

MUSTERERKENNUNG

Vorlesung im Sommersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 6. März 2017

Teil XI

Klassifikatorübergreifende Verfahren

Aufgabenstellungen und Lösungsprinzipien

... die allen/vielen Klassifikationsverfahren gemeinsam sind

Offene Fragen

Alle Lernverfahren implementiert · alle Musterdaten gesammelt:

- Detektion von Ausnahmesituationen
 - *Musterklassen ohne Lerndaten*
- Kostenstruktur und Gütevergleich bei Zweiklassenproblemen
- Reduktion von Mehrklassenproblemen auf $K = 2$
- Ensembles kooperierender Klassifikatoren
 - *Expertenkomitees:* $\left\{ \begin{array}{l} \text{Zusammensetzung} \\ \text{Entscheidungsstrategien} \end{array} \right\}$
- Kombinatorische Auswahl des optimalen Merkmalsatzes
- Schritthaltendes Lernen und Klassifikatoradaption
- Lernen mit teilweise etikettierten Trainingsdaten

Ein-Klassen-Szenarien

Zwei-Klassen-Szenarien

Mehr-Klassen-Szenarien

Ensemblemethoden — Bagging & Boosting

Kombinatorische Merkmalauswahl

Inkrementelles Lernen

Transduktion

Unterscheidung zweier Klassen Ω_1, Ω_0

Keine Lernstichprobe ω_0 für Ω_0 verfügbar!

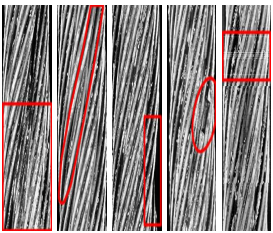
Fehlerdetektion

Ω_1 : intakte Muster
 Ω_0 : defekte Muster

Fehler sind ziemlich
selten.

Die Fehlerklasse ist
inhomogen.

Beisp. „Seilbahn“



Statistischer Klassifikator

Prüfgröße \hookrightarrow Laplace-Prinzip
 (vom unzureichenden Grunde)

$$u_1(\mathbf{x}) = p_1 \cdot \mathcal{N}(\mathbf{x} \mid \mu_1, \mathbf{S}_1)$$

$$u_0(\mathbf{x}) = p_0 \cdot \text{const}$$

Entscheidungsregel

$$\delta(\mathbf{x}) = \begin{cases} \Omega_1 & \mathcal{N}(\mathbf{x} \mid \mu_1, \mathbf{S}_1) \geq \theta \\ \Omega_0 & \text{sonst} \end{cases}$$

Phantomdaten

(diskriminativ/nichtparametrisch)

Auswürfeln einer Surrogatprobe $\tilde{\omega}_0 \subset \Omega_0$

Zweiklassenszenarium

Verifikationsaufgabe — $\Omega_1/\Omega_0 = \begin{cases} \text{positive} \\ \text{negative} \end{cases}$ Rückmeldung

Kostenmatrix

$$\mathbf{R} = \begin{pmatrix} r_{00} & r_{01} \\ r_{10} & r_{11} \end{pmatrix} = \begin{pmatrix} 0 & \rho_{FA} \\ \rho_{FR} & 0 \end{pmatrix}, \quad \rho_{TA} = 0 = \rho_{TR}$$

Prüfgrößen

$$u_\lambda(\mathbf{x}) = \begin{cases} \rho_{FR} \cdot P(\Omega_1, \mathbf{x}) & \lambda = 1 \\ \rho_{FA} \cdot P(\Omega_0, \mathbf{x}) & \lambda = 0 \end{cases}$$

$$\Delta(\mathbf{x}) \stackrel{\text{def}}{=} \log \frac{u_1(\mathbf{x})}{u_0(\mathbf{x})} = \log \frac{\rho_{FR}}{\rho_{FA}} + \log \text{odds}(\mathbf{x})$$

Entscheidungsregel

$$\lambda^*(\mathbf{x}) = \begin{cases} 1 & \log \text{odds}(\mathbf{x}) \geq \theta \\ 0 & \log \text{odds}(\mathbf{x}) \leq \theta \end{cases}, \quad \theta \stackrel{\text{def}}{=} \log \frac{\rho_{FA}}{\rho_{FR}}$$

Ein-Klassen-Szenarien

Zwei-Klassen-Szenarien

Mehr-Klassen-Szenarien

Ensemblemethoden — Bagging & Boosting

Kombinatorische Merkmalauswahl

Inkrementelles Lernen

Transduktion

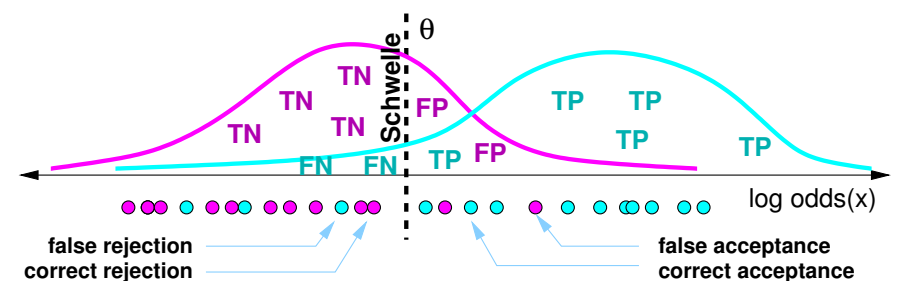
Kalibrierung der Verifikationsregel

Arbeitspunkt

Die Einstellung der Prüfgrößenschwelle θ in der Verifikationsregel

$$\log \text{odds}(\mathbf{x}) \stackrel{?}{\geq} \theta$$

heißt Arbeitspunkt. Er bestimmt, wieviele Muster akzeptiert werden und wieviele zurückgewiesen.



Ein Zoo von Fehlermaßen ...

... bei festem Arbeitspunkt θ

Verifikationsergebnis

θ	accept	reject	
good guys	N_{tt}	N_{tf}	$N_{t.}$
bad guys	N_{ft}	N_{ff}	$N_{f.}$
	$N_{.t}$	$N_{.f}$	

Referenzbezogene Raten

$N_{tt}/N_{t.}$ TPR, sensitivity, recall
 $N_{tf}/N_{t.}$ FNR, loss = $1 - \text{recall}$
 $N_{ft}/N_{f.}$ FPR, fallout
 $N_{ff}/N_{f.}$ TNR, vigilance = $1 - \text{fallout}$

Populationsbezogene Raten

$(N_{tt} + N_{ff})/N_{..}$ accuracy
 $(N_{tf} + N_{ft})/N_{..}$ error

Gewichteter F-Value

$$\left(\frac{\beta}{\text{PRE}} + \frac{1-\beta}{\text{REC}} \right)^{-1} = \frac{N_{tt}}{2N_{tt} + N_{tf} + N_{ft}}$$

Hypothesenbezogene Raten

$N_{tt}/N_{t.}$ precision (Ausbeute)
 $N_{tf}/N_{f.}$ waste
 $N_{ft}/N_{t.}$ garbage = $1 - \text{precision}$
 $N_{ff}/N_{f.}$ care = $1 - \text{waste}$

R Programmcode

code:receiver-operator

```
load ("data/twoclass.rda")

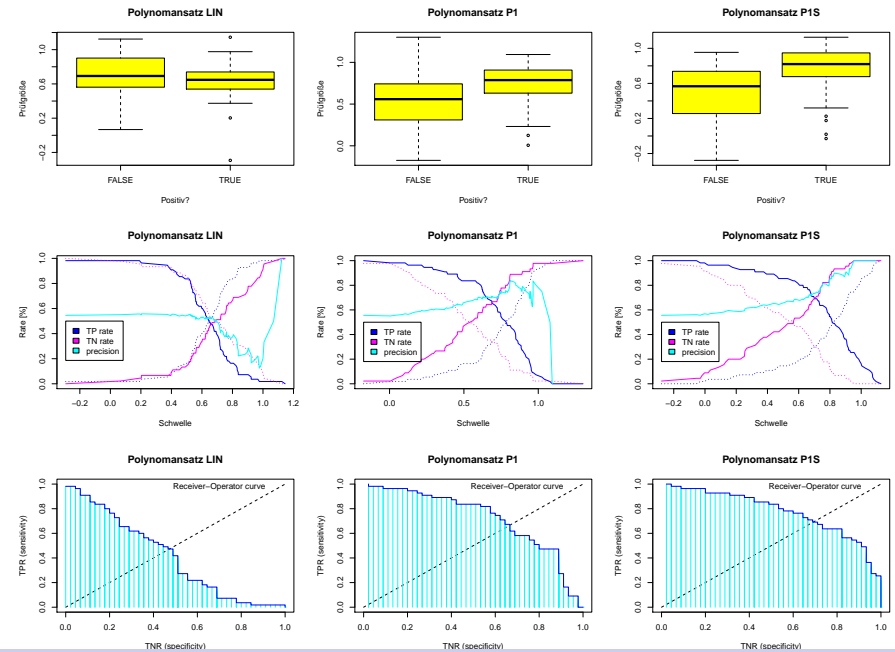
plot.scorebox <- function (u, ...)
  plot (u~factor(names(u)), xlab="Positiv?", ylab="Prüfgröße", ...)

plot.recall <- function (u, copa=c("blue","magenta","cyan"), ...) {
  u <- sort(u)
  hit <- names(u) == "TRUE"
  fnr <- cumsum(hit) / sum(hit)
  tpr <- 1-fnr
  tnr <- cumsum(!hit) / sum(!hit)
  fpr <- 1-tnr
  pre <- (sum(hit)-cumsum(hit)) / (length(hit):1-1)
  plot (u, tpr, type="n", xlab="Schwelle", ylab="Rate [%]", ...)
  lines (u,tpr,col=copa[1],type="l")
  lines (u,fnr,col=copa[1],type="l",lty="dotted")
  lines (u,tnr,col=copa[2],type="l")
  lines (u,fpr,col=copa[2],type="l",lty="dotted")
  lines (u,pre,col=copa[3],type="l")
  legend (min(u), 0.5,
    legend=c("TP rate","TN rate","precision"), fill=copa)
}

plot.ROC <- function (u, ...) {
  u <- sort(u)
  hit <- names(u) == "TRUE"
  fnr <- cumsum(hit) / sum(hit)
  tpr <- 1-fnr
  tnr <- cumsum(!hit) / sum(!hit)
  fpr <- 1-tnr
  plot (0:1, 0:1, type="l", lty="dashed",
    legend=c("TP rate","TN rate","precision"), fill=copa)
```

ROC — Receiver-Operator-Charakteristik

Diabetesdaten der Pima-Indianer mit verschiedenen Polynomklassifikatoren



Vergleich von Klassifikatoren

Problem (Kostenmatrix unbekannt)

Ohne Kenntnis der Kostenmatrix läßt sich der **tatsächliche Arbeitspunkt** θ der Entscheidungsregel nicht fixieren.

Problem (Pareto-Phänomen)

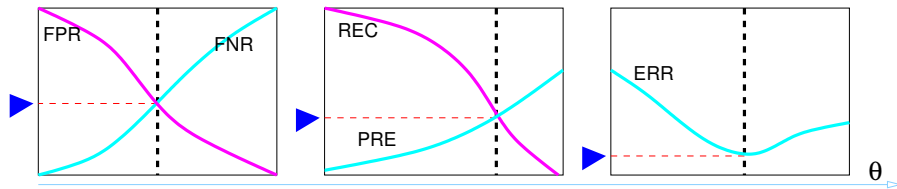
Die individuellen arbeitspunktabhängigen Fehlerkurven zweier wettbewerbender Klassifikatorprüfgrößen sind **nicht notwendig vergleichbar**.

Lösung

- Bestimmung eines **optimalen Arbeitspunktes** θ^* und Vergleich der korrespondierenden Gütewerte (BEP, EER, 45°-Tangente ...)
- Definition einer geeigneten **kumulativen Kenngröße** für die Gütekurve (AUC, mittlere Ausbeute)

Fehlermaße mit gleitendem Arbeitspunkt

Minimale und balancierte Fehlerraten



Equal Error Rate

Balancepunkt von Fehlrückweisungsrate und Fehlakzeptanzrate

$$\text{EER} \stackrel{\text{def}}{=} \text{FNR}_{\theta^*}, \quad \theta^* = \underset{\theta}{\operatorname{argmin}} |\text{FNR}_{\theta} - \text{FPR}_{\theta}|$$

Break-Even Point

Schnittpunkt der **recall**-Kurve mit der **precision**-Kurve; impliziert $N_{t.}^{\theta} = N_{.t}^{\theta}$

$$\text{BEP} \stackrel{\text{def}}{=} \text{REC}_{\theta^*}, \quad \theta^* = \underset{\theta}{\operatorname{argmin}} |\text{REC}_{\theta} - \text{PRE}_{\theta}|$$

Ein-Klassen-Szenarien

Zwei-Klassen-Szenarien

Mehr-Klassen-Szenarien

Ensemblemethoden — Bagging & Boosting

Kombinatorische Merkmalauswahl

Inkrementelles Lernen

Transduktion

AUC — Area Under (ROC) Curve

Wilcoxon-Mann-Whitney Statistik

Definition

Den Flächeninhalt unter der Receiver-Operator-Kurve

$$\text{ROC} : \begin{cases} [0, 1] & \rightarrow [0, 1] \\ \text{TNR}(\theta) & \mapsto \text{TPR}(\theta) \end{cases}$$

bezeichnen wir als **AUC-Wert**.

Bemerkungen

- Empirisch werden einfach alle TPR (Recallwerte) gemittelt, die sich an den Schwellen θ der Negativmuster ergeben.
- Ähnlich definiert ist die AVP ('average precision') des Information Retrieval.
- Der AUC-Wert ist äquivalent zur **Wilcoxon-Rangsummenstatistik**

$$\frac{1}{N_t \cdot N_f} \cdot \sum_{x \in \omega_1} \sum_{y \in \omega_0} \mathbb{I}_{u(x) > u(y)},$$

welche die Wahrscheinlichkeit $P(\mathbb{U}_1 > \mathbb{U}_0)$ dafür schätzt, daß die Prüfgröße eines Positivmusters größer ist als die Prüfgröße eines Negativmusters.

mehr Information

Zwei-Klassen-Diskriminanten

Mächtige Formalismen zur Zwei-Klassen-Unterscheidung

- Supportvektormaschinen
- Logistische Regression
- Mehrschichtenperzeptron
- Polynomklassifikator (Importvektormaschine)

Nutzung von Zwei-Klassen-Diskriminanten

- Kontrastive Dichteschätzung
- Mehr-Klassen-Entscheidung \blacklozenge „Jeder gegen alle“
- Mehr-Klassen-Entscheidung \blacklozenge „Jeder gegen jeden“
- Mehr-Klassen-Entscheidung \blacklozenge Fehlerkorrigierende Codes

Logarithmierte Gewinnquote

Zwei Klassen $\{\Omega_1, \Omega_0\}$ — eine Prüfgröße $\log \text{odds}(x)$

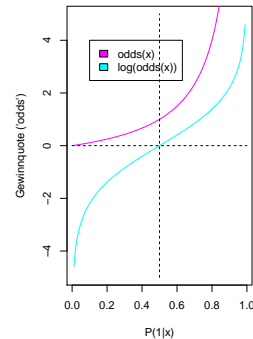
$$\kappa^*(x) = \operatorname{argmax}_{\lambda \in \{1,0\}} \frac{p_\lambda \cdot f_\lambda(x)}{p_1 f_1(x) + p_0 f_0(x)} = \begin{cases} 1 & p_1 \cdot f_1(x) \geq p_0 \cdot f_0(x) \\ 0 & p_1 \cdot f_1(x) < p_0 \cdot f_0(x) \end{cases}$$

Definition

Die a posteriori Klassenverteilung kann in der Zweiklassenwelt in einer einzigen Größe, der **Gewinnquote** oder den **Odds**

$$\text{odds}(x) \stackrel{\text{def}}{=} \frac{P(1|x)}{P(0|x)} = \frac{p_1 \cdot f_1(x)}{p_0 \cdot f_0(x)}$$

subsumiert werden.



Gewinnquote als Dichteverhältnis

Lemma

Aus der Chancenfunktion $\text{odds}(\cdot)$ sind die a posteriori Klassenwahrscheinlichkeiten via

$$P(1|x) = \frac{\text{odds}(x)}{1 + \text{odds}(x)} \quad \text{bzw.} \quad P(0|x) = \frac{1}{1 + \text{odds}(x)}$$

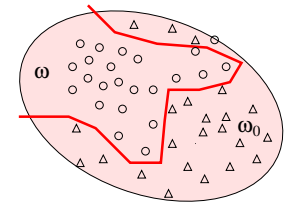
exakt reproduzierbar.

Zwischen den beiden Klassenverteilungen f_1 und f_0 für „positive“ ($\kappa = 1$) und „negative“ ($\kappa = 0$) Muster besteht der wichtige Zusammenhang

$$f_1(x) = \text{odds}(x) \cdot \frac{p_0}{p_1} \cdot f_0(x) \quad (\forall x \in \Omega).$$

Bemerkung

Ist also die Negativverteilung f_0 bekannt sowie die Klassendiskriminante $\text{odds} : \Omega \rightarrow \mathbb{R}^+$, so ist auch die Positivverteilung f_1 schon festgelegt.



Kontrastive Verteilungsdichteschätzung

(Algorithmus)

1 KONTRASTMODELL

Wähle eine Verteilungsfunktion $f_0 : \Omega \rightarrow \mathbb{R}$.

2 AUSWÜRFELN DER KONTRASTPROBE

Erzeuge eine Beispieldatensammlung $\omega_0 = \{y_1, \dots, y_S\}$ für $f_0(\cdot)$.

3 DISKRIMINANTE LERNEN

Lerne unter Verwendung der Positiv- und Negativmuster in ω bzw. ω_0 einen binären Klassifikator mit der Diskriminanzfunktion $\widehat{\text{odds}}(\cdot)$ an.

4 DICHTESCHÄTZUNG

Reproduziere die gesuchte Verteilung näherungsweise durch

$$\hat{f}(x) \stackrel{\text{def}}{=} \widehat{\text{odds}}(x) \cdot \frac{|\omega_0|}{|\omega|} \cdot f_0(x).$$

(sumf0h0f1A)

Bemerkung

Statt Setzen der Kontrastdichte und Auswürfeln der Kontrastdaten können auch vorliegende Daten als Kontrast verwendet werden; die Parameter der Kontrastdichte werden daraus geschätzt.

Klassifikatoren für $K > 2$ Klassen

Standardtechniken: $\tilde{\omega}_\kappa = \bigcup_\lambda \omega_\lambda$ oder $\tilde{\omega}_\kappa = \bigcup_{\lambda \neq \kappa} \omega_\lambda$

Klassenweise Kontrastdiskriminanten

Hypothetische Kontrastverteilungen h_1, \dots, h_K :

$$\text{odds}_\lambda(x) = \frac{q_\lambda \cdot f_\lambda(x)}{(1 - q_\lambda) \cdot h_\lambda(x)}, \quad \lambda \in \{1, \dots, K\}$$

(Kontrastdaten im Anzahlverhältnis $(1 - q_\kappa)$ zu q_κ)

Klassenweise Kontrastschätzung

$$\hat{f}_\lambda(x) \stackrel{\text{def}}{=} \widehat{\text{odds}}_\lambda(x) \cdot \frac{1 - q_\lambda}{q_\lambda} \cdot h_\lambda(x), \quad \lambda \in \{1, \dots, K\}$$

Reproduzierte Bayesregel

$$\hat{P}(\lambda|x) \propto T_\lambda \cdot \widehat{\text{odds}}_\lambda(x) \cdot \frac{S_\lambda}{T_\lambda} \cdot h_\lambda(x) \propto S_\lambda \cdot \widehat{\text{odds}}_\lambda(x)$$

(Die letzte Proportionalität gilt nur, wenn alle Kontrastverteilungen identisch sind!)

Kreuzklassifizierende Diskriminanten

Jeder gegen jeden — aber welches ist die wahrscheinlichste Klasse?

Ideale Kreuzdiskriminanten

$$\text{odds}_{\kappa\lambda}(x) = \frac{q_{\kappa} \cdot f_{\kappa}(x)}{q_{\lambda} \cdot f_{\lambda}(x)} \quad (\forall \lambda, \kappa, x)$$

Für ein festes $x \in \Omega$ sei $\mathbf{R}(x)$ die quadratische Matrix mit Einträgen

$$r_{\kappa\lambda} \stackrel{\text{def}}{=} \log \text{odds}_{\kappa\lambda}(x) .$$

Fakt

Die Matrix \mathbf{R} hat einen Rang ≤ 2 und besitzt wegen

$$r_{\kappa\lambda} = \underbrace{\log(q_{\kappa} \cdot f_{\kappa}(x))}_{=: u_{\kappa}} - \underbrace{\log(q_{\lambda} \cdot f_{\lambda}(x))}_{=: u_{\lambda}} ,$$

die Darstellung $(\mathbf{u}\mathbf{1}^{\top} - \mathbf{1}\mathbf{u}^{\top})$ mit den K -Klassen-Diskriminanten $u_{\lambda}(\cdot)$, $\lambda = 1, \dots, K$.

Ein-Klassen-Szenarien

Zwei-Klassen-Szenarien

Mehr-Klassen-Szenarien

Ensemblemethoden — Bagging & Boosting

Kombinatorische Merkmalauswahl

Inkrementelles Lernen

Transduktion

Kreuzklassifizierende Diskriminanten

Näherungsweise Rekonstruktion der originalen Bayesprüfgrößen

Minimierung des mittleren quadratischen Fehlers

$$\varepsilon(\mathbf{u}) \stackrel{\text{def}}{=} \left\| \hat{\mathbf{R}} - \mathbf{u}\mathbf{1}^{\top} + \mathbf{1}\mathbf{u}^{\top} \right\|_{\mathcal{F}}^2$$

Eindeutige Lösung

$$\hat{u}_{\kappa} = \frac{1}{K} \sum_{\lambda=1}^K \hat{r}_{\kappa\lambda} = \frac{1}{K} \sum_{\lambda=1}^K \log \widehat{\text{odds}}_{\kappa\lambda}(x)$$

für alle Klassen $\kappa \in \{1, \dots, K\}$

Reproduzierte Bayesregel

$$\kappa^*(x) = \underset{\kappa}{\operatorname{argmax}} \prod_{\lambda=1}^K \widehat{\text{odds}}_{\kappa\lambda}(x)$$

Ensemblemethoden

Klassifikationsentscheidungen — im Expertenteam gefällt

Welche Wissensquellen werden kombiniert?

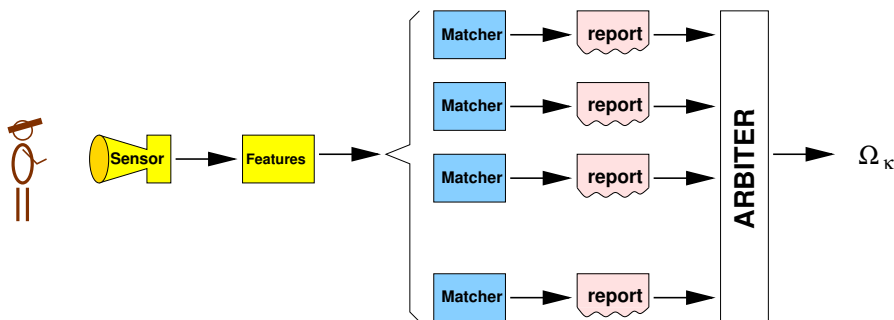
- **Orthogonale Sensoren**
- **Wiederholte Messungen**
- **Konkurrierende Klassifikatortypen**
- **Streuende Parameterschätzwerte**
- **Spezialisierte Entscheider**

Wie werden die Quellen kombiniert?

- **Sensordatenfusion** (Mittel- und Produktbildung)
- **Interaktionsmodelle** (Zeitreihen; Konsortiale Klassifikation)
- **Prüfgrößenkopplung** (Metaklassifikation)
- **Votierungsverfahren** (Mehrheit und Präferenz)

Komitee-Maschinen treffen Mehrheitsentscheidung

Ein Sensor · eine Messung · M unabhängige Experten



Society of Experts

Konkurrierende
Klassifikationsverfahren
oder konkurrierende
Modellkapazitäten

Bagging

Bootstrap Aggregation
Konkurrierende
ausgewürfelte
Lerndatensammlungen

Randomisierung

Konkurrierende
Parametersätze
(zufällige Startwerte
iterativer
Lernverfahren)

Bootstrappingverfahren zum Lernen und Testen

Simuliertes Auswürfeln von lerndatenbezogenen a posteriori-Verteilungen

Kreuzvalidierung

- Anlernen des Klassifikators mit $\tilde{\omega}$
- Testen des Klassifikators mit ω

$$\hat{p}_\varepsilon(\omega, \tilde{\omega}) \stackrel{\text{def}}{=} (1 - \frac{1}{e}) \cdot p_\varepsilon(\tilde{\omega} \mid \delta_{\tilde{\omega}}) + \frac{1}{e} \cdot p_\varepsilon(\omega \setminus \tilde{\omega} \mid \delta_{\tilde{\omega}})$$

Fehlerratenstreuung

Bestimme Mittelwert und Varianz der Schätzungen $\hat{p}_\varepsilon(\omega, \tilde{\omega}^{(m)})$, $m = 1..M$.

Parameterschätzung

Mittlere Parameterschätzwerte $\theta^{(m)}$ einer Reihe von Bootstrappproben $\tilde{\omega}^{(m)}$.

Bootstrap Aggregation

Mittlere Entscheidungen $\delta^{(m)}(x)$ einer Reihe von Bootstrappproben $\tilde{\omega}^{(m)}$.

Bootstrap-Stichproben

T -mal zufälliges Ziehen aus ω mit Zurücklegen

Lemma

Es sei $\omega = \{\mathbf{x}_1, \dots, \mathbf{x}_T\} \subset \Omega$ ein Datensatz und

$$\tilde{\omega} = \{\mathbf{x}_{\tau(1)}, \mathbf{x}_{\tau(2)}, \mathbf{x}_{\tau(3)}, \dots, \mathbf{x}_{\tau(T)}, \}, \quad \tau: [1..T] \rightarrow [1..T]$$

eine **Bootstrapprobe** der Größe T daraus.

1. Die Wahrscheinlichkeit dafür, beim Bootstrap übergangen zu werden, beträgt asymptotisch

$$\lim_{T \rightarrow \infty} P(\mathbf{x}_t \notin \tilde{\omega}) = \lim_{T \rightarrow \infty} (1 - \frac{1}{T})^T = \frac{1}{e} \approx 0.368.$$

2. Ist $s: \Omega^T \rightarrow \mathbb{R}$ eine beliebige Stichprobenstatistik, so ist ihr Erwartungswert für Originalprobe und Bootstrapprobe gleich:

$$\mathcal{E}_\omega[s(\mathbb{X}_1, \dots, \mathbb{X}_T)] = \mathcal{E}_{\tilde{\omega}}[s(\mathbb{X}_{\tau(1)}, \dots, \mathbb{X}_{\tau(T)})]$$

Boosting — Konsultation komplementärer Klassifikatoren

Ein Klassifikortyp · Parameter nach sequentieller Probenselektion

(Algorithmus)

- 1 INITIALISIERUNG
Setze konstante Gewichte $\beta_t^{(1)} = 1/T$ und für alle $m = 1, 2, \dots$:
- 2 ENTSCHEIDUNGSREGEL LERNEN
Lerne neuen Klassifikator $\delta^{(m)} = \delta(\omega, \beta^{(m)})$.
- 3 REKLASSIFIKATIONSTEST
Berechne eine neue (gewichtete) Fehlerrate
 $\varepsilon_m := p_\varepsilon(\omega, \beta^{(m)} \mid \delta^{(m)})$.
- 4 TERMINIERUNG (falls $\varepsilon_m \notin (0, \frac{1}{2})$)
- 5 ERFOLGSGEWICHTE AKTUALISIEREN

$$\beta_t^{(m+1)} \propto \begin{cases} \beta_t^{(m)} & \mathbf{x}_t \text{ falsch klassifiziert} \\ \beta_t^{(m)} \cdot \varepsilon_m / (1 - \varepsilon_m) & \mathbf{x}_t \text{ korrekt klassifiziert} \end{cases}$$

Nächster AdaBoost.M1-Iterationsschritt \rightsquigarrow 2

(summiert)

Boosting — aufmerksamkeitsgesteuertes Lernen

Sequentielles Lernen schwer klassifizierbarer Muster durch „Nachsitzen“

Klassifikationsphase

Expertenentscheidungen gewichtet summieren:

$$u_{\kappa}(\mathbf{x}) = \sum_{m=1}^M -\log \frac{\varepsilon_m}{1 - \varepsilon_m} \cdot \delta_{\kappa}^{(m)}(\mathbf{x})$$

Bemerkungen

1. AdaBoost endet wenn $p_{\varepsilon}^{(m)} \geq 1/2$ — kein Problem für $K = 2$.
2. **Gewichtetes Lernen** von Klassifikatoren:
 - Integration in das Parameterlernverfahren (ML, QM)
 - Simulation durch gewichtetes Auswürfeln einer Bootstrap-Lernprobe.
 - Vergrößerte Simulation durch Probe mit $C \cdot T \cdot \beta_t^{(m)}$ vielen \mathbf{x}_t -Kopien
3. Je geringer die Rate $\varepsilon_m = p_{\varepsilon}^{(m)}$, desto größer das Abstimmungsgewicht.

Ein-Klassen-Szenarien

Zwei-Klassen-Szenarien

Mehr-Klassen-Szenarien

Ensemblemethoden — Bagging & Boosting

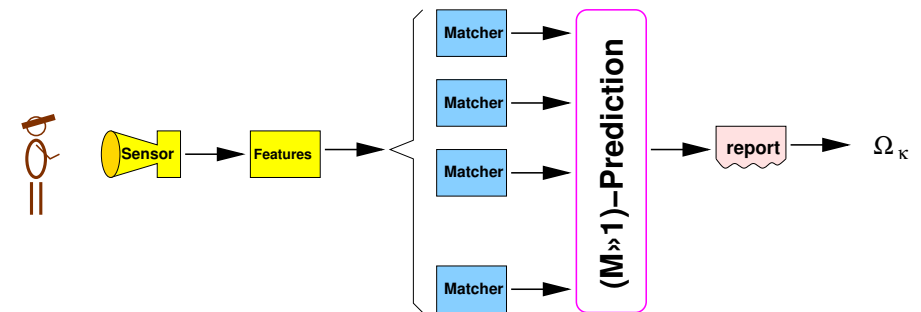
Kombinatorische Merkmalauswahl

Inkrementelles Lernen

Transduktion

Prüfgrößeninterpolation

Frühe Integration quantitativer Expertenurteile



Lineare Interpolation

Die **Mischungskoeffizienten** $\{\pi_1, \dots, \pi_M\}$ des Vorhersagemodells

$$\mathbf{u}(\mathbf{x}) = \sum_{m=1}^M \pi_m \cdot \mathbf{u}^{(m)}(\mathbf{x}) \quad \approx \quad \delta^*(\mathbf{x}), \quad \sum_{m=1}^M \pi_m = 1$$

werden mit einer Variante des EM-Algorithmus bestimmt (*held-out*-Daten!)

Merkmalgewinnung

Komplexer Arbeitsprozeß beim Entwurf eines ME-Systems

1. **Domänenwissen** ad hoc Merkmalsatz entwerfen
2. **Vergleichbarkeit** Merkmale normieren
3. **Interaktionen** Produkte/Konjunktionen dazu
4. **Hilfsvariablen** Linearkombinationen/Disjunktionen dazu
5. **Variablenreihung** Nützlichkeit der Merkmale quantifizieren
6. **Ausreißerdetektion** Kriterium sind die höchstrangigen Merkmale
7. **Teilmengenauswahl** schnelles, robustes Kriterium
8. **Wettbewerb austragen** Klassifikations- und Selektionstechniken
9. **Endauswertung** stabile Lösung mit Bootstrap/Kreuzvalidierung

Merkmalauswahl

Welche Teilmenge des Merkmalsatzes $\{x_1, \dots, x_D\}$ wird verwendet?

Entwurfsziele

geringere Fehlerrate
schnellere Verarbeitung
Strukturaufdeckung

Kombinatorische Explosion

Es gibt insgesamt 2^D
viele Teilmengen
 $\mathcal{S} \subset \{x_1, \dots, x_D\}$.

Standalone-Gütekriterien

Relevanz eines Merkmals

Redundanz eines Merkmals

Ranker

1. Bewerte alle Merkmale x_d nach Relevanz
2. Verwende nur die M Bestbewerteten

Filter

1. Bewerte Merkmale nach Relevanz
2. Analysiere Abhängigkeiten untereinander
3. Treffe eine kompatible Auswahl

Wrapper

1. Wähle Basisklassifikationsverfahren
2. Bewerte Teilmengen \mathcal{S} durch Evaluation
3. Suche die höchstperformante Teilmenge

Ranker — Variablenreihung nach Relevanz

Wie groß ist der individuelle Beitrag von x_d zur Klassenentscheidung?

Pearsons Korrelationskoeffizient

Abhängigkeitsmaß für stetige Zufallsvariable (schwierig für $K \neq 2$)

$$r(d) = \frac{\text{Cov}[\mathbb{X}_d, Y]}{\sqrt{\text{Var}[\mathbb{X}_d] \cdot \text{Var}[Y]}}$$

Transinformation

Divergenz zwischen gemeinsamer Verteilung und Produktverteilung

$$r(d) = \mathcal{D}(f_{\mathbb{X}_d Y} \| f_{\mathbb{X}_d} f_Y) = \int \int f_{\mathbb{X}_d Y}(x_d, y) \cdot \log \frac{f_{\mathbb{X}_d Y}(x_d, y)}{f_{\mathbb{X}_d}(x_d) \cdot f_Y(y)} dx_d dy$$

Ein-Merkmal-Klassifikatoren

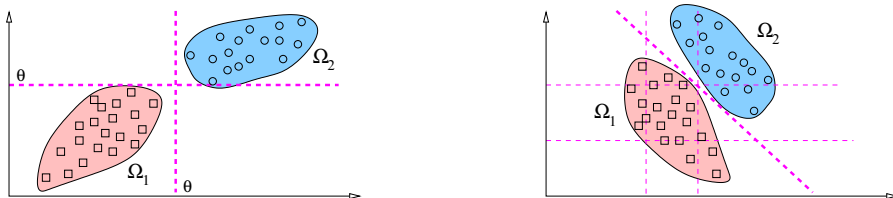
Betreibe ein Standardklassifikationssystem mit Mustern $\mathbf{x} = (x_d) \in \mathbb{R}^1$ und ermittle z.B. die Held-out-Erkennungsrate $r(d) = 1 - p_\epsilon(\delta(\omega^{(d)}))$.

Ranker — Variablenreihung nach Relevanz

$$\mathcal{S}_\theta \stackrel{\text{def}}{=} \{x_d \mid r(d) \geq \theta\}, \quad \theta \in \mathbb{R}_+$$

PRO

Extrem schnelles und einfaches Verfahren



CONTRA

Die Welt ist komplizierter, als der Ranker denkt:

- Ein **hoch relevantes** Merkmal ist u.U. **völlig nutzlos**, wenn es den bereits selektierten keine Musterinformation hinzufügt.
- Zwei oder mehrere **nicht relevante** Merkmale besitzen u.U. **hohen Nutzwert**, wenn sie der Auswahl **gemeinsam** hinzugefügt werden.

FCBF — Fast Correlation-Based Filtering

Relevanz eines Merkmals · Redundanz eines Merkmals

Relevanz von \mathbb{X}_i

Symmetrische Unsicherheit

$\tilde{\mathfrak{S}}(\mathbb{X}_i, Y)$ zwischen Merkmal und Klassenvariable, wobei

$$\tilde{\mathfrak{S}}(\mathbb{X}, Y) = \frac{2 \cdot \mathfrak{S}(\mathbb{X}, Y)}{\mathcal{H}(\mathbb{X}) + \mathcal{H}(Y)}$$

mit der **Transinformation**

$$\mathfrak{S}(\mathbb{X}, Y) \stackrel{\text{def}}{=} \mathcal{H}(\mathbb{X}) + \mathcal{H}(Y) - \mathcal{H}(\mathbb{X}, Y)$$

Redundanz von \mathbb{X}_i

Für zwei relevante ($\text{SU} \geq \theta$)

Merkmale $\mathbb{X}_i, \mathbb{X}_j$ heie \mathbb{X}_j **ebenbürtig** zu \mathbb{X}_i , falls

$$\tilde{\mathfrak{S}}(\mathbb{X}_i, \mathbb{X}_j) \geq \tilde{\mathfrak{S}}(\mathbb{X}_i, Y)$$

gilt. Merkmal \mathbb{X}_i ist **redundant**, wenn es ein ebenbürtigen Wettbewerber \mathbb{X}_j gibt mit:

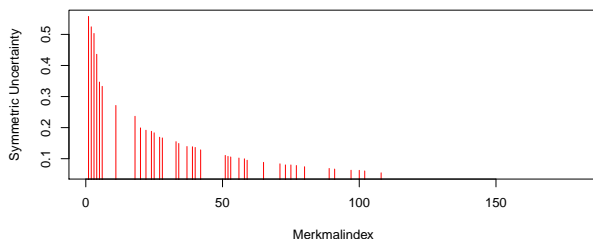
$$\tilde{\mathfrak{S}}(\mathbb{X}_j, Y) \geq \tilde{\mathfrak{S}}(\mathbb{X}_i, Y)$$

Problem

Die Redundanzeigenschaft ist nur **relativ** zu einem willkürlichen Abschnitt $\mathcal{S}_\theta = \{x_i \mid \tilde{\mathfrak{S}}(\mathbb{X}_i, Y) \geq \theta\}$ relevanter Merkmale definiert!

FCBF — Anwendungsbeispiel

DNA-Datensatz · 3 Klassen · 180 %₁-Merkmale · 3186 Muster



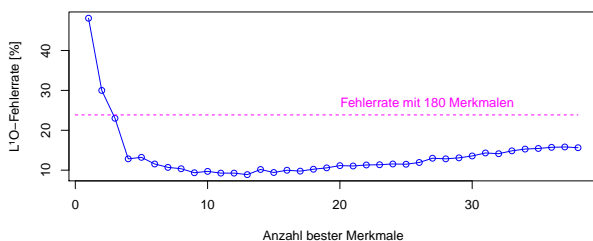
FCBF wählt
maximal ($\delta = 0$)
40 von 180
Merkmalen aus.

Fehlerraten

23.8%	180
09.3%	9
12.8%	4
23.0%	3

Erfolgsbilanz

Faktor 20 Zeit
Faktor 2.5 Fehler



Bemerkung

$180^3 = 5.8 \cdot 10^6$
Kombinationen,
um
(x₉₀, x₈₅, x₉₃) zu
finden!

Wrapper-Methoden

Gütebewertung durch Klassifikationstest

Mit Hilfe eines *robusten* und *CV-fähigen* Klassifikationsverfahrens wird

$$J(S) \stackrel{\text{def}}{=} 1 - \hat{p}_e^{\text{CV}}(\delta(\omega^S)) , \quad \omega^S = \{\mathbf{x}^S \mid \mathbf{x} \in \omega\}$$

definiert.

Kombinatorische Optimierung

Wähle die bestbewertete Merkmalteilmenge

$$S^* = \underset{S \subseteq \mathcal{F}}{\operatorname{argmax}} J(S) , \quad \mathcal{F} \stackrel{\text{def}}{=} \{x_1, \dots, x_D\}$$

Suchverfahren

Vollständige Suche

Gierige Suche

Zulässige Suche

Approximative Suche

Aufwand $O(2^D)$
Vorwärts-/Rückwärtsauswahl
A*-Algorithmus, Branch&Bound
Prophet statt Orakel

SFS — Sequential Forward Selection

Gierige Bottom-up Auswahlstrategie

1 INITIALISIERUNG

Setze $S \leftarrow \emptyset$.

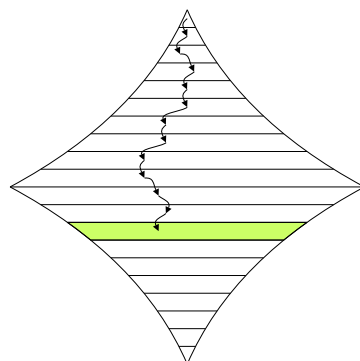
2 AUSWAHL AUFWÄRTS

Bestimme das beste zusätzliche
Merkmal

$$x_d = \operatorname{argmax} \{J(S, x_d) \mid x_d \notin S\}$$

3 TERMINIERUNG

Wenn $J(S, x_d) \leq J(S) \rightsquigarrow$ Ende.
Sonst $S \leftarrow S \cup \{x_d\}$ und \rightsquigarrow 2.



Bemerkung

SFS trifft voreilige Entscheidungen
(Horizont=1) und verfehlt i.a. die
Optimallösung.

$$S^{(1)} \subset S^{(2)} \subset \dots \subset S^{(i)} \subset \dots$$

↗ Redundanz · ↘ Kombination

SBE — Sequential Backward Elimination

Gierige Top-down Auswahlstrategie

1 INITIALISIERUNG

Setze $S \leftarrow \mathcal{F} = \{x_1, \dots, x_D\}$.

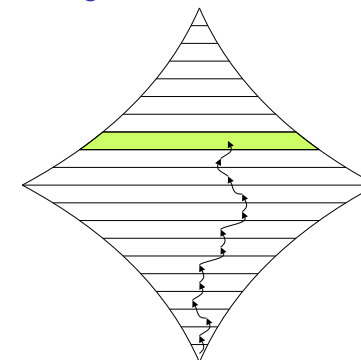
2 AUSWAHL ABWÄRTS

Bestimme das verzichtbarste
Merkmal

$$x_d = \operatorname{argmax} \{J(S \setminus x_d) \mid x_d \in S\}$$

3 TERMINIERUNG

Wenn $J(S \setminus x_d) \leq J(S) \rightsquigarrow$ Ende.
Sonst $S \leftarrow S \setminus \{x_d\}$ und \rightsquigarrow 2.



Bemerkung

SBE aufwändiger als SFS:
· Start mit umfangreicheren S
· Längerer Weg zum Ziel

$$O(n \cdot (D - n)^2) \gg O(n^2 \cdot (D - n))$$

↗ Redundanz · ↘ Kombination

PTA(k, ℓ) — Plus k Take Away ℓ

Gierige bidirektionale Auswahlstrategie

- 1 INITIALISIERUNG · $S' \leftarrow S \leftarrow \emptyset$
- 2 MEHRFACHAUSWAHL AUFWÄRTS
Bestimme k -mal nacheinander das beste zusätzliche Merkmal

$$x_d = \operatorname{argmax} \{J(S, x_d) \mid x_d \notin S\}$$

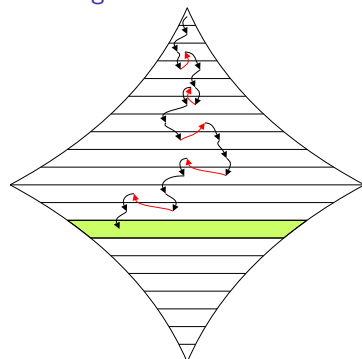
und setze $S' \leftarrow S' \cup \{x_d\}$.

- 3 MEHRFACHAUSWAHL ABWÄRTS
Bestimme ℓ -mal nacheinander das verzichtbarste Merkmal

$$x_d = \operatorname{argmax} \{J(S \setminus x_d) \mid x_d \in S\}$$

und setze $S' \leftarrow S' \setminus \{x_d\}$.

- 4 TERMINIERUNG
Wenn $J(S') \leq J(S) \rightsquigarrow$ Ende.
Sonst $S \leftarrow S'$ und \rightsquigarrow 2.



Bemerkung

Redundant gewordene Merkmale können jetzt wieder eliminiert werden.

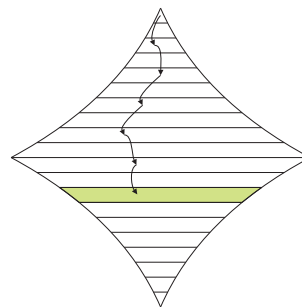
Mehraufwand gegenüber SFS:

$$(k+\ell) / (k-\ell)$$

⬇ Redundanz · ⬆ Kombination

Verallgemeinerte sequentielle Verfahren

Gierige blockweise Auswahlstrategien · Polynomialer Aufwand $O(n^{m+1} \cdot (D - n))$



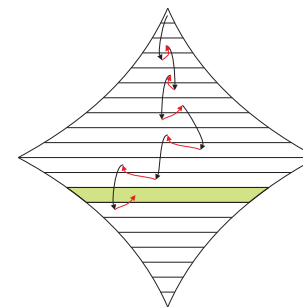
GSFS(m)

Je Aufwärtsschritt wird die performanteste Erweiterungsmenge

$$S_{\oplus}^* = \{x_{d_1}, x_{d_2}, \dots, x_{d_m}\}$$

ermittelt und ggf. hinzugefügt.

⬇ Redundanz · ⬆ Kombination



GPTA(k, ℓ)

In den Akkumulations- und Eliminationsschritten werden optimale Auswahlen S_{\oplus}^* , S_{\ominus}^* mit

$$|S_{\oplus}^*| = k \quad \text{bzw.} \quad |S_{\ominus}^*| = \ell$$

getroffen.

SFFS — Sequential Forward Floating Search

Pulsierende bidirektionale Suche (Pudil, 1994)

- 1 INITIALISIERUNG
Setze $S \leftarrow \emptyset$, $n = 0$, $\iota_0 = J(\emptyset)$.

- 2 AUSWAHL AUFWÄRTS
Bestimme Bestmerkmal

$$x_d = \operatorname{argmax} \{J(S, x_d) \mid x_d \notin S\}$$

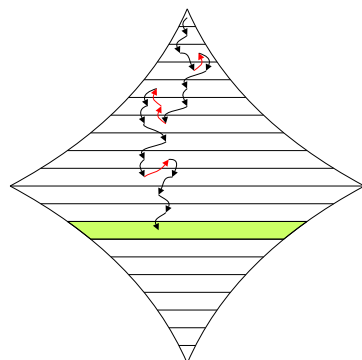
und setze $S \leftarrow S \cup \{x_d\}$, $n++$,
 $\iota_n \leftarrow J(S)$.

- 3 AUSWAHL ABWÄRTS
Bestimme Schlechtmerkmal

$$x_d = \operatorname{argmax} \{J(S \setminus x_d) \mid x_d \in S\}$$

Wenn $J(S \setminus x_d) \leq \iota_{n-1}$, dann \rightsquigarrow 2.
Sonst $S \leftarrow S \setminus \{x_d\}$, $n--$, $\iota_n \leftarrow J(S)$, \rightsquigarrow 3.

- 4 TERMINIERUNG
Zielkardinalität n^* oder ι_n -Schranke.



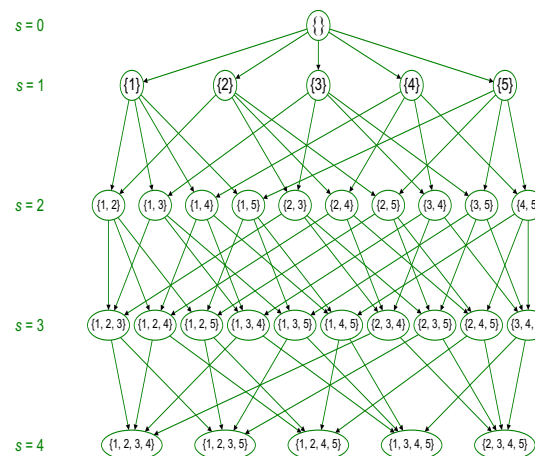
Bemerkung

Effizient und effektiv.

Liefert näherungsweise Bestlösungen $|S_n| = n$ im Kern des traversierten Bereichs
 $0 \leq n \leq n_{\max}$.

A*-Algorithmus

Suche im Nachbarschaftsgraphen aller $S \subset \mathcal{F}$



Merkmalmengen und Nachbarschaftsrelation

Pro

A* findet die optimale Lösung.
A* findet die n -besten Lösungen.

Contra

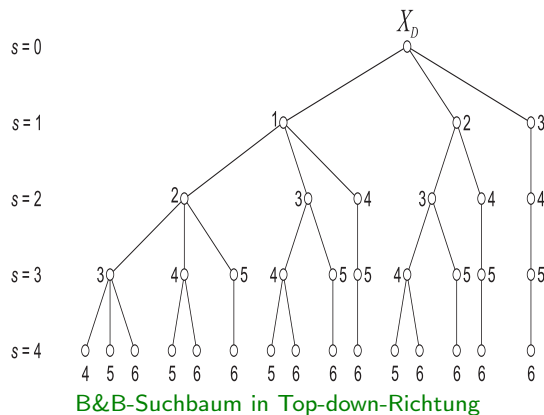
Je nach Schärfe der heuristischen Restschätzung eine aufwändige Suche.

Ausweg

Adaptives Lernen eines **Propheten** für unausgewertete $J(S)$.

Branch&Bound-Algorithmus

Trivialvariante des A*-Algorithmus · Restschätzung ist $\hat{h}(\cdot) \equiv 0$



Zulässigkeit

Nur garantiert, wenn die Gütefunktion **monoton** ist:

$$S_1 \subseteq S_2 \Rightarrow J(S_1) \leq J(S_2)$$

Effizienz

B&B ist nur **Top-down** anwendbar.



Hoher $J(\cdot)$ -Aufwand, späte Resultate.

Dualisierte Quadratmittelaufgabe $\|y - Xa\|^2 \xrightarrow{!} \text{MIN}$

Speicheraufwand $O(T^2)$ und Rechenaufwand $O(T^3)$

Duale Lösungsdarstellung

als Linearkombination der Objektvektoren:

$$a = X^T b = \sum_{t=1}^T b_t \cdot x_t, \quad b \in \mathbb{R}^T$$

Duale Regressionsfehlerformel

in Abhängigkeit vom Vektor b der Lösungskoeffizienten:

$$\varepsilon(b) = \|y - X \cdot X^T b\|^2 \xrightarrow{!} \text{MIN}$$

Duale Gauß'sche Normalgleichungen

Lineares Gleichungssystem (Dimension $T \times T$) mit Gram'scher Matrix:

$$G^2 \cdot b = G \cdot y, \quad G = X \cdot X^T, \quad G_{st} = \langle x_s, x_t \rangle$$

Kombinatorische Regression

Ausgleichsrechnung für Funktionen $f : \mathcal{P}\mathcal{F} \rightarrow \mathbb{R}$ bzw. $f : 2^{\{1,2,\dots,D\}} \rightarrow \mathbb{R}$

Aufgabenstellung

Güteschätzung $\tilde{J}(S)$ für Merkmaltelmengen $S \subseteq \mathcal{F} = \{x_1, \dots, x_D\}$

Termexpansion

Binärattribute: Merkmal- k -Subsets

$$\phi : \mathcal{P}\mathcal{F} \rightarrow \{0, 1\}^{\binom{D}{k}}$$

mit

$$\phi_U(S) = \begin{cases} 1 & U \subseteq S \\ 0 & U \not\subseteq S \end{cases}$$

Quadratmittelaufgabe

der Dimension $\binom{D}{k}$ (bzw. 2^D)

Duale Quadratmittelaufgabe

Gramsche n^2 -Matrix mit Einträgen

$$K_k(S_i, S_j) = \langle \phi(S_i), \phi(S_j) \rangle$$

Kombinat. Kernoperator

$$\begin{aligned} K_k(S, S') &= \sum_{|U|=k} \phi_U(S) \cdot \phi_U(S') \\ &= \sum_{|U|=k} \phi_U(S \cap S') \\ &= \binom{|S \cap S'|}{k} \quad \text{bzw. } 0 \end{aligned}$$

Dto. für die kumulativen Kerne

$$K^{\leq k} := K_0 + K_1 + K_2 + \dots + K_k$$

Ein-Klassen-Szenarien

Zwei-Klassen-Szenarien

Mehr-Klassen-Szenarien

Ensemblemethoden — Bagging & Boosting

Kombinatorische Merkmalauswahl

Inkrementelles Lernen

Transduktion

Inkrementelles Lernen

$$\{\mathbf{x}_1, \dots, \mathbf{x}_T\} \rightarrow \{\mathbf{x}_1, \dots, \mathbf{x}_{T+1}\}$$

Gegeben:

Ein Klassifikator mit den Parametern $\theta^{(T)}$ auf Grundlage der Lerndaten $\omega^{(T)} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$

Gesucht:

Ein schnelles Berechnungsverfahren für die Parameter $\theta^{(T+1)}$ zur erweiterten Probe $\omega^{(T+1)} = \omega^{(T)} \cup \{\mathbf{x}_{T+1}\}$

Inverse Kovarianzmatrix:

Aktualisierungsformeln für NVK

(Aufzählung der \mathbf{x}_t für ein ausgewähltes Ω_κ)

Mittelwertvektor:

$$\mu^{(T+1)} = \frac{T}{T+1} \cdot \mu^{(T)} + \frac{1}{T+1} \cdot \mathbf{x}_{T+1}$$

Kovarianzmatrix:

$$R^{(T+1)} = \frac{T}{T+1} \cdot R^{(T)} + \frac{\mathbf{x}_{T+1} \cdot \mathbf{x}_{T+1}^\top}{T+1}$$

(via Sherman-Morrison-Woodbury-Formel)

$$(\mathbf{A} - \mathbf{u}\mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} + \frac{1}{1 - \mathbf{v}^\top \mathbf{A}^{-1} \mathbf{u}} \cdot \mathbf{A}^{-1} \mathbf{u} \mathbf{v}^\top \mathbf{A}^{-1}$$

Gewichtetes inkrementelles Lernen

Gegeben:

Ein Klassifikator mit den Parametern $\theta^{(T)}$ auf Grundlage der Lerndaten

$$\omega^{(T)} = \{(\mathbf{x}_t, \beta_t)\}_{t=1}^T$$

Gesucht:

Ein schnelles Berechnungsverfahren für die Parameter $\theta^{(T+1)}$ zur erweiterten Probe

$$\omega^{(T+1)} = \{(\mathbf{x}_t, \tilde{\beta}_t)\}_{t=1}^{T+1}$$

Schnelle Aktualisierungsformeln

für „selbstähnliche“ Gewichtstrukturen

$$\tilde{\beta}_t = \begin{cases} \alpha_t \cdot \beta_t & t \leq T \\ 1 - \alpha_t & t = T+1 \end{cases}, \quad \alpha_t \in (0, 1)$$

Gewichteter Mittelwertvektor

$$\mu^{(T+1)} = \alpha_t \cdot \mu^{(T)} + (1 - \alpha_t) \cdot \mathbf{x}_{T+1}$$

Beispielgewichtungen

Ungewichtetes Mittel:

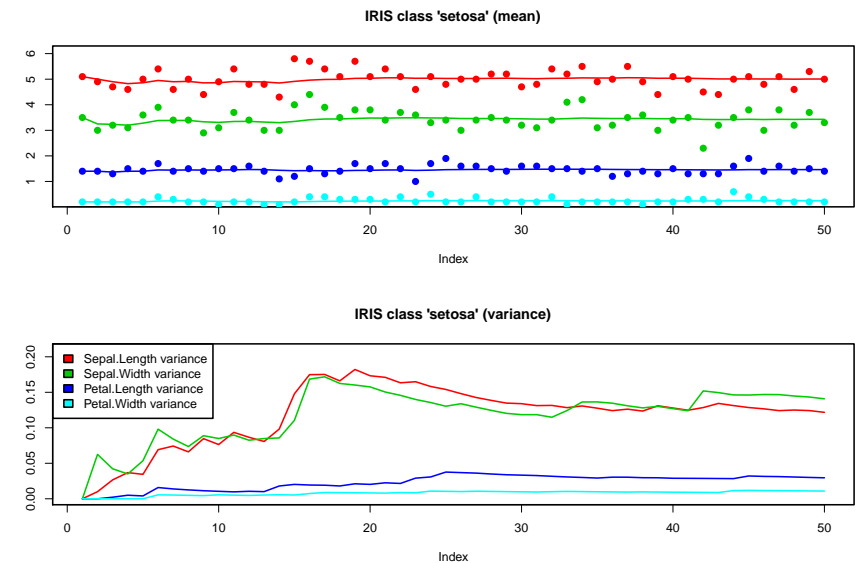
$$\beta_t = 1/T, \quad \tilde{\beta}_t = 1/(T+1), \quad \alpha_t \equiv T/(T+1)$$

Exponentielles Vergessen:

$$\beta_t = (1-\alpha) \cdot \alpha^{T-t}, \quad \tilde{\beta}_t = (1-\alpha) \cdot \alpha^{(T+1)-t}$$

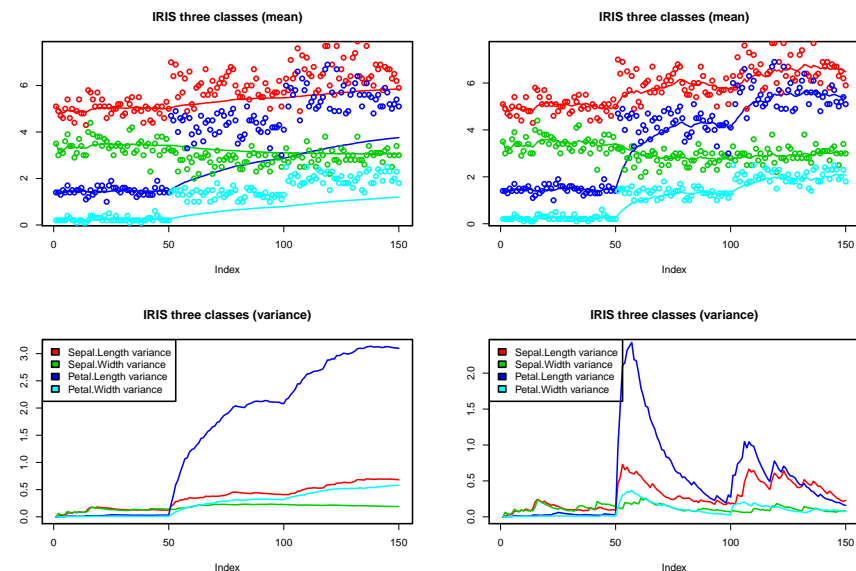
IRIS-Datensatz · Klasse „setosa“ [1:50]

Inkrementelle Mittelwerte und Varianzen



IRIS-Datensatz · 3 Klassen [1:50:100:150]

Mittelwerte und Varianzen · Uniforme und exponentielle Gewichte



Ein-Klassen-Szenarien

Zwei-Klassen-Szenarien

Mehr-Klassen-Szenarien

Ensemblemethoden — Bagging & Boosting

Kombinatorische Merkmalauswahl

Inkrementelles Lernen

Transduktion

Zusammenfassung (11)

1. Im Fall eines **Zweiklassenproblems** ist die Gestaltung der **Kostenmatrix** äquivalent zur Festsetzung eines Schwellenwertes für die **logodds-Prüfgröße**.
2. Eine kostenunabhängige Leistungscharakteristik ist die **Receiver-Operator-Kurve**, das einschlägige **skalare** Gütemaß ist die **Wilcoxon-Mann-Whitney-** oder **AUC**-Statistik.
3. **Mehrklassenentscheidungsregeln** lassen sich auf Zweiklassenregeln zurückführen; wir unterscheiden die **Jeder-gegen-jeden-** und die **Jeder-gegen-den-Rest**-Strategie.
4. Das **Bagging**-Ensemble mittelt die Entscheidungen von Klassifikatorinstanzen mit unterschiedlich ausgewürfelten **Bootstrap**-Lernstichproben.
5. Das **Boosting**-Ensemble staffelt seine Entscheidungsregeln nach **erfolgsgewichteten** Lernstichproben.
6. Wettbewerbende Klassifikatorinstanzen werden auf **Prüfgrößenbasis** durch **EM-Interpolation** fusioniert.
7. **Ranker** wählen Merkmale ausschließlich auf Grundlage ihrer **Relevanz** aus.
8. **Filter** streben gleichzeitig die Unterdrückung **redundanter** Merkmalkombinationen an.
9. **Wrapper** suchen die fehlerminimale unter **allen** Merkmalkombinationen; diese NP-harte Aufgabe wird mit heuristischen **Graphsuchverfahren** angegriffen.
10. Rekursionsformeln für **inkrementelle Lernschritte** erlauben dem Klassifikator **lebenslanges Lernen** und **schritthaltende Anpassung** an gleitende Umgebungsbedingungen.