sind speicherintensiv
- der **induktive Bias** ist: die Funktion ist glatt, variiert lokal wenig
- die Methoden sind lokal, falls $k < N$
- kritisch ist die Wahl der Metrik \Rightarrow Datamining, selbst ein Lernproblem
- es gibt Algorithmen um die nächsten Nachbarn effizient zu finden
- Gewichtungen werden erst zur Laufzeit bekannt

Anmerkung: Es scheint, dass viele "deep learning" Methoden effektiv die nächste-Nachbarn Lösung berechnen, d.h. gut als "starke Kompression mit lokaler Generalisierung" verstanden werden können.

Der Fluch der Dimensionalität

Problem: Der hochdimensionale Raum ist quasi "leer".

Beispiel (Volumen einer D -dim Kugel):

Sei $r = 1$. Wie Groß ist der Anteil des Volumens der Kugel, der zwischen $r = 1$ und $r = 1 - \epsilon$ liegt?

$$V_D(r) = K_D r^D$$

$$\frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - \underbrace{(1 - \epsilon)^D}_{\substack{< 1 \\ \xrightarrow{D \rightarrow \infty} 0}} \xrightarrow{D \rightarrow \infty} 1$$

\Rightarrow das gesamte Volumen der Einheitskugel konzentriert sich für $D \rightarrow \infty$ in der Kugelschale

Kernel regression

- nutze Kernel Funktionen, die von der Distanz abhängen $K_i(d(x_i, x_q))$
- verwendet Gaussfunktionen $K_\sigma(x_q - x_i)$

$$K_\sigma(x_q - x_i) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\|x_i - x_q\|^2}{2\sigma^2}}$$

- durch Substitution erhalten wir den **Nadaraya-Watson kernel regression estimator**

$$f(x_q) = \sum_i y_i \frac{K_\sigma(x_q - x_i)}{\sum_j K_\sigma(x_q - x_j)}$$

- verwendet formal alle Trainingsdaten, praktisch sind aber nur wenige relevant (da die Gaussfunktion lokal ist)

3 Unüberwachtes Lernen (unsupervised learning)

3.1 Prototypen und Clustering

- Clustering zielt darauf, Gruppen von "ähnlichen" Datenpunkten zu finden (und gegebenenfalls zu repräsentieren)
- in einfacher Form ordnet ein Clusteringalgorithmus einen Punkt jeweils einem Prototypen zu (auch genannt: representative, code word, center)
- "Ähnlichkeit, Nähe" wird durch eine Metrik gemessen
- Metrik ist definiert durch die folgenden drei Eigenschaften:
 - Symmetrie: $d(a, b) = d(b, a)$
 - Positive Definitheit: $d(a, b) \geq 0$
 - Dreiecksungleichung: $d(a, b) + d(b, c) \geq d(a, c)$

Metriken

- Euklidische Metrik:

$$d(a, b) = \sqrt{\sum_{i=1}^D (a_i - b_i)^2}$$

- Mahalanobis Distanz:

$$d(a, b) = \sqrt{(a - b)^T \Sigma^{-1} (a - b)}, \text{ mit } \Sigma^{-1} = \langle (x - \langle x \rangle) (x - \langle x \rangle)^T \rangle$$

Vektorquantifizierung

- ein sehr einfacher Clusteralgorithmus
- gegeben die Daten x_n
- initialisiere und optimiere K Prototypen w_k , $k = 1, \dots, K$
 - Zuordnung von Daten zum nächstliegenden Prototypen bzgl. Metrik $d()$

$$x_n \mapsto w_c, \text{ mit } c = \operatorname{argmin}_k d(x_n, w_k)$$

- Lernziel: wähle die Position der Prototypen um den Quantisierungsfehler zu minimieren:

$$E_{VQ} = \frac{1}{N} \sum_{n=1}^N d(x_n, w_c(x_n))$$

Einfacher VQ-Algorithmus

1. Wähle # der Prototypen K (Modellselektion)
2. initialisiere die Position der K Prototypen
3. iteriere über die Daten:
 - a) wähle zufällig x_k
 - b) finde den nächstliegenden Prototypen w_c
 - c) bewege den Prototypen in Richtung des Datenpunktes

$$w_c^{new} = \epsilon(w_c^{old} - x_k), \text{ mit Lernrate } \epsilon \text{ (die über die Zeit geringer wird)}$$

- d) Abbruch, wenn der Fehler konstant bleibt

k-means

- ein einfacher und effektiver Clusteralgorithmus
- K verschiedene Cluster (Modellselektion)
- K Prototypen w_k , $k = 1, \dots, K$
- N Datenpunkte x_n , $n = 1, \dots, N$
- ordne jedem Punkt genau einem Cluster zu durch eine Variable

$$r_{nk} = \begin{cases} 1 & \text{wenn } \|x_n - w_k\|^2 < \|x_n - w_l\|^2 \forall l \neq k, k = 1, \dots, K \\ 0 & \text{sonst} \end{cases}$$

- minimiere den Quantisierungsfehler J (für euklidische Metrik)

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - w_k\|^2$$

- direkte simultane Optimierung von J bzgl. r_{nk} und w_k unmöglich
- Ansatz: Wechsel zwischen zwei Mengen von Variablen

Algorithm 3 k-means

Require: Number of prototypes K **Require:** Initialize prototypes $\{w_k\}$, $k = 1, \dots, K$

- 1: **repeat**
 - 2: **E step.** Minimize J with respect to r_{nk} while keeping the w_k fixed
 - 3: **M step.** Minimize J with respect to w_k while keeping r_{nk} fixed
 - 4: **until** convergence criterion is fulfilled
-

- Wechsel zwischen sogenannten E-step und M-step (vereinfachte Version des viel allgemeineren E(xpectation)-M(aximization Algorithmus)
- E-step: Zuordnung von Punkten zu Clustern für feste Prototypen

$$r_{nk} = \begin{cases} 1 & \text{wenn } \|x_n - w_k\|^2 < \|x_n - w_l\|^2 \ \forall l \neq k, \ k = 1, \dots, K \\ 0 & \text{sonst} \end{cases}$$

- M-step: Update der Clusterzentren w_k für feste Zuordnung

$$w_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}}$$

(damit liegt der neue Prototyp im Zentrum der zugeordneten Punkte)

- k-means konvergiert!

LGB

- wie findet (oder automatisch verbessert) man die richtige Anzahl an Clustern
- Vorgehen: Iterativ erhöhe Anzahl von Clustern bis
 - Quatisierungsfehler unter einem Schwellwert ist
 - oder: Maximum an Clustern erreicht ist

Algorithm 4 LGB**Require:** Maximal number of prototypes L_{max} **Require:** Threshold e_{min} for minimal quantization error E_{VQ} **Require:** Initial codebook comprising a single prototype $w_1 = \frac{1}{N} \sum_{n=1}^N x_n$

- 1: **repeat**
- 2: **for all** prototypes $w_l, l = 1, \dots, L$ (or for the M prototypes with largest quantization errors) **do**
- 3: Create two new prototypes

$$c_1 = w_l + \epsilon$$

$$c_2 = w_l - \epsilon$$

where $\epsilon \in \mathbb{R}^D$ is an random vector of small length

- 4: $L = 2L$ (or $L = L + M$)
- 5: Apply an arbitrary clustering algorithm with L prototypes starting from the current codebook until convergence, e.g. k-means with $k = L$
- 6: **until** quantization error $E_{VQ} < e_{min}$ or size of codebook $|\{w_k\}| = L_{max}$

Gaussian Mixture Modell (GMM)

- k-means gruppiert "hart" jeden Punkt zu einer Gruppe, GMM kann "soft" clustering

Soft Clustering

- Wahrscheinlichkeit für Cluster
- probabilistischer Ansatz

Modellselektion: Wie viele Gauß Funktionen?**Parameteroptimierung:** Mittelwerte, Varianz(matrizen), a-Priori Wahrscheinlichkeiten**Datenmodell:**

$$P(x) = \sum_{k=1}^K \underbrace{\pi_k}_{\text{prior}} \underbrace{N(x|w_k, \Sigma_k)}_{\text{Likelihood}}, \quad \sum_k \pi_k = 1$$

$$\text{Posterior: } P_k(x_n) = \frac{\pi_k N(x_n|w_k, \Sigma_k)}{\sum_{i=1}^K \pi_i N(x_n|w_i, \Sigma_i)}$$

Posterior: Wahrscheinlichkeit, dass Punkt x_n zu Klasse k gehört, entspricht einer weichen Zuordnung.

Problem: Likelihood hängt nun von allen Parametern ab.

- Parameter können nicht direkt berechnet werden

Lösung: E(xpectation)-M(aximization)

- Init: Wähle Startparameter
- Wiederhole:
 - E: Berechne Posterior aller Punkte
 - M: Berechne (optimiere) Parameter

EM für GMM

Gegeben ein GMM, Ziel ist es den Likelihood zu maximieren.

1. Initialisiere die Mittelwert μ_k , die Kovarianzen Σ_k und die Priors π_k und berechne den Anfangswert des log Likelihoods
2. **E step.** Berechne die Posteriors (responsibilities) mit den momentanen Parametern

$$\gamma_k(x_n) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{i=1}^K \pi_i N(x_n | \mu_i, \Sigma_i)}$$

3. **M step.** Optimierte die Parameter

$$\begin{aligned} N_k &= \sum_{n=1}^N \gamma_k(x_n) \\ \mu_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma_k(x_n) x_n \\ \Sigma_k &= \frac{1}{N_k} \sum_{n=1}^K \gamma_k(x_n) (x_n - \mu_k)(x_n - \mu_k)^T \\ \pi_k &= \frac{N_k}{N} \end{aligned}$$

4. Berechne den log Likelihood

$$\ln p(X | \mu, \Sigma) = \sum_{n=1}^N \ln \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k)$$

Konvergieren die Parameter oder der log Likelihood? Wenn nicht gehe zu 2..

Anmerkungen

- Likelihood
 - Erhöht sich mit jeder Iteration
 - Algorithmus konvergiert
- Initialisierung
 - Kann zufällig erfolgen
 - k-means oder LGB zur Initialisierung der Mittelwerte
 - Methoden zur automatischen Wahr der Anzahl existiert
- Update Schritte
 - Geschlossene Form für Gauß-Verteilungen
 - Für generelle Mixtures muss das nicht sein