

MUSTERERKENNUNG

Vorlesung im Sommersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 30. Mai 2017

Teil VI

Numerische Klassifikation

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

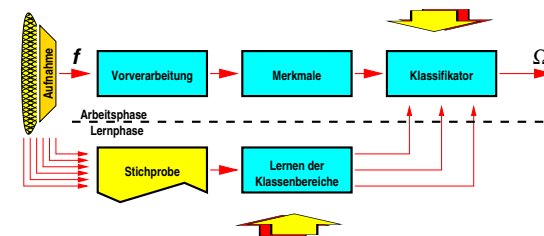
Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Numerische Klassifikationsaufgabe

Raten der wahren Musterklasse Ω_κ von $\mathbf{f} \in \Omega$



Musterrepräsentation

ein Vektor mit numerischen Merkmalen

$$\mathbf{x} = (x_1, \dots, x_D)^\top \in \mathbb{R}^D$$

Zielsetzung

Die **geratene** Klasse $\hat{\kappa}$ von \mathbf{x} bzw. \mathbf{f} ist (möglichst oft) gleich der **wahren** Klasse κ .

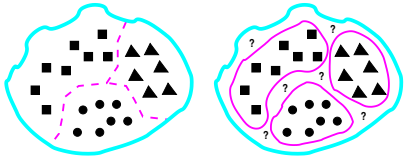
Klassifikationsergebnis

die vermutete Klassenzugehörigkeit

$$\hat{\kappa}(\mathbf{x}) \in \{1, 2, \dots, K\}$$

Entscheidungsgrundlage eines Klassifikators

Modelle für eine harte oder eine weiche Klasseneinteilung

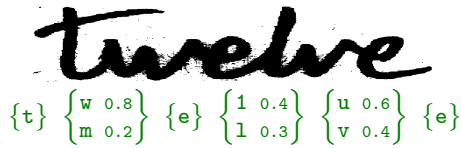


Entscheidungsfunktion

$$\delta : \begin{cases} \mathbb{R}^D & \rightarrow \{1, \dots, K\} \\ \mathbf{x} & \mapsto \hat{\kappa} = \delta(\mathbf{x}) \end{cases}$$

Probabilistisches Modell

Der Klassifikator kennt
Zugehörigkeitsgrade $\hat{P}(\kappa|\mathbf{x})$
der Muster zu den Klassen Ω_{κ} .



Geometrisches Modell

Der Klassifikator kennt nur eine
Klassenpartition $\{\hat{\Omega}_1, \dots, \hat{\Omega}_K\}$
des Merkmalraums oder ihre
Grenzen.



Maschinelles Lernen eines Klassifikators

Lehrer $\hat{=}$ endliche Menge von Beispielmustern

Induktion

Der Klassifikator ist aus einer **endlichen Probe** $\omega \subset \Omega$ des Merkmalraums zu „lernen“.



Lernmethode

- alle Muster aus ω etikettiert
- kein Muster aus ω etikettiert
- einige Muster aus ω etikettiert

überwacht

unüberwacht

teilüberwacht

Das Etikettieren von Mustern ist i.a. *viel teurer* als das Beschaffen selbst!

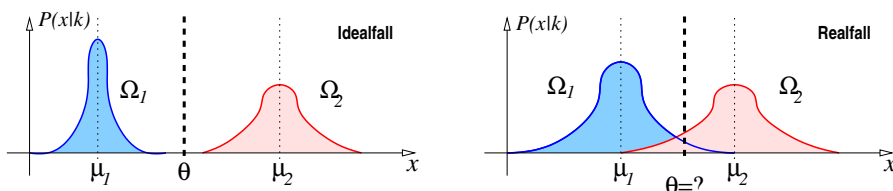
Klassifikatorleistung

Erwartete und beobachtete Fehlerhäufigkeit

Klassifikatortest

Die Fehlerrate einer Entscheidungsregel δ wird mit Hilfe einer (etikettierten) **Teststichprobe** ω' geschätzt:

$$\hat{p}_{\varepsilon}(\{\omega'_{\kappa}\}|\delta) \stackrel{\text{def}}{=} \frac{\text{falsch klassifiziert}}{\text{alle Muster}} = \frac{\sum_{\kappa=1}^K \#\{\mathbf{x} \in \omega'_{\kappa} \mid \delta(\mathbf{x}) \neq \kappa\}}{\sum_{\kappa=1}^K |\omega'_{\kappa}|}$$



Bemerkung

Überlagern sich die Klassengebiete im gewählten Merkmalraum, so können nicht *alle* Muster im Überlappungsbereich korrekt klassifiziert werden!

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Statistische Entscheidungstheorie

Viel Wind um einen Papiertiger?

Voraussetzungen

- Die Wahrscheinlichkeitsverteilung des **mustererzeugenden Prozesses** ist bekannt.
- Ein **Gütekriterium** (Kostenfunktion) für mögliche Klassifikationsentscheidungen liegt vor.

Resultat

- Die **Bayesregel** ist der optimale (kostenminimale) Klassifikator.
- Die Bayesregel ist sogar optimal unter allen unscharfen Entscheidungsregeln.

Problem

Die Bayesregel ist leider nicht effektiv berechenbar ...

Qualitätskriterium

Kostenmodell für die (Fehl-)klassifikation eines Musters

Kostenmatrix

$$\begin{pmatrix} r_{1,1} & r_{1,2} & \dots & r_{1,K} \\ \vdots & \dots & \ddots & \vdots \\ r_{K,1} & r_{K,2} & \dots & r_{K,K} \end{pmatrix} \quad \text{bzw.} \quad \begin{pmatrix} r_{1,0} & r_{1,1} & r_{1,2} & \dots & r_{1,K} \\ \vdots & \dots & \dots & \ddots & \vdots \\ r_{K,0} & r_{K,1} & r_{K,2} & \dots & r_{K,K} \end{pmatrix}$$

Eine Entscheidung für die Klasse Ω_λ bei Vorlage eines Musters aus Ω_κ verursacht die **Kosten** $r_{\kappa\lambda}$.

- Der Eintrag $r_{\kappa,0}$ bezieht die **Rückweisungskosten**.
- Plausible Forderung ($\lambda \neq \kappa$): $0 \leq r_{\kappa\kappa} < r_{\kappa 0} < r_{\kappa\lambda}$
- 0/1-Kostenmatrix zur Zählung der Klassifikationsfehler:

$$r_{\kappa\lambda} = \begin{cases} 0 & \lambda = \kappa \\ 1 & \lambda \neq \kappa \end{cases}$$

Generatives Modell

Wahrscheinlichkeitsverteilung für die Mustererzeugung („1 + K Urnen“)

Gemeinsame Verteilung

Ein Muster ist durch seine Klasse $\kappa \in \{1..K\}$ und den Vektor $\mathbf{x} \in \mathbb{R}^D$ seiner Merkmale charakterisiert:

$$P(\kappa) \cdot P(\mathbf{x}|\kappa) = P(\mathbf{x}, \kappa) = P(\mathbf{x}) \cdot P(\kappa|\mathbf{x})$$

Randverteilungen

A priori
Klassenwahrscheinlichkeiten

$$P(\kappa) = \int P(\mathbf{x}, \kappa) d\mathbf{x}$$

Marginale Merkmalvektordichte

$$P(\mathbf{x}) = \sum_{\kappa=1}^K P(\mathbf{x}, \kappa)$$

Bedingte Verteilungen

Klassenbedingte Merkmalvektordichte

$$P(\mathbf{x}|\kappa) = P(\mathbf{x}, \kappa) / P(\kappa)$$

A posteriori
Klassenwahrscheinlichkeiten

$$P(\kappa|\mathbf{x}) = P(\mathbf{x}, \kappa) / P(\mathbf{x})$$

Entscheidungsregel

Scharfe und randomisierte Entscheidungsregeln

Definition

Für einen Problemkreis mit K Musterklassen und dem Merkmalraum \mathbb{R}^D heißt eine Abbildung

$$\delta : \begin{cases} \mathbb{R}^D & \rightarrow \mathbb{R}^{K+1} \\ \mathbf{x} & \mapsto \delta(\mathbf{x}) = (\delta_0(\mathbf{x}), \delta_1(\mathbf{x}), \dots, \delta_K(\mathbf{x}))^\top \end{cases}$$

Entscheidungsregel, falls die Normierungsbedingung

$$\sum_{\lambda=0}^K \delta_\lambda(\mathbf{x}) = 1, \quad (\kappa = 1, \dots, K)$$

gilt. Die Regel heißt **scharf** (hart), wenn alle $\delta(\mathbf{x})$ Einheitsvektoren sind; andernfalls heißt die Regel **unscharf** (weich, randomisiert).

Eine unscharfe Regel liefert zum Beispiel Aussagen der Form $\delta(\mathbf{x}) = (1/10, 2/10, 7/10, 0)^\top$ oder $\delta(\mathbf{x}) = (1/2, 1/8, 1/8, 1/4)^\top$

Risikoformel

Zu erwartende Klassifikationskosten einer Entscheidungsregel

Lemma

Die zu erwartenden Kosten (das **Risiko**) der Klassifikation von Mustern durch die Entscheidungsregel

$$\delta(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^{K+1}$$

auf Grundlage der Kostenmatrix $[r_{\kappa\lambda}]$ betragen

$$\mathfrak{R}(\delta) = \sum_{\kappa=1}^K p_{\kappa} \sum_{\lambda=0}^K r_{\kappa\lambda} \int_{\mathbb{R}^D} P(\mathbf{x}|\Omega_{\kappa}) \cdot \delta_{\lambda}(\mathbf{x}) d\mathbf{x}$$

Im Falle der 0/1-Kostenmatrix lautet die Risikoformel

$$\mathfrak{R}(\delta) = \sum_{\kappa=1}^K p_{\kappa} \int_{\mathbb{R}^D} P(\mathbf{x}|\Omega_{\kappa}) \cdot \sum_{\lambda \neq \kappa} \delta_{\lambda}(\mathbf{x}) d\mathbf{x} = \sum_{\kappa=1}^K p_{\kappa} \int_{\mathbb{R}^D} P(\mathbf{x}|\Omega_{\kappa}) \cdot (1 - \delta_{\kappa}(\mathbf{x})) d\mathbf{x}$$

und quantifiziert die **Fehlerwahrscheinlichkeit** der Entscheidungsregel $\delta(\cdot)$.

Beweis.

Das Risiko $\mathfrak{R}(\delta)$ setzt sich zusammen als gewichtete Summe der klassenweisen Einzelrisiken:

$$\mathfrak{R}(\delta) = \sum_{\kappa=1}^K p_{\kappa} \cdot \mathfrak{R}_{\kappa}(\delta)$$

$\mathfrak{R}_{\kappa}(\delta)$ wiederum setzt sich aus den einzelnen Wagnisarten zusammen, die sich aus der jeweiligen Klassifikationsentscheidung ergeben:

$$\mathfrak{R}_{\kappa}(\delta) = \sum_{\lambda=0}^K r_{\kappa\lambda} \cdot P(\text{klassifiziert nach } \Omega_{\lambda} \mid \text{stammt aus } \Omega_{\kappa})$$

Den bedingten Wahrscheinlichkeitsausdruck gewinnen wir durch Marginalisierung:

$$P(\Omega_{\lambda}|\Omega_{\kappa}) = \int_{\mathbb{R}^D} P(\Omega_{\lambda}, \mathbf{x} \mid \Omega_{\kappa}) d\mathbf{x}$$

Mehrmalige Anwendung der Definition bedingter Wahrscheinlichkeiten nebst Erweiterung ergibt:

$$P(\Omega_{\lambda}, \mathbf{x} \mid \Omega_{\kappa}) = \frac{P(\Omega_{\lambda}, \mathbf{x}, \Omega_{\kappa})}{P(\Omega_{\kappa})} \cdot \frac{P(\mathbf{x}, \Omega_{\kappa})}{P(\mathbf{x}, \Omega_{\kappa})} = P(\Omega_{\lambda} \mid \mathbf{x}, \Omega_{\kappa}) \cdot P(\mathbf{x} \mid \Omega_{\kappa}) = \delta_{\lambda}(\mathbf{x}) \cdot P(\mathbf{x} \mid \Omega_{\kappa})$$

□

Optimale Entscheidungsregel

Allgemeine Kostenmatrix · Rückweisungsmöglichkeit

Satz

Die optimale (randomisierte) Entscheidungsregel δ^* , welche das Risiko $\mathfrak{R}(\delta)$ hinsichtlich einer gegebenen Kostenmatrix $[r_{\kappa\lambda}]$ minimiert, klassifiziert die Merkmalvektoren gemäß

$$\delta_{\kappa}^*(\mathbf{x}) = \begin{cases} 1 & u_{\kappa}(\mathbf{x}) = \min_{\lambda} u_{\lambda}(\mathbf{x}) \\ 0 & \text{sonst} \end{cases}, \quad \mathbf{x} \in \mathbb{R}^D$$

mit der **Prüfgröße**

$$u_{\lambda}(\mathbf{x}) \stackrel{\text{def}}{=} \sum_{\kappa=1}^K r_{\kappa\lambda} \cdot p_{\kappa} \cdot P(\mathbf{x}|\Omega_{\kappa}), \quad \lambda = 0, 1, \dots, K.$$

Bemerkung

Diese optimale Entscheidungsregel ist *scharf*!

Beweis.

Mit obiger Prüfgrößendefinition gilt die Darstellung

$$\mathfrak{R}(\delta) = \int_{\mathbb{R}^D} \underbrace{\sum_{\lambda=0}^K \underbrace{\left[\sum_{\kappa=1}^K r_{\kappa\lambda} p_{\kappa} P(\mathbf{x}|\Omega_{\kappa}) \right]}_{u_{\lambda}(\mathbf{x})}}_{I(\mathbf{x})} \delta_{\lambda}(\mathbf{x}) d\mathbf{x}$$

Der Integrand $I(\mathbf{x})$ gehorcht der Abschätzung

$$\begin{aligned} I(\mathbf{x}) &= \sum_{\lambda=0}^K u_{\lambda}(\mathbf{x}) \cdot \delta_{\lambda}(\mathbf{x}) \\ &\geq \sum_{\lambda=0}^K u_{\min}(\mathbf{x}) \cdot \delta_{\lambda}(\mathbf{x}) \\ &= u_{\min}(\mathbf{x}) \sum_{\lambda=0}^K \delta_{\lambda}(\mathbf{x}) = u_{\min}(\mathbf{x}) \end{aligned}$$

Dieser Minimalwert $u_{\min}(\mathbf{x}) = \min_{\lambda} u_{\lambda}(\mathbf{x})$ läßt sich durch die Wahl

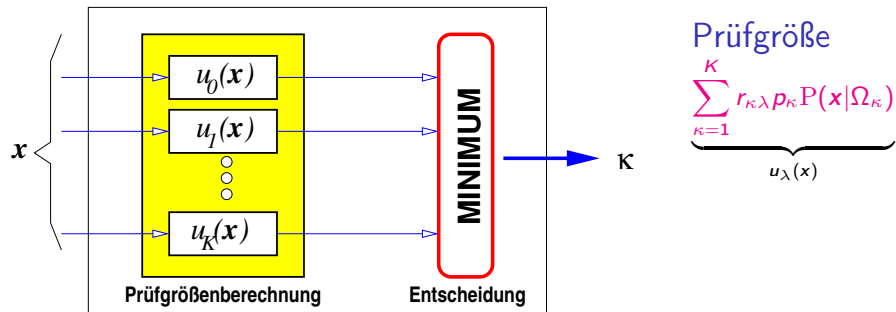
$$\delta^*(\mathbf{x}) = \kappa^* \text{-ter Einheitsvektor}$$

mit $u_{\min}(\mathbf{x}) = u_{\kappa^*}(\mathbf{x})$ systematisch erzwingen.

□

Der optimale Klassifikator

Allgemeine Kosten · Vereinfachte Kosten



Vereinfachte Kostenmatrix

$$r_{\kappa\lambda} = \begin{cases} r_w & \kappa = \lambda \\ r_z & \lambda = 0 \\ r_f & \text{sonst} \end{cases} \Rightarrow \begin{aligned} u_0(x) &= r_z \cdot \sum_{\kappa=1}^K p_{\kappa} \cdot P(x|\Omega_{\kappa}) = r_z \cdot P(x) \\ u_{\lambda}(x) &= r_f \cdot P(x) + (r_w - r_f) \cdot p_{\lambda} \cdot P(x|\Omega_{\lambda}) \end{aligned}$$

Beweis.

Die vereinfachte Kostenmatrix mit den drei Parametern r_w , r_f , r_z wird bei Klassifikatoren mit erzwungener Entscheidung o.B.d.A. zu

$$r_{\kappa\lambda} = \begin{cases} 0 & \kappa = \lambda \\ 1 & \kappa \neq \lambda \end{cases}$$

Das Risiko $\mathfrak{R}(\delta)$ entspricht dann der *Fehlerwahrscheinlichkeit* des Klassifikators und seine Prüfgrößen lauten

$$u_{\lambda}(x) = \sum_{\kappa \neq \lambda} p_{\kappa} \cdot P(x|\Omega_{\kappa}) = \sum_{\kappa \neq \lambda} P(x, \Omega_{\kappa}) = P(x) - P(x, \Omega_{\lambda})$$

Statt diese Prüfgröße zu *minimieren*, können wir auch die gemeinsame Wahrscheinlichkeit $P(x, \Omega_{\lambda})$ von Muster und Klasse *maximieren*, denn der Randdichtewert $P(x)$ ist konstant bezüglich der Wahl von Ω_{λ} . \square

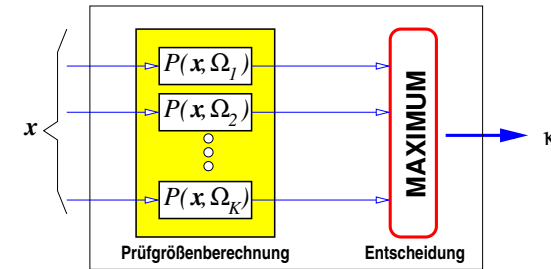
Bemerkung

Herzlichen Dank für dieses interessante *theoretische* Resultat. Aber wie sieht es denn nun eigentlich aus mit der *praktischen* Berechenbarkeit von $\hat{\kappa}(x)$ nach der Bayesregel?

1. Die klassenbedingten Musterwahrscheinlichkeiten $P(x|\Omega_{\kappa})$ sind nicht bekannt und nur sehr unzureichend und/oder aufwändig zu modellieren.
2. Die a priori Klassenwahrscheinlichkeiten $P(\Omega_{\kappa})$ sind ebenfalls unbekannt, jedoch sind i.a. leicht robuste Schätzwerte zu beschaffen.
3. Der Randverteilungsdichtewert $P(x)$ ist ebenfalls unbekannt, aber bei Vorliegen der obengenannten Größen unproblematisch durch Aufsummieren der $P(x, \Omega_{\lambda})$ zu berechnen.
4. Zur Klassifikation selbst wird der Wert $P(x)$ allerdings garnicht benötigt, denn er ist ja unabhängig vom Klassenindex κ .

Bayesregel

Klassifikator mit der minimalen Fehlerrate



Satz

Die Bayesregel (*Maximum a posteriori-Klassifikator*)

$$\delta(x) = \underset{\lambda}{\operatorname{argmax}} P(\Omega_{\lambda}|x) = \underset{\lambda}{\operatorname{argmax}} \frac{p_{\lambda} \cdot P(x|\Omega_{\lambda})}{P(x)}$$

ist derjenige Klassifikator mit der kleinstmöglichen Fehlerrate p_{ϵ}^{BA} (**Bayesfehlerrate**).

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Bayesregel

Garantiert minimale Fehlerrate — praktisch nicht verwertbar

Klassifiziere ein neues Muster $\mathbf{x} \in \mathbb{R}^D$ gemäß

$$\hat{\kappa}(\mathbf{x}) = \operatorname{argmax}_{\lambda=1..K} u_{\lambda}(\mathbf{x})$$

mit der Prüfgröße

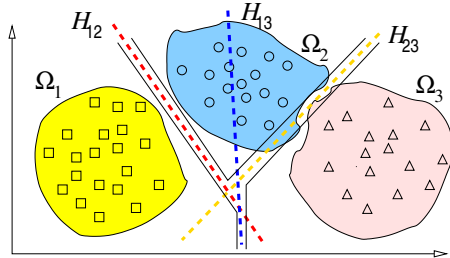
$$u_{\lambda}(\mathbf{x}) = P(\Omega_{\lambda}|\mathbf{x}) = \frac{P(\Omega_{\lambda}) \cdot P(\mathbf{x}|\Omega_{\lambda})}{P(\mathbf{x})}$$

Typen realer Klassifikatoren

1. Schätze die K Verteilungen $P(\cdot|\Omega_{\lambda})$ statistisch
2. Schätze die K Trennfunktionen $u_{\lambda}(\cdot)$ diskriminativ
3. Rate „auf Zuruf“ die Zugehörigkeit $\hat{\kappa}(\cdot)$ partitionierend

Diskriminativer Klassifikator

Schätze Klassenzugehörigkeiten oder -trennfunktionen für jedes Muster



Modellinformation

- beste Klasse je Muster
- Zugehörigkeitsmaß je Muster
- explizite Klassengrenzen
- Verteilungsmodell je Klasse

Vorgehensweise

Approximiere geeignete Trennfunktionen

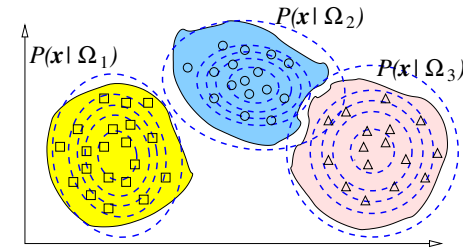
$$h_{\kappa\lambda}(\mathbf{x}) = r \text{ mit } \begin{cases} r \geq 0 & \mathbf{x} \in \omega_{\kappa} \\ r < 0 & \mathbf{x} \in \omega_{\lambda} \end{cases}$$

zur Trennung der Klassengebiete Ω_{κ} , Ω_{λ} und setze

$$\delta(\mathbf{x}) = \kappa^*, \quad \text{falls } h_{\kappa^*\lambda}(\mathbf{x}) \geq 0 \text{ für alle } \lambda \neq \kappa^*$$

Statistischer Klassifikator

Schätze Wahrscheinlichkeitsverteilung $f_{\kappa}(\mathbf{x})$ für jede Musterklasse



Modellinformation

- beste Klasse je Muster
- Zugehörigkeitsmaß je Muster
- Verteilungsmodell je Klasse
- explizite Klassengrenzen

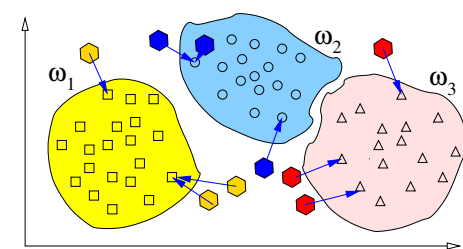
Vorgehensweise

Schätze die wahren klassenbedingten Verteilungsdichten $P(\mathbf{x}|\Omega_{\kappa})$ mit Hilfe eines (parametrischen?) Verteilungsmodells $f_{\kappa}(\mathbf{x}) = f(\mathbf{x}|\theta_{\kappa})$ und entscheide nach der näherungsweisen Bayesregel

$$\delta(\mathbf{x}) = \operatorname{argmax}_{\kappa=1..K} \hat{P}(\Omega_{\kappa}|\mathbf{x}) = \operatorname{argmax}_{\kappa=1..K} \frac{p_{\kappa} \cdot f_{\kappa}(\mathbf{x})}{\sum_{\lambda=1}^K p_{\lambda} \cdot f_{\lambda}(\mathbf{x})}$$

Partitionierender Klassifikator

Rate bestpassenden Klassenindex $\hat{\kappa}(\mathbf{x})$ für jedes Muster



Modellinformation

- beste Klasse je Muster
- Zugehörigkeitsmaß je Muster
- explizite Klassengrenzen
- Verteilungsmodell je Klasse

Vorgehensweise

Ermittle zum Eingabemuster \mathbf{x} den *nächsten Nachbarn*

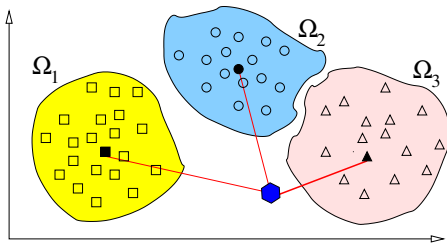
$$\mathbf{x}_{\text{NN}} = \operatorname{argmin}_{\mathbf{z} \in \bigcup \omega_{\kappa}} \|\mathbf{z} - \mathbf{x}\|$$

in der Lernstichprobe und setze

$$\delta(\mathbf{x}) = \kappa \quad \text{für } \mathbf{x}_{\text{NN}} \in \omega_{\kappa}$$

Abstandsmessender Klassifikator

Entscheide auf Grundlage der Ähnlichkeiten zwischen Muster und allen Klassenprototypen



Modellinformation

- beste Klasse je Muster
- Zugehörigkeitsmaß je Muster
- Klassengrenzen (Voronoizellen)
- Verteilungsmodell je Klasse

Vorgehensweise

Wähle geeignete Klassenprototypen μ_κ und setze

$$\delta(\mathbf{x}) = \underset{\kappa}{\operatorname{argmin}} \|\mathbf{x} - \mu_\kappa\|_q$$

Definition der Prototypen?

Definition des Abstandsmaßes?

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Informationshierarchie

Stärkere und schwächere Klassifikatormodelle

Verteilungsfunktionen

Wie sind die Muster einer Klasse über den Raum hinweg gestreut?

???

nicht aus den $u_\kappa(\mathbf{x})$ reproduzierbar, da $P(\mathbf{x})$ fehlt

Trennfunktionen

Wie verlaufen die Grenzen zwischen je zwei Klassen?

$$h_{\kappa\lambda}(\mathbf{x}) = \log \frac{u_\kappa(\mathbf{x})}{u_\lambda(\mathbf{x})}$$

Zugehörigkeiten

Wie verteilt sich die Paßfähigkeit eines Musters auf die Klassen?

$$u_\kappa(\mathbf{x}) = P(\Omega_\kappa|\mathbf{x}) \propto P(\Omega_\kappa) \cdot P(\mathbf{x}|\Omega_\kappa)$$

Sicherheit u/o Alternativen nicht aus $\hat{\kappa}(\mathbf{x})$ erschließbar

Partitionen

Welche Muster gehören zu welchen Klassen?

$$\hat{\kappa}(\mathbf{x}) = \underset{\lambda=1..K}{\operatorname{argmax}} u_\lambda(\mathbf{x})$$

nicht aus $h_{\kappa\lambda}(\mathbf{x})$ wg. „Bermuda“-Syndrom

Naiver Bayesklassifikator

Idealisierende Annahme klassenweiser Merkmalunabhängigkeit

Lemma

Sind die Merkmale x_1, \dots, x_D der Muster jeder Klasse Ω_κ **statistisch unabhängig**, so lautet die („naive“) Bayes-Entscheidungsregel:

$$\hat{\kappa}(\mathbf{x}) = \underset{\lambda=1..K}{\operatorname{argmax}} \left(P(\Omega_\lambda) \cdot \prod_{d=1}^D P(x_d|\Omega_\lambda) \right)$$

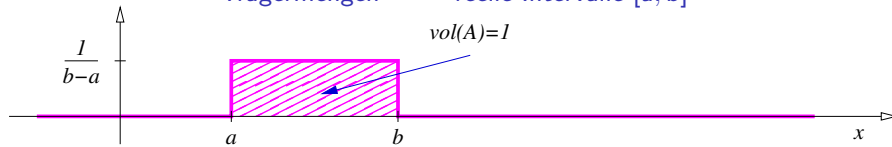
Bemerkung

Es sind $K \cdot D$ univariate Dichtefunktionen $f_{\lambda,d}(\cdot)$ aus den Lerndaten zu schätzen.

- Gleichverteilung (uniforme Dichte)
- Normalverteilung (Gaußdichte)
- Laplace- oder Cauchydicke (super-/subnormal)
- Parzen- oder Histogrammdichte (parameterfrei)

Univariate Gleichverteilungsdichten

Trägersmengen $\hat{=}$ reelle Intervalle $[a, b]$



Definition

Ist $\mathcal{A} \subset \mathbb{R}^D$ eine D -dimensionale Punktmenge mit endlichem Hypervolumen $vol(\mathcal{A})$, so heißt

$$f_{\mathcal{A}}^{\text{unif}} : \mathbf{x} \mapsto \begin{cases} 1/vol(\mathcal{A}) & \mathbf{x} \in \mathcal{A} \\ 0 & \mathbf{x} \notin \mathcal{A} \end{cases}$$

die **Gleichverteilungs-** oder **uniforme Dichte** zum Träger(ereignis) \mathcal{A} .

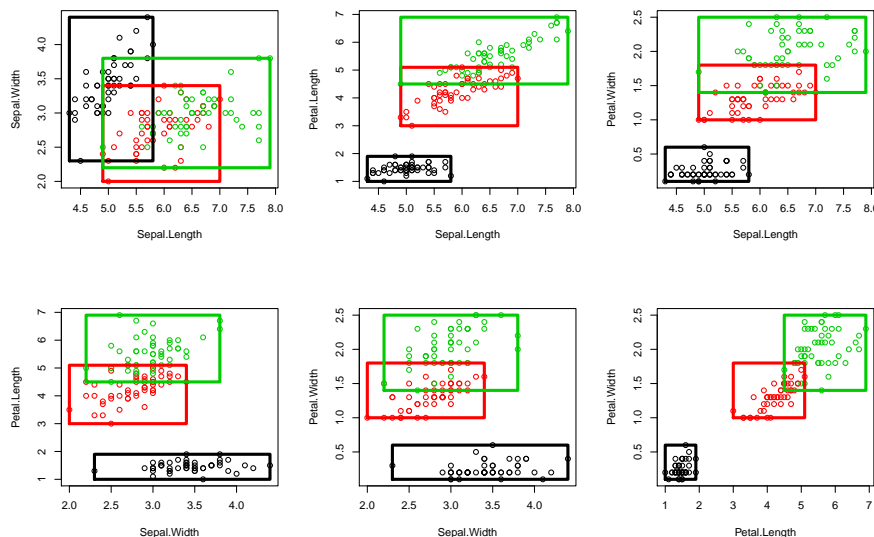
Der univariate Spezialfall

$$f_{[a,b]}^{\text{unif}} : x \mapsto \begin{cases} 1/(b-a) & x \in [a, b] \\ 0 & x \notin [a, b] \end{cases}, \quad a, b \in \mathbb{R}, \quad a < b$$

heißt **uniforme Intervalldichte** über $[a, b]$.

Univariate Gleichverteilungsdichten

Beispiel: Iris-Datensatz, 2D-Träger für alle (x_i, x_j) -Kombinationen



Univariate Gleichverteilungsdichten

Maximum-Likelihood-Schätzung der Verteilungsparameter

Lemma

Für eine Lernprobe $\omega = \{z_1, \dots, z_T\}$ nimmt die (unlogarithmierte) Maximum-Likelihood-Zielfunktion

$$\ell_{\text{ML}}(a, b) = f_{[a,b]}^{\text{unif}}(\omega) = \prod_{t=1}^T f_{[a,b]}^{\text{unif}}(z_t)$$

ihren maximalen Wert bei den Intervallgrenzen

$$a^* \stackrel{\text{def}}{=} \min_{t=1..T} z_t, \quad b^* \stackrel{\text{def}}{=} \max_{t=1..T} z_t$$

an.

Bemerkung

Es ist $\ell_{\text{ML}}(a, b) = 1 / (b - a)^T$, sofern alle z_t in $[a, b]$ liegen und Null sonst.

UNB — Entscheidungsregel

... für den naiven Bayesklassifikator mit Gleichverteilung

Prüfgröße

$$u_{\kappa}(\mathbf{x}) = \hat{P}(\Omega_{\kappa}) \cdot \hat{f}_{\kappa}^{\text{unif}}(\mathbf{x}) = \frac{T_{\kappa}}{T} \cdot \frac{1}{\prod_d (b_{\kappa,d} - a_{\kappa,d})} \cdot \underbrace{\mathbb{I}_{a_{\kappa} \leq \mathbf{x} \leq b_{\kappa}}}_{\mathbb{I}_{\mathbf{x} \in \mathcal{H}_{\kappa}}}$$

Entscheidungsregel

$$\hat{\kappa}(\mathbf{x}) = \operatorname{argmax}_{\lambda=1..K} u_{\kappa}(\mathbf{x}) = \operatorname{argmax}_{\lambda | \mathbf{x} \in \mathcal{H}_{\lambda}} \frac{T_{\kappa}}{vol(\mathcal{H}_{\lambda})}$$

Bemerkung

- Die Muster außerhalb von $\bigcup_{\lambda} \mathcal{H}_{\lambda}$ können keiner Klasse zugeordnet werden.
- Klassen mit einem Ausreißer werden *extrem* benachteiligt ($vol(\mathcal{H}_{\lambda}) \rightarrow \infty$).
- Die UNB-Regel verallgemeinert miserabel gegenüber „neuen“ Mustern $\mathbf{x} \notin \omega$.

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Approximationskriterium

In welchem Sinne ist $P(\Omega_\kappa|\mathbf{x})$ durch $h_\kappa(\mathbf{x})$ anzunähern?

Minimaler Klassifikationsfehler

$$\#_{\kappa, \mathbf{x}} \left(h_\kappa(\mathbf{x}) \neq \max_{\lambda} h_{\lambda}(\mathbf{x}) \right) \xrightarrow{!} \text{MIN}$$

Maximale Rückschlußwahrscheinlichkeit

logit-Modell oder probit-Modell; Optimierung mit linearen N.B.

$$\prod_{\kappa, \mathbf{x}} h_\kappa(\mathbf{x}) \xrightarrow{!} \text{MAX} \quad \text{wobei} \quad \sum_{\lambda=1}^K h_{\lambda}(\mathbf{x}) = 1 \quad (\forall \mathbf{x} \in \Omega)$$

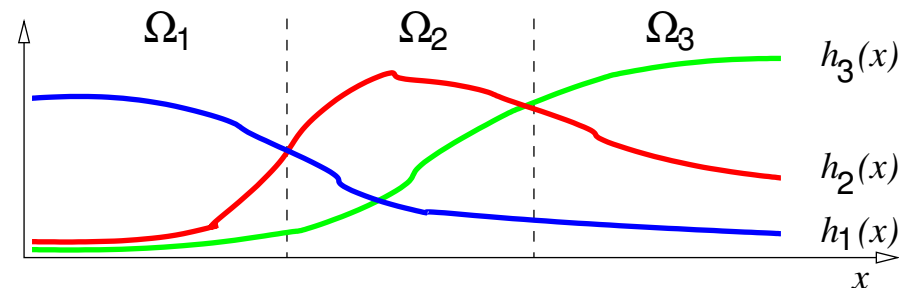
Minimaler Rekonstruktionsfehler

Geometr. Abstand zur **idealen** Trennfunktion; Ausgleichsrechnung

$$\sum_{\kappa, \mathbf{x}} \|h_{\lambda}(\mathbf{x}) - \mathbf{1}_{\lambda=\kappa}\|^2 \xrightarrow{!} \text{MIN} \quad (\forall \lambda = 1..K)$$

Verteilungsfreie Klassifikatoren

Direkte Approximation der a posteriori-Klassenwahrscheinlichkeit



Trenn- oder Diskriminantenfunktionen

Für jede Klasse ist die wahre Rückschlußwahrscheinlichkeit durch eine geeignete Modellfunktion anzunähern:

$$h_\kappa(\mathbf{x}) \approx P(\Omega_\kappa|\mathbf{x}) = \frac{P(\Omega_\kappa) \cdot P(\mathbf{x}|\Omega_\kappa)}{P(\mathbf{x})}$$

Quadratmittelklassifikator

Die ideale Trennfunktion sagt 1/0 zur richtigen/falschen Klasse

Definition

Die erwartete quadratische Abweichung

$$\mathcal{E}_{\mathbb{X}, \mathbb{K}}[\|\mathbf{h}(\mathbb{X}) - \mathbf{e}^{(\mathbb{K})}\|^2] = \sum_{\lambda=1}^K \sum_{\kappa=1}^K \int P(\mathbf{x}, \Omega_\kappa) \cdot (h_{\lambda}(\mathbf{x}) - \mathbf{1}_{\kappa=\lambda})^2 d\mathbf{x}$$

heißt **Rekonstruktionsfehler** der idealen Trennfunktion einer Klassifikationsaufgabe.

Bemerkungen

1. Die Approximation der idealen Trennfunktion ($\mathbf{h}(\mathbf{x})$ liefert 1 für die „richtige“ Klasse) fordert viel mehr als unbedingt nötig.
2. Trotzdem gilt, bei uneingeschränkter Optimierung: $h_{\lambda}(\mathbf{x}) = P(\Omega_{\lambda}|\mathbf{x})$ besitzt den minimalen Rekonstruktionsfehler.
3. In praxi wird der Rekonstruktionsfehler auf einer klassifizierten Lernstichprobe ermittelt.

Linearer Quadratmittelklassifikator

Diskriminanten sind linear in den Merkmalvektorkomponenten

$$\mathbf{x} \in \mathbb{R}^D \Rightarrow \mathbf{h}(\mathbf{x}) = \begin{pmatrix} h_1 \\ h_2 \\ \dots \\ h_K \end{pmatrix} = \begin{pmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \mathbf{a}_2^\top \mathbf{x} \\ \dots \\ \mathbf{a}_K^\top \mathbf{x} \end{pmatrix} \Rightarrow \kappa^*(\mathbf{x}) = \underset{\lambda}{\operatorname{argmax}} h_\lambda(\mathbf{x})$$

Lemma

Der Klassifikator mit den linearen Prüfgrößen

$$h_\lambda(\mathbf{x}) = \sum_{d=1}^D a_{\lambda,d} \cdot x_d, \quad \lambda = 1..K \quad (\text{kürzer: } \mathbf{h}(\mathbf{x}) = \mathbf{A}^\top \mathbf{x})$$

besitzt auf der etikettierten Stichprobe $\{(\mathbf{x}_t, \kappa_t) \mid t = 1, 2, \dots, T\}$ den (empirischen) Rekonstruktionsfehler

$$\varepsilon(\mathbf{A}) = \sum_{\lambda=1}^K \varepsilon(\mathbf{a}_\lambda) = \sum_{\lambda=1}^K \sum_{t=1}^T (a_{\lambda}^\top \mathbf{x}_t - \mathbf{1}_{\kappa_t=\lambda})^2 = \sum_{\lambda=1}^K \|\mathbf{X} \mathbf{a}_\lambda - \mathbf{y}_\lambda\|^2$$

mit der Datenmatrix \mathbf{X} und den Klassenindikatoren \mathbf{y}_λ , $\lambda = 1..K$.

Beweis.

Der Gesamtfehler $\varepsilon(\mathbf{A})$ zerfällt in eine Summe ungekoppelter, klassenbezogener Fehlerkomponenten

$$\varepsilon(\mathbf{a}_\lambda) = \mathbf{a}_\lambda^\top \mathbf{X}^\top \mathbf{X} \mathbf{a}_\lambda - \mathbf{a}_\lambda^\top \mathbf{X}^\top \mathbf{y}_\lambda - \mathbf{y}_\lambda^\top \mathbf{X} \mathbf{a}_\lambda + \mathbf{y}_\lambda^\top \mathbf{y}_\lambda,$$

die separat optimiert werden können; dazu leiten wir partiell ab:

$$\nabla_{\mathbf{a}_\lambda} \varepsilon(\mathbf{a}_\lambda) = 2 \cdot \mathbf{X}^\top \mathbf{X} \mathbf{a}_\lambda - 2 \cdot \mathbf{X}^\top \mathbf{y}_\lambda = 2 \cdot (T \cdot \mathbf{R}) \mathbf{a}_\lambda - 2 \cdot (T_\lambda \cdot \hat{\boldsymbol{\mu}}_\lambda)$$

Das Resultat setzen wir (komponentenweise) gleich Null, kürzen durch $2 \cdot T$ und erhalten die Gaußschen Normalgleichungen. □

Im Satz und im Beweis wurden folgende Hilfsgrößen verwendet:

$$\mathbf{R} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top, \quad \boldsymbol{\mu}_\lambda = \frac{1}{T_\lambda} \sum_{t|\kappa_t=\lambda} \mathbf{x}_t, \quad T_\lambda = \sum_{t|\kappa_t=\lambda} 1, \quad \hat{\boldsymbol{\mu}}_\lambda = \frac{T_\lambda}{T}$$

Optimaler linearer Quadratmittelklassifikator

Gaußsches Normalgleichungssystem

Satz

Die Diskriminantenparameter \mathbf{a}_λ , $\lambda = 1..K$, für den linearen Quadratmittelklassifikator mit minimalem Rekonstruktionsfehler $\varepsilon(\mathbf{A})$ ergeben sich aus den **Gaußschen Normalgleichungen**

$$\mathbf{R} \cdot \mathbf{a}_\lambda = \mathbf{m}_\lambda, \quad \lambda = 1, \dots, K$$

mit den Stichprobenstatistiken

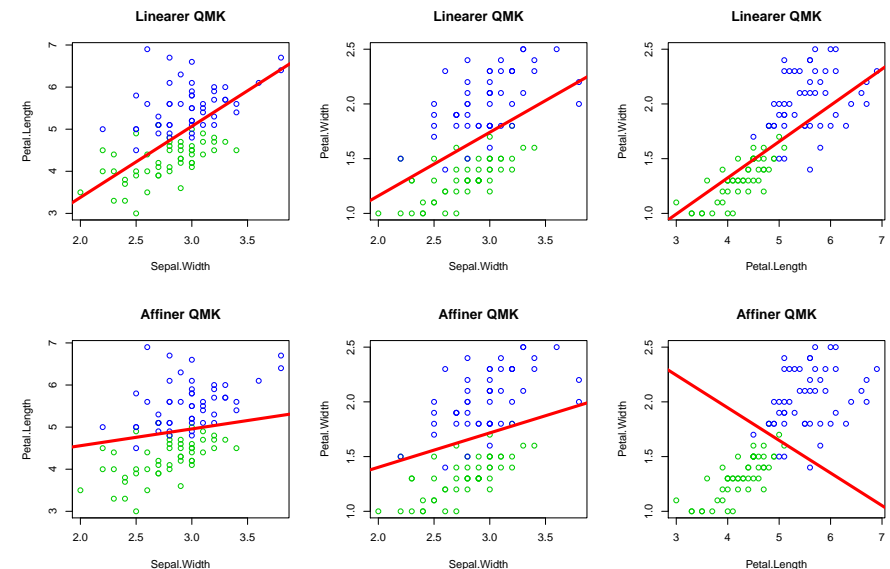
$$\mathbf{R} = \hat{\mathbf{S}} + \hat{\boldsymbol{\mu}} \hat{\boldsymbol{\mu}}^\top \quad \text{und} \quad \mathbf{m}_\lambda = \hat{\rho}_\lambda \cdot \hat{\boldsymbol{\mu}}_\lambda.$$

Bemerkung

Ist die $(T \times T)$ -Matrix \mathbf{R} der zweiten Stichprobenmomente invertierbar, so gilt $\mathbf{A} = \mathbf{R}^{-1} \cdot \mathbf{M}$ mit $\mathbf{M} = (\mathbf{m}_1, \dots, \mathbf{m}_K)$.

Linearer vs. affiner Quadratmittelklassifikator

Beispiel: Iris-Datensatz, 2D-Träger für einige (x_i, x_j) -Kombinationen



Termlinearer Quadratmittelklassifikator

Diskriminanten sind linear in den Termkoeffizienten

$$\mathbf{x} \in \mathbb{R}^D \Rightarrow \underbrace{\phi(\mathbf{x}) = \begin{pmatrix} \phi_1(\mathbf{x}) \\ \phi_2(\mathbf{x}) \\ \vdots \\ \phi_L(\mathbf{x}) \end{pmatrix}}_{\mathbf{z} \in \mathbb{R}^L} \Rightarrow \mathbf{h}(\mathbf{z}) = \mathbf{g}^\phi(\mathbf{x}) = \underbrace{\begin{pmatrix} \mathbf{a}_1^\top \mathbf{z} \\ \mathbf{a}_2^\top \mathbf{z} \\ \vdots \\ \mathbf{a}_K^\top \mathbf{z} \end{pmatrix}}_{\mathbf{A}^\top \mathbf{z}} \Rightarrow \kappa^*(\mathbf{x})$$

1 TERMEXPANSION

Erweitere die Datenmatrix $\mathbf{X} \in \mathbb{R}^{T \times D}$ durch Termberechnung zu $\mathbf{X}^\phi \in \mathbb{R}^{T \times L}$.

2 STICHPROBENSTATISTIKEN

Berechne die Stichprobenmomente \hat{p}_λ , $\hat{\mu}_\lambda^\phi$ und \mathbf{R}^ϕ für alle Klassen $\lambda = 1..K$.

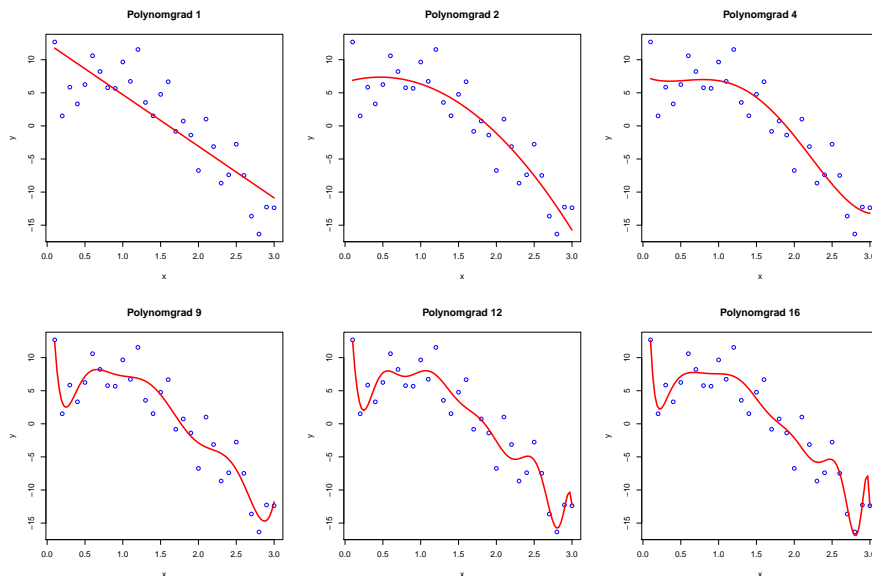
3 GAUßSCHE NORMALENGLEICHUNGEN

Invertiere \mathbf{R}^ϕ und berechne alle Koeffizientenvektoren $\mathbf{a}_\lambda \in \mathbb{R}^L$.

(Algorithmus)

Überanpassungseffekt bei Ausgleichspolynomen

Weiß verrauschte Daten zur Kurve $y = 7 + 2x - 3x^2$



Quadratmittelklassifikator mit Polynomtermen

Riskantes Spiel mit dem Fluch der Dimension

Affine Terme

$$L = D + 1$$

$$a_{\lambda,0} + \sum_{i=1}^D a_{\lambda,i} \cdot x_i$$

Quadratische Terme

$$L = (D+1)(D+2)/2$$

$$a_{\lambda,0} + \sum_{i=1}^D a_{\lambda,i} \cdot x_i + \sum_{i=1}^D \sum_{j \geq i} b_{\lambda,i,j} \cdot x_i x_j$$

Kubische Terme

$$L = (D+1)(D+2)(D+3)/6$$

$$a_{\lambda,0} + \sum_{i=1}^D a_{\lambda,i} \cdot x_i + \sum_{i=1}^D \sum_{j \geq i} b_{\lambda,i,j} \cdot x_i x_j + \sum_{i=1}^D \sum_{j \geq i} \sum_{k \geq j} c_{\lambda,i,j,k} \cdot x_i x_j x_k$$

Polynomterme n -ten Grades

$$L = \binom{D+n}{n}$$

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Nichtparametrische Dichteschätzung

Problem

Die Näherung der Verteilungen $P(\mathbf{x}|\Omega_\kappa)$ durch parametrische Dichtefunktionen (z.B. Normal- oder Gleichverteilung) birgt das Risiko einer **fehlerhaften Modellierung** der Lerndaten ω_κ in sich.

Empirische Verteilung der Daten?

Die Auszählung der Ereignishäufigkeiten zur Formulierung einer kanonischen Verteilung

$$f^{\text{emp}}(\mathbf{x}|\omega_\kappa) \propto \begin{cases} 1 & \mathbf{x} \in \omega_\kappa \\ 0 & \mathbf{x} \notin \omega_\kappa \end{cases}$$

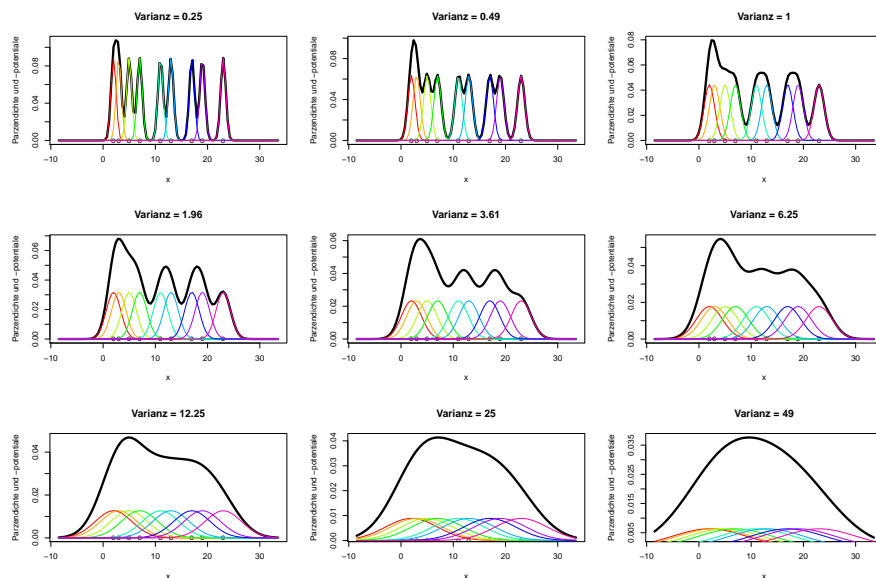
ist nur in diskreten Ereignisräumen opportun.

Lösung

Die (infinitesimal hohen und schmalen) Dirac-Gipfel der empirischen Datenverteilung werden durch geeignete Kern- oder Potentialfunktionen verflacht.

Parzen-Dichteschätzung

Potenzialfunktionen $\hat{=}$ Gaußkerne mit verschiedenen Streuungen σ^2



Parzen-Dichteschätzung

Potenzialfunktionen & Mittelung ihrer Verschiebungsinstanzen

Definition

Wir bezeichnen eine stetige Abbildung $g : \mathbb{R}^D \rightarrow \mathbb{R}_0^+$ als **Potenzialfunktion**, wenn sie flächennormiert ist und ihre Masse sich um den Ursprung konzentriert:

$$\int_{\mathbb{R}^D} g(\xi) d\xi = 1 \quad \text{und} \quad \int_{\mathbb{R}^D} \|\xi\|^2 \cdot g(\xi) d\xi < \infty$$

Die Funktion

$$f_{\omega}^{\text{parzen}}(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T g(\mathbf{x} - \mathbf{z}_t) \quad \text{bzw.} \quad f_{\omega, \sigma}^{\text{parzen}}(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \frac{1}{\sigma^D} \cdot g\left(\frac{\mathbf{x} - \mathbf{z}_t}{\sigma}\right)$$

heißt **Parzendichte** der Datenprobe $\omega = \{\mathbf{z}_1, \dots, \mathbf{z}_T\}$ mit Potentialfunktion $g(\cdot)$ (und Konzentration $1/\sigma$).

Bemerkungen

1. $g(\cdot)$ ist Gaußglocke, Dreieck/Kegel, uniformes Rechteck/Hyperwürfel.
2. Mit $g(\mathbf{x})$ ist auch jede verschobene und skalierte Instanz normiert und \mathbf{z} -konzentriert.
3. Die oben definierte Summe ist offensichtlich nichtnegativ und flächennormiert, repräsentiert also eine Verteilungsdichte.

Parzen-Dichteschätzung

Theoretisches Resultat: Asymptotische Konvergenz

Satz

Die multivariate Zufallsvariable $\mathbb{X} \in \mathbb{R}^D$ sei gemäß der stetigen Dichte $f_{\mathbb{X}}$ verteilt und $\omega = \{\mathbf{z}_n | n \in \mathbb{N}\}$ sei eine Folge zufällig und unabhängig gezogener Stichprobenelemente.

Genügt die normierte und beschränkte Potentialfunktion $g : \mathbb{R}^D \rightarrow \mathbb{R}_0^+$ der Bedingung

$$\lim_{|\mathbf{x}| \rightarrow \infty} g(\mathbf{x}) \cdot \prod_{d=1}^D x_d = 0$$

und ist $[h_\nu]$ eine positive reelle Zahlenfolge mit den Eigenschaften

$$\lim_{\nu \rightarrow \infty} h_\nu^q = 0 \quad \text{und} \quad \lim_{\nu \rightarrow \infty} \nu \cdot h_\nu^q = \infty, \quad (\forall q \in \mathbb{N})$$

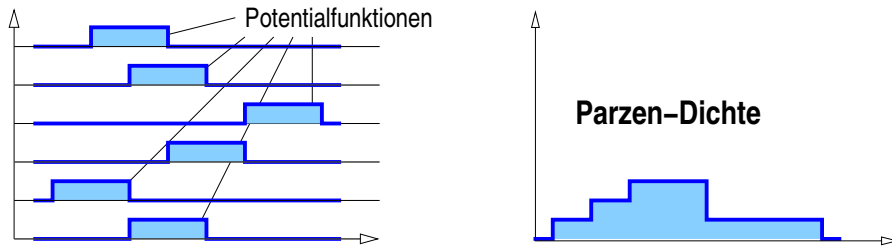
dann konvergiert die Parzen-Schätzung

$$\hat{P}_N(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h_N^D} \cdot g\left(\frac{\mathbf{x} - \mathbf{z}_n}{h_N}\right)$$

im quadratischen Mittel gegen die wahre Dichte $f_{\mathbb{X}}$.

Uniforme Parzenschätzung

Potenzialfunktion $\hat{=}$ Gleichverteilung auf Nullpunktumgebung



Definition

Eine Parzenschätzung mit

$$g(\mathbf{x}) = \begin{cases} 1/V & \mathbf{x} \in \mathcal{U} \\ 0 & \mathbf{x} \notin \mathcal{U} \end{cases}$$

für eine Nullpunktumgebung $\mathcal{U} \subset \mathbb{R}^D$ mit Hypervolumen $vol(\mathcal{U}) = V$ heißt **uniforme Parzenschätzung**.

Histogramm-Klassifikator

UPS-ähnlicher Klassifikator — und mit ganz ähnlichen Problemen

1 PARZELLIERUNG

Zerlege Merkmalraum \mathbb{R}^D in Hyperquader mit Volumina $V = \Delta^D$.

2 TREFFERQUOTEN SPEICHERN

Je Quader \mathcal{Q} setze $m_\kappa(\mathcal{Q}) = \#(\omega_\kappa \cap \mathcal{Q})$.

3 ENTSCHEIDUNGSREGEL

Klassifiziere \mathbf{x} nach maximalem $m_\kappa(\mathcal{Q}_\mathbf{x})$ in dem Quader mit $\mathbf{x} \in \mathcal{Q}_\mathbf{x}$.

Beispiel

- Merkmalraum ist \mathbb{R}^{20}
- Stichprobenumfang $N = |\omega| = 10^6$ Muster
- Ca. 10 Intervalle je Koordinate

Von den 10^{20} Zellen sind mindestens $10^{20} - 10^6$ leer ...

Fluch der Dimension

Fast alle Zellen ($\mathcal{Q}_\mathbf{x}$, $\mathcal{U}(\mathbf{x})$) sind leer!

Uniforme Parzenschätzung

UPS-Klassifikator — mit symmetrischer Nullumgebung

Klassendichte

$$\hat{P}(\mathbf{x}|\Omega_\kappa) = \frac{m_\kappa^{(V)}(\mathbf{x})}{T_\kappa \cdot V}, \quad m_\kappa^{(V)}(\mathbf{x}) \text{ ist die Anzahl der } \mathbf{z} \in \omega_\kappa \text{ mit } \mathbf{x} \in \mathcal{U}(\mathbf{z})$$

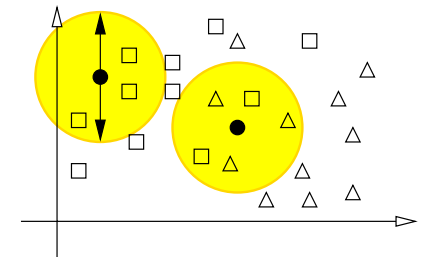
Prüfgröße

$$u_\kappa(\mathbf{x}) = \hat{p}_\kappa \cdot \hat{P}(\mathbf{x}|\Omega_\kappa) = \frac{T_\kappa}{T} \cdot \frac{m_\kappa^{(V)}(\mathbf{x})}{T_\kappa \cdot V} = \frac{m_\kappa^{(V)}(\mathbf{x})}{T \cdot V}$$

Entscheidungsregel

$$\delta(\mathbf{x}) = \operatorname{argmax}_\kappa m_\kappa^{(V)}(\mathbf{x})$$

(maximale Anzahl der ω_κ -Befunde in $\mathcal{U}(\mathbf{x})$)



Hypersphären-Schätzung und -klassifikator

Wie UPS — aber Befundzahl fixieren und Volumen variieren

Klassendichte

$$\hat{P}(\mathbf{x}|\Omega_\kappa) = \frac{m}{T_\kappa \cdot V_\kappa^{(m)}(\mathbf{x})}, \quad \text{mit } V_\kappa^{(m)}(\mathbf{x}) = \min_{\rho} \{vol(\mathcal{U}_\rho(\mathbf{x})) \mid |\mathcal{U}_\rho(\mathbf{x}) \cap \omega_\kappa| \geq m\}$$

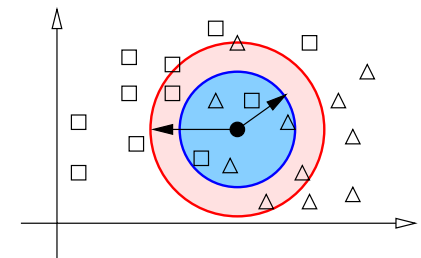
Prüfgröße

$$u_\kappa(\mathbf{x}) = \hat{p}_\kappa \cdot \hat{P}(\mathbf{x}|\Omega_\kappa) = \frac{T_\kappa}{T} \cdot \frac{m}{T_\kappa \cdot V_\kappa^{(m)}(\mathbf{x})} = \frac{m}{T \cdot V_\kappa^{(m)}(\mathbf{x})}$$

Entscheidungsregel

$$\delta(\mathbf{x}) = \operatorname{argmin}_\kappa V_\kappa^{(m)}(\mathbf{x})$$

(ω_κ versammelt m Muster auf kleinstem Raum um \mathbf{x})



k-Nächste-Nachbarn-Regel

Wie HSK — aber gesucht sind Befunde *beliebiger* Klasse

Prüfgröße

UPS der Rückschlußwahrscheinlichkeiten

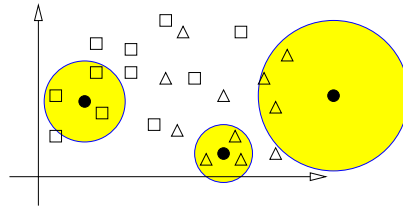
$$\hat{P}(\Omega_\kappa | \mathbf{x}) = \frac{\hat{P}(\Omega_\kappa) \cdot \hat{P}(\mathbf{x} | \Omega_\kappa)}{\hat{P}(\mathbf{x})} = \frac{\frac{T_\kappa}{T} \cdot \frac{m_\kappa^{(k)}(\mathbf{x})}{T_\kappa \cdot V}}{\frac{k}{T \cdot V}} = \frac{m_\kappa^{(k)}(\mathbf{x})}{k}$$

Unter den k nächsten Nachbarn von \mathbf{x} in ω (Hypervolumen = V) seien jeweils genau $m_\kappa^{(k)}(\mathbf{x})$ aus der Klasse κ gewesen.

Entscheidungsregel

$$\delta(\mathbf{x}) = \operatorname{argmax}_{\kappa} m_\kappa^{(k)}(\mathbf{x})$$

(die Mehrheitsfraktion unter den k nächsten \mathbf{x} -Nachbarn)



Nichtparametrische Klassifikatoren

Vorteile

- Keine anfechtbaren Vorannahmen über die Verteilung oder die Trennflächen der Musterklassen
- Keine aufwändige Lernprozedur

Nachteile

- Irrsinnig hoher Klassifikationsaufwand: $O(T \cdot D)$ je Muster
- Irrsinnig hoher Speicherbedarf: $O(T \cdot D)$ für ω

Offene Fragen

1. Welches V bzw. Δ für UPS- und Histogramm-Klassifikatoren?
2. Welche Befundraten m bzw. k für Hypersphären- und k NN-Klassifikatoren?
3. Alternative Distanzmaße für Hypersphären- und k NN-Klassifikatoren?

Nächster-Nachbar-Regel

NNR \triangleq Spezialfall von HSK ($m = 1$) und auch von k NN ($k = 1$)

Prüfgröße

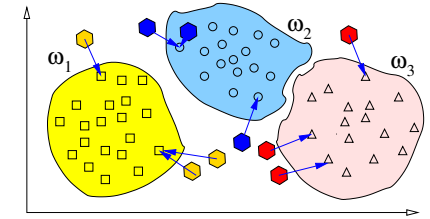
$$u_\kappa(\mathbf{x}) = \frac{1}{T \cdot V_\kappa^{(1)}(\mathbf{x})} \quad \text{bzw.} \quad u_\kappa(\mathbf{x}) = \frac{m_\kappa^{(1)}(\mathbf{x})}{1}$$

$$u_\kappa(\mathbf{x}) = \frac{1}{T \cdot \rho_\kappa^D}, \quad \rho_\kappa = \min_{\mathbf{z} \in \omega_\kappa} \|\mathbf{x} - \mathbf{z}\|$$

Entscheidungsregel

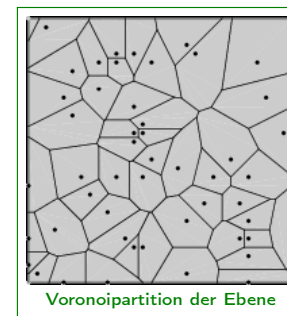
$$\delta(\mathbf{x}) = \operatorname{argmin}_{\kappa} \min_{\mathbf{z} \in \omega_\kappa} \|\mathbf{x} - \mathbf{z}\|$$

(Klassenindex des nächsten Nachbarn von \mathbf{x} in ω)



Voronoizellen

Nächster-Nachbar-Regel mit euklidischem Abstand



Definition

Für eine Teilmenge $\omega = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ des Raumes \mathbb{R}^D heißen die Mengen

$$\mathcal{V}_t \stackrel{\text{def}}{=} \{\mathbf{z} \mid \|\mathbf{z} - \mathbf{x}_t\|_2 \leq \|\mathbf{z} - \mathbf{x}_s\|_2, s \neq t\}$$

die **Voronoizellen** von ω .

Bemerkungen

1. Die Voronoizellen \mathcal{V}_t , $t = 1..T$, besitzen paarweise leeren Durchschnitt oder einen Durchschnitt vom Volumen Null; das Mengensystem wird deshalb oft als **Voronoipartition** bezeichnet.
2. Nur im Inneren der Zellen ist die 1NN-Regel eindeutig.

Voronoi-Partition und Delauney-Triangulierung

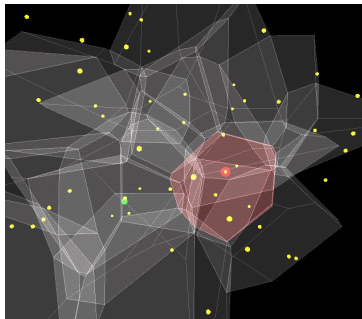
Topologische Repräsentation für schnelle Suchverfahren [mehr Information](#)

Definition

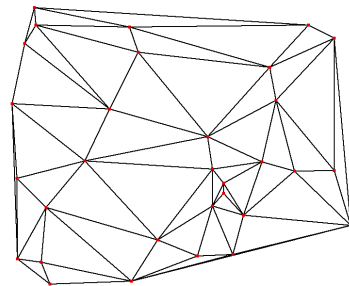
Für eine Teilmenge $\omega = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ des Raumes \mathbb{R}^D bezeichnen wir den ungerichteten Graphen $(\mathcal{X}, \mathcal{E})$ mit $\mathcal{X} = \{1, \dots, T\}$ und

$$\{s, t\} \in \mathcal{E} \Leftrightarrow \mathcal{V}_s \cap \mathcal{V}_t \neq \emptyset$$

als **Delauney-Triangulierung** der Punktmenge.



Voronoi-Zellen im \mathbb{R}^3 [Wolfram](#)



Delauneygraph

Klassifikator versus Klassifikationsverfahren

Konkreter Klassifikator

Entscheidungsregel im Raum Ω für K Klassen:

$$\delta : \begin{cases} \Omega & \rightarrow \{0, 1\}^K \text{ oder } [0, 1]^K \\ \mathbf{x} & \mapsto (\delta_1(\mathbf{x}), \dots, \delta_K(\mathbf{x})) \end{cases}, \quad \sum_{\lambda} \delta_{\lambda}(\mathbf{x}) = 1$$

Etikettierte Stichprobe

K -dimensionales Feld von Mustersequenzen (Multimengen)

$$\omega : \{1, 2, \dots, K\} \rightarrow \Omega^*$$

Abstraktes Klassifikationsverfahren

Aus einer endlichen, etikettierten Lernstichprobe wird eine Entscheidungsregel für K Klassen gelernt:

$$\mathfrak{A} : \begin{cases} (\Omega^*)^K & \rightarrow ([0, 1]^K)^{\Omega} \\ \omega = [\omega_k] & \mapsto \delta_{\omega, \mathfrak{A}} \end{cases}$$

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Reklassifikationsfehlerrate

Lernstichprobe $\hat{=}$ Teststichprobe

$$\Omega \supset \omega \xrightarrow{\text{Lernen}} \delta_{\omega} \xrightarrow{\text{Testen}} \hat{p}_{\epsilon}(\omega | \delta_{\omega})$$

Definition

Für ein Klassifikationsverfahren \mathfrak{A} bezeichnen wir die relative Häufigkeit

$$\hat{p}_{\epsilon}(\omega | \delta_{\omega, \mathfrak{A}})$$

von Fehlklassifikationen der gelernten Entscheidungsregel auf ihren eigenen Lerndaten $\omega = (\omega_1, \dots, \omega_K)$ als **Reklassifikationsfehlerrate**.

Bemerkungen

1. Die Reklassifikationsrate ist nicht eindeutig für das Verfahren \mathfrak{A} .
2. Die Reklassifikationsrate unterschätzt den Fehler bei Anwendung der Entscheidungsregel auf *neue* Muster.
3. Der Trend zur *optimistischen* Bewertung ist besonders kraß, wenn ω geringen Umfang besitzt.

Reklassifikationsfehlerrate

Iris-Datensatz · 3 Klassen · je 50 Muster des \mathbb{R}^4

Uniform-NB

5.3%	set	ver	vir
set	50	0	0
ver	0	50	0
vir	0	8	42

Affin-QMK

15.3%	set	ver	vir
set	50	0	0
ver	0	34	16
vir	0	7	43

1NN-Regel

0.0%	set	ver	vir
set	50	0	0
ver	0	50	0
vir	0	0	50

Beispiele

Der Fehler der 1NN-Regel beträgt praktisch Null.
Die UNB-Regel patzt u.U. in den Überlappungsgebieten.
Der QMK tendiert für hohen Polynomgrad zum Fehler Null.

Held-out-Fehlerrate

Iris-Datensatz · 50/100 Muster zum Lernen/Testen

Uniform-NB

16%	set	ver	vir
set	34	0	0
ver	5	28	0
vir	10	1	22

Affin-QMK

21%	set	ver	vir
set	34	0	0
ver	0	16	17
vir	0	4	29

1NN-Regel

4%	set	ver	vir
set	34	0	0
ver	0	31	2
vir	0	2	31

Beispiele

Der Fehler der 1NN-Regel beträgt 4% statt 0%.
Die UNB-Regel verdreifacht ihren Fehler.
Beim QMK erhöht sich der Fehler um ein Drittel.

Die Diskrepanz der Fehlerraten korrespondiert mit der Anzahl der Freiheitsgrade.

Neuklassifikationsfehlerrate

Held-out: Lernstichprobe ! = Teststichprobe

$$\Omega \supset \omega_L \xrightarrow{\text{Lernen}} \delta_{\omega_L} \xrightarrow{\text{Testen}} \hat{p}_\varepsilon(\omega_T | \delta_{\omega_L})$$

Definition

Für ein Klassifikationsverfahren \mathfrak{A} bezeichnen wir die relative Häufigkeit

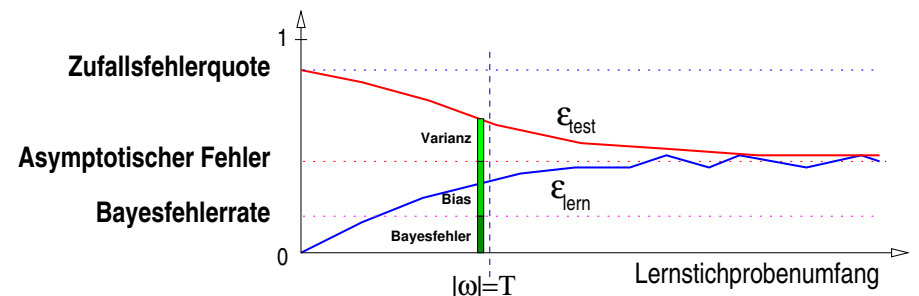
$$\hat{p}_\varepsilon(\omega_T | \delta_{\omega_L, \mathfrak{A}})$$

von Fehlklassifikationen der gelernten Entscheidungsregel auf einer zu den Lerndaten ω_L disjunkten Testdatenprobe ω_T als **Held-out-Fehlerrate**.

Bemerkungen

1. Die Held-out-Fehlerrate auf ω_T ist ein Näherungswert der Fehlerrate, die sich bei Anwendung von $\delta_{\omega_L, \mathfrak{A}}$ auf *neue* Muster ergibt.
2. Die Held-out-Fehlerrate von $\delta_{\omega_L, \mathfrak{A}}$ unterschätzt die Leistungsfähigkeit des Verfahrens \mathfrak{A} , da die Entscheidungsregel nur bezüglich einer endlichen Probe ω_L optimiert wurde.
3. Der Trend zur *pessimistischen* Bewertung ist besonders kraß, wenn ω_L geringen Umfang besitzt.

Fehlerrate, Überanpassung und Unteranpassung



- **Bayesfehler** — was theoretisch herauszuholen ist (Bayesregel)
- **Grenzfehler** — das Verfahren \mathfrak{A} mit infiniten Lernprobe
- **Zufallsfehler** — die „blinde“ Entscheidung: häufigste Klasse
- **Induktiver Bias** — schuld ist das unzureichende Datenmodell
- **Varianz** — schuld ist unsere begrenzte Lernstichprobe

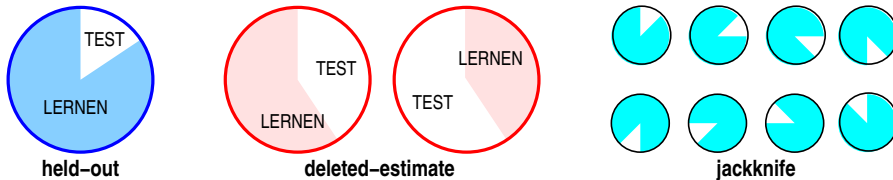
Rotationsverfahren

Robuste und ökonomische Klassifikatorentwurfstechnik

Fakt

Der Stichprobenumfang $|\omega^{\text{lern}}|$ ist verantwortlich für die Güte des angelernten Klassifikators.

Der Stichprobenumfang $|\omega^{\text{test}}|$ ist verantwortlich für die Genauigkeit seiner geschätzten Fehlerrate.



Rotations-Fehlerrate

Iris-Datensatz · 10-fache Kreuzvalidierung

Uniform-NB

11.3%	set	ver	vir
set	50	0	0
ver	3	43	4
vir	3	7	40

Affin-QMK

16%	set	ver	vir
set	49	1	0
ver	0	35	15
vir	0	8	42

1NN-Regel

4%	set	ver	vir
set	50	0	0
ver	0	47	3
vir	0	3	47

Beispiele

Die beiden parametrischen Klassifikatoren profitieren vom erhöhten Umfang der Lerndatenprobe (135 statt 50 Muster).

Die 1NN-Regel war offenbar bereits bei 50 Lernmustern in den Sättigungsbereich gelangt.

Rotationsverfahren

Benutze (fast) alle Muster zum Lernen UND zum Testen!

(Algorithmus)

- 1 ZERLEGE die Daten DISJUNKT: $\omega = \omega^{(1)} \cup \omega^{(2)} \cup \dots \cup \omega^{(R)}$.
- 2 LERNE die Klassifikatoren $\delta^{(r)} = \delta_{\omega^{(r)}}$, $\omega^{(r)} = \omega \setminus \omega^r$, $r = 1..R$.
- 3 KLASSIFIZIERE alle Muster $\mathbf{x} \in \omega^r$ mit der Regel $\delta^{(r)}$.
- 4 KUMULIERE den absoluten und relativen Gesamtfehler.

(summiertog(A))

Bemerkungen

1. Deleted-estimate ($R = 2$)
Jackknife (R -fache Kreuzvalidierung)
Leave-one-out ($R = |\omega|$).
2. Datenpartitionierung zufällig oder systematisch?
(Klassenzugehörigkeit, Doubletten, Zeitreihen)

Genauigkeit der Fehlerschätzung

Wie viele Testmuster werden für eine seriöse Fehlervoraussage benötigt?

Fakt

Ein Klassifikationstest entspricht dem Ziehen aus einer Urne mit einem (unbekannten) Anteil p_ϵ von „Nieten“.

Lösung

Für einen Klassifikator δ mit wahrer Fehlerrate p_ϵ gehorcht die diskrete Zufallsvariable

$$\mathbb{M} = \text{Anzahl falschklassifizierter unter } N \text{ Mustern}$$

der **Binomialverteilung**:

$$P(\mathbb{M} = n_f) = \mathcal{B}(n_f | p_\epsilon, N) = \binom{N}{n_f} \cdot p_\epsilon^{n_f} \cdot (1 - p_\epsilon)^{N - n_f}$$

Genauigkeit der Fehlerschätzung

Große Teststichprobe → Fehleranzahl ist normalverteilt

Satz (Binomialer Grenzwertsatz)

Ist die Zufallsvariable \mathbb{X}_N für jedes $N \in \mathbb{N}$ binomialverteilt gemäß $\mathcal{B}(\cdot | N, p)$, so gilt die asymptotische Näherungsformel

$$\lim_{N \rightarrow \infty} P \left(\frac{\mathbb{X}_N - \mathcal{E}[\mathbb{X}_N]}{\sqrt{\text{Var}[\mathbb{X}_N]}} \leq x \right) = \frac{1}{\sqrt{2\pi}} \cdot \int_{-\infty}^x e^{-\frac{\xi^2}{2}} d\xi$$

mit den Mittelwerten $\mu = \mathcal{E}[\mathbb{X}_N] = N \cdot p$ und Varianzen $\sigma^2 = \text{Var}[\mathbb{X}_N] = N \cdot p \cdot (1 - p)$ der Binomialverteilung.

Beispiel

Für den Kreuzvalidierungstest des QMK auf den Iris-Daten gilt:

$$N = 150, \quad \hat{p}_\varepsilon = 0.16, \quad \mu = 24, \quad \sigma^2 = 20.16, \quad \sigma \approx 4.5$$

Folglich liegt mit 95% Sicherheit die Zahl der Nieten im Intervall $[15, 33]$, die wahre Fehlerrate also zwischen 10% und 22% — ein katastrophales Paradebeispiel indiskutabler Evaluierungsbedingungen!

Analytische Aussagen zur Bayesfehlerrate?

Satz (Fehlerrate der NN-Regel)

Ist $p_\varepsilon^{\text{Bayes}}$ die Bayesfehlerrate eines Klassifikationsproblems $[\Omega_\kappa]$ und $p_\varepsilon^{\text{NN}}$ die asymptotische Fehlerrate ($|\omega| \rightarrow \infty$) des Nächster-Nachbar-Klassifikators, so gilt

$$p_\varepsilon^{\text{Bayes}} \leq p_\varepsilon^{\text{NN}} \leq p_\varepsilon^{\text{Bayes}} \cdot \left(2 - \frac{K}{K-1} \cdot p_\varepsilon^{\text{Bayes}} \right).$$

Für hinreichend kleine Fehlerraten folgt daraus

$$p_\varepsilon^{\text{NN}}/2 \leq p_\varepsilon^{\text{Bayes}} \leq p_\varepsilon^{\text{NN}}.$$

Bemerkungen

- Die simple $L^1 O$ -Evaluierung eines Datensatzes liefert eine gute Schätzung der nicht-asymptotischen Fehlerrate $p_\varepsilon^{\text{NN}}(\omega)$.
 - Für sehr große Proben ω können wir schließen, daß die Bayesfehlerrate oberhalb $\hat{p}_\varepsilon^{\text{NN}}(\omega) / 2$ liegt.
- Wackeliges Kriterium für „Klassifikatorentwicklung erfolgreich abgeschlossen“

Genauigkeit der Fehlerschätzung

Signifikanzniveau & Konfidenzintervall

Faustregel

Wenn wir auf ein Signifikanzniveau von $\mu \pm C\sigma$ orientieren und gegenüber der unbekannten Fehlerrate p_ε einen relativen Fehler von $\Delta > 0$ unterschreiten möchten, so werden mindestens

$$N = (1-p_\varepsilon)/p_\varepsilon \cdot C^2/\Delta^2$$

Testmuster benötigt

Bemerkungen

- Zur Signifikanzschwelle $C = 3$ (99%) und relativer Abweichung $\Delta = 10\%$ sind 2700, 17000 oder 90000 Muster gefordert für Fehlerraten um 25%, 5% oder 1%.
- Die Übertragung der Fehlerschranken für geschätzte Rate bezüglich wahrer Rate auf Fehlerschranken für die wahre Rate bezüglich des Schätzwerts ist natürlich strenggenommen ein (unzulässiger) *Abduktionsschluß*.

Beweis.

Alle Abschätzungen der Fehlerraten (0/1-Risiko) analysieren die punktwisen Fehlerwahrscheinlichkeiten $p_\varepsilon(\mathbf{x})$ für Bayesregel und 1NN-Regel und bilden abschließend die Erwartungswerte.

Es bezeichne $p_\lambda(\mathbf{x})$ die wahre a posteriori Wahrscheinlichkeit, daß Muster \mathbf{x} zur Klasse Ω_λ gehört.

Die Bayesregel klassifiziert \mathbf{x} nach maximalem $p_\lambda(\mathbf{x})$; das Risiko beträgt also

$$p_\varepsilon^{\text{Bayes}}(\mathbf{x}) = 1 - \max_\lambda p_\lambda(\mathbf{x}).$$

Die asymptotische 1NN-Regel findet zu \mathbf{x} einen nächsten Nachbarn $\tilde{\mathbf{x}} \in \omega$, der wegen $|\omega| \rightarrow \infty$ infinitesimal nahe bei \mathbf{x} liegt. Folglich besitzt $\tilde{\mathbf{x}}$ dieselbe a posteriori Verteilung wie \mathbf{x} selbst und wird von 1NN gemäß $p_\lambda(\mathbf{x})$ klassifiziert. Das Risiko beträgt also

$$p_\varepsilon^{\text{NN}}(\mathbf{x}) = \sum_\lambda p_\lambda(\mathbf{x}) \cdot (1 - p_\lambda(\mathbf{x})) = 1 - \sum_\lambda p_\lambda^2(\mathbf{x}).$$

Die punktweise Abschätzung von $p_\varepsilon^{\text{NN}}(\mathbf{x})$ nach unten und oben ergibt sich durch mehr oder weniger langwierige Standardumformungen ... □

Univariate Gleichverteilungsdichten

Leave-One-Out-Schätzung der Verteilungsparameter

$$a^* = \min_{t=1..T} z_t - \alpha, \quad b^* = \max_{t=1..T} z_t + \beta$$

Lemma

Für eine Lernprobe $\omega = \{z_1, \dots, z_T\}$ nimmt die (unlogarithmierte) Leave-One-Out-Zielfunktion

$$\ell_{L^1O}(\alpha, \beta) = \prod_{t=1}^T f_{[a^{(t)}, b^{(t)}]}^{unif}(z_t)$$

mit den Intervallgrenzschätzern

$$a^{(t)} \stackrel{\text{def}}{=} \min_{s \neq t} z_s - \alpha, \quad b^{(t)} \stackrel{\text{def}}{=} \max_{s \neq t} z_s + \beta$$

ihren maximalen Wert bei den Dilatationsparametern

$$\alpha^* \stackrel{\text{def}}{=} z_{(2)} - z_{(1)}, \quad \beta^* \stackrel{\text{def}}{=} z_{(T)} - z_{(T-1)}$$

an, d.h., $a^* = 2z_{(1)} - z_{(2)}$ und $b^* = 2z_{(T)} - z_{(T-1)}$.

Axiomatische Wahrscheinlichkeitstheorie

Zufallsexperiment · Elementarereignisse · Ereignisse

Definition

Es sei \mathcal{U} der Raum aller möglichen Ausgänge eines „Zufallsexperiments“. Ein Mengensystem $\mathfrak{E} \subseteq \mathfrak{P}\mathcal{U}$ heißt **Ereignisraum** über \mathcal{U} , wenn es die folgenden Eigenschaften einer σ -**Algebra** besitzt:

- (1) $\mathcal{U} \in \mathfrak{E}$
- (2) $A \in \mathfrak{E} \Rightarrow A^c = \mathcal{U} \setminus A \in \mathfrak{E}$
- (3) $A_n \in \mathfrak{E} (\forall n \in \mathbb{N}) \Rightarrow \bigcup_{n \in \mathbb{N}} A_n \in \mathfrak{E}$

Die Elemente von \mathfrak{E} heißen **Ereignisse**, die Elemente von \mathcal{U} heißen **Elementarereignisse**. Wir bezeichnen ferner $\mathcal{U}, \emptyset \in \mathfrak{E}$ als das **sichere** bzw. das **unmögliche** Ereignis.

Aufgabenstellung

Statistische Entscheidungstheorie

Klassifikatortypen

Uniformer naiver Bayesklassifikator

Linearer Quadratmittelklassifikator

Nichtparametrische Klassifikatoren

Verallgemeinerung auf neue Muster

Mathematische Hilfsmittel

Axiomatische Wahrscheinlichkeitstheorie

Eigenschaften und Beispiele von Ereignisräumen

Lemma

Ein Ereignisraum \mathfrak{E} über \mathcal{U} ist insbesondere eine Boolesche Algebra:

- (4) $\emptyset \in \mathfrak{E}$
- (5) $A_1, A_2 \in \mathfrak{E} \Rightarrow A_1 \cup A_2 \in \mathfrak{E}$
- (6) $A_1, A_2 \in \mathfrak{E} \Rightarrow A_1 \cap A_2 \in \mathfrak{E}$

Beispiele

1. $\mathcal{U} = \{1, 2, 3, 4, 5, 6\}$ (Würfel), $\mathfrak{E} = \mathfrak{P}\mathcal{U}$, $|\mathfrak{E}| = 2^6 = 64$
2. $\mathcal{U} = \{(a, b, c) \mid a, b, c \in \{K, Z\}\}$ (drei Münzen), $|\mathfrak{P}\mathcal{U}| = 2^8$
3. $\mathcal{U} = \mathbb{N}_0$ (Verkehrstote Deutschland 1984)
 $\mathfrak{E} = \{A_k \mid k \in \mathbb{N}_0\}$ mit $A_k \stackrel{\text{def}}{=} \{k\}$ oder $A_k \stackrel{\text{def}}{=} \{1, \dots, k\}$
4. $\mathcal{U} = \mathbb{R}_0^+$ (Lebensdauer einer Glühbirne)
 $\mathfrak{E} = \{A_{r,\delta} \mid r, \delta \in \mathbb{R}_0^+\}$ mit $A_{r,\delta} \stackrel{\text{def}}{=} \{x \mid r \leq x \leq r + \delta\}$
5. $\mathcal{U} = \mathbb{R}^5$ (Ertrag von fünf Weizenarten), z.B. $A_{3,1} = \{x \mid x_3 \geq 2x_1\}$

Axiomatische Wahrscheinlichkeitstheorie

Kolmogorov-Axiome für Wahrscheinlichkeitsräume

Definition (Kolmogorov)

Ist \mathfrak{E} ein Ereignisraum über \mathfrak{U} , so heißt die Abbildung

$$P : \mathfrak{E} \rightarrow \mathbb{R}$$

eine **Wahrscheinlichkeitsfunktion** über \mathfrak{U} falls gilt:

- (1) $P(A) \geq 0$ für alle $A \in \mathfrak{E}$
- (2) $P(\mathfrak{U}) = 1$
- (3) $P\left(\bigcup_{n \in \mathbb{N}} A_n\right) = \sum_{n \in \mathbb{N}} P(A_n)$

Dabei seien die A_n, A_m paarweise disjunkt.

Die Eigenschaft (3) heißt **σ -Additivität** von P . Das Tripel $(\mathfrak{U}, \mathfrak{E}, P)$ heißt dann **Wahrscheinlichkeitsraum** über \mathfrak{U} .

Axiomatische Wahrscheinlichkeitstheorie

Endliche Wahrscheinlichkeitsräume · Gleichverteilung

Bemerkung

Ein Wahrscheinlichkeitsraum $(\mathfrak{U}, \mathfrak{E}, P)$ mit endlicher Menge $\mathfrak{U} = \{\varepsilon_1, \dots, \varepsilon_N\}$ von Elementarereignissen ist vollständig durch die Wahrscheinlichkeitsmassen $p_j = P(\{\varepsilon_j\})$, $j = 1, \dots, N$ charakterisiert, denn für jedes Ereignis $A \in \mathfrak{E} = \mathfrak{P}\mathfrak{U}$ gilt dann:

$$P(A) = \sum_{\varepsilon \in A} P(\{\varepsilon\}) = \sum_{\varepsilon_j \in A} p_j$$

Eine endliche Verteilung

$$P(\{\varepsilon_j\}) = \frac{1}{N}, \quad j = 1, \dots, N$$

heißt **Gleichverteilung** oder **uniforme** Verteilung.

Lemma

In einem gleichverteilten Wahrscheinlichkeitsraum $(\mathfrak{U}, \mathfrak{E}, P)$ gilt:

$$A \in \mathfrak{E} \Rightarrow P(A) = \frac{n_A}{N} = \frac{|A|}{|\mathfrak{U}|} = \frac{\#(\text{„günstige Fälle“})}{\#(\text{„alle Fälle“})}$$

Axiomatische Wahrscheinlichkeitstheorie

Gleichungen und Ungleichungen für Wahrscheinlichkeitsräume

Satz

In einem Wahrscheinlichkeitsraum $(\mathfrak{U}, \mathfrak{E}, P)$ gelten für alle $A, B, A_n \in \mathfrak{E}$ die Eigenschaften:

- (1) $P : \mathfrak{E} \rightarrow [0, 1]$
- (2) $P(\emptyset) = 0$
- (3) $P(A^c) = 1 - P(A)$
- (4) $P(A) = P(A, B) + P(AB^c)$
- (5) $P(A \setminus B) = P(AB^c) = P(A) - P(A, B)$
- (6) $P(A \cup B) = P(A) + P(B) - P(A, B)$
- (7) $A \subseteq B \Rightarrow P(A) \leq P(B)$
- (8) $P\left(\bigcup_{n \in \mathbb{N}} A_n\right) \leq \sum_{n \in \mathbb{N}} P(A_n)$

Dabei sei die Kurzschreibweise $A, B \stackrel{\text{def}}{=} A \cap B$ vereinbart.

Prinzip vom unzureichenden Grunde

Principle of insufficient reason

Laplace-Prinzip

Haben wir keinen Anlaß, etwas anderes (Klügeres?) anzunehmen, so betrachten wir die Elementarereignisse eines W.-Raumes als gleichwahrscheinlich.

Klassische Wahrscheinlichkeitstheorie

Probabilistisches Schließen heißt **Zählen!**

Mathematisches Hilfsmittel ist die **Kombinatorik**.

Problem

Wie sieht eine unendliche Gleichverteilung aus?

Wie groß ist die Wahrscheinlichkeit, aus \mathbb{N} eine gerade Zahl zu ziehen?

$$\mathbb{N} = \{1, 2, 3, 4, 6, 5, 8, 10, 12, 7, 14, 16, 18, 20, 9, 22, 24, 26, 28, 30, 11, 32, 34, 36, \dots\}$$

Axiomatische Wahrscheinlichkeitstheorie

Bedingte Wahrscheinlichkeiten & Teilraumbildung

Definition

Sei $(\mathcal{U}, \mathfrak{E}, P)$ ein Wahrscheinlichkeitsraum und $A, B \in \mathfrak{E}$ mit $P(B) \neq 0$. Dann heißt

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

die **bedingte Wahrscheinlichkeit** des Ereignisses A unter der Bedingung B .

Satz

Sei $(\mathcal{U}, \mathfrak{E}, P)$ ein Wahrscheinlichkeitsraum und $B \in \mathfrak{E}$ mit $P(B) \neq 0$. Dann ist auch $(\mathcal{U}, \mathfrak{E}, P^B)$ mit

$$P^B(A) \stackrel{\text{def}}{=} P(A|B) \quad \text{für alle } A \in \mathfrak{E}$$

ein Wahrscheinlichkeitsraum.

Axiomatische Wahrscheinlichkeitstheorie

Marginalisierung und totale Wahrscheinlichkeit

Lemma

Im Wahrscheinlichkeitsraum $(\mathcal{U}, \mathfrak{E}, P)$ gelte für \mathcal{U} die paarweise disjunkte Zerlegung $\mathcal{U} = \bigcup_{i=1}^n B_i$. Dann gilt für jedes Ereignis A die Gleichung

$$P(A) = \sum_{i=1}^n P(A, B_i) = \sum_{i=1}^n P(A|B_i) \cdot P(B_i)$$

der **totalen Wahrscheinlichkeit** sowie für jedes $j = 1, \dots, n$ die **Mehrwege-Bayesformel**

$$P(B_j|A) = \frac{P(A|B_j) \cdot P(B_j)}{\sum_{i=1}^n P(A|B_i) \cdot P(B_i)}.$$

Axiomatische Wahrscheinlichkeitstheorie

Bayesformel & Kettenregel

Lemma (Bayesformel)

Für Ereignisse $A, B \in \mathfrak{E}$ mit $P(A) \neq 0 \neq P(B)$ gilt die Gleichung

$$P(B|A) = \frac{P(B, A)}{P(A)} = \frac{P(A, B)}{P(A)} = \frac{P(A|B) \cdot P(B)}{P(A)}$$

Handelt es sich bei B um eine Ursache und bei A um eine Wirkung, so bezeichnen wir $P(B|A)$ auch als **a posteriori Wahrscheinlichkeit** (von B unter A).

Lemma (Kettenregel)

Für die Ereignisse A_1, \dots, A_n eines Wahrscheinlichkeitsraumes sei $P(A_1 \dots A_{n-1}) \neq 0$. Dann gilt die Faktorzerlegung

$$P(A_1 \dots A_n) = P(A_1) \cdot P(A_2|A_1) \cdot P(A_3|A_1 A_2) \cdot \dots \cdot P(A_n|A_1 \dots A_{n-1})$$

Axiomatische Wahrscheinlichkeitstheorie

Statistische und paarweise statistische Unabhängigkeit

Definition

Zwei Ereignisse A, B heißen **paarweise** (statistisch) **unabhängig** ($A \not\sim B$), wenn eine der folgenden Bedingungen erfüllt ist:

- (1) $P(AB) = P(A) \cdot P(B)$
- (2) $P(A|B) = P(A)$ und $P(B) \neq 0$
- (3) $P(B|A) = P(B)$ und $P(A) \neq 0$

Die Ereignisse $A_1, \dots, A_n \in \mathfrak{E}$ heißen statistisch **unabhängig** genau dann, wenn für alle denkbaren Indexkombinationen gilt:

$$\begin{aligned} P(A_i A_j) &= P(A_i) \cdot P(A_j) \\ P(A_i A_j A_k) &= P(A_i) \cdot P(A_j) \cdot P(A_k) \\ &\vdots = \vdots \\ P\left(\bigcap_{\ell} A_{\ell}\right) &= \prod_{\ell} P(A_{\ell}) \end{aligned}$$

Beispiele

- Paarweise statistische Unabhängigkeit (Würfelpaar-Versuch):

A = „die Augensumme ist ungerade“

B = „der erste Wurf war eine '6'“

C = „die Augensumme beträgt sieben“

Dann gilt $A \not\sim B$, $B \not\sim C$, aber $A \sim C$.

- **Achtung** — aus der paarweisen Unabhängigkeit folgt im allgemeinen **nicht** die (allgemeine) Unabhängigkeit. Die Ereignisse

A_1 = „der erste Wurf war ungerade“

A_2 = „der zweite Wurf war ungerade“

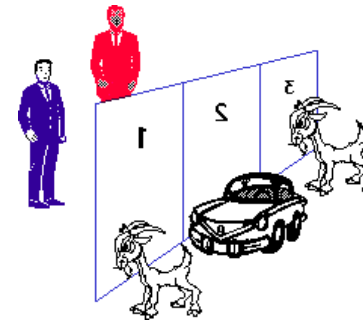
A_3 = „die Summe war ungerade“

sind zwar paarweise unabhängig, doch es gilt

$$P(A_1 A_2 A_3) = 0 \neq \frac{1}{8} = P(A_1) \cdot P(A_2) \cdot P(A_3)$$

Das Monty-Hall-Paradoxon

Der Fachmann staunt ... und der Laie wundert sich!



Fragestellung

Der Moderator einer Spielshow zeigt dem Kandidaten 3 Türen. „**Hinter einer der drei Türen steht der Hauptgewinn, ein Volvo. Hinter den anderen beiden Türen stehen Ziegen. Welche Tür wählen Sie?**“ Nachdem der Kandidat sich entschieden hat, öffnet der Moderator eine der beiden anderen Türen — mit einer stinkenden Ziege dahinter! „ **Bleiben Sie bei Ihrer Wahl oder möchten Sie noch einmal umwählen?**“

Zusammenfassung (6)

1. Ziel des Klassifikatorentwurfs ist die **Minimierung** der zu erwartenden **Fehlerrate**.
2. Der Klassifikator wird mit Methoden des **maschinellen Lernens** aus einer etikettierten **Lernstichprobe** gewonnen.
3. Laut statistischer **Entscheidungstheorie** wird die minimale Fehlerrate durch die **Bayesregel** (MAP-Regel) garantiert.
4. Die Bayesregel ist **nicht praktikabel**, denn sie erfordert die exakte Kenntnis des **Mustererzeugungsprozesses**.
5. **Statistische Klassifikatoren** wie der uniform-naive Bayesklassifikator schätzen ein Näherungsmodell der Musterverteilung.
6. **Diskriminative Klassifikatoren** wie die lineare Quadratmittel-Trennfunktion approximieren die a posteriori Klassenwahrscheinlichkeiten.
7. **Nichtparametrische Klassifikatoren** wie die kNN-Regel verzichten auf explizite Strukturannahmen; der Klassifikationsaufwand steigt dann mit dem Lernprobenumfang.
8. Eine faire Schätzung der **Klassifikatorleistung** erfordert strikte **Trennung von Lern- und Testdaten**, z.B. mittels Kreuzvalidierung (L^1O).