

MASCHINELLES LERNEN & DATAMINING

Vorlesung im Wintersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 13. Oktober 2017

Teil IV

Vorhersage und Kategorisierung

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Vorhersage und statistische Abhängigkeit

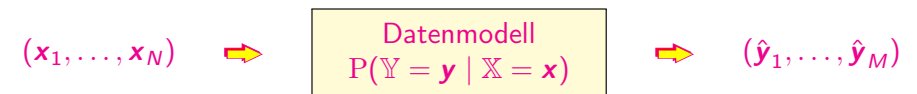
Charakterisierung der statistischen Unabhängigkeit

Zwei Variablenmengen $\mathbb{X} = (\mathbb{X}_1, \dots, \mathbb{X}_N)$ und $\mathbb{Y} = (\mathbb{Y}_1, \dots, \mathbb{Y}_M)$ heißen **statistisch unabhängig** voneinander gdw. gilt:

$$(\forall \mathbf{x})(\forall \mathbf{y}) \quad P(\mathbb{X} = \mathbf{x}, \mathbb{Y} = \mathbf{y}) = P(\mathbb{X} = \mathbf{x}) \cdot P(\mathbb{Y} = \mathbf{y})$$

Für Tupel \mathbf{x} mit $P(\mathbf{x}) \neq 0$ ist das äquivalent zu:

$$P(\mathbb{Y} = \mathbf{y} \mid \mathbb{X} = \mathbf{x}) = P(\mathbb{Y} = \mathbf{y})$$



Fakt

Im Fall statistischer Abhängigkeit besteht eine Chance, die Werte der **endogenen** Variablen \mathbb{Y}_m aus den Werten der **exogenen** Variablen \mathbb{X}_n zu „erraten“.

Statistische Prädiktion von Einzelvariablen

Quellvariable $(\mathbb{X}_1, \dots, \mathbb{X}_N)$ → Zielvariable $\mathbb{Y}_1 =: \mathbb{Y}$

Maschinelles Lernen eines Vorhersagemodells

Entscheidungsfunktion: $f : \mathcal{X}_1 \times \dots \times \mathcal{X}_N \rightarrow \mathcal{Y}$

Kostenfunktion („loss“): $\mathcal{L}(\mathbf{x}, y, \hat{y})$ mit $\hat{y} = f(\mathbf{x})$

Risiko (zu minimieren): $\mathfrak{R}(f) := \mathcal{E}_{P(\mathbf{x}, y)}[\mathcal{L}(\mathbb{X}, \mathbb{Y}, f(\mathbb{Y}))]$

\mathbb{Y} nominal

$\mathcal{L}(\mathbf{x}, y, \hat{y}) = c_{y\hat{y}}$
Kostenmatrix \mathbf{C} mit
 $c_{\kappa\kappa} \leq c_{\kappa\lambda}$

Spezialfall

(Fehlerrate)
 $c_{\kappa\lambda} = \begin{cases} 0 & \kappa = \lambda \\ 1 & \kappa \neq \lambda \end{cases}$

\mathbb{Y} ordinal

$\mathcal{L}(\mathbf{x}, y, \hat{y}) = c_{y\hat{y}}$
Diskrepanzmatrix \mathbf{C}
mit $c_{k\ell} \leq c_{k'\ell'}$ für
 $k' \leq k \leq \ell \leq \ell'$

Spezialfall

(Linearskala)
 $c_{k\ell} = |z_k - z_\ell|$

\mathbb{Y} kardinal

$\mathcal{L}(\mathbf{x}, y, \hat{y}) = d(y, \hat{y})$
metrische Distanzmaße
 $d(y, \hat{y}) = |y - \hat{y}|^p, p \geq 0$

Spezialfall

(Quadratmittel)
 $d(y, \hat{y}) = (y - \hat{y})^2$

Klassifikationsverfahren

Welche(n) Skalentyp(en) besitzen die exogenen Variablen ?

Numerisch

NV-Klassifikator
Polynomklassifikator
Multilayer-Perzeptron
Supportvektormaschine

Metrisch

Nächste-Nachbar-Regeln
SVM + Kerneltrick
MDS + **numerisch**

Diskret

Versionenraumverfahren
Kanonische+naive Bayesregel
Markovnetze

Numerisch & diskret

Entscheidungsbäume
Loglinearmodelle
Bayesnetze
(Konversion)

Optimale Prädiktion in den Spezialkonfigurationen

$$\mathfrak{R}(f) = \mathcal{E}[\mathcal{L}(\mathbb{X}, \mathbb{Y}, f(\mathbb{Y}))] = \int \sum_y P(\mathbf{x}, y) \cdot c_{y, f(\mathbf{x})} dx$$

Klassifikation (Bayesregel)

\mathbb{Y} ist **nominal**

Modus

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_{\kappa \in \Omega_{\mathbb{Y}}} P(\mathbb{Y} = \kappa \mid \mathbf{x})$$

Ordinale Klassifikation

\mathbb{Y} ist **ordinal**

Median

$$\hat{y}(\mathbf{x}) = \operatorname{median}_{\ell \in \Omega_{\mathbb{Y}}} P(\mathbb{Y} = \ell \mid \mathbf{x})$$

Quadratmittel-Regression

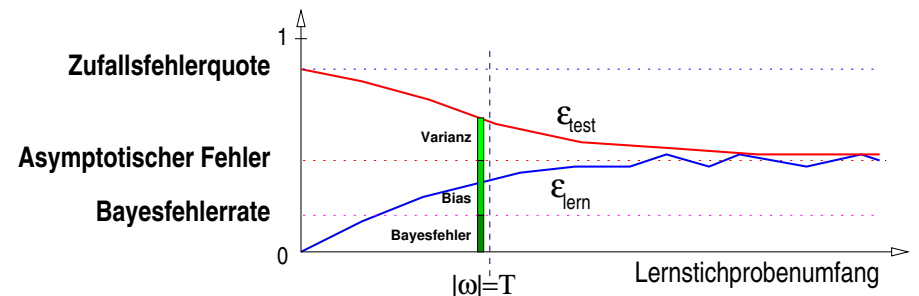
\mathbb{Y} ist **kardinal**

Mean

$$\hat{y}(\mathbf{x}) = \mathcal{E}_{\mathbb{Y}|\mathbf{x}}[\mathbb{Y}] = \int_{\mathbb{R}} P(y|\mathbf{x}) \cdot y dy$$

Fehlerrate, Überanpassung & Unteranpassung

Was wir schon in der Vorlesung „Mustererkennung“ über das Lernen gelernt haben



Lernstichprobe des Klassifikationsverfahrens

$$\omega = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$$

Fehlerrate auf den Lerndaten

ϵ_{lern}

Fehlerrate auf den Testdaten (\approx Fehlerwahrscheinlichkeit)

ϵ_{test}

Bayesfehler — weniger geht nicht

Zufallsfehler — mehr muss nicht

Grenzfehler — Daten! Daten!!

Bias

Datenmodell ↓

Varianz

Lernprobe ↓

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Kardinalitätskonflikt

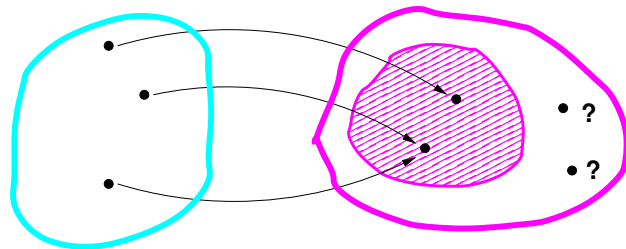
„Worüber man nicht reden kann, darüber soll man schweigen.“

Intension

$\{\phi \mid \phi : \Omega \rightarrow \{0, 1\} \text{ terminierender Algorithmus}\}$

Extension

Potenzmenge $\mathfrak{P}\Omega$



abzählbar unendlich

$\langle 0 \rangle, \langle 1 \rangle, \langle 00 \rangle, \langle 01 \rangle, \langle 10 \rangle, \langle 11 \rangle, \langle 000 \rangle, \langle 001 \rangle, \langle 010 \rangle, \dots, \langle 0000 \rangle, \langle 0001 \rangle, \langle 0010 \rangle, \dots, \langle 00000 \rangle, \dots, \langle 000000 \rangle, \dots$

überabzählbar unendlich

$\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}, \dots, \{x_1, x_2\}, \{x_1, x_3\}, \dots, \{x_2, x_3\}, \dots, \{x_1, x_2, x_3\}, \{x_1, x_2, x_4\}, \dots, \{x_1, x_2, x_3, x_4\}, \{x_1, x_2, x_3, x_5\}, \dots$

$$\phi \mapsto \Omega_\phi \stackrel{\text{def}}{=} \{x \in \Omega \mid \phi(x) = 1\} \in \mathfrak{P}\Omega$$

Begriffe (Konzepte)

Intensionaler Zugehörigkeitstest: $\phi(x) = \phi(x_1, \dots, x_N) = 1 ?$



Scharf oder unscharf?

Ist dieses Element aus

$$\Omega = \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_N$$

$$\text{ein } \left\{ \begin{array}{l} \text{Tiger ?} \\ \text{Liger ?} \\ \text{Töwe ?} \\ \text{Löwe ?} \end{array} \right\}$$

Extension eines Begriffs

Mengentheoret. charakterisiert

$$\mathcal{C} \subseteq \Omega$$

$x \in \mathcal{C} \iff x \text{ „ist“ ein } \mathcal{C}$

Intension eines Begriffs

Algorithmisch charakterisiert

$$\phi : \Omega \rightarrow \{0, 1\}$$

„Parser“ ϕ entscheidbare Funkt.
(\rightsquigarrow endlich aufschreibbar)

Abstraktion

Konzentration auf die Objekteigenschaften zu Lasten der Objektidentität

Gegeben

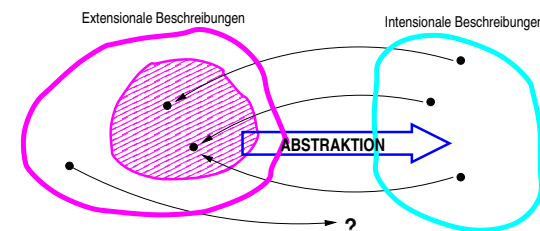
ist ein Konzept
 $\mathcal{C} \subseteq \Omega$ (extensional)

Gesucht

ist eine „krispe“ intensionale Beschreibung
 $\phi : \Omega \rightarrow \{0, 1\}$ mit der Eigenschaft

$$x \in \mathcal{C} \iff \phi(x) = 1$$

(\rightsquigarrow Kompatibilität: $\Omega_\phi = \mathcal{C}$)



Problem

Nicht jedes Konzept ist abstrahierbar!

Induktion

Verallgemeinerung oder „Lernen aus Beispielen“

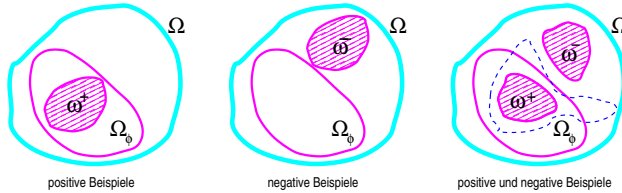
Gegeben

Positivbeispiele

$$\omega^+ \subseteq \Omega_\phi$$

Negativbeispiele

$$\omega^- \subseteq \Omega \setminus \Omega_\phi$$



Gesucht

kompatible intensionale Beschreibung ψ :

$$(\forall \mathbf{x} \in \Omega) \quad \mathbf{x} \in \omega^+ \Rightarrow \psi(\mathbf{x}) = 1$$

$$(\forall \mathbf{x} \in \Omega) \quad \mathbf{x} \in \omega^- \Rightarrow \psi(\mathbf{x}) = 0$$

Potentielle Kandidaten sind alle ψ mit

$$\omega^+ \subseteq \Omega_\psi \subseteq \Omega \setminus \omega^-$$

Lernverfahren

- Hypothesenraum
- Lösungsvielfalt
- Auswahlkriterium

Kategorisierung von Objekten

(a.k.a. „Klassifikation“)

Extensionale Charakt. vs.

Mengenpartition

$\mathcal{C}_1, \dots, \mathcal{C}_K \subseteq \Omega$ mit

$$\bigcup_{\kappa=1}^K \mathcal{C}_\kappa = \Omega$$

und für alle $\kappa \neq \lambda$:

$$\mathcal{C}_\kappa \cap \mathcal{C}_\lambda = \emptyset$$

Spezialfall $K = 2$

Konzept $\mathcal{C}_1 = \mathcal{C}$ und sein

Komplement $\mathcal{C}_2 = \Omega \setminus \mathcal{C}$

Intensionale Charakterisierung

keine wirklich zwingende

Verallgemeinerung:

1. Charakteristische Fkt.

$$\phi : \Omega \rightarrow \{1, \dots, K\} \subset \mathbb{R}$$

2. Konzepttupel

$$\phi : \Omega \rightarrow \{0, 1\}^K \text{ mit } \Omega_{\phi_\kappa} \stackrel{!}{=} \mathcal{C}_\kappa$$

3. Diskriminanten

$$\phi : \Omega \rightarrow \mathbb{R}^K \text{ mit } \mathbf{x} \in \mathcal{C}_\kappa \text{ gdw.}$$

$$\phi_\kappa(\mathbf{x}) \geq \phi_\lambda(\mathbf{x}) \quad (\forall \lambda)$$

4. Nominale Regression

$$P : \Omega \times \{\xi_1, \dots, \xi_K\} \rightarrow \mathbb{R}$$

plus Bayesregel

Induktionsproblematik

Lerndaten? Hypothesen? Verfahren?

Übergeneralisierung

Es werden *Oberbegriffe* von ϕ gelernt.

Überspezialisierung

Es werden *Unterbegriffe* von ϕ gelernt.

Fehlgranulation

Überanpassung oder *Unteranpassung*

Natürlichkeit, Fortsetzbarkeit

Gelernte Verallgemeinerung versagt bei Wiederabruf

Abhilfe

Negativbeispiele bereitstellen

Abhilfe

Repräsentative

Positivstichprobe

Abhilfe

Adäquates Sortiment von Hypothesen

Abhilfe

Occam's razor: einfache Erklärung

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Aussagenlogisches Lernen

Begriffe lernen · Klassifikation · Gruppierung

Segel-Szenarium

\mathcal{X}_1 sky	\mathcal{X}_2 air	\mathcal{X}_3 humidity	\mathcal{X}_4 wind	\mathcal{X}_5 water	\mathcal{X}_6 forecast
$\left\{ \begin{array}{l} \text{sunny} \\ \text{rainy} \\ \text{cloudy} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{warm} \\ \text{cool} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{normal} \\ \text{high} \\ \text{low} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{strong} \\ \text{weak} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{warm} \\ \text{cool} \end{array} \right\}$	$\left\{ \begin{array}{l} \text{same} \\ \text{change} \end{array} \right\}$

Single Representation Trick

z.B. Hypothesen als unvollständige Attributwertspezifikationen

- **Objekte** \triangleq Attributbelegungen

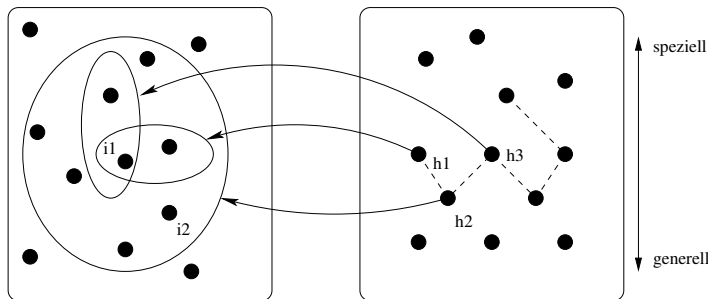
$$(\text{sunny}, \text{warm}, \text{normal}, \text{strong}, \text{warm}, \text{same}) \in \Omega$$

- **Hypothesen** \triangleq partielle Attributbelegungen

$$(\text{sunny}, ?, ?, \text{strong}, ?, ?) \in \mathcal{H}$$

Hypothesen und Objektmengen

$$\Omega(h) \triangleq \{x \in \Omega \mid h \models x\}$$



Segel-Szenarium

$$\begin{aligned} i_1 : & (\text{sunny}, \text{warm}, \text{high}, \text{strong}, \text{cool}, \text{same}) \\ i_2 : & (\text{sunny}, \text{warm}, \text{high}, \text{light}, \text{warm}, \text{same}) \end{aligned}$$

$$\begin{aligned} h_1 : & (\text{sunny}, ?, ?, \text{strong}, ?, ?) \\ h_2 : & (\text{sunny}, ?, ?, ?, ?, ?) \\ h_3 : & (\text{sunny}, ?, ?, ?, \text{cool}, ?) \end{aligned} \quad \begin{aligned} h_2 & \supseteq h_1, h_3 \\ & \text{bzw.} \\ h_1, h_3 & \Rightarrow h_2 \end{aligned}$$

Hypothesenraum

Konjunktionen positiver Literale (KPL)

Definition

Es sei $\Omega = \mathcal{X}_1 \times \dots \times \mathcal{X}_N$ ein Objektraum. Dann heißen die Elemente aus

$$\mathcal{H} = (\mathcal{X}_1 \cup \{?\}) \times \dots \times (\mathcal{X}_N \cup \{?\})$$

KPL-Hypothesen über Ω . Die Menge \mathcal{H} heißt **KPL-Hypothesenraum** über Ω .

Ein Beispielobjekt $x \in \Omega$ genügt der Hypothese $h \in \mathcal{H}$ (x **erfüllt** h bzw. $h \models x$) genau dann, wenn gilt:

$$\forall i = 1, \dots, N : (h_i = ?) \vee (h_i = x_i)$$

Bemerkung

Vollständige KPL \triangleq Objekte
 Leere KPL \triangleq Konzept $\mathcal{C} = \Omega$
 Definiere $h_\emptyset \triangleq$ Leerkonzept
 $\mathcal{C} = \emptyset$

Segel-Szenarium

Objektraum: $|\Omega| = 144$
 KPL-Hypothesenraum: $|\mathcal{H}| = 1296$
 Konzeptraum:
 $|\mathfrak{P}\Omega| = 2^{144} \approx 1000^{14.4} \approx 10^{43}$

Der Verband aller KPL-Hypothesen

Definition

Für jede Hypothese $h \in \mathcal{H}$ sei $\Omega(h) \stackrel{\text{def}}{=} \{x \in \Omega \mid h \models x\}$ (Extension) definiert. Die Menge \mathcal{H} erbt von $\mathfrak{P}\Omega$ die **Inklusionsrelation** (h ist „allgemeiner“ oder „genereller“ als h'):

$$h \supseteq h' \iff \forall x \in \Omega : (h' \models x \Rightarrow h \models x)$$

Der Raum aller DNF-Hypothesen (disjunktive Normalform) ist die Boolesche Algebra $(\mathfrak{P}\Omega, \subseteq)$.

Lemma

Der Raum (\mathcal{H}, \subseteq) bildet eine Halbordnung.

$$\left\{ \begin{array}{l} \text{reflexiv} \\ \text{transitiv} \\ \text{antisymmetrisch} \end{array} \right\}$$

Die KPL-Hypothesen sind abgeschlossen gegenüber Durchschnittsbildung.

Die KPL-Hypothesen sind nicht abgeschlossen gegenüber der Mengenvereinigung, es existiert das Supremum je zweier Hypothesen:

$$(h \vee h')_n \stackrel{\text{def}}{=} \begin{cases} v & (\exists v \in \mathcal{X}_n) \ h_n = v = h'_n \\ ? & h_n \neq h'_n \\ ? & h_n = ? = h'_n \end{cases}$$

Sukzessiver Generalisierungsalgorithmus

Definition

Eine Hypothese $h \in \mathcal{H}$ heißt **konsistent** mit den Lerndaten (ω^+, ω^-) genau dann wenn gilt:

$$\begin{aligned} x \in \omega^+ &\Rightarrow h \models x \\ x \in \omega^- &\Rightarrow h \not\models x \end{aligned}$$

Bemerkungen

1. Konsistenz falls $\omega^+ \subseteq \Omega_h \subseteq \Omega \setminus \omega^-$
2. Jedes h ist konsistent mit (\emptyset, \emptyset) .
3. Kein h ist konsistent wenn $\omega^+ \cap \omega^- \neq \emptyset$.
4. Auch für disjunkte (ω^+, ω^-) enthält \mathcal{H} nicht notwendig eine konsistente Hypothese!

1 INITIALISIERUNG

Setze $h \leftarrow h_\emptyset$.

2 GENERALISIERUNG

Setze für alle $x \in \omega^+$:

$$h \leftarrow h \vee x$$

(„speziellste Erweiterung“ von h um x)

3 TERMINIERUNG

Das Ergebnis ist h .

Keine $x \in \omega^-$ verwendet.

Resultat genügt allen $x \in \omega^+$.
 h ist minimal mit dieser Eigenschaft.

Wenn konsistente Hypothese existiert, wird sie gefunden.

Nur für KPL (Supremum!) realisierbar.

Für $\mathcal{H} = \mathfrak{P}\Omega$ ist SGA trivial.

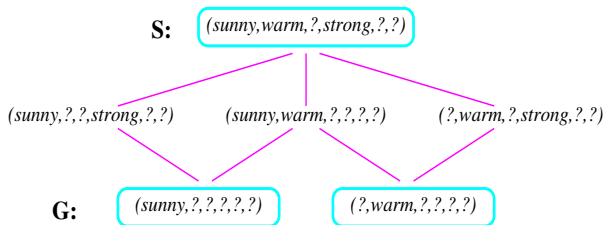
Der Versionsraum

Definition

Die Menge der mit den Lernbeispielen konsistenten Hypothesen

$$\{h \in \mathcal{H} \mid h \text{ konsistent mit } (\omega^+, \omega^-)\}$$

heißt **Versionenraum** von (ω^+, ω^-) bezüglich \mathcal{H} und wird mit $\mathfrak{V}(\mathcal{H}, \omega^+, \omega^-)$ (oder \mathfrak{V}) bezeichnet.



Beispiel

Versionenraum mit 6 Hypothesen

1x minimal
2x maximal
3x weder/noch

Minimale und maximale VR-Elemente

$$\mathfrak{V}_S \stackrel{\text{def}}{=} \{h \in \mathfrak{V} \mid \forall h' \in \mathfrak{V} : h' \subseteq h \Rightarrow h' = h\}$$

$$\mathfrak{V}_G \stackrel{\text{def}}{=} \{h \in \mathfrak{V} \mid \forall h' \in \mathfrak{V} : h \subseteq h' \Rightarrow h' = h\}$$

Kandidateneliminationsalgorithmus

Suche operiert auf („Kandidaten“-) Mengen von Hypothesen

1 INITIALISIERUNG

Setze $H \leftarrow \mathcal{H}$

2 GENERALISIERUNG / SPEZIALISIERUNG

Eliminiere für alle $x \in \omega^+ \cup \omega^-$

- a Fall $x \in \omega^+$: alle $h \in H$ mit $h \not\models x$
- b Fall $x \in \omega^-$: alle $h \in H$ mit $h \models x$

3 TERMINIERUNG

Das Ergebnis ist h , falls $H = \{h\}$ ist.

Am Ende enthält die Kandidatenmenge genau die konsistenten Hypothesen aus \mathcal{H} .

Es gibt keine, eine oder mehrere Lösungen.

Das Verfahren ist aus Aufwandsgründen impraktikabel!

Versionenräume als Halbordnungsintervalle

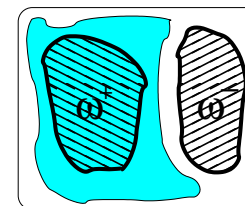
Beispiel

Im vollständigen Hypothesenraum $\mathcal{H} = \mathfrak{P}\Omega$ sind die minimalen und die maximalen VR-Elemente eindeutig:

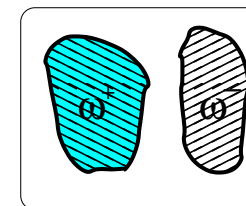
$$\mathfrak{V}_S = \{\omega^+\} \quad \text{und} \quad \mathfrak{V}_G = \{\Omega \setminus \omega^-\}$$

Versionenräume besitzen die Gestalt einer **Intervalldarstellung**:

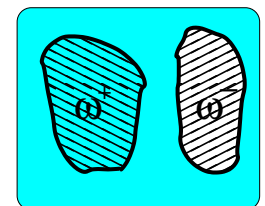
$$\mathfrak{V}(\mathcal{H}, \omega^+, \omega^-) = \{h \in \mathcal{H} \mid \omega^+ \subseteq h \subseteq \Omega \setminus \omega^-\} = [\mathfrak{V}_S, \mathfrak{V}_G]_{\mathcal{H}}$$



eine VR-Hypothese



die kleinste VR-Hypothese



die größte VR-Hypothese

Der Versionenraum-Darstellungssatz

Die Intervalldarstellung gilt in allen beliebigen Hypothesenräumen

Definition

Sei $\mathcal{H} \subseteq \mathfrak{H}\Omega$; ein **einfaches HO-Intervall** in \mathcal{H} hat die Form:

$$[h_u, h_o]_{\mathcal{H}} \stackrel{\text{def}}{=} \{h \in \mathcal{H} \mid h_u \subseteq h \subseteq h_o\}$$

Ein **verallgemeinertes HO-Intervall** in \mathcal{H} hat die Form:

$$[\mathcal{H}_u, \mathcal{H}_o]_{\mathcal{H}} = \{h \in \mathcal{H} \mid \exists h_u \in \mathcal{H}_u, \exists h_o \in \mathcal{H}_o : h_u \subseteq h \subseteq h_o\}$$

Satz

Für den Versionenraum \mathfrak{V} der Beispieldaten ω^+ und ω^- bezüglich \mathcal{H} gilt eine Intervalldarstellung:

$$\mathfrak{V}(\mathcal{H}, \omega^+, \omega^-) = [\mathfrak{V}_S, \mathfrak{V}_G]_{\mathcal{H}}$$

Dabei sind \mathfrak{V}_S und \mathfrak{V}_G die Mengen der \subseteq -minimalen (\subseteq -maximalen) Elemente des Versionenraums \mathfrak{V} .

$$[\mathcal{H}_u, \mathcal{H}_o]_{\mathcal{H}} = \bigcup_{h_u \in \mathcal{H}_u} \bigcup_{h_o \in \mathcal{H}_o} [h_u, h_o]_{\mathcal{H}}$$

Beweis.

Inklusionsrichtung \subseteq :

Sei $h \in \mathfrak{V}$. Sei $G(h) := \{h' \in \mathfrak{V} \mid h' \supseteq h\}$.

Wegen $h \in G(h)$ ist $G(h) \neq \emptyset$.

- Sei h_G ein maximales Element aus $G(h)$.
Dann ist $h_G \in \mathfrak{V}_G$ und $h_G \supseteq h$.

(Die Existenz eines $h_S \in \mathfrak{V}_S$ zeigt man/frau analog.)

Inklusionsrichtung \supseteq :

Sei $h \in \mathcal{H}$ mit $h \subseteq h_G \in \mathfrak{V}_G$ und $h \supseteq h_S \in \mathfrak{V}_S$.

Zu zeigen: h ist konsistent mit (ω^+, ω^-) , d.h. $h \in \mathfrak{V}$.

1. Sei $x \in \omega^+$.
Wegen $h_S \in \mathfrak{V}_S \subseteq \mathfrak{V}$ gilt $h_S \models x$.
Wegen $h \supseteq h_S$ gilt auch $h \models x$.
2. Sei $x \in \omega^-$.
Wegen $h_G \in \mathfrak{V}_G \subseteq \mathfrak{V}$ gilt $h_G \not\models x$.
Wegen $h \subseteq h_G$ gilt auch $h \not\models x$.

□

Versionenraum-Kandidateneliminationsalgorithmus

(Algorithmus)

1 INITIALISIERUNG

Setze $G \leftarrow \{\Omega\}$ und $S \leftarrow \{\emptyset\}$.

2⁺ POSITIVE BEISPIELE

Für alle $x \in \omega^+$:

- Entferne alle $h \in G$ mit $h \not\models x$
- Für alle $h \in S$:

Generalisiere h zu h' mit $h' \models x$

Behalte $h' \in S$, falls h' spezieller als G

- Entferne alle nichtminimalen $h \in S$

2⁻ NEGATIVE BEISPIELE

Für alle $x \in \omega^-$:

- Entferne alle $h \in S$ mit $h \models x$ • Für alle $h \in G$:

Spezialisiere h zu h' mit $h' \not\models x$

Behalte $h' \in G$, falls h' allgemeiner als S

- Entferne alle nichtmaximalen $h \in G$

3 TERMINIERUNG

Das Ergebnis ist h , falls $G = \{h\} = S$ ist.

(zumf3hogg1A)

Bemerkungen

1. Grundidee: alle Versionenräume werden als „Intervalle“ $[S, G]$ abgespeichert, und auch die Hypothesenelimination geschieht auf S, G und nicht auf \mathfrak{V} .
2. Es gilt natürlich $\mathcal{H} = [\emptyset, \Omega]_{\mathcal{H}}$.
3. Wenn es geeignete Hypothesen mit $\Omega(h_{\emptyset}) = \emptyset$ und $\Omega(h_{\Omega}) = \Omega$ gibt, kann entsprechend initialisiert werden.
4. Hypothesen $h \in G$, die einem Positivbeispiel $x \in \omega^+$ nicht genügen, dürfen ohne weiteres eliminiert werden, da jegliche Spezialisierung von h ebenfalls an x scheitern würde. Dasselbe gilt für $h \in S$, $x \in \omega^-$ mit $h \models x$.
5. Gilt jedoch für $x \in \omega^+$ und ein $h \in S$ die Aussage $h \not\models x$, so darf h wegen der Gefährdung des Teilraums $[h, G]$ nicht einfach gelöscht werden!
6. Von allen Generalisierungen h' von h mit $h' \models x$ interessieren natürlich nur diejenigen mit $[h', G] \neq \emptyset$, und die auch minimal sind in S mit dieser Eigenschaft.
7. Am Ende sind alle Hypothesen aus S und aus G und auch aus $[S, G]$ konsistent mit den Beispieldaten, und $[S, G]$ ist auch diesbezüglich vollständig.

Beispiel (Segeln im KPL-Hypothesenraum)

Versionenraum nach VRE-Algorithmus

	sky	air	humidity	wind	water	forecast
$h_1 \in S$	sunny	warm	?	strong	?	?
h_2	sunny	?	?	strong	?	?
h_3	sunny	warm	?	?	?	?
h_4	?	warm	?	strong	?	?
$h_5 \in G$	sunny	?	?	?	?	?
$h_6 \in G$?	warm	?	?	?	?

Unbeobachtete („neue“) Objekte

Vorhersage des Konzepts „go_sailing“:

							S					G	G
	sky	air	hum	wind	water	fore	h_1	h_2	h_3	h_4		h_5	h_6
x_1	sunny	warm	norm	strong	cool	change	1	1	1	1		1	1
x_2	rainy	cold	norm	weak	warm	same	0	0	0	0		0	0
x_3	sunny	warm	norm	weak	warm	same	0	0	1	0		1	1

Die induktive Hülle

Definition

Wir bezeichnen die Menge

$$\overline{\omega^+} \stackrel{\text{def}}{=} \{x \in \Omega \mid h \in \mathfrak{H}(\omega^+, \omega^-) \Rightarrow h \models x\}$$

als **induktive Hülle** der Positivbeispiele ω^+ und die Menge

$$\overline{\omega^-} \stackrel{\text{def}}{=} \{x \in \Omega \mid h \in \mathfrak{H}(\omega^+, \omega^-) \Rightarrow h \not\models x\}$$

als **induktive Hülle** der Negativbeispiele ω^- . Die Elemente aus

$$\omega^? \stackrel{\text{def}}{=} \Omega \setminus (\overline{\omega^+} \cup \overline{\omega^-})$$

heißen **ambige Objekte** von Ω bezüglich \mathcal{H} , ω^+ und ω^- .

Lemma

Die Operatoren $\omega^+ \mapsto \overline{\omega^+}$ und $\omega^- \mapsto \overline{\omega^-}$ sind tatsächlich Hüllenoperatoren:

- $\omega_1^+ \subseteq \omega_2^+ \Rightarrow \overline{\omega_1^+} \subseteq \overline{\omega_2^+}$ (Monotonie)
- $\omega^+ \subseteq \overline{\omega^+}$ und $\omega^- \subseteq \overline{\omega^-}$ (Inklusion)
- $\overline{\omega^+} = \overline{\overline{\omega^+}}$ und $\overline{\omega^-} = \overline{\overline{\omega^-}}$ (Involution)

Parlamentarischer Alltag im Versionenraum

Positiver Konsens

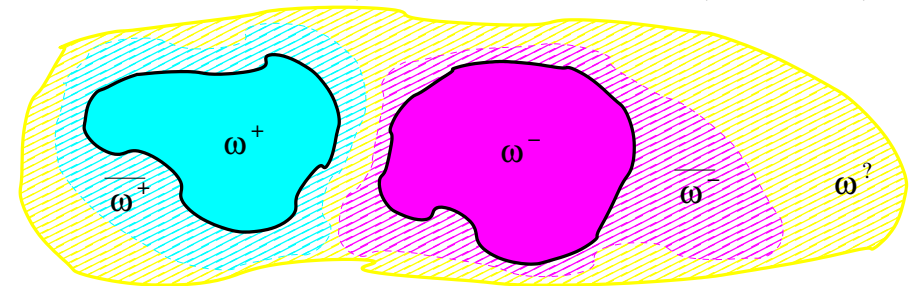
Für alle $h \in \mathfrak{H}$ gilt
 $h \models x$

Negativer Konsens

Für alle $h \in \mathfrak{H}$ gilt
 $h \not\models x$

Ambiges Votum

Ex. $h_+, h_- \in \mathfrak{H}$ mit
 $h_+ \models x$ und $h_- \not\models x$



Bemerkungen

- Alle $h \in \mathfrak{H}$ sind konsistent $\Rightarrow \left\{ \begin{array}{l} \text{alle } x \in \omega^+ \text{ werden einstimmig akzeptiert} \\ \text{alle } x \in \omega^- \text{ werden einstimmig abgewiesen} \end{array} \right\}$.
- Für $\mathfrak{H} = \emptyset$ folgt $\overline{\omega^+} \cap \overline{\omega^-} = \emptyset$.
- Ist $x \in \Omega$ ambig, so ex. Hypothesen $h^+ \in \mathfrak{H}_G$, $h^- \in \mathfrak{H}_S$ mit $\left\{ \begin{array}{l} h^+ \models x \\ h^- \not\models x \end{array} \right\}$.

Der induktive Bias

Großartiger Lernerfolg durch mangelhafte Ausdrucksfähigkeit

Induktives Schließen

steht und fällt mit dem **Ausdrucksdefizit** des Hypothesenraums \mathcal{H} .

ω^+	sunny	warm	normal	strong	warm	same
ω^+	rainy	warm	normal	weak	warm	same
\mathfrak{H}_S	?	warm	normal	?	warm	same
$\omega^?$	sunny	warm	normal	weak	warm	same

$\Rightarrow x$

Aussagenlogisch orientierte Hypothesenräume

- Konjunktion positiver Literale KPL: $x_i = \xi$ oder ?
- Konjunktion positiver und negativer Literale $x_i = \xi$ oder $x_i \neq \xi$ oder ?
- Konjunktion disjunktiver Komplexe $x_i \in \mathcal{X}^+$
- Disjunktion positiver (und negativer) Literale (dual zu oben)
- Disjunktion von Konjunktionen positiver Literale $\mathcal{H} = \mathfrak{P}\Omega$

Lernen einelementiger Versionenräume

Zur Auswahl neuer Lernbeispiele

- Erweiterung von ω^+ um Beispiele aus $\overline{\omega^+}$ ist überflüssig.
Erweiterung von ω^- um Beispiele aus $\overline{\omega^-}$ ist überflüssig.
- Erweiterung von ω^+ um Beispiele aus $\overline{\omega^-}$ bewirkt Inkonsistenz.
Erweiterung von ω^- um Beispiele aus $\overline{\omega^+}$ bewirkt Inkonsistenz.
- Nur die Erweiterung von (ω^+, ω^-) um ambige Beispiele $x \in \omega^?$ ist zugleich **konsistent** und **produktiv**!

Exploratives Lernen

Sukzessives Akquirieren produktiver neuer Beispiele, bis

- der Versionenraum \mathfrak{V} nur noch ein h enthält oder
- der Versionenraum \mathfrak{V} leergelaufen ist.

Ambiguität und Rückweisung

Votierungstechniken für die Entscheidungsphase

Faules Lernen

Fallbasiertes Schließen

$$x \mapsto \begin{cases} \Omega^+ & x \in \omega^+ \\ \Omega^- & x \in \omega^- \\ \Omega^? & \text{sonst} \end{cases}$$

(keine Verallgemeinerung)

Fleißiges Lernen

einer Hypothese ($\mathfrak{V} = \{h^*\}$)

$$x \mapsto \begin{cases} \Omega^+ & h^* \models x \\ \Omega^- & h^* \not\models x \end{cases}$$

Orakel

Occam's Razor (MDL, BIC, AIC)

Wahrscheinlichkeiten ...

Einstimmigkeit

$$x \mapsto \begin{cases} \Omega^+ & h \in \mathfrak{V}_S \Rightarrow h \models x \\ \Omega^- & h \in \mathfrak{V}_G \Rightarrow h \not\models x \\ \Omega^? & \text{sonst} \end{cases}$$

Generalkonsens

$$x \mapsto \begin{cases} \Omega^+ & h \in \mathfrak{V}_G \Rightarrow h \models x \\ \Omega^- & \text{sonst} \end{cases}$$

Mehrheitsvotum

$$x \mapsto \begin{cases} \Omega^+ & |\{h \in \mathfrak{V} \mid h \models x\}| > |\mathfrak{V}|/2 \\ \Omega^- & |\{h \in \mathfrak{V} \mid h \models x\}| < |\mathfrak{V}|/2 \\ \Omega^? & |\{h \in \mathfrak{V} \mid h \models x\}| = |\mathfrak{V}|/2 \end{cases}$$

Gibbs-Sampling

Auswürfeln von $h^* \in \mathfrak{V}$ und

$$x \mapsto \begin{cases} \Omega^+ & h^* \models x \\ \Omega^- & h^* \not\models x \end{cases}$$

Lernen mit Orakel

(Algorithmus)

1 INITIALISIERUNG

Setze $G \leftarrow \{\Omega\}$ und $S \leftarrow \{\emptyset\}$ und $\omega^? \leftarrow \Omega$.

2 EXPLORATIONSSCHRITT

Solange $\omega^? \neq \emptyset$ gilt:

- Wähle ein Beispiel $x \in \omega^?$ aus
- Befrage das Orakel nach $x \in \mathcal{C}$
- Modifiziere den Versionenraum vermöge

$$\mathfrak{V} \leftarrow \mathfrak{V}(\omega^+ \cup \{x\}, \omega^-)$$

im Fall einer positiven Antwort und vermöge

$$\mathfrak{V} \leftarrow \mathfrak{V}(\omega^+, \omega^- \cup \{x\})$$

im Fall einer negativen Antwort des Orakels.

- Aktualisiere die Menge $\omega^?$ der ambigen Objekte

3 TERMINIERUNG

Das Ergebnis ist h , falls $G = \{h\} = S$ gilt.

(zumTheogIA)

ILP — Induktive logische Programmierung

Hypothesenraumbias wird explizit durch eine logische Theorie \mathcal{B} vorgegeben

Gegeben

Hypothesen $h \in \mathcal{H}$ sind **prädikatenlogische** Formeln

Objekte $x \in \Omega$ als Singleton-Hypothesen h_x

Positive und negative Lerndatensätze ω^+, ω^-

p.l. Formelmenge \mathcal{B} als expliziter Bias („Sachbereichstheorie“)

Gesucht

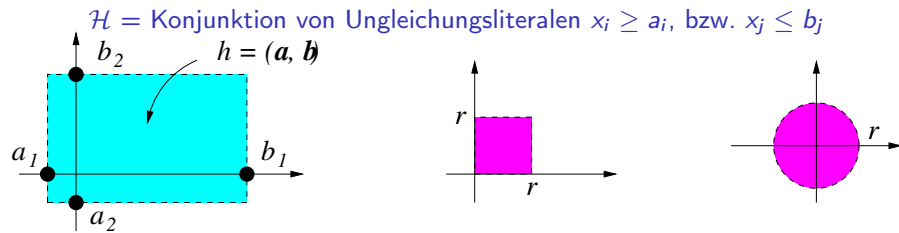
Eine Hypothese $h \in \mathcal{H}$ mit den Eigenschaften

- Vollständigkeit** $\mathcal{B}, h \models \omega^+$
- Korrektheit** für alle $x \in \omega^-$ gilt $\mathcal{B}, h \not\models x$
- Konsistenz** $\mathcal{B}, h, \omega^+, \omega^- \not\models \text{false}$

die zudem ein Gütekriterium $\left\{ \begin{smallmatrix} \text{speziell} \\ \text{generell} \\ \text{interessant} \\ \text{kurz} \end{smallmatrix} \right\}$ optimiert.

➡ nicht entscheidbar in der Prädikatenlogik erster Stufe

Numerisches Beispiel ($\Omega = \mathbb{R}^N$)



Hypothesenraum

$$h = \{x \mid a_i \leq x_i \leq b_i \text{ für alle } i\}$$

Konzeptraum I

$$C_r = \{x \mid \max_i x_i \leq r\}$$

Konzeptraum II

$$C_r = \{x \mid \|x\| \leq r\}$$

Computational Learning Theory

- Was wird (asymptotisch) gelernt?

{ korrektes Quadrat
u.U. nichts }

- Wie groß ist der erwartete Klassifikationsfehler?

{ 0%
22% }

Sterne und ihre Vereinigung

Lemma

Für jede Hypothese $h \in \mathcal{S}(x|\omega^-)$ gilt:

- h wird von x erfüllt. $h \models x$
- h wird von keinem $y \in \omega^-$ erfüllt. $(\forall y \in \omega^-) h \not\models y$
- h ist maximal mit diesen Eigenschaften, d.h., es gilt:

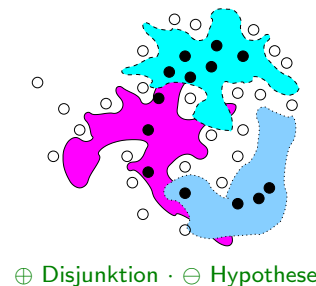
$$h' \supset h \Rightarrow h' \not\models x \text{ oder ex. } y \in \omega^- : h' \models y$$

Lemma

Aus **nichtleeren** Sternen lassen sich konsistente Disjunktionen konstruieren, d.h. die Vereinigungsmenge

$$h^* = \bigcup_{x \in \omega^+} \mathcal{S}(x|\omega^-)$$

ist konsistent mit (ω^+, ω^-) .

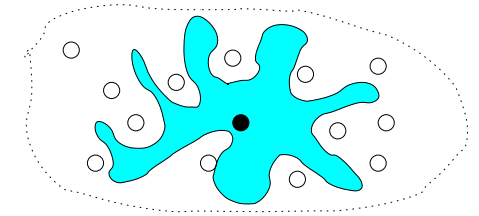
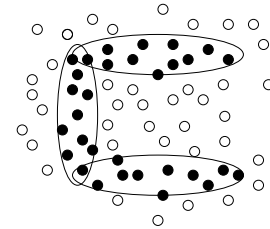


Michalskis Stern

Abgrenzung eines Positivbeispiels gegen alle Negativbeispiele

Problem

Die Menge ω^+ ist schwer gegen ω^- abgrenzbar.
 ω^+ zerfällt jedoch in einfacher strukturierte Teilmengen.



Definition

Es seien $\omega^+ \subset \Omega$, $\omega^- \subset \Omega$ und $x \in \omega^+$ ein Positivbeispiel. Die Hypothesenmenge

$$\mathcal{S}(x|\omega^-) \stackrel{\text{def}}{=} \mathfrak{V}_G(\mathcal{H}, \{x\}, \omega^-)$$

heißt **Stern** von x gegen ω^- .

Achtung!

Der Stern ist keine Hypothese, sondern ein Intervall.

Sternerzeugungsalgorithmus

- Wähle zufällig ein $x^* \in \omega^+$.
- Erzeuge den Stern
 $\mathcal{S}(x^*|\omega^-) = \mathfrak{V}_G(\mathcal{H}, \{x^*\}, \omega^-)$.
- Wähle eine Vorzugshypothese $h^* \in \mathcal{S}(x^*|\omega^-)$ mit maximaler Präferenz $\gamma(h^*)$.
- Wenn $h^* \models x$ für alle $x \in \omega^+$, so \rightsquigarrow 6.
- Tilge alle $x \in \omega^+$ mit $h^* \models x$ und \rightsquigarrow 1.
- Bilde die logische Disjunktion

$$h_{\text{dis}} \stackrel{\text{def}}{=} h_1 \vee \dots \vee h_r$$

aller bislang erzeugten Hypothesen.

Gegeben

\mathcal{H} , ω^+ , ω^- und eine **Präferenzfunktion**
 $\gamma : \mathcal{H} \mapsto \mathbb{R}$

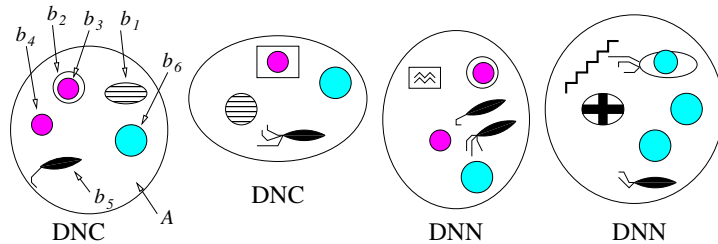
Gesucht

Eine disjunktive Beschreibung
 $h_1 \vee \dots \vee h_r$,
 $h_i \in \mathcal{H}$, die konsistent mit (ω^+, ω^-) ist.

Beispiel — Konzeptualisierung von Krebszellen

Aufgabenstellung

Unterscheide Krebszellen (DNC) von gesunden Zellen (DNN) auf Grundlage numerischer, kategorialer und struktureller Zellmerkmale.



Mensch-Maschine-Mensch-Zyklus

- (1) Definiere relevante Deskriptoren · etikettiere (ω^+ , ω^-)
- (2) Lerne induktiv passende Hypothesen für $\mathcal{C} = \mathcal{C}_{\text{DNC}}$.
- (3) Evaluere, analysiere und modifiziere das Szenarium.

Beispiel — Konzeptualisierung von Krebszellen

Objektbeschreibung der ersten DNC-Zelle

$$\begin{aligned} & \text{contains}(c, b_1, \dots, b_6) \wedge \text{circ}(c) = 8 \wedge \text{pplasm}(c) = A \\ & \wedge \text{shape}(b_1) = \text{ellipse} \wedge \text{texture}(b_1) = \text{stripes} \wedge \text{weight}(b_1) = 4 \\ & \wedge \text{orient}(b_1) = \text{NW} \wedge \text{shape}(b_2) = \text{circle} \wedge \text{contains}(b_2, b_3) \\ & \wedge \text{texture}(b_2) = \text{blank} \wedge \text{weight}(b_2) = 3 \wedge \dots \\ & \wedge \text{shape}(b_6) = \text{circle} \wedge \text{texture}(b_6) = \text{shaded} \wedge \text{weight}(b_6) = 5 \end{aligned}$$

DNC-Charakterisierung durch prädikatenlogische Formel

$$\begin{aligned} & \exists_1 b (\text{weight}(b) = 5) \\ & \exists_1 b (\text{shape}(b) = \text{circle} \wedge \text{texture}(b) = \text{shaded} \wedge \text{weight}(b) \geq 3) \\ & \exists b_1 \exists b_2 (\text{contains}(b_1, b_2) \wedge \text{shape}(b_1) = \text{circle} \wedge \text{shape}(b_2) = \text{circle}) \\ & \dots \wedge \dots \wedge \dots \end{aligned}$$

Beispiel — Konzeptualisierung von Krebszellen

Globale (zellbezogene) Merkmale

1. $\text{circ} \in \{1, 2, \dots, 10\}$ (Anzahl der Zellsegmente)
2. $\text{pplasm} \in \{A, B, C, D\}$ (Protoplasmatyp der Zelle)

Lokale (segmentbezogene) Merkmale

- $\text{shape}(i) \in \{\text{triangle}, \text{circle}, \text{ellipse}, \text{heptagon}, \text{square}, \text{boat}, \text{spring}\}$
(bzw. eine Baumstruktur dieser Formklassen)
- $\text{texture}(i) \in \{\text{blank}, \text{shaded}, \text{black}, \text{grey}, \text{stripes}, \text{crossed}, \text{wavy}\}$
- $\text{weight}(i) \in \{1, 2, 3, 4, 5\}$
- $\text{orient}(i) \in \{N, NE, E, SE, S, SW, W, NW\}$
- $\text{contains}(c, b_1, b_2, \dots) \in \{T, F\}$
- $\text{hastails}(c, b_1, b_2, \dots) \in \{T, F\}$

Konzeptuelle Klassifikation

Gegeben

Klassenspezifische Lernstichproben

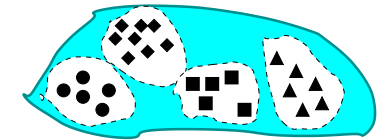
$$\omega_\kappa \subseteq C_\kappa \subseteq \Omega, \quad \kappa = 1, \dots, K$$

für die Konzepte $C_1, \dots, C_K \in \mathcal{C}$ mit $C_\kappa \cap C_\lambda = \emptyset$ für $\kappa \neq \lambda$.

Gesucht

Ein konsistentes System $h_1, \dots, h_K \in \mathcal{H}$, d.h. für alle $1 \leq \kappa \leq K$ gilt:

$$\begin{cases} h_\kappa \models \mathbf{x} & \mathbf{x} \in \omega_\kappa \\ h_\kappa \not\models \mathbf{x} & \mathbf{x} \in \omega_\lambda, \lambda \neq \kappa \end{cases} \quad \text{und „quodlibet“ sonst}$$



Versionenraum-Methode

Berechne für jede Objektklasse κ einen **diskriminativen VR**

$$\mathfrak{V}_\kappa = \mathfrak{V}(\mathcal{H}, \omega_\kappa, \bigcup_{\lambda \neq \kappa} \omega_\lambda)$$

Stern-Methode

Berechne für jedes κ eine **Sterndisjunktion**

$$h_\kappa^* = \bigcup_{\mathbf{x} \in \omega_\kappa} \mathcal{S}(\mathbf{x} \mid \omega \setminus \omega_\kappa)$$

Votierung beim K -Klassen-Problem

$$\beta(\mathbf{x}) = (\beta_1, \dots, \beta_K) \in \{1, 0, ?\}^K$$

Stern-Methode

Es gibt keine Fehlanzeigen.
Aber es gibt u.U. Konflikte.
Und es gibt u.U. Leerrunden.

Versionenraum-Methode

Es gibt Konflikte & Leerrunden.
Es gibt auch Fehlanzeigen:

$$\left\{ \begin{array}{l} \text{eine FA statt PRO} \\ \text{weniger als } K - 1 \text{ CONs} \end{array} \right\}$$

Definition

Sei $\mathcal{A} \subseteq \Omega$. Das Hypothesensystem (h_1, \dots, h_K) heißt **konzeptuelle Partition** von \mathcal{A} , wenn es für jedes $\mathbf{x} \in \mathcal{A}$ einen Klassenindex κ gibt mit

$$\forall \lambda = 1, \dots, K : (h_\lambda \models \mathbf{x} \iff \lambda = \kappa)$$

Bemerkungen

1. Ein konsistentes System (h_1, \dots, h_K) ist konzeptuelle Partition seiner Lerndaten $\bigcup_{\kappa} \omega_{\kappa}$.
2. Läßt sich jedes konsistente System zu einer konzeptuellen Partition von Ω erweitern?

Die Bayesregel

... ist der theoretisch optimale Klassifikator.

Satz

Ist für den Objektraum \mathcal{X} und das Klasseninventar $\mathcal{K} = \{1, \dots, K\}$ die wahre Verbundverteilung $P(\kappa, \mathbf{x})$ bekannt, so liefert die **Bayesentscheidungsregel** (MAP-Regel)

$$\begin{aligned} \delta(\mathbf{x}) &= \operatorname{argmax}_{\kappa \in \mathcal{K}} P(\kappa | \mathbf{x}) \\ &= \operatorname{argmax}_{\kappa \in \mathcal{K}} \frac{P(\kappa) \cdot P(\mathbf{x} | \kappa)}{P(\mathbf{x})} \end{aligned}$$

die minimale erwartete Klassifikationsfehlerrate.

Bemerkung

Die Aussage gilt natürlich nur, wenn das **korrekte** Wahrscheinlichkeitsmodell verwendet wird.

Marginal 1

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

Marginal 2

$$P(\kappa) = \sum_{\mathbf{x} \in \mathcal{X}} P(\kappa, \mathbf{x})$$

Bedingte Vtl.

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

Posterior Vtl.

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

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Die Bayesregel für diskrete Attribute

Kanonische multivariat-diskrete Verteilung („Hypertabelle“)

Lemma

Der Objektraum $\mathcal{X} = \mathcal{X}_1 \times \dots \times \mathcal{X}_N$ enthalte ausschließlich **diskrete** Attribute mit Wertebereichen \mathcal{X}_n der Größe $L_n = |\mathcal{X}_n|$.

1. Die gemeinsame Verteilung $P(\kappa, \mathbf{x})$ ist durch die $K \cdot L_1 \cdot \dots \cdot L_N$ Einträge

$$p_{\kappa, \mathbf{x}_1, \dots, \mathbf{x}_N} = P(\kappa, \mathbf{x}), \quad \kappa \in \mathcal{K}, \mathbf{x} \in \mathcal{X}$$

eines $(1 + N)$ -dimensionalen Hyperwürfels $\mathbf{P} \in [0, 1]^{K \times L_1 \times \dots \times L_N}$ charakterisiert.

2. Für einen etikettierten Lerndatensatz $\{(\kappa_t, \mathbf{x}_t) \mid t = 1..T\}$ mit den absoluten Häufigkeiten $T_{\kappa, \mathbf{x}}, (\kappa, \mathbf{x}) \in \mathcal{K} \times \mathcal{X}$, lauten die Maximum-Likelihood-Parameter

$$\hat{b}_{\kappa, \mathbf{x}_1, \dots, \mathbf{x}_N} = T_{\kappa, \mathbf{x}_1, \dots, \mathbf{x}_N} / T.$$

3. Die Bayesentscheidungsregel lautet $\delta(\mathbf{x}) = \operatorname{argmax}_{\kappa \in \mathcal{K}} \hat{b}_{\kappa, \mathbf{x}_1, \dots, \mathbf{x}_N}$.

Beispiel — kanonische Bayesregel

... mit ML-geschätzten Verteilungsparametern

Lerndatensammlung „Tennis“

ω	outlook	temp	humid	wind	Tennis?
ω_1	sunny	hot	high	weak	no
ω_2	sunny	hot	high	strong	no
ω_3	overcast	hot	high	weak	yes
ω_4	rain	mild	high	weak	yes
ω_5	rain	cool	normal	weak	yes
ω_6	rain	cool	normal	strong	no
ω_7	overcast	cool	normal	strong	yes
ω_8	sunny	mild	high	weak	no
ω_9	sunny	cool	normal	weak	yes
ω_{10}	rain	mild	normal	weak	yes
ω_{11}	sunny	mild	normal	strong	yes
ω_{12}	overcast	mild	high	strong	yes
ω_{13}	overcast	hot	normal	weak	yes
ω_{14}	rain	mild	high	strong	no
ω_{neu}	sunny	cool	high	strong	?

Parameter

$$2 \cdot 3^2 \cdot 2^2 = 72$$

Einträge

14 Einsen

58 Nullen

Neuzugang

Nulleintrag bei
(yes, ω_{neu}) und
(no, ω_{neu}).

Nennerausdruck
 $\hat{P}(\omega_{neu}) = 0$

Dann gilt für die a posteriori Wahrscheinlichkeit:

$$P(\text{no} \mid (\text{sunny}, \text{cool}, \text{high}, \text{strong})^\top) = \text{undef.}$$

Beispiel — naive Bayesregel

... mit ML-geschätzten Verteilungsparametern

Attribut „outlook“

	sunny	over	rain	Σ
yes	2	4	3	9
no	3	0	2	5
Σ	5	4	5	14

Attribut „humidity“

	high	normal	Σ
yes	3	6	9
no	4	1	5
Σ	7	7	14

Attribut „temp“

	hot	mild	cool	Σ
yes	2	4	3	9
no	2	2	1	5
Σ	4	6	4	14

Attribut „wind“

	weak	strong	Σ
yes	6	3	9
no	2	3	5
Σ	8	6	14

Parametertabelle und Neuklassifikation (6 + 6 + 4 + 4 = 20 Einträge)

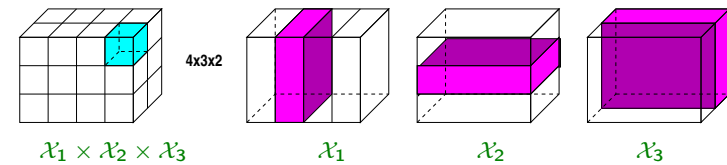
$$P(\text{no}, \text{sunny}, \text{cool}, \text{high}, \text{strong}) = \frac{5}{14} \cdot \frac{3}{5} \cdot \frac{1}{5} \cdot \frac{4}{5} \cdot \frac{3}{5} = \frac{180}{8750} = 0.02057$$

$$P(\text{yes}, \text{sunny}, \text{cool}, \text{high}, \text{strong}) = \frac{9}{14} \cdot \frac{2}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} \cdot \frac{3}{9} = \frac{486}{91854} = 0.005291$$

$$P(\text{no} \mid (\text{sunny}, \text{cool}, \text{high}, \text{strong})^\top) = \frac{0.02057}{0.02057 + 0.005291} = 0.7954$$

Die naive Bayesregel

Klassenbedingte statistische Unabhängigkeit zwischen allen Objektattributen



NBK-Entscheidungsregel

$$\delta(\mathbf{x}) = \operatorname{argmax}_{\kappa \in \mathcal{K}} P(\kappa, \mathbf{x}) = \operatorname{argmax}_{\kappa \in \mathcal{K}} \left\{ P(\kappa) \cdot \prod_{n=1}^N P(x_n \mid \kappa) \right\}$$

(maximale faktorisierte Verbundwahrscheinlichkeit)

Modellparameter und ihre ML-Schätzwerte

$$\hat{a}_\kappa = \frac{T_\kappa}{T}, \quad \hat{b}_{\xi \mid \kappa, n} = \frac{T_{\kappa, n, \xi}}{T_\kappa}, \quad T_\kappa = \sum_{\xi \in \mathcal{X}_1} T_{\kappa, 1, \xi}, \quad \begin{cases} \kappa = 1..K \\ n = 1..N \\ \xi = 1..L_n \end{cases}$$

Das sind $K \cdot \sum_n L_n$ Parameter statt $K \cdot \prod_n L_n$ Parameter!

NTF — Nichtnegative Tensorfaktorisierung

Mischung naiver Verbundverteilungen von $N \in \{2, 3\}$ nominalen Attributen

Matrix ($\Omega_1 \times \Omega_2$)

Verteilungsparameter

$$P(i, j) =: x_{ij}$$

Naive Faktorisierung

$$P(i, j) = p_1(i) \cdot p_2(j)$$

Mischungsmodell

$$P(i, j) = \sum_{m=1}^M \underbrace{\pi_m}_{v_{im}} \cdot \underbrace{p_1^{(m)}(i) \cdot p_2^{(m)}(j)}_{a_{jm}}$$

Reduktion

$$L_1 \cdot L_2 \rightarrow M \cdot (1 + L_1 + L_2)$$

Würfel ($\Omega_1 \times \Omega_2 \times \Omega_3$)

Verteilungsparameter

$$P(i, j, k) =: x_{ijk}$$

Naive Faktorisierung

$$P(i, j, k) = p_1(i) \cdot p_2(j) \cdot p_3(k)$$

Mischungsmodell

$$P(i, j, k) = \sum_{m=1}^M \pi_m \cdot p_1^{(m)}(i) \cdot p_2^{(m)}(j) \cdot p_3^{(m)}(k)$$

Reduktion

$$L_1 \cdot L_2 \cdot L_3 \rightarrow M \cdot (1 + L_1 + L_2 + L_3)$$

NTF — Nichtnegative Tensorfaktorisierung

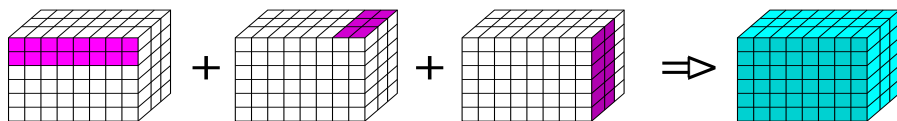
Mischung naiver Verbundverteilungen nominaler Attribute

Wahrscheinlichkeitshyperwürfel ($\Omega_1 \times \Omega_2 \times \dots \times \Omega_N$)

Naive Mischung

$$P(x_1, \dots, x_N) = \sum_{m=1}^M \pi_m \cdot \prod_{n=1}^N p_n^{(m)}(x_n)$$

Parameter lernen nach EM-Prinzip (*expectation-maximization*)



Reduktion

$$L_1 \cdot \dots \cdot L_N \Rightarrow M \cdot (1 + L_1 + \dots + L_N)$$

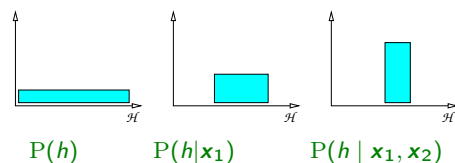
$$\Pi \rightsquigarrow M \cdot \Sigma$$

Parametrische (modellgetriebene) Bayesregel

Optimale Entscheidungsregel unter der Annahme $C^* \in \mathcal{H} \neq \mathcal{C}$

Problem

Die Auswahl und exklusive Nutzung der bestpassenden Hypothese aus \mathcal{H} ist willkürlich und wegen $|\omega| < \infty$ suboptimal.



Bayestheorie

1. Hypothesen sind nicht nur wahr oder falsch
2. Hypothesen treffen „weiche“ Entscheidungen
3. Jedes Lernbeispiel erhöht/vermindert inkrementell die Hypothesenwahrscheinlichkeiten
4. Vorwissen lässt sich mit den Lerndaten verzahnen
5. Mathematisch abgesicherter Votierungsmechanismus

EM-Algorithmus für das NTF-Modell

1 Initialisierung · 2 E-Schritt · 3 M-Schritt · 4 Abbruch

A posteriori Wahrscheinlichkeiten der Komponentenauswahl

Für jedes Lerndatensatzobjekt $\mathbf{x}_1, \dots, \mathbf{x}_T$ berechne

$$\gamma_t(m) \stackrel{\text{def}}{=} P(\mathbb{M} = m \mid \mathbf{x}_t, \theta^{\text{alt}}) = \pi_m^{\text{alt}} \cdot \prod_{n=1}^N \theta_{m,n,x_{tn}}^{\text{alt}} / P^{\text{alt}}(\mathbf{x}_t)$$

Neuschätzung durch a posteriori Erwartungswerte

$$\hat{\pi}_m = \sum_{t=1}^T \gamma_t(m) / T \quad \text{und} \quad \hat{\theta}_{m,n,\xi} = \sum_{t=1}^T \gamma_t(m) \cdot \mathbf{1}_{x_{tn}=\xi} / \sum_{t=1}^T \gamma_t(m)$$

Startparameter

zufällig · wiederholt · lokale Optima

Rechenaufwand

$$O(I_{\max} \cdot T \cdot M \cdot (N + \sum_n L_n))$$

Hypothesen- und Objektwahrscheinlichkeiten

Lemma (A priori und a posteriori Hypothesenwahsch'keit)

Für alle Hypothesen $h \in \mathcal{H}$ und gegebene Daten $\omega \subset \Omega$ gilt die folgende Darstellung der **a posteriori** Hypothesenwahrscheinlichkeit:

$$P(h|\omega) = \frac{P(\omega|h) \cdot P(h)}{P(\omega)}$$

„Deduktion“ $\Uparrow P(h)$ · „Abduktion“ $\Downarrow P(\omega|h)$ · „explaining-away“ $\Downarrow P(\omega)$

Definition

Für ein K -Klassenproblem über Ω mit Hypothesenraum \mathcal{H} und den Lerndaten ω heißt

$$\delta_{\text{Bayes}}(\mathbf{x}) = \operatorname{argmax}_{\kappa} \sum_{h \in \mathcal{H}} P(\mathbf{x} \in C_{\kappa} \mid h) \cdot P(h|\omega)$$

die **Bayes-Entscheidungsregel** für das Objekt \mathbf{x} .

Bemerkung

Diese Entscheidungsregel realisiert die **minimale asymptotische Fehlerrate**.

Hypothesenauswahltechniken

MDL-Prinzip

Minimum Description Length (Rissanen'87)

$$\begin{aligned} h_{\text{MAP}} &= \operatorname{argmax}_{h \in \mathcal{H}} P(\omega|h) \cdot P(h) \\ &= \operatorname{argmin}_{h \in \mathcal{H}} \{-\log_2 P(\omega|h) - \log_2 P(h)\} \end{aligned}$$

ML-Schätzer

$$h_{\text{ML}} = \operatorname{argmax}_{h \in \mathcal{H}} P(\omega|h)$$

MAP-Schätzer

$$h_{\text{MAP}} = \operatorname{argmax}_{h \in \mathcal{H}} P(h|\omega)$$

Gibbs-Sampler

$$h_{\text{GS}} \sim P(h|\omega)$$

Versionenraum

$$h_{\text{VR}} \in \mathfrak{V}(\mathcal{H}, \omega)$$

- 1 Wähle eine Codierung \mathcal{C}_1 für die Hypothesen
- 2 Wähle eine bedingte Codierung \mathcal{C}_2 für die Klassen
- 3 Berechne die „kürzeste Erklärung“

$$h_{\text{MDL}} \stackrel{\text{def}}{=} \operatorname{argmin}_{h \in \mathcal{H}} \{\ell_{\mathcal{C}_1}(h) + \ell_{\mathcal{C}_2}(\omega|h)\}$$

der Lerndaten.

Algorithmus

zumitragla

Bayesregel für gemischte Attributskalen

$$\mathbf{z} = (\mathbf{x}, \mathbf{y}) \in \underbrace{\mathbb{R} \times \dots \times \mathbb{R}}_{\Omega' = \mathbb{R}^{N'}} \times \underbrace{\mathcal{X}_1 \times \dots \times \mathcal{X}_{N''}}_{\Omega''}$$

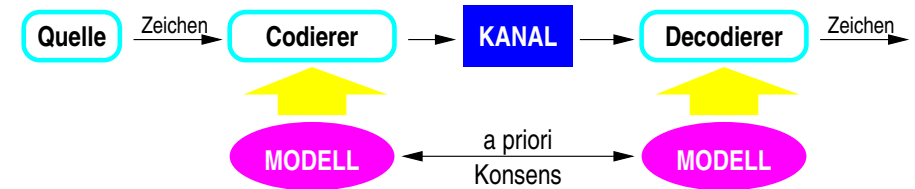
Normalzerlegung

$$P(\mathbf{z}) = P(\mathbf{y}) \cdot P(\mathbf{x}|\mathbf{y}) = P(y_1, \dots, y_{N''}) \cdot \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_{\mathbf{y}}, \mathbf{S}_{\mathbf{y}})$$

- $L^\times = \prod_{n=1}^{N''} L_n$ kanonische W'keitsparameter
- $N' + \binom{N'}{2}$ Dichteparameter je NV-Dichte
- insgesamt $O(K \cdot L^\times \cdot N'^2)$ Parameter
- Unabhängigkeitsannahme für Ω'' bringt wenig Vorteile.
- Unabhängigkeit in Ω' reduziert auf $O(K \cdot L^\times \cdot N')$ Parameter.

Noiseless Source Coding Theorem

Quellenentropie \triangleq minimale Bitanzahl nach optimaler Datenkompression



Satz (Shannon 1949)

Ein Zufallsprozeß erzeuge Zeichenfolgen über dem Alphabet $\{s_1, \dots, s_L\}$ mit den Wahrscheinlichkeiten q_1, \dots, q_L .

1. Die **optimale** Codierung dieser Quelle verwendet für jedes Zeichen s_i ein Codewort der Länge $-\log_2 q_i$ Bit.
2. Ihre mittlere Codewortlänge beträgt $\mathcal{H}(q_1, \dots, q_L)$ Bit.

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Diskriminative Klassifikatoren

$$\kappa(\mathbf{x}) = \operatorname{argmax}_{\lambda} h_{\lambda}(\mathbf{x})$$

Definition

Es sei C_1, \dots, C_K ein K -Klassen-Problem über $\Omega = \mathbb{R}^N$. Die Elemente von $\mathbf{h} = (h_1, \dots, h_K)^{\top}$ mit

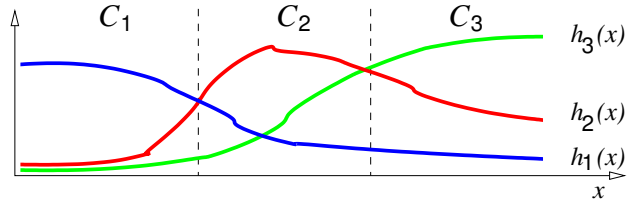
$$h_{\kappa} : \mathbb{R}^N \rightarrow \mathbb{R}, \quad \kappa = 1, \dots, K$$

heißen **Trennfunktionen** der Klassen $\kappa = 1, \dots, K$.

Die Abbildungen $\mathbf{d} = (d_1, \dots, d_K)^{\top}$ mit

$$d_{\kappa} : \mathbf{x} \mapsto \begin{cases} 1 & \mathbf{x} \in C_{\kappa} \\ 0 & \mathbf{x} \notin C_{\kappa} \end{cases}$$

heißen **ideale Trennfunktionen** des Problems.



Lernen als skalare Regressionsaufgabe

Zerlegung in Zweiklassenprobleme

Für jedes κ ergibt sich das QM-Approximationsproblem

$$h_{\kappa} \approx d_{\kappa} : \mathbf{x} \mapsto \begin{cases} 1 & \mathbf{x} \in \omega^+, \omega^+ = \omega_{\kappa} \\ 0 & \mathbf{x} \in \omega^-, \omega^- = \bigcup_{\lambda \neq \kappa} \omega_{\lambda} \end{cases}$$

Skalares Regressionsproblem

Für die Daten $\{(\mathbf{x}_t, y_t) \in \mathbb{R}^N \times \mathbb{R} \mid t = 1, \dots, T\}$ finde die Regressionsfunktion $h \in \mathcal{H}$ mit minimalem Fehler

$$\varepsilon(h) \stackrel{\text{def}}{=} \sum_{t=1}^T (y_t - h(\mathbf{x}_t))^2$$

$$\left\{ \begin{array}{l} \text{quadratisches Fehlermaß} \\ \{0, 1\}\text{-Zielgröße} \end{array} \right\} \Rightarrow \left\{ \begin{array}{l} \text{entkoppelte Diskriminantfunktionen} \\ \text{kein Problem mit } K > 2 \end{array} \right\}$$

Quadratmittelklassifikator

Willkürlicher Zielausdruck — willkürliches Straffunktional

Definition

Es sei C_1, \dots, C_K ein K -Klassen-Problem über $\Omega = \mathbb{R}^N$ und \mathcal{H} eine Menge von Trennfunktionen. Die Trennfunktion $\mathbf{h} \in \mathcal{H}$ mit minimalem erwarteten quadratischen Fehler

$$\varepsilon(\mathbf{h}) \stackrel{\text{def}}{=} \mathcal{E}[\|\mathbf{h}(\mathbb{X}) - \mathbf{d}(\mathbb{X})\|^2]$$

heißt **Quadratmitteldiskriminante**, der zugehörige Klassifikator heißt **Quadratmittelklassifikator**.

Sind ferner die Lerndaten $\omega_1, \dots, \omega_K$ gegeben, so heißt der Klassifikator mit minimalem Fehler

$$\varepsilon(\mathbf{h}, \{\omega_{\kappa}\}) \stackrel{\text{def}}{=} \sum_{\kappa=1}^K \sum_{\mathbf{x} \in \omega_{\kappa}} \|\mathbf{h}(\mathbf{x}) - \mathbf{e}^{(\kappa)}\|^2$$

empirischer QMK. Dabei bezeichne $\mathbf{e}^{(\kappa)}$ den κ -ten Einheitsvektor.

Multivariate lineare Regression

Linearer Ansatz

$$h(\mathbf{x}) = \sum_n a_n \cdot x_n = \mathbf{a}^{\top} \mathbf{x}, \quad \mathbf{x}, \mathbf{a} \in \mathbb{R}^N$$

Affiner Ansatz

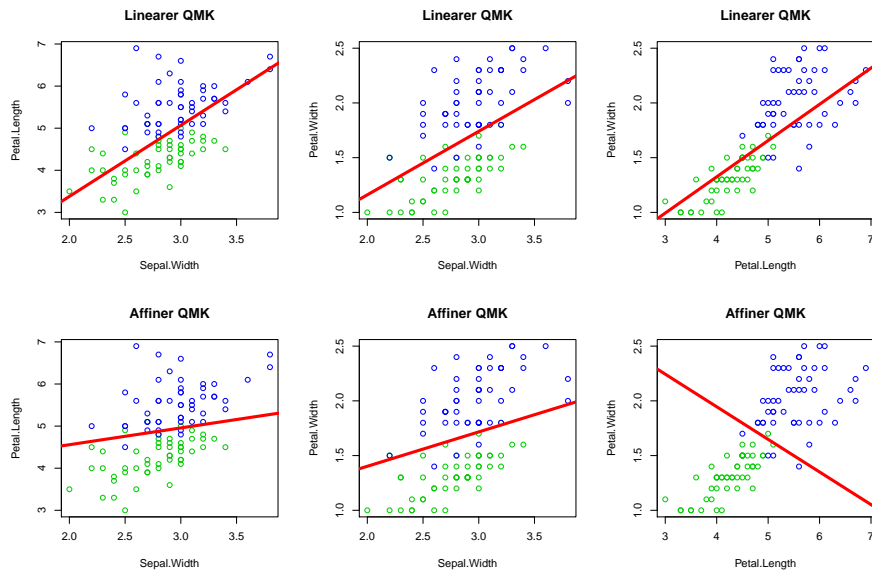
$$h(\mathbf{x}) = a_0 + \sum_n a_n \cdot x_n = \mathbf{a}^{\top} \mathbf{x}', \quad \left\{ \begin{array}{l} \mathbf{a} \in \mathbb{R}^{N+1} \\ \mathbf{x}' \stackrel{\text{def}}{=} (1, \mathbf{x}^{\top})^{\top} \end{array} \right.$$

Was heißt hier eigentlich „Regressionsproblem“ ?

- Datenmodell $\mathbb{Y} = h(\mathbb{X}_1, \dots, \mathbb{X}_N) + \mathbb{E}$ mit Störterm $\mathbb{E} \sim \mathcal{N}(0, \sigma^2)$
- Datenprobe $\{\mathbf{x}_t, y_t\}_1^T$ bzw. gemeinsame Datenverteilung $f_{\mathbb{X}, \mathbb{Y}}(\cdot, \cdot)$
- ➔ Posterior-Erwartungswerte $\hat{h}(\mathbf{x}) = \mathcal{E}_{\mathbb{Y}|\mathbf{x}}[\mathbb{Y}]$, also $\hat{h}(\mathbf{x}) = \int f(y|\mathbf{x}) \cdot y \, dy$

Linearer vs. affiner Quadratmittelklassifikator

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



Linearer versus affiner Ansatz

Lineare Funktionen allein beschreiben wegen $h(\mathbf{0}) = 0$ ausschließlich Hyperebenen, die durch den Koordinatenursprung verlaufen und sind als Regressionsmodell unzureichend. Affine Funktionen verfügen zusätzlich über den y -Schnittpunkt a_0 (*intercept*); affine Regression kann aber leicht auf lineare Regression zurückgeführt werden. Wir verwenden wieder die Notation \mathbf{X} , \mathbf{y} für Datenmatrix und Zielwertvektor und betrachten den Abweichungsvektor $\mathbf{y} - a_0\mathbf{1} - \mathbf{X}\mathbf{a}$ des affinen Modells sowie den resultierenden quadratischen Fehler:

$$\begin{aligned}\varepsilon(a_0, \mathbf{a}) &= \|\mathbf{y} - a_0\mathbf{1} - \mathbf{X}\mathbf{a}\|^2 \\ &= (\mathbf{y} - a_0\mathbf{1} - \mathbf{X}\mathbf{a})^\top \cdot (\mathbf{y} - a_0\mathbf{1} - \mathbf{X}\mathbf{a}) \\ &= (\mathbf{y}^\top - a_0\mathbf{1}^\top - \mathbf{a}^\top \mathbf{X}^\top) \cdot (\mathbf{y} - a_0\mathbf{1} - \mathbf{X}\mathbf{a}) \\ &= \underbrace{\|\mathbf{y}\|^2}_c + \underbrace{a_0^2 T - 2a_0 T \mu_y}_{\varepsilon(a_0)} + \underbrace{\mathbf{a}^\top \mathbf{X}^\top \mathbf{X} \mathbf{a} - 2\mathbf{a}^\top \mathbf{X}^\top \mathbf{y}}_{\varepsilon(\mathbf{a})} + \underbrace{2a_0 T \mu_x^\top \mathbf{a}}_0\end{aligned}$$

Wenn wir o.B.d.A. mittelwertfreie Vektordaten annehmen ($\mu_x = \mathbf{0}$), so verschwindet der Kopplungsterm und wir dürfen a_0 und \mathbf{a} separat optimieren. Für a_0 ergibt sich nach Nullsetzen der Ableitung

$$\partial \varepsilon(a_0) / \partial a_0 = 2Ta_0 - 2T\mu_y$$

der Minimalwert $a_0 = \mu_y$. Der Fehler $\varepsilon(\mathbf{a})$ und die Konstante $\|\mathbf{y}\|^2$ ergeben zusammen den Minimierungsausdruck der linearen Regressionsaufgabe ...

Multivariate lineare Regression

Lösen des Systems der Gaußschen Normalgleichungen

Satz

Es seien die Regressionsdaten $(\mathbf{x}_t, y_t) \in \mathbb{R}^N \times \mathbb{R}$, $t = 1, \dots, T$ in der Matrixnotation

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)^\top, \quad \mathbf{y} = (y_1, \dots, y_T)^\top$$

dargestellt, und es sei $h: \mathbf{x} \mapsto \mathbf{a}^\top \mathbf{x}$ linear. Dann lautet der quadratische Regressionsfehler

$$\varepsilon(h) = \varepsilon(\mathbf{a}) = \|\mathbf{y} - \mathbf{X}\mathbf{a}\|^2$$

und wird durch jede Lösung \mathbf{a} der **Gaußschen Normalgleichungen**

$$\mathbf{X}^\top \mathbf{X} \mathbf{a} = \mathbf{X}^\top \mathbf{y}$$

minimiert.

Beweis.

Zur Lösung des Quadratmittelpblems setzen wir die partiellen Ableitungen der Koeffizienten a_1, \dots, a_N gleich Null — wir verwenden die Gradientenvektorschreibweise:

$$\begin{aligned}\nabla_{\mathbf{a}} \varepsilon(\mathbf{a}) &= \nabla_{\mathbf{a}} \left\{ \|\mathbf{y}\|^2 + \mathbf{a}^\top \mathbf{X}^\top \mathbf{X} \mathbf{a} - 2\mathbf{y}^\top \mathbf{X} \mathbf{a} \right\} \\ &= \mathbf{0} + 2\mathbf{X}^\top \mathbf{X} \mathbf{a} - 2\mathbf{X}^\top \mathbf{y} \\ &= 2 \cdot (\mathbf{X}^\top \mathbf{X} \mathbf{a} - \mathbf{X}^\top \mathbf{y}) \\ &\stackrel{!}{=} \mathbf{0}\end{aligned}$$

□

Bemerkung

Das LGS $\mathbf{X}^\top \mathbf{X} \mathbf{a} = \mathbf{X}^\top \mathbf{y}$ heißt System der *Gaußschen Normalgleichungen*. Wir schreiben auch kürzer $\mathbf{R}\mathbf{a} = \mathbf{m}$; dabei ist \mathbf{R} wieder die unzentrierte, unnormierte Kovarianzmatrix der Vektordaten.

Ausgleichsrechnung und Lineare Gleichungssysteme

Was Sie schon immer über lineare Algebra wissen wollten, aber nie zu fragen wagten

$$\mathbf{X} \cdot \mathbf{a} \stackrel{!}{=} \mathbf{y} \quad \text{mit dem Fehlervektor } \mathbf{e} := \mathbf{X}\mathbf{a} - \mathbf{y}$$

LGS eindeutig

Matrix \mathbf{X} ist quadratisch und vollrangig.



$f: \mathbf{z} \mapsto \mathbf{X}\mathbf{z}$ bijektiv

$\mathbf{a} = \mathbf{X}^{-1}\mathbf{y}$ ist die eindeutige Lösung mit Fehler $\mathbf{e} = \mathbf{0}$.

Bemerkung

Das Gaußsche Normalgleichungssystem ist entweder eindeutig lösbar oder besitzt unendlich viele Lösungen.

... überbestimmt

Matrix \mathbf{X} hat den vollen Spaltenrang.



$f: \mathbf{z} \mapsto \mathbf{X}\mathbf{z}$ ist injektiv
 $\mathbf{X}^\top \mathbf{X}$ ist regulär

$\mathbf{a} = (\mathbf{X}^\top \mathbf{X})^{-1} \cdot \mathbf{X}^\top \mathbf{y}$ ist eine Lösung mit minimalem Gesamtfehler $\|\mathbf{e}\|$.

... unterbestimmt

Matrix \mathbf{X} hat den vollen Zeilenrang.



$f: \mathbf{z} \mapsto \mathbf{X}\mathbf{z}$ ist surjektiv
 $\mathbf{X}\mathbf{X}^\top$ ist regulär

$\mathbf{a} = \mathbf{X}^\top (\mathbf{X}\mathbf{X}^\top)^{-1} \cdot \mathbf{y}$ ist eine Lösung mit Fehler $\mathbf{e} = \mathbf{0}$ und minimaler Länge $\|\mathbf{a}\|$.

Minimalnormlösung des GNG-Systems

Lemma

Es sei $\mathbf{R} \cdot \mathbf{a} = \mathbf{m}$ bzw. $\mathbf{X}^\top \mathbf{X} \cdot \mathbf{a} = \mathbf{X}^\top \mathbf{y}$ das GNG-System einer linearen Regressionsaufgabe.

1. Die Matrix \mathbf{R} ist symmetrisch und positiv-semidefinit.
2. Das Gleichungssystem hat stets mindestens eine Lösung.
3. Ist \mathbf{R} invertierbar, so existiert eine eindeutige Lösung:

$$\mathbf{a}^* = \mathbf{R}^{-1} \cdot \mathbf{m}$$

4. Ist \mathbf{X}^+ die Pseudoinverse der Datenmatrix, so löst

$$\mathbf{a}^+ = \mathbf{X}^+ \cdot \mathbf{y}$$

das Gleichungssystem und besitzt unter allen Lösungen die minimale Norm $\|\mathbf{a}^+\|$.

Die Berechnung der Minimalnormlösung ist **nicht praktikabel!**

System der Gaußschen Normalgleichungen

Linearer Quadratmittelklassifikator ($K = 2$)

$$\mathbf{R} \cdot \mathbf{a} = \mathbf{m}$$

$$\mathbf{R} = \frac{1}{T} \cdot \mathbf{X}^\top \mathbf{X} = \mathbf{S} + \boldsymbol{\mu} \boldsymbol{\mu}^\top$$

$$\mathbf{m} = \frac{1}{T} \cdot \mathbf{X}^\top \mathbf{y} = \frac{1}{T} \cdot \sum_{\omega^+} \mathbf{x}_t = \mathbf{p}^+ \cdot \boldsymbol{\mu}^+$$

$\mathbf{p}^+ = T^+/T$, $\boldsymbol{\mu}^+$ Positivstatistiken; \mathbf{R} Momentenmatrix der Gesamtprobe.

Linearer Quadratmittelklassifikator ($K > 2$)

$$\mathbf{R} \cdot \mathbf{a}_1 = \mathbf{m}_1$$

$$\vdots = \vdots$$

$$, \quad \mathbf{m}_K = \frac{1}{T} \cdot \sum_{\omega_K} \mathbf{x}_t = \mathbf{p}_K \cdot \boldsymbol{\mu}_K$$

$$\mathbf{R} \cdot \mathbf{a}_K = \mathbf{m}_K$$

Kompaktschreibweise: $\mathbf{R} \cdot \mathbf{A} = \mathbf{M}$ mit $\mathbf{M} = (\mathbf{p}_1 \boldsymbol{\mu}_1, \dots, \mathbf{p}_K \boldsymbol{\mu}_K)$.

Beweis.

Für eine beliebige Rechteckmatrix mit der SV-Zerlegung $\mathbf{X} = \mathbf{V}\mathbf{D}\mathbf{U}^\top$ heißt die Matrix

$$\mathbf{X}^+ = \mathbf{U}\mathbf{D}^+ \mathbf{V}^\top$$

die **Moore-Penrose-Inverse** oder **Pseudoinverse**. Die Pseudoinverse \mathbf{D}^+ einer Diagonalmatrix \mathbf{D} wiederum enthält auf ihrer Diagonalen die Pseudo-Reziproken:

$$d_n^+ = \begin{cases} 1/d_n & d_n \neq 0 \\ 0 & d_n = 0 \end{cases}, \quad n = 1, \dots, N$$

Diese Pseudoinverse gehorcht der **Moore-Penrose-Gleichung**, denn es gilt:

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} \mathbf{X}^+ &= \mathbf{U}\mathbf{D}\mathbf{V}^\top \cdot \mathbf{V}\mathbf{D}\mathbf{U}^\top \cdot \mathbf{U}\mathbf{D}^+ \mathbf{V}^\top \\ &= \mathbf{U}\mathbf{D}^2 \mathbf{D}^+ \mathbf{V}^\top = \mathbf{U}\mathbf{D}\mathbf{V}^\top = \mathbf{X}^\top \end{aligned}$$

Folglich löst $\mathbf{a}^+ = \mathbf{X}^+ \mathbf{y}$ auch die Gaußschen Normalgleichungen:

$$\mathbf{T} \cdot \mathbf{R} \mathbf{a}^+ = \mathbf{X}^\top \mathbf{X} \cdot \mathbf{X}^+ \mathbf{y} = \mathbf{X}^\top \cdot \mathbf{y} = \mathbf{T} \cdot \mathbf{m}$$

Der Beweis der Minimaleigenschaft erfordert einen Lagrange-Ansatz:

$$\frac{1}{2} \cdot \|\mathbf{a}\|^2 + \lambda \cdot \left\| \mathbf{X}^\top \mathbf{X} \mathbf{a} - \mathbf{X}^\top \mathbf{y} \right\|^2 \stackrel{!}{\rightarrow} \text{MIN}$$

Gratregularisierung

Lemma

Der regularisierte quadratische Regressionsfehler

$$\varepsilon_\lambda(\mathbf{a}) = \|\mathbf{y} - \mathbf{X}\mathbf{a}\|^2 + \lambda \cdot \|\mathbf{a}\|^2$$

($\lambda > 0$) wird durch die (eindeutige) Lösung

$$\mathbf{a}_\lambda^* = (\mathbf{R}_\lambda)^{-1} \cdot \mathbf{m}, \quad \mathbf{R}_\lambda \stackrel{\text{def}}{=} \mathbf{R} + \lambda \mathbf{E}$$

minimiert.

Beweis.

Der regularisierte Quadratmittelfehler besitzt den Gradientenvektor

$$\nabla_{\mathbf{a}} \varepsilon_\lambda(\mathbf{a}) = \nabla_{\mathbf{a}} \varepsilon(\mathbf{a}) + \lambda \cdot \nabla_{\mathbf{a}} \|\mathbf{a}\|^2 = 2 \cdot (\mathbf{R}\mathbf{a} + \lambda \mathbf{a} - \mathbf{m}) = 2 \cdot (\mathbf{R}_\lambda \mathbf{a} - \mathbf{m})$$

Die Gratregularisierungsmatrix ist für $\lambda \neq 0$ stets invertierbar, denn wegen

$$\mathbf{R}_\lambda = \mathbf{R} + \lambda \mathbf{E} = \mathbf{U} \mathbf{D}^2 \mathbf{U}^\top + \lambda \mathbf{U} \mathbf{E} \mathbf{U}^\top = \mathbf{U} \cdot (\mathbf{D}^2 + \lambda \mathbf{E}) \cdot \mathbf{U}^\top = \mathbf{U} \cdot (\mathbf{D}^2)_\lambda \cdot \mathbf{U}^\top$$

besitzen alle Eigenwerte von \mathbf{R}_λ die Form $d_n^2 + \lambda > 0$. \square

Lineare versus nichtlineare QMK

Angriffspunkt: 1.Quellvariable 2.Berechnungsweg 3.Zielvariable

Termexpansion

GNGS für alle Koeffizienten

- Linear & affin $O(N^1)$
- Quadratisch $O(N^2)$
- Kubisch $O(N^3)$
- Polynomansatz $O(\binom{N+p}{p})$

Nominale Attribute

Kontrastmatrizen ($L_n - 1$)

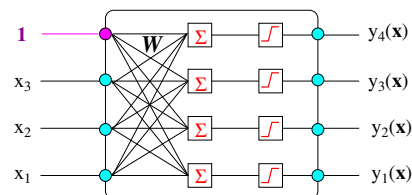
- ohne Interaktionsterme $O(\ell)$, $\ell = \sum_n L_n$
- einfache Interaktionsterme $O(\ell^2)$

Neuronale

Berechnungsmodelle

Error Backpropagation

- Mehrschichtenperzeptron
- Radiale Basisfunktionen
- Time-Delay Neural Network



Gelenkfunktion $\phi(\mathbf{y}) = \mathbf{x}^\top \mathbf{a}$

Generalized Linear Model

Gewichtete & nichtquadratische Regression

Historische Wurzeln des IRLS: „Iteratively Reweighted Least Squares“

Quadratischer Fehler	$\sum_t (\mathbf{x}_t^\top \mathbf{a} - y_t)^2$	$=$	$\ \mathbf{X}\mathbf{a} - \mathbf{y}\ _2^2$
Allgemeiner L_p -Fehler	$\sum_t \mathbf{x}_t^\top \mathbf{a} - y_t ^p$	$=$	$\ \mathbf{X}\mathbf{a} - \mathbf{y}\ _p^p$
Gewichteter Fehler	$\sum_t w_t^2 \cdot (\mathbf{x}_t^\top \mathbf{a} - y_t)^2$	$=$	$\ \mathbf{W} \cdot (\mathbf{X}\mathbf{a} - \mathbf{y})\ _2^2$

$$\mathbf{W} = \text{diag}(w_1, \dots, w_T)$$

Gewichtete Ausgleichsrechnung

Wegen $\|\mathbf{W} \cdot (\mathbf{X}\mathbf{a} - \mathbf{y})\|_2^2 = \|\mathbf{W}\mathbf{X}\mathbf{a} - \mathbf{W}\mathbf{y}\|_2^2 = \|\tilde{\mathbf{X}}\mathbf{a} - \tilde{\mathbf{y}}\|_2^2$ lautet der

Lösungskoeffizientenvektor $\mathbf{a} = (\mathbf{X}^\top \mathbf{W}^\top \mathbf{W} \mathbf{X})^{-1} \cdot \mathbf{X}^\top \mathbf{W}^\top \mathbf{W} \mathbf{y}$

Ausgleichsrechnung in der L_p -Fehlernorm (Betrag/Minimum)

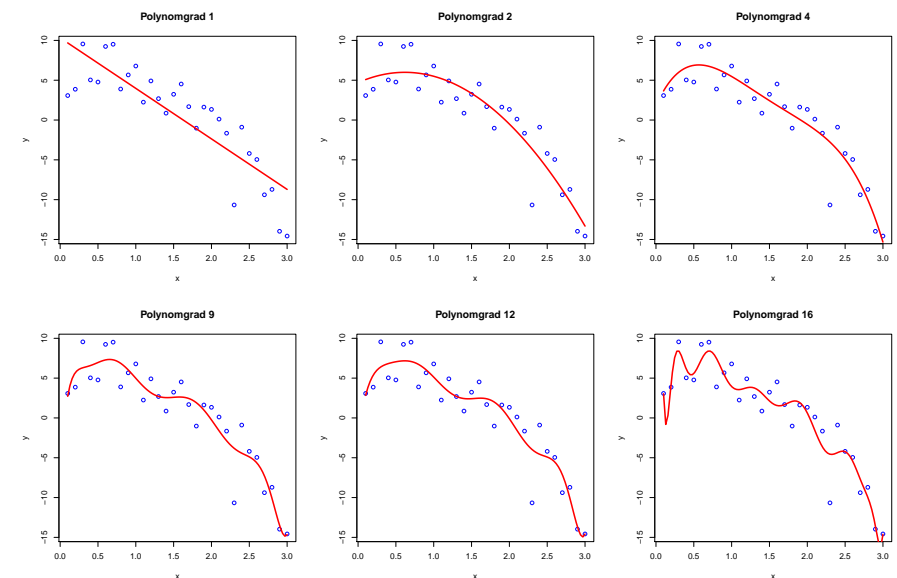
Die Fehlerminimierung kann wegen

$$\|\mathbf{e}\|_p^p = \sum_t |e_t|^p = \sum_t |e_t|^{p-2} |e_t|^2 = \sum_t w_t^2 e_t^2$$

auf IRLS mit Gewichten $w_t = |e_t|^{(p-2)/2}$ zurückgeführt werden.

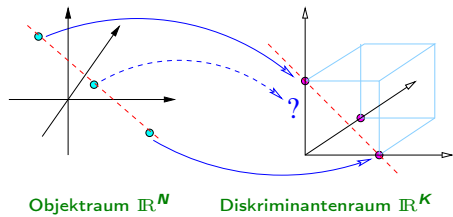
Überanpassungseffekt bei Ausgleichspolynomen

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



Maskierungseffekt

Lineare Quadratmitteldiskriminanten in Mehrklassensituationen



Kollineare
Klassenzentren

$$\mu_{\kappa} = \beta_{\kappa} \mathbf{a} + \mathbf{b}$$

↪ kollineare

Diskriminatenvektoren

$$\mathbf{h}(\mu_{\kappa}) = \beta_{\kappa} \tilde{\mathbf{a}} + \tilde{\mathbf{b}}$$

(ideal = κ -te Einheitsvektoren)

Quadratmittel

Minimiere den
quadratischen
Vorhersagefehler

$$\sum_t (y_t - \mathbf{a}^T \mathbf{x}_t)^2$$

- ⊖ negative Werte!
- ⊖ nicht normiert!

Logit

Maximiere
Datenwahrsc'hkeit

$$\prod_t P(y_t | \mathbf{x}_t)$$

mit Posterior-Wahrscheinlichkeiten der
logistischen Form $P(\Omega_1 | \mathbf{x}) \propto e^{\mathbf{a}^T \mathbf{x}}$

Probit

Maximiere
Wahrsc'hkeitssumme

$$\sum_t P(y_t | \mathbf{x}_t)$$

Lineare logistische Regression

Zweiklassenmodell

Lineares Vorhersagemodell für die **log-odds**

$$\log \frac{P(\Omega_1 | \mathbf{x})}{P(\Omega_0 | \mathbf{x})} \stackrel{!}{=} h(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$$

Mehrklassenmodell

$K - 1$ Modelle für logarithmierte
Kontrastwahrscheinlichkeiten

$$\log \frac{P(\Omega_{\lambda} | \mathbf{x})}{P(\Omega_K | \mathbf{x})} \stackrel{!}{=} h_{\lambda}(\mathbf{x}) = \mathbf{a}_{\lambda}^T \mathbf{x}$$

für alle $1 \leq \lambda < K$.

Konsistente W'keiten

Alle $P(\Omega_{\lambda} | \mathbf{x}) \in [0, 1]$.

Alle Odds $\in [0, +\infty]$.

Log-odds $\in [-\infty, +\infty]$.

Umkehrformeln

Alle Klassen $\lambda \neq K$

$$p_{\lambda}(\mathbf{x}) = \frac{\exp(\mathbf{a}_{\lambda}^T \mathbf{x})}{1 + \sum_{\kappa \neq K} \exp(\mathbf{a}_{\kappa}^T \mathbf{x})}$$

Referenzklasse $\lambda = K$

$$p_{\lambda}(\mathbf{x}) = \frac{1}{1 + \sum_{\kappa \neq K} \exp(\mathbf{a}_{\kappa}^T \mathbf{x})}$$

bzw. $p_1(\mathbf{x}) =$

$$1 - p_0(\mathbf{x}) = \frac{e^{\mathbf{a}^T \mathbf{x}}}{(1 + e^{\mathbf{a}^T \mathbf{x}})}.$$

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Maximum-Likelihood-Schätzung

Vereinfachter Fall: $K = 2$

Lemma

Für das binäre logistische Modell $p_1(\mathbf{x}) \propto \exp(\mathbf{a}^T \mathbf{x})$ mit den Lerndaten $(\mathbf{x}_t, y_t) \in \mathbb{R}^N \times \{1, 0\}$, $t = 1, \dots, T$ gilt:

1. Die ML-Zielgröße besitzt die Darstellung

$$\ell(\mathbf{a}) = \log \prod_{t=1}^T p_{y_t}(\mathbf{x}_t) = \sum_{t=1}^T \left\{ y_t \cdot \mathbf{a}^T \mathbf{x}_t - \log [1 + \exp(\mathbf{a}^T \mathbf{x}_t)] \right\}.$$

2. Für ihren Gradientenvektor der partiellen Ableitungen gilt:

$$\nabla_{\mathbf{a}} = \frac{\partial \ell(\mathbf{a})}{\partial \mathbf{a}} = \sum_{t=1}^T \mathbf{x}_t \cdot (y_t - p_1(\mathbf{x}_t)) = \mathbf{X}^T \cdot (\mathbf{y} - \mathbf{p})$$

3. Für ihre Hessematrix der gemischten partiellen Ableitungen gilt

$$\mathbf{H}_{\mathbf{a}} = \frac{\partial^2 \ell(\mathbf{a})}{\partial \mathbf{a} \partial \mathbf{a}^T} = - \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^T \cdot p_1(\mathbf{x}_t) \cdot (1 - p_1(\mathbf{x}_t)) = -\mathbf{X}^T \mathbf{W} \mathbf{X},$$

wobei $\mathbf{W} = \text{diag}(w_1, \dots, w_T)$ und $w_t = p_1(\mathbf{x}_t) \cdot (1 - p_1(\mathbf{x}_t))$ bezeichne.

Der IRLS-Algorithmus

„Iteratively Reweighted Least Squares“

(Algorithmus)

1 INITIALISIERUNG $\mathbf{a} \leftarrow \mathbf{0}$

2 NEWTON-RAPHSON-SCHRITT

$$\mathbf{a} \leftarrow \mathbf{a} + (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \cdot \mathbf{X}^\top \cdot (\mathbf{y} - \mathbf{p})$$

Die diagonale Skalierungsmatrix $\mathbf{W} \in \mathbb{R}^{T \times T}$ hat Einträge $w_{tt} = p_t \cdot (1 - p_t)$, $p_t = \hat{p}_1(\mathbf{x}_t)$.

3 TERMINIERUNG

Prüfe Abbruchbedingung; gehe \rightsquigarrow 2 oder ENDE.

(sumfthog1A)

Newton-Raphson-Optimierungsschritt

Gradientenaufstieg mit quadratisch berechneter Schrittweite

$$\mathbf{a} \leftarrow \mathbf{a} - \mathbf{H}_a^{-1} \cdot \nabla_a$$

Maximum-Likelihood-Schätzung

Allgemeiner Fall: $K \geq 2$

Lemma

Für das logistische Modell $p_\lambda(\mathbf{x}) \propto \exp(\mathbf{a}_\lambda^\top \mathbf{x})$ mit den Lerndaten $(\mathbf{x}_t, g_t) \in \mathbb{R}^N \times \{1, \dots, K\}$, $t = 1, \dots, T$ gilt:

1. Die ML-Zielgröße besitzt die Darstellung

$$\ell(\mathbf{A}) = \log \prod_{t=1}^T p_{g_t}(\mathbf{x}_t) = \sum_{t=1}^T \left\{ \mathbf{a}_{g_t}^\top \mathbf{x}_t - \log \sum_{\nu} e^{\mathbf{a}_\nu^\top \mathbf{x}_t} \right\}, \quad \mathbf{a}_K = \mathbf{0} \in \mathbb{R}^N.$$

2. Für die $K \cdot N$ partiellen Ableitungen ihrer Gradientenmatrix gilt:

$$\frac{\partial \ell(\mathbf{A})}{\partial \mathbf{a}_{\lambda,i}} = \sum_{t=1}^T \mathbf{x}_{t,i} \cdot (y_{t,\lambda} - p_\lambda(\mathbf{x}_t))$$

3. Für die $K^2 \cdot N^2$ gemischten partiellen Ableitungen ihres Hestetensors gilt:

$$\frac{\partial^2 \ell(\mathbf{A})}{\partial \mathbf{a}_{\lambda,i} \cdot \partial \mathbf{a}_{\kappa,j}} = \sum_{t=1}^T \mathbf{x}_{t,i} \mathbf{x}_{t,j} \cdot (p_\lambda(\mathbf{x}_t) \cdot p_\kappa(\mathbf{x}_t) - \delta_{\lambda,\kappa} \cdot p_\lambda(\mathbf{x}_t))$$

Es bezeichne $y_{t,\lambda} = \delta_{g_t,\lambda}$ die Klassenindikatorfunktion der Lerndaten.

IRLS und Regularisierung

Was heißt eigentlich „wiederholte Neugewichtung“?

Newtonschritt = gewichtete lineare Regression

$$\begin{aligned} \mathbf{a}^* &\leftarrow \mathbf{a} + (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \cdot \mathbf{X}^\top \cdot (\mathbf{y} - \mathbf{p}) \\ &= (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \cdot \mathbf{X}^\top \mathbf{W}^{1/2} \cdot \mathbf{W}^{1/2} \cdot \underbrace{(\mathbf{X} \mathbf{a} + \mathbf{W}^{-1}(\mathbf{y} - \mathbf{p}))}_z \end{aligned}$$

\mathbf{a}^* ist Lösung der gewichteten Regressionsaufgabe $\|\mathbf{W}^{1/2} \cdot (\mathbf{z} - \mathbf{X} \mathbf{a})\|^2 \xrightarrow{!} \text{MIN}$

WLS-Regularisierung in jedem Newtonschritt

Löse gewichtetes GNG-System $\mathbf{R}_W \mathbf{a} = \mathbf{m}_W$ mit $\mathbf{R}_W = \mathbf{X}^\top \mathbf{W} \mathbf{X}$ und $\mathbf{m}_W = \mathbf{X}^\top \mathbf{W} \mathbf{z}$ mittels regularisierter Koeffizientenmatrix:

$$\mathbf{R}_{W,\lambda} \stackrel{\text{def}}{=} (\mathbf{R}_W)_\lambda = \mathbf{X}^\top \mathbf{W} \mathbf{X} + \lambda \cdot \mathbf{E}$$

Das Probit-Modell ($K = 2$)

Logistisches Wahrscheinlichkeitsmodell

mit einer **additiven** Zielfunktion

$$p_y(\mathbf{x}) = \frac{e^{y \cdot \mathbf{a}^\top \mathbf{x}}}{1 + e^{\mathbf{a}^\top \mathbf{x}}}, \quad \ell(\mathbf{a}) = \sum_{t=1}^T p_{y_t}(\mathbf{x}_t) \xrightarrow{!} \text{MAX}$$

Gradientenvektor

$$\nabla_a \ell(\mathbf{a}) = \sum_{t=1}^T p_{y_t}(\mathbf{x}_t) \cdot \{y_t - p_1(\mathbf{x}_t)\} \cdot \mathbf{x}_t = \mathbf{X}^\top \mathbf{Q}(\mathbf{y} - \mathbf{p})$$

mit $\mathbf{Q} = \text{diag}(\{p_{y_t}(\mathbf{x}_t)\}_t)$

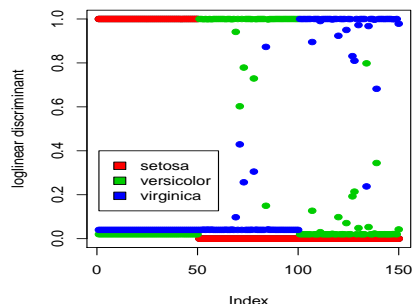
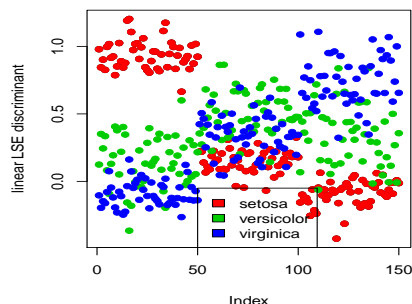
Hessematrix

$$\mathbf{H}_a = \sum_{t=1}^T p_{y_t}(\mathbf{x}_t) \cdot \{(y_t - p_1(\mathbf{x}_t))^2 + p_1^2(\mathbf{x}_t) - p_1(\mathbf{x}_t)\} \cdot \mathbf{x}_t \mathbf{x}_t^\top = -\mathbf{X}^\top \mathbf{W} \mathbf{P} \mathbf{X}$$

mit $\mathbf{W} = \text{diag}(\{p_1(\mathbf{x}_t) \cdot p_0(\mathbf{x}_t)\}_t)$ und $\mathbf{P} = \text{diag}(\{p_{y_t}(\mathbf{x}_t) - p_{1-y_t}(\mathbf{x}_t)\}_t)$

Reklassifikationsexperiment — Irisblüten-Datensatz

3 Klassen · 4 numerische Attribute · 50+50+50 Objekte



Quadratmittellmodell

3 affine Prädiktoren $\mathbf{a}_\lambda \in \mathbb{R}^5$
Starke Schwankung um 1 und 0
Vertauschungen $\left\{ \begin{array}{l} \text{'versicolor'} \\ \text{'virginica'} \end{array} \right\}$

Loglinearmodell

2 affine Prädiktoren $\mathbf{a}_\lambda \in \mathbb{R}^5$
Fast alle Wahrsch'keiten bei $\{0, 1\}$
Fast perfekte Klassenidentifikation

Dualisierung der Regressionsaufgabe

Der Schlüssel zum Kerneltrick für Quadratmittel- & logistische Prädiktoren

Definition

Für eine Datenmatrix $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)^\top$ des \mathbb{R}^N bezeichne $\text{Lin}(\mathbf{X})$ die **lineare Hülle** der Vektoren und $\text{Lin}(\mathbf{X}^\perp)$ ihren **Orthogonalraum**.

Lineare Hülle

Die Menge aller Linearkombinationen der Matrixzeilen \mathbf{x}_t ; sie bildet den kleinsten Untervektorraum von \mathbb{R}^N , der alle \mathbf{x}_t enthält.

$$\text{Lin}(\mathbf{X}) = \{\mathbf{X}^\top \mathbf{a} \mid \mathbf{a} \in \mathbb{R}^T\}$$

Orthogonalraum

Die Menge aller Vektoren, die auf allen Matrixzeilen \mathbf{x}_t senkrecht stehen; sie bildet den größten Untervektorraum von \mathbb{R}^N , der keines der \mathbf{x}_t enthält.

$$\text{Lin}(\mathbf{X}^\perp) = \{\mathbf{z} \in \mathbb{R}^N \mid \mathbf{X}\mathbf{z} = \mathbf{0}\}$$

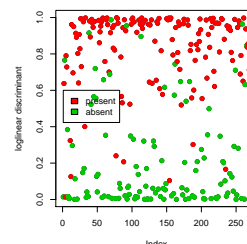
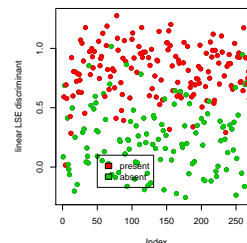
Lemma

Lineare Hülle und Orthogonalraum spannen stets den Gesamttraum auf:

$$\mathbb{R}^N = \text{Lin}(\mathbf{X}) \oplus \text{Lin}(\mathbf{X}^\perp)$$

Reklassifikationsexperiment — Herzkrankheiten-Datensatz

2 Klassen · 13 diskrete & numerische Attribute · 270 Objekte



Auszug Datenfriedhof

```
70.0 1.0 4.0 130.0 322.0 0.0 2.0 109.0 0.0 2.4 2.0 3.0 3.0 2
67.0 0.0 3.0 115.0 564.0 0.0 2.0 160.0 0.0 1.6 2.0 0.0 7.0 1
57.0 1.0 2.0 124.0 261.0 0.0 0.0 141.0 0.0 0.3 1.0 0.0 7.0 2
64.0 1.0 4.0 128.0 263.0 0.0 0.0 105.0 1.0 0.2 2.0 1.0 7.0 1
74.0 0.0 2.0 120.0 269.0 0.0 2.0 121.0 1.0 0.2 1.0 1.0 3.0 1
65.0 1.0 4.0 120.0 177.0 0.0 0.0 140.0 0.0 0.4 1.0 0.0 7.0 1
56.0 1.0 3.0 130.0 256.0 1.0 2.0 142.0 1.0 0.6 2.0 1.0 6.0 2
... ..
```

Attribute, Skalen, Werte

- age (IR) -0.02511018
- sex {male, female} 1.89901910
- chest pain {A, B, C, D} 1.741, 0.784, 2.748
- blood pressure (IR) 0.03110868
- serum cholestoral (IR) 0.00655756
- fasting blood sugar {T, F} -0.37604461
-
- thal {normal, fixed, defect} -0.318, 1.468
- intercept -7.68704469

Darstellungssatz für QM-Lösungen

Endlichdimensionaler Spezialfall des Satzes von Kimeldorf & Wahba (1971)

Satz

Die regulisierten (unregulisierten) und gewichteten (ungewichteten) Quadratmittelaufgaben mit den Normalengleichungen

$$\begin{aligned} \mathbf{R} \cdot \mathbf{a} &= \mathbf{m} & (LSE) \\ \mathbf{R}_\lambda \cdot \mathbf{a} &= \mathbf{m} & (RLSE) \\ \mathbf{R}_w \cdot \mathbf{a} &= \mathbf{m}_w & (WLSE) \\ \mathbf{R}_{w,\lambda} \cdot \mathbf{a} &= \mathbf{m}_w & (RWLSE) \end{aligned}$$

besitzen jeweils mindestens eine Lösung, die sogar als Linearkombination der Datenvektoren $\mathbf{x}_1, \dots, \mathbf{x}_T$ darstellbar ist, d.h. es gilt:

$$\mathbf{a}^* \in \text{Lin}(\mathbf{X})$$

Bezeichnungen

für die nicht normierten und unzentrierten Momente:

$$\begin{aligned} \mathbf{m} &= \mathbf{X}^\top \mathbf{y} \\ \mathbf{R} &= \mathbf{X}^\top \mathbf{X} \\ \mathbf{R}_\lambda &= \mathbf{R} + \lambda \mathbf{E} \\ \mathbf{m}_w &= \mathbf{X}^\top \mathbf{W} \mathbf{z} \\ \mathbf{R}_w &= \mathbf{X}^\top \mathbf{W} \mathbf{X} \\ \mathbf{R}_{w,\lambda} &= \mathbf{R}_w + \lambda \mathbf{E} \end{aligned}$$

- REPRÄSENTATION FÜR LSE-LÖSUNG

Ist $\mathbf{a} = \mathbf{a}_0 + \mathbf{a}_\perp$ mit $\mathbf{a}_0 \in \text{Lin}(\mathbf{X})$ und $\mathbf{a}_\perp \in \text{Lin}(\mathbf{X}^\perp)$ eine Lösung der GNG $\mathbf{R}\mathbf{a} = \mathbf{m}$, so gilt:

$$\mathbf{m} = \mathbf{R}\mathbf{a} = \mathbf{X}^\top \mathbf{X} \mathbf{a}_0 + \mathbf{X}^\top \mathbf{X} \mathbf{a}_\perp = \mathbf{R} \mathbf{a}_0$$

Wir können folglich auch eine Lösung in $\text{Lin}(\mathbf{X})$ finden.

- REPRÄSENTATION FÜR RLSE-LÖSUNG

Ist $\mathbf{a} = \mathbf{a}_0 + \mathbf{a}_\perp$ eine Lösung der GNG $\mathbf{R}_\lambda \mathbf{a} = \mathbf{m}$, so gilt:

$$\mathbf{m} = \mathbf{R}_\lambda \mathbf{a} = \mathbf{X}^\top \mathbf{X} \mathbf{a} + \lambda \mathbf{a}_0 + \lambda \mathbf{a}_\perp$$

Da sowohl $\mathbf{m} = \mathbf{X}^\top \mathbf{y}$ als auch $\mathbf{X}^\top \mathbf{X} \mathbf{a}$ und $\lambda \mathbf{a}_0$ offensichtlich aus $\text{Lin}(\mathbf{X})$ sind, ist das auch für den verbleibenden Ausdruck $\lambda \mathbf{a}_\perp$ der Fall. Wegen $\lambda > 0$ folgt $\mathbf{a}_\perp = \mathbf{0}$, also ist $\mathbf{a} = \mathbf{a}_0$ zwingend aus der linearen Hülle von \mathbf{X} .

- REPRÄSENTATION FÜR WLSE-LÖSUNG

Im IRLS-Schritt sei $\mathbf{a} = \mathbf{a}_0 + \mathbf{a}_\perp$ eine Lösung der GNG $\mathbf{R}_w \mathbf{a} = \mathbf{m}_w$. Dann gilt:

$$\mathbf{m}_w = \mathbf{R}_w \mathbf{a} = \mathbf{X}^\top \mathbf{W} \mathbf{X} \cdot \mathbf{a}_0 + \mathbf{X}^\top \mathbf{W} \mathbf{X} \cdot \mathbf{a}_\perp = \mathbf{R}_w \mathbf{a}_0$$

- REPRÄSENTATION FÜR RWLSE-LÖSUNG

Für die Lösung $\mathbf{a} = \mathbf{a}_0 + \mathbf{a}_\perp$ im regularisierten IRLS-Schritt gilt wie bei RLSE:

$$\mathbf{m}_w = \mathbf{R}_{w,\lambda} \mathbf{a} = (\mathbf{X}^\top \mathbf{W} \mathbf{X} + \lambda \mathbf{E}) \cdot (\mathbf{a}_0 + \mathbf{a}_\perp) = \mathbf{X}^\top \mathbf{W} \mathbf{X} \mathbf{a} + \lambda \mathbf{a}_0 + \lambda \mathbf{a}_\perp$$

Regularisierung dualisierter QM-Aufgaben

Ungewichteter und gewichteter Fall

Lemma

Die Lösungen der dualisierten **LSE-Aufgabe** lauten je nach Regularisierungstechnik:

$$\begin{aligned} \mathbf{b}^* &= \mathbf{G}^{-1} \cdot \mathbf{y} & \varepsilon(\mathbf{b}) &= \|\mathbf{y} - \mathbf{G}\mathbf{b}\|^2 \\ \mathbf{b}^* &= (\mathbf{G} + \lambda \mathbf{E})^{-1} \cdot \mathbf{y} & \varepsilon_\lambda(\mathbf{b}) &= \|\mathbf{y} - \mathbf{G}\mathbf{b}\|^2 + \lambda \cdot \|\mathbf{X}^\top \mathbf{b}\|^2 \\ \mathbf{b}^* &= (\mathbf{G}^2 + \lambda \mathbf{E})^{-1} \cdot \mathbf{G}\mathbf{y} & \varepsilon'_\lambda(\mathbf{b}) &= \|\mathbf{y} - \mathbf{G}\mathbf{b}\|^2 + \lambda \cdot \|\mathbf{b}\|^2 \end{aligned}$$

Lemma

Die Lösungen der dualisierten **WLSE-Aufgabe** lauten je nach Regularisierungstechnik:

$$\begin{aligned} \mathbf{b}^* &= \mathbf{G}^{-1} \cdot \mathbf{z} & \varepsilon_w(\mathbf{b}) &= \|\mathbf{z} - \mathbf{G}\mathbf{b}\|_w^2 \\ \mathbf{b}^* &= (\mathbf{W}\mathbf{G} + \lambda \mathbf{E})^{-1} \cdot \mathbf{W}\mathbf{z} & \varepsilon_{w,\lambda}(\mathbf{b}) &= \|\mathbf{z} - \mathbf{G}\mathbf{b}\|_w^2 + \lambda \cdot \|\mathbf{X}^\top \mathbf{b}\|^2 \\ \mathbf{b}^* &= (\mathbf{G}\mathbf{W}\mathbf{G} + \lambda \mathbf{E})^{-1} \cdot \mathbf{G}\mathbf{W}\mathbf{z} & \varepsilon'_{w,\lambda}(\mathbf{b}) &= \|\mathbf{z} - \mathbf{G}\mathbf{b}\|_w^2 + \lambda \cdot \|\mathbf{b}\|^2 \end{aligned}$$

MSE[⊥] — die dualisierte Quadratmittelaufgabe

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

Duale Lösungsdarstellung

als Linearkombination der Objektvektoren:

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

Duale Regressionsfehlerformel

in Abhängigkeit vom Vektor \mathbf{b} der Lösungskoeffizienten:

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

Duale Gauß'sche Normalgleichungen

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

$$\mathbf{G}^2 \cdot \mathbf{b} = \mathbf{G} \cdot \mathbf{y}, \quad \mathbf{G} = \mathbf{X} \cdot \mathbf{X}^\top$$

Beweis.

- UNREGULARISIERTE LÖSUNG:
Der Gradientenvektor der Zielgröße

$$\varepsilon(\mathbf{b}) = \|\mathbf{y} - \mathbf{G}\mathbf{b}\|^2 = \mathbf{y}^\top \mathbf{y} - 2 \cdot \mathbf{b}^\top \mathbf{G}\mathbf{y} + \mathbf{b}^\top \mathbf{G}^2 \mathbf{b}$$

lautet

$$\nabla_{\mathbf{b}} \varepsilon(\mathbf{b}) = \mathbf{0} - 2 \cdot \mathbf{G}\mathbf{y} + 2 \cdot \mathbf{G}^2 \mathbf{b}.$$

Nullsetzen ergibt die GNG. Unter der Annahme einer regulären Gramschen Matrix ergibt sich die Lösung durch Multiplikation beider Gleichungsseiten mit \mathbf{G}^{-2} .

- REGULARISIERTE LÖSUNG I:
Wir regularisieren im Vektorraum \mathbb{R}^N ; der Fehlerterm besitzt den Gradientenvektor

$$\nabla_{\mathbf{b}} \varepsilon_\lambda(\mathbf{b}) = -2\mathbf{G}\mathbf{y} + 2\mathbf{G}^2 \mathbf{b} + 2\lambda \cdot \mathbf{G}\mathbf{b} = -2\mathbf{G} \cdot (\mathbf{y} - (\mathbf{G} + \lambda \mathbf{E}) \cdot \mathbf{b})$$

Da \mathbf{G}_λ regulär ist für $\lambda > 0$ liefert $\mathbf{b} = \mathbf{G}_\lambda^{-1} \mathbf{y}$ eine Lösung.

- REGULARISIERTE LÖSUNG II:
Wir regularisieren im Vektorraum \mathbb{R}^T ; der Fehlerterm besitzt den Gradientenvektor

$$\nabla_{\mathbf{b}} \varepsilon'_\lambda(\mathbf{b}) = -2\mathbf{G}\mathbf{y} + 2\mathbf{G}^2 \mathbf{b} + 2\lambda \cdot \mathbf{b} = -2 \cdot (\mathbf{G}\mathbf{y} - (\mathbf{G}^2 + \lambda \mathbf{E}) \cdot \mathbf{b})$$

Da auch $(\mathbf{G}^2)_\lambda$ regulär ist für $\lambda > 0$ liefert $\mathbf{b} = (\mathbf{G}^2)_\lambda^{-1} \cdot \mathbf{G}\mathbf{y}$ eine Lösung.

Beweis.

Das zweite Lemma dient der schrittweisen Berechnung und Regularisierung im IRLS-Algorithmus für loglineare Modelle.

An Stelle des Fehlerfunktional $\|\mathbf{y} - \mathbf{G}\mathbf{b}\|^2$ wird

$$\|\mathbf{z} - \mathbf{G}\mathbf{b}\|_{\mathbf{W}}^2 \stackrel{\text{def}}{=} (\mathbf{z} - \mathbf{G}\mathbf{b})^\top \cdot \mathbf{W} \cdot (\mathbf{z} - \mathbf{G}\mathbf{b})$$

minimiert. Wir unterscheiden wieder zwischen der Regularisierung im Raum \mathbb{R}^N und im Raum \mathbb{R}^T .

- Ist \mathbf{G} invertierbar, so hängt die Lösung $\mathbf{b}^* = \mathbf{G}^{-1}\mathbf{b}$ nicht von der (diagonalen) Gewichtmatrix \mathbf{W} ab, denn $\mathbf{z} \approx \mathbf{G}\mathbf{b}$ wird ja mit exakter Gleichheit erfüllt:

$$\nabla \varepsilon_{\mathbf{W}}(\mathbf{b}) = -2\mathbf{G}\mathbf{W}\mathbf{z} + 2\mathbf{G}\mathbf{W}\mathbf{G}\mathbf{b} = -2\mathbf{G}\mathbf{W} \cdot (\mathbf{z} - \mathbf{G}\mathbf{b})$$

- Bei Regularisierung im Raum \mathbb{R}^N ergibt sich:

$$\nabla \varepsilon_{\mathbf{W},\lambda}(\mathbf{b}) = -2\mathbf{G}\mathbf{W}\mathbf{z} + 2\mathbf{G}\mathbf{W}\mathbf{G}\mathbf{b} + 2\lambda\mathbf{G}\mathbf{b} = -2\mathbf{G} \cdot (\mathbf{W}\mathbf{z} - (\mathbf{W}\mathbf{G})_\lambda \cdot \mathbf{b})$$

- Bei Regularisierung im Raum \mathbb{R}^T ergibt sich:

$$\nabla \varepsilon'_{\mathbf{W},\lambda}(\mathbf{b}) = -2\mathbf{G}\mathbf{W}\mathbf{z} + 2\mathbf{G}\mathbf{W}\mathbf{G}\mathbf{b} + 2\lambda\mathbf{b} = -2 \cdot (\mathbf{G}\mathbf{W}\mathbf{z} - (\mathbf{G}\mathbf{W}\mathbf{G})_\lambda \cdot \mathbf{b})$$

□

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Kombinatorische Regression

Aufgabenstellung

Klassifikation von Texten
 $\mathbf{v} \in \Omega = \mathcal{V}^*$ über Wortschatz \mathcal{V}

Termexpansion

Binärattribute:
 Wort- m -Tupel oder
 Wort- m -Subsets

$$\phi : \Omega \rightarrow \{0, 1\}^{\mathcal{V}^m}$$

$$\text{mit } \phi_{\mathbf{u}}(\mathbf{x}) = \begin{cases} 1 & \mathbf{u} \in \mathbf{x} \\ 0 & \mathbf{u} \notin \mathbf{x} \end{cases}$$

Loglinearmodell

der Dimension L^m bzw. $\binom{L}{m}$

Duales Loglinearmodell

Gramsche T^2 -Matrix mit Einträgen

$$K(\mathbf{x}_s, \mathbf{x}_t) = \langle \phi(\mathbf{x}_s), \phi(\mathbf{x}_t) \rangle$$

Kombinat. Kernoperator

$$\begin{aligned} K(\mathbf{x}, \mathbf{y}) &= \sum_{\mathbf{u} \in \mathcal{V}^m} \phi_{\mathbf{u}}(\mathbf{x}) \cdot \phi_{\mathbf{u}}(\mathbf{y}) \\ &= \begin{cases} |V_{\mathbf{x}}^m \cap V_{\mathbf{y}}^m| & \text{Tupel} \\ \binom{|V_{\mathbf{x}}^m \cap V_{\mathbf{y}}^m|}{m} & \text{Subsets} \end{cases} \end{aligned}$$

Die Zählaufgaben $|V_{\mathbf{x}}^m \cap V_{\mathbf{y}}^m|$ sind sehr effizient zu bewältigen.

Ordinale Regression

Reelle Quellattribute $\mathbb{X}_1, \dots, \mathbb{X}_N \rightarrow$ geordnetes Zielattribut $\mathbb{Y} \in \{1, \dots, L\}$

Nominales Attribut

A posteriori Verteilung

$$p_{\ell}(\mathbf{x}) \stackrel{\text{def}}{=} P(\mathbb{Y} = \ell \mid \mathbb{X} = \mathbf{x})$$

Normierungsbedingung:

$$\sum_{\ell} p_{\ell}(\mathbf{x}) = 1$$

Nominale Beispiele

RedGreenBlue-Skala:

$$\mathbf{p} = (\frac{1}{2}, \frac{1}{6}, \frac{1}{3})$$

Unfairer Würfel:

$$\mathbf{p} = (0, 0, 0, \frac{1}{4}, \frac{1}{4}, \frac{1}{2})$$

Müssen ordinale Verteilungen zwangsläufig „unimodal“ sein?

Ordinales Attribut

Kumulative a post. Verteilung

$$q_{\ell}(\mathbf{x}) \stackrel{\text{def}}{=} P(\mathbb{Y} \leq \ell \mid \mathbb{X} = \mathbf{x})$$

Skalenbindung:

$$\mathbf{x} \rightsquigarrow \mathbf{z}(\mathbf{x}) \in J_{\ell} \subset \mathbb{R}$$

Ordinale Beispiele

HighMediumLow-Skala:

$$\mathbf{p} = (\frac{1}{2}, \frac{1}{6}, \frac{1}{3})$$

Zensurenskala:

$$\mathbf{p} = (\frac{1}{2}, 0, 0, 0, \frac{1}{2}, 0)$$

Postulat der verborgenen dichten Qualitätsskala

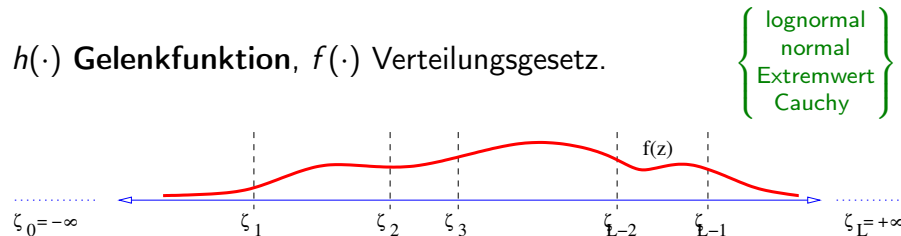
„Cumulative link model“ — Agresti 2002

Kumulatives Gelenkunktionsmodell

Latente Variable \mathbb{Z} auf der Skala $-\infty = \zeta_0 < \zeta_1 < \dots < \zeta_L = +\infty$
mit

$$\mathbb{Y} = \ell \Leftrightarrow \mathbb{Z} \in (\zeta_{\ell-1}, \zeta_\ell] \quad \text{und} \quad \mathbb{Z} \sim f(\mu = h(\mathbf{x}), \sigma^2 = 1)$$

$h(\cdot)$ Gelenkunktion, $f(\cdot)$ Verteilungsgesetz.

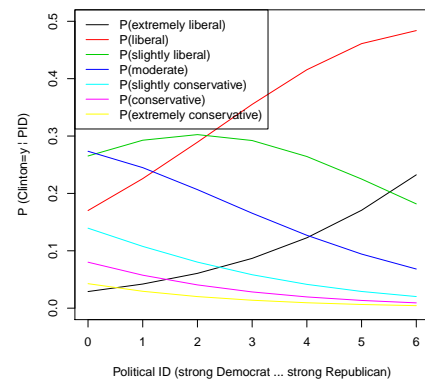


Bemerkung

Es gilt $q_\ell(\mathbf{x}) = P(\mathbb{Y} \leq \ell \mid \mathbb{X} = \mathbf{x}) = P(\mathbb{Z} \leq \zeta_\ell \mid \mathbb{X} = \mathbf{x}) = F(\zeta_\ell - h(\mathbf{x}))$.

Beispiel — Präsidentschaftswahlen USA'96

<http://www.stat.washington.edu/quinn/classes/536/data/n96r.dat>



Datensatz

944 Versuchspersonen
11 Attribute, u.a.:

- Pol/ID Clinton*
- Alter
- Bildungsgrad
- Einkommen
- Stimme für ...
- TV-News/Woche
- Pol/ID selbst*
- Pol/ID Dole*

*) Pol/IDs in 7 Stufen

POLR-Datenanalyse

Zielattribut $\hat{=}$ Einschätzung von Clintons politischer Haltung

Quellattribut $\hat{=}$ politische Selbsteinschätzung

Fixiert: 3 TV/Woche, 44 Jahre, 12 Schuljahre, 35–40 Kilodollar

Proportional Odds Linear Regression

Lineares Binomialmodell für die Gelenkunktion

POLR-Modell

Lineare Vorhersage der logarithmierten Chancenfunktionen:

$$\log \text{odds}_\ell(\mathbf{x}) \stackrel{\text{def}}{=} \log \frac{q_\ell(\mathbf{x})}{1 - q_\ell(\mathbf{x})} = \log \frac{P(\mathbb{Y} \leq \ell \mid \mathbb{X} = \mathbf{x})}{P(\mathbb{Y} > \ell \mid \mathbb{X} = \mathbf{x})} \stackrel{!}{=} \mathbf{a}^\top \mathbf{x} + \zeta_\ell$$

Bemerkungen

1. Normierung

$$\sum_\ell p_\ell(\mathbf{x}) = \sum_\ell (q_\ell(\mathbf{x}) - q_{\ell-1}(\mathbf{x})) = q_L(\mathbf{x}) - q_0(\mathbf{x}) = 1 - 0$$

2. Monotonie

$$k \leq \ell \Rightarrow \zeta_k \leq \zeta_\ell \Rightarrow \log_k(\mathbf{x}) \leq \log_\ell(\mathbf{x}) \Rightarrow q_k(\mathbf{x}) \leq q_\ell(\mathbf{x})$$

3. Proportionale Chancen

$$\log \frac{\text{odds}_\ell(\mathbf{x})}{\text{odds}_\ell(\mathbf{x}')} = \mathbf{a}^\top (\mathbf{x} - \mathbf{x}') \quad \text{ist unabhängig von } \ell$$

Lernen von Präferenzrelationen

Objektive Präferenz

Aus einer Serie gewonnener, verlorener oder unentschiedener Partien $(\mathbf{x}_t, \mathbf{y}_t) \in \Omega \times \Omega$ ist eine passende Qualitätsrelation (Ω, \preceq) zu lernen.



„Tourniermetapher“

Geschlossene Welten

Objektraum Ω und/oder Subjektraum \mathfrak{S} bilden ein endliches **Inventar**.
(Nominalattribut)

Subjektive Präferenz

Aus einer Serie persönlicher Nennungen, Wertungen oder Reihungen $(\mathbf{s}_t, \mathbf{x}_t) \in \mathfrak{S} \times \Omega$ ist eine **Schar** passender Qualitätsrelationen $(\Omega, \preceq_s)_{s \in \mathfrak{S}}$ zu lernen.



„Jurorenmetapher“

Offene Welten

Objekte u/o Subjekte sind durch ihre **Eigenschaften** charakterisiert.
(Attributvektoren)

Objektive Präferenz durch logistische Regression

Bilaterales Ereignismodell

$\mathbb{X} \triangleq$ Objekt #1 (Herausforderer)
 $\mathbb{Y} \triangleq$ Objekt #2 (Gegner)
 $\mathbb{Z} \triangleq$ Resultat \pm , „Sieg“ oder \pm , „Tor“ ...

Logistisch-lineares Erfolgsmodell

$$\log \text{odds}(\mathbf{x}, \mathbf{y}) = \underbrace{\mathbf{a}^\top \mathbf{x} + \mathbf{b}^\top \mathbf{y} + \zeta}_{g(\mathbf{x}) - h(\mathbf{y})}$$

Präferenzinterpretation

\mathbf{x} hat immer dann bessere Gewinnchancen als \mathbf{y} wenn $g(\mathbf{x}) > h(\mathbf{y})$ gilt.

Intervallordnung ?

Es gilt $p(\mathbf{x}, \mathbf{x}) \leq \frac{1}{2} \iff g_{\mathbf{x}} \leq h_{\mathbf{x}}$.

$$\mathbf{x} \succ \mathbf{y} \text{ gdw. } [g_{\mathbf{x}}, h_{\mathbf{x}}] \supset [g_{\mathbf{y}}, h_{\mathbf{y}}]$$

Fußballturnier

GER : BRA 3:1
 USA : LBY 0:1
 UK : IRAN 2:2

Punktestandbezogen

x_{GER}	x_{BRA}	+
x_{BRA}	x_{GER}	-
x_{USA}	x_{LBY}	-
x_{LBY}	x_{USA}	+
x_{UK}	x_{IRAN}	-
x_{IRAN}	x_{UK}	-

Torstandbezogen

x_{GER}	x_{BRA}	+	3
x_{GER}	x_{BRA}	-	87
x_{BRA}	x_{GER}	+	1
x_{BRA}	x_{GER}	-	89
...	...		
x_{UK}	x_{IRAN}	+	2
x_{UK}	x_{IRAN}	-	88
x_{IRAN}	x_{UK}	+	2
x_{IRAN}	x_{UK}	-	88

Spezielle Form des POLR-Modells

$$\log \text{odds}_{\ell}(\mathbf{x}, \mathbf{y}) = \mathbf{a}^\top (\mathbf{x} - \mathbf{y}) + \begin{cases} -\infty & \ell = 0 \\ -\zeta & \ell = 1 \\ +\zeta & \ell = 2 \\ +\infty & \ell = 3 \end{cases}$$

Beweis.

Aus der strukturellen Symmetrie

$P(> | \mathbf{x}, \mathbf{y}) = P(< | \mathbf{y}, \mathbf{x})$ folgt:

$$\begin{aligned} \Rightarrow p_1(\mathbf{x}, \mathbf{y}) &= p_3(\mathbf{y}, \mathbf{x}) \\ \Rightarrow q_1(\mathbf{x}, \mathbf{y}) - 0 &= 1 - q_2(\mathbf{y}, \mathbf{x}) \\ \Rightarrow \frac{q_1(\mathbf{x}, \mathbf{y})}{1 - q_1(\mathbf{x}, \mathbf{y})} &= \frac{1 - q_2(\mathbf{y}, \mathbf{x})}{q_2(\mathbf{y}, \mathbf{x})} \\ \Rightarrow \text{odds}_1(\mathbf{x}, \mathbf{y}) &= \text{odds}_2^{-1}(\mathbf{y}, \mathbf{x}) \\ \Rightarrow +\log \text{odds}_1(\mathbf{x}, \mathbf{y}) &= -\log \text{odds}_2(\mathbf{y}, \mathbf{x}) \\ \Rightarrow 0 &= \mathbf{a}^\top \mathbf{x} + \mathbf{b}^\top \mathbf{y} + \zeta_1 + \mathbf{a}^\top \mathbf{y} + \mathbf{b}^\top \mathbf{x} + \zeta_2 \\ \Rightarrow 0 &= (\mathbf{a} + \mathbf{b})^\top (\mathbf{x} + \mathbf{y}) + (\zeta_1 + \zeta_2) \\ \Rightarrow \mathbf{b} &= -\mathbf{a} \text{ und } \zeta_1 = -\zeta_2 \end{aligned}$$

□

Objektive Präferenz durch Proportional-Odds Regression

Trilaterales Ereignismodell

$\mathbb{X} \triangleq$ Objekt #1 (Herausforderer)
 $\mathbb{Y} \triangleq$ Objekt #2 (Gegner)
 $\mathbb{Z} \triangleq$ Resultat aus $\{\xi_1, \xi_2, \xi_3\} = \{>, \dot{=}, <\}$

POLR Erfolgsmodell

$$\log \text{odds}_{\ell}(\mathbf{x}, \mathbf{y}) = \underbrace{\mathbf{a}^\top \mathbf{x} + \mathbf{b}^\top \mathbf{y} + \zeta_{\ell}}_{\mathbf{a}^\top (\mathbf{x} - \mathbf{y}) \pm \zeta}$$

Präferenzinterpretation

\mathbf{x} hat immer dann bessere Gewinnchancen als \mathbf{y} wenn $\log \text{odds}_1(\mathbf{x}, \mathbf{y}) > 0$ gilt, also

$$g_{\mathbf{x}} := \mathbf{a}^\top \mathbf{x} > \mathbf{a}^\top \mathbf{y} + \zeta =: h_{\mathbf{y}}$$

Semi-Ordnung !

$$\mathbf{x} \succ \mathbf{y} \text{ gdw. } [g_{\mathbf{x}}, g_{\mathbf{x}} + \zeta] \supset [g_{\mathbf{y}}, g_{\mathbf{y}} + \zeta]$$

Fußballturnier

GER : BRA 3:1
 USA : LBY 0:1
 UK : IRAN 2:2

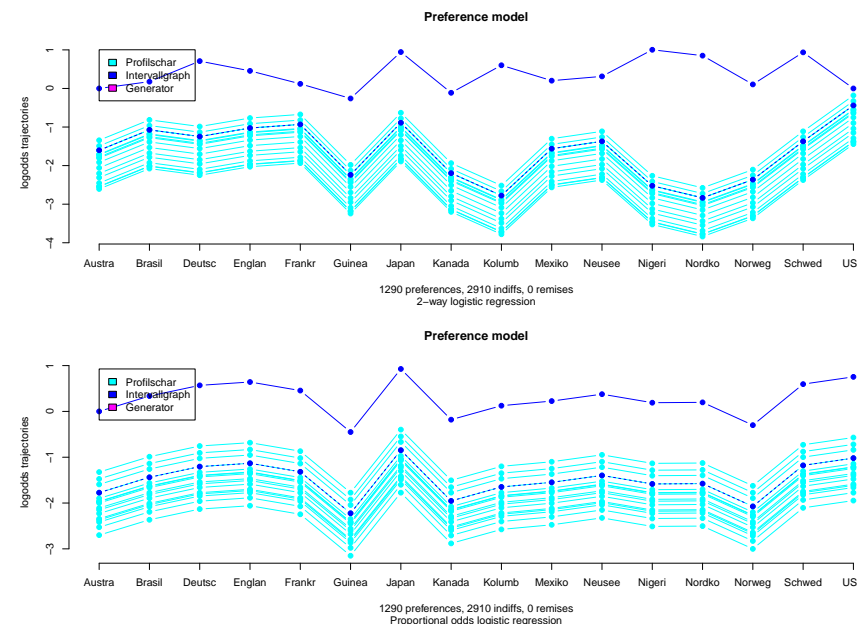
Punktestandbezogen

x_{GER}	x_{BRA}	\succ
x_{BRA}	x_{GER}	\prec
x_{USA}	x_{LBY}	\prec
x_{LBY}	x_{USA}	\succ
x_{UK}	x_{IRAN}	$\dot{=}$
x_{IRAN}	x_{UK}	$\dot{=}$

Torstandbezogen

x_{GER}	x_{BRA}	\succ	3
x_{GER}	x_{BRA}	$\dot{=}$	86
x_{GER}	x_{BRA}	\prec	1
x_{BRA}	x_{GER}	\succ	1
x_{BRA}	x_{GER}	$\dot{=}$	86
x_{BRA}	x_{GER}	\prec	3
...	...		

Beispiel — Frauenfußball-WM 2011



Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

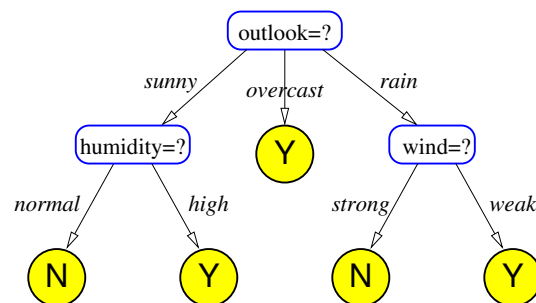
Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Entscheidungsbaum

Hierarchie sequentieller Auswahlfragen („multiple choice“)



Sportwetterempfehlungen

Vier nominale Wetterlagevariablen gegeben

Klassifikationsziel:

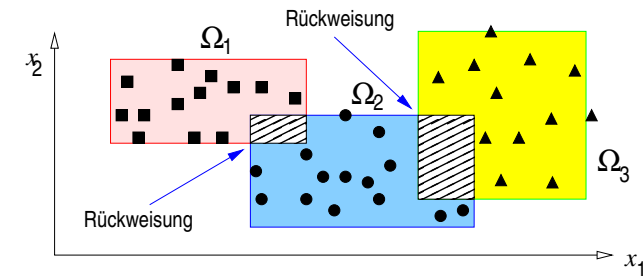
„Ist dieses Wetter zum Tennisspielen geeignet?“

Fragetypus

Wertverzweigung

Parallelepiped-Klassifikator

Vollständige Konjunktion je zweier Literale $x_n \geq a_n$, $x_n \leq b_n$, $n = 1, \dots, N$



Vorteile

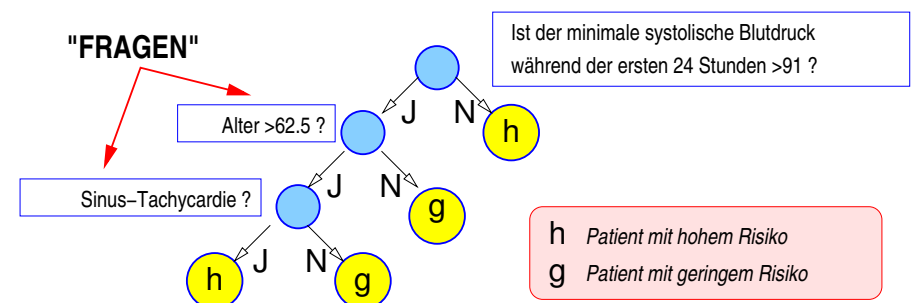
Extrem schnelle Lernphase
Effiziente Abrufphase
Klassengebiete intuitiv zu deuten
Nominalattribute handhabbar[€]

Nachteile

Achsenparallele Grenzen
Unimodale Klassengebiete
Ausgedehnte Rückweisungszonen
Keine Rückschlußwahrsh'keiten

Binärer Entscheidungsbaum

Hierarchie sequentieller Ja/Nein-Fragen



Diagnose für Herzinfarktpatienten

19 Attribute gemessen bzw. erfragt
Patienten 30 Tage unter klinischer Beobachtung

Klassifikationsziel:

„Ist ein zweiter, diesmal tödlicher Infarkt eingetreten?“

Fragetypus

Schwellwertdichotomie
Zielwertdichotomie

Struktur eines Entscheidungsbaumes

Knoten $\hat{=}$ Fragen

\mathcal{B} bezeichnet die Menge aller Knoten.
Innere Knoten $\beta \in \mathcal{B}$ beherbergen eine **Entscheidungsfrage**:

$$Q(\beta) : \Omega \rightarrow \{1, \dots, L\}$$

Wurzelknoten $\hat{=}$ Startposition

$\beta_{\Delta} \in \mathcal{B}$ besitzt keinen Vorgänger.
In β_{Δ} beginnt die Befragung des Objekts.

Kanten $\hat{=}$ Antworten

Für $\beta \in \mathcal{B}$ ist β^{\uparrow} der **Vorgängerknoten** und $\beta^{(1)}, \dots, \beta^{(L)}$ sind die unmittelbaren **Nachfolger**.

Blattknoten $\hat{=}$ Ergebnis

Die $\beta \in \mathcal{B}_{\ell}$ besitzen keine Nachfolger, aber eine **Klassenmarkierung**:

$$\delta_{\ell} : \mathcal{B}_{\ell} \rightarrow \{1, \dots, K\}$$

Klassifikation eines Objekts

Hierarchisches Interview — „Durchschleusen“ bis zum Blattknoten

1 INITIALISIERUNG

Setze $\beta = \beta_{\Delta}(\mathcal{B})$.

2 BEFRAGUNG

Reiche \mathbf{x} gemäß $Q(\beta)$ an einen Kindknoten weiter:

$$\beta \leftarrow \beta^{(i)}, \quad i = Q(\beta)(\mathbf{x})$$

3 TERMINIERUNG

Ist β ein Blattknoten, so lautet das Resultat:

$$\delta(\mathbf{x}) = \delta_{\ell}(\beta)$$

Andernfalls \rightsquigarrow 2.

Befragung der Attributwerte

Dichotomien („Yin-Yang“-Fragen) und Wertverzweigungen

Attribut-Wert-Gleichungen

$$x_i = low$$

bei nominalen Merkmalen
(das negative Literal $x_i \neq low$ ist dazu dual)

Attribut-Wert-Ungleichungen

$$x_i \leq 3.14$$

bei ordinalen Merkmalen
(auch $x_i \geq 17$ oder Intervalle $18 \leq x_i \leq 65$ denkbar)

Wertverzweigungen

$$x_i = red|blue|green$$

bei Attributen mit kleinem $|\mathcal{X}_n|$
(eine Nachfolgerkante je Attributwert)

Teilmengenzugehörigkeit

$$x_i \in \{cloudy, rainy\}$$

bei nominalen Attributen

Reguläre Ausdrücke

$$x_i = ababb * c * ba$$

bei Wort- oder Zeichenketten

Simultanes Schleusen einer Objektmenge

Definition

Ist $(\mathcal{B}, Q, \delta_{\ell})$ ein Entscheidungsbaum über Ω und $\omega \subseteq \Omega$ ein Datensatz, so definieren wir die **assoziierten Objektmengen** ω_{β} induktiv durch:

$$\omega(\beta) \stackrel{\text{def}}{=} \begin{cases} \omega & \beta = \beta_{\Delta} \\ \{\mathbf{x} \in \omega_{\beta} \mid Q(\beta)(\mathbf{x}) = j\} & \beta = \beta^{(j)} \end{cases}$$

Lemma

Ist $(\mathcal{B}, Q, \delta_{\ell})$ ein Entscheidungsbaum über Ω , so gilt:

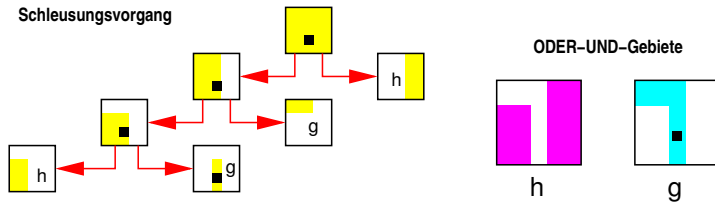
$$\Omega = \bigsqcup_{\beta \in \mathcal{B}_{\ell}} \Omega(\beta)$$

Der Entscheidungsbaum definiert ferner eine vollständige Zerlegung von Ω in Klassengebiete:

$$\Omega = \bigsqcup_{\kappa=1}^K \Omega_{\kappa}, \quad \Omega_{\kappa} \stackrel{\text{def}}{=} \bigcup_{\delta_{\ell}(\beta)=\kappa} \Omega(\beta)$$

Entscheidungsbäume als Hypothesen

Disjunktionen von Literalkonjunktionen



Beispiel

Intuitiv interpretierbare Klassenentscheidungen:

$$\Omega_h = (\{x_b \not\leq 91\} \wedge \{x_a \not\leq 62.5\} \wedge \{x_s \not\leq 0\}) \vee (\{x_b \leq 91\})$$

$$\Omega_g = (\{x_b \not\leq 91\} \wedge \{x_a \not\leq 62.5\} \wedge \{x_s \leq 0\}) \vee (\{x_b \leq 91\} \wedge \{x_a \leq 62.5\})$$

mit den Variablen (Merkmalen)

x_b = minimaler systolischer Blutdruck
 x_a = Alter des Patienten
 x_s = Sinus-Tachycardie? (0 oder 1)

Lernen eines Entscheidungsbaumes

aus klassenetikettierten Beispielobjekten: $\omega = \omega_1 \uplus \omega_2 \uplus \dots \uplus \omega_K \subset \Omega$

Trennschärfe

Der Baum soll die Beispiele möglichst *korrekt klassifizieren*.

➡ **Konsistenz**

Induktionskraft

Der Baum soll die Beispiele in geeigneter Weise *verallgemeinern*.

➡ **geringe Knotenzahl**

Hypothesenraum

Welche Größe? Welche Form?

Welches Attribut? Welche Frage?

➡ gigantische Auswahl an E-Bäumen

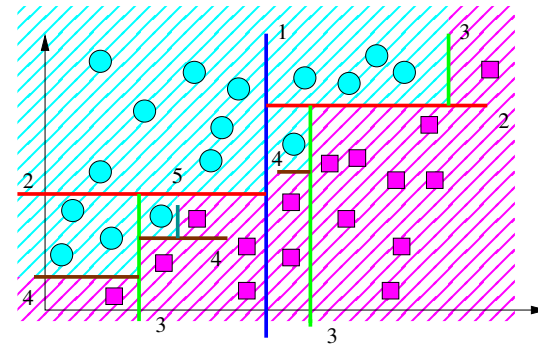
Lernen $\hat{=}$ Suche

Vollständige Suche ist NP-hart.

➡ lokal optimierende Suche (Bergsteiger-Algorithmus, „divide-and-conquer“)

Entscheidungsbäume für numerische Attribute?

Rekursive Halbraumbildung nach sukzessiven Schwellwertabfragen $x_n \leq \theta$



Vom Entscheidungsbaum induzierte Klassengebiete

$$\hat{\Omega}_K = \bigcup_{m=1}^{M_K} \hat{\Omega}_{K,m}, \quad \hat{\Omega}_{K,m} = \bigcap_{l=1}^{M_{K,m}} H_{K,m,l} = \text{Halbraum} \begin{cases} x_d \leq \theta \\ \text{oder} \\ x_d > \theta \end{cases}$$

TDI-Lernalgorithmus

Gierige Top-Down Induktion von Entscheidungsbaumen

(Algorithmus)

1 INITIALISIERUNG

Erzeuge einen Wurzelknoten $\beta = \beta_{\Delta}$ mit den assoziierten Stichproben $\omega_1, \dots, \omega_K$.

2 STOPPTTEST

Ist β hinreichend **reinklassig**, so beende die lokale Konstruktion mit der Blattmarkierung

$$\delta_{\ell}(\beta) = \operatorname{argmax}_K |\omega_K(\beta)|.$$

3 FRAGEAUSWAHL

Wähle eine Frage $Q(\beta)$ mit maximaler Reduktion der **Entscheidungsunsicherheit**.

4 EXPANSION

Bilde die Nachfolgerknoten $\beta^{(1)}, \dots, \beta^{(L)}$ bezüglich $Q(\beta)$ und ihre assoziierten Stichproben

$$\omega_K(\beta^{(l)}), \quad l = 1, \dots, L.$$

5 REKURSION

Fahre mit den Nachfolgern $\beta^{(l)}$ von β bei Schritt 2 fort

(zum nächsten A)

Stoppkriterium

Wann endet der Züchtungsvorgang?

Lokale Stoppkriterien

Wann endet die Knotenexpansion in einem Blatt?

- Wenn $\omega(\beta)$ nur noch einen Datenvektor enthält.
- Wenn $\omega(\beta)$ nur noch Daten einer Klasse enthält. ➡ Konsistenz
- Wenn $|\omega(\beta)|$ eine gegebene Schranke unterschreitet.
- Wenn $|\omega(\beta)| - \max_{\lambda} |\omega_{\lambda}(\beta)|$ eine Schranke unterschreitet.

Überanpassung an die Lernbeispiele

Gefährlich in großen Bäumen durch *Zersplitterung* von ω auf die Blattknoten.

- Ist es wirklich weise, einen konsistenten Baum zu konstruieren?
- ➡ globale a posteriori Stoppkriterien
a.k.a. Baumbeschneidungstechniken, „pruning“

Auswahlregel für die „beste“ nächste Frage

Vorgehensweise

Welches ist die (lokal) zielführendste Frage?

1. Definiere Entscheidungsunsicherheit einer Häufigkeitsverteilung
2. Definiere Entscheidungsunsicherheit eines Baumknotens
3. Definiere Entscheidungsunsicherheit einer Frage (in β)
4. Definiere Entscheidungsunsicherheit eines Teilbaums (unter β)

Relative Klassenhäufigkeit

in der Teilstichprobe ω_{β} zum Knoten $\beta \in \mathcal{B}$:

$$\hat{p}_{\kappa}(\beta) = \frac{\text{Anzahl der } \Omega_{\kappa}\text{-Muster in } \beta}{\text{Anzahl aller Muster in } \beta} = \frac{|\omega_{\kappa}(\beta)|}{\sum_{\lambda=1}^K |\omega_{\lambda}(\beta)|}$$

➡ ML-Schätzwert für $P(\mathbf{x} \in \Omega_{\kappa} \mid \mathbf{x} \in \Omega_{\beta})$

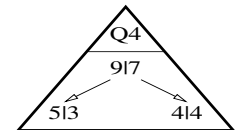
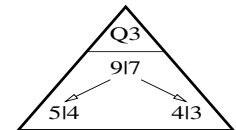
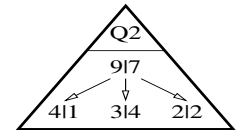
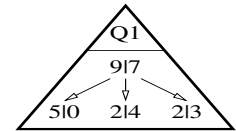
Die Frage nach der richtigen Frage

... bei Yuichiro Anzai im Autohaus ... (Beispiel)

Japanische Gebrauchtfahrzeuge

und ihre Veräußerungschancen am Markt

Objekt	cm ³	Türen	Autom.	Farbe	$\mathbf{x} \stackrel{?}{\in} C$
\mathbf{x}_1	2000	2T	ja	hell	+
\mathbf{x}_2	2800	4T	ja	hell	+
\mathbf{x}_3	2000	2T	nein	dunkel	−
\mathbf{x}_4	1600	4T	ja	dunkel	−
\mathbf{x}_5	1600	4T	ja	hell	−
\mathbf{x}_6	2800	4T	ja	dunkel	+
\mathbf{x}_7	2000	4T	ja	hell	+
\mathbf{x}_8	2000	5T	nein	hell	−
\mathbf{x}_9	1600	2T	nein	hell	+
\mathbf{x}_{10}	2800	5T	ja	hell	+
\mathbf{x}_{11}	2800	5T	nein	dunkel	+
\mathbf{x}_{12}	2000	4T	ja	dunkel	−
\mathbf{x}_{13}	1600	2T	nein	dunkel	+
\mathbf{x}_{14}	2800	2T	nein	dunkel	+
\mathbf{x}_{15}	1600	4T	nein	hell	−
\mathbf{x}_{16}	2000	5T	ja	dunkel	−



Entscheidungsunsicherheit

Gütemaß für den Entmischungsgrad einer Verteilung

Definition

Es sei $K \in \mathbb{N}$ und $\{p_{\kappa} \mid \kappa = 1, \dots, K\}$ eine diskrete Wahrscheinlichkeitsverteilung. Eine Abbildung

$$\mathfrak{S} : \{p_1, \dots, p_K\} \mapsto u \in \mathbb{R}$$

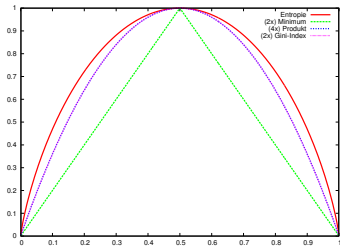
heißt Maß für die **Entscheidungsunsicherheit** (Homogenität, „impurity“), falls gilt:

1. Die Größe $\mathfrak{S}(\cdot)$ ist nichtnegativ.
2. $\mathfrak{S}(\cdot)$ ist maximal für die Gleichverteilung $p_{\kappa} \equiv 1/K$
3. $\mathfrak{S}(\cdot)$ ist minimal für die definiten Verteilungen

$$\mathbf{e}_{\lambda} = (\underbrace{0, \dots, 0}_{\lambda-1}, \underbrace{1, 0, \dots, 0}_{K-\lambda}), \quad \lambda \in \{1, \dots, K\}$$

Homogenitätsmaße

Minimum · Produkt · Gini-Index · Entropie



Extremalwerte

	min	max
\mathfrak{S}_m	0	$1/K$
\mathfrak{S}_p	0	$1/K^K$
\mathfrak{S}_g	0	$1 - 1/K$
\mathfrak{S}_e	0	$\log_2 K$

Lemma

Die folgenden Abbildungen sind (für festes $K \in \mathbb{N}$) Beispiele für Homogenitätsmaße:

$$\mathfrak{S}_m(\mathbf{p}) \stackrel{\text{def}}{=} \min_{\kappa} p_{\kappa}$$

$$\mathfrak{S}_g(\mathbf{p}) \stackrel{\text{def}}{=} \sum_{\lambda \neq \kappa} p_{\lambda} \cdot p_{\kappa}$$

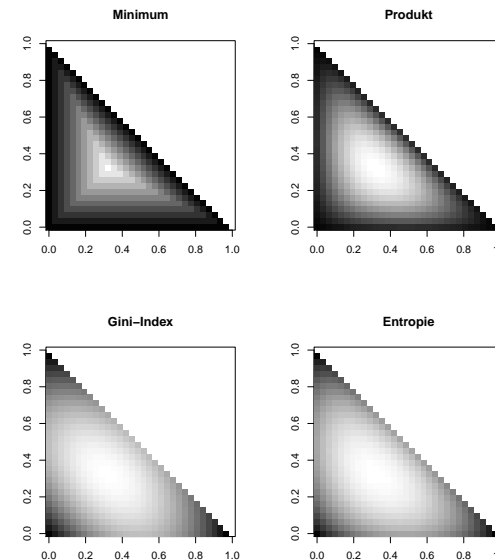
$$\mathfrak{S}_p(\mathbf{p}) \stackrel{\text{def}}{=} \prod_{\kappa} p_{\kappa}$$

$$\mathfrak{S}_e(\mathbf{p}) \stackrel{\text{def}}{=} - \sum_{\kappa} p_{\kappa} \cdot \log_2 p_{\kappa}$$

Für den Gini-Index gilt $\mathfrak{S}_g(\mathbf{p}) = 1 - \|\mathbf{p}\|^2$.

Homogenitätsmaße

Drei Ereignisse — Darstellung in der (p_1, p_2) -Ebene



Minimum/Produkt

Geringe Homogenität (Unsicherheit) wird bereits dann signalisiert, wenn nur eines der drei Ereignisse unwahrscheinlich ist.

➔ **unbrauchbar**

Entropie/Gini

Grundverschiedene Formeln, aber kaum unterschiedliche Funktionswerte.

➔ **praktisch äquivalent**

Rechenbeispiel (Gini-Index)

Ausgangsknoten β

Der Knoten β beherbergt die Verteilung $\mathbf{p} = (0.5, 0.3, 0.2)$, also gilt

$$\mathfrak{S}_{\text{Gini}}(\mathbf{p}) = 1 - 0.25 - 0.09 - 0.04 = 0.62$$

Erste Frage Q_1

Die Entscheidungsunsicherheiten der Q_1 -Nachfolger lauten

$$\mathfrak{S}_{\text{Gini}}(\beta^{(1)} | Q_1) = 1 - 0.64 - 0.04 = 0.32$$

$$\mathfrak{S}_{\text{Gini}}(\beta^{(2)} | Q_1) = 1 - 0.04 - 0.36 - 0.04 = 0.56$$

Die mittlere E.U. der Nachfolger bzw. ihre Reduktion betragen also

$$\mathfrak{S}_{\text{Gini}}(\beta | Q_1) = 0.5 \cdot 0.32 + 0.5 \cdot 0.56 = 0.44$$

$$\Delta_{Q_1} \mathfrak{S}_{\text{Gini}}(\beta) = 0.62 - 0.44 = 0.18$$

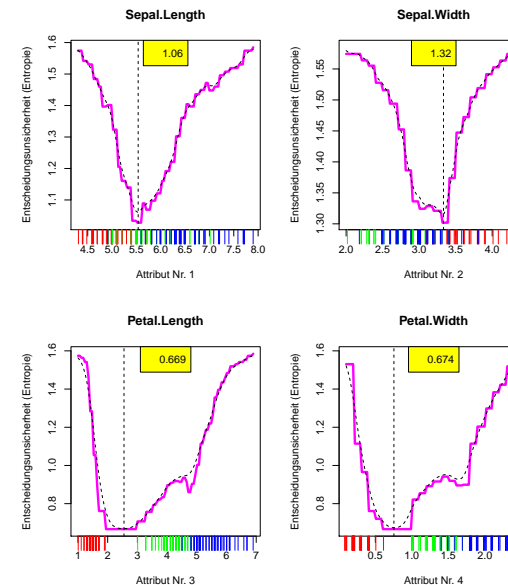
Zweite Frage Q_2

Auf dieselbe Weise errechnet sich für die konkurrierende Frage der Wert

$$\Delta_{Q_2} \mathfrak{S}_{\text{Gini}}(\beta) = 0.62 - 0.6 \cdot 0.5 - 0.4 \cdot 0.5 = 0.12$$

Folglich ist Q_1 der Frage Q_2 vorzuziehen.

Rechenbeispiel (Entropiemaß)



IRIS-Datensatz

150 Objekte

4 Attribute

3 Kategorien $\left\{ \begin{matrix} 50 \\ 50 \\ 50 \end{matrix} \right\}$

Wurzelknoten

Berechne für jedes Attribut x_n den EU-minimalen Schwellenwert θ_n

$$Q(\beta) : x_3 \stackrel{?}{\leq} 2.65$$

Reduktion der Entscheidungsunsicherheit

Entscheidungsunsicherheit im Knoten β

$$\mathfrak{S}(\beta) \stackrel{\text{def}}{=} \mathfrak{S}(\hat{\boldsymbol{\rho}}^{(\beta)}), \quad \hat{\rho}_{\kappa} \stackrel{\text{def}}{=} \frac{|\omega_{\kappa}(\beta)|}{|\omega(\beta)|}$$

Verzweigungswahrscheinlichkeiten der Frage Q in β

$$\hat{P}(\beta^{(i)}|\beta) \stackrel{\text{def}}{=} \frac{|\omega(\beta^{(i)})|}{|\omega(\beta)|}, \quad i = 1, \dots, L$$

Entscheidungsunsicherheit nach der Frage Q in β

$$\mathfrak{S}(\beta|Q) \stackrel{\text{def}}{=} \sum_i \hat{P}(\beta^{(i)}|\beta) \cdot \mathfrak{S}(\beta^{(i)})$$

Reduktion der Entscheidungsunsicherheit durch Q

$$\Delta_Q(\beta) \stackrel{\text{def}}{=} \mathfrak{S}(\beta) - \mathfrak{S}(\beta|Q)$$

Aufspüren und Tilgen nutzloser Teilbäume

CART Pruning

Züchtung eines überangepaßten Baumes

Sukzessive Vergrößerung (Entfernen schwacher Äste)

Auswahl des besten Teilbaumes

Lerndaten ω
Modellstrafterm
Validierungsdaten $\tilde{\omega}$

Lokaler Resubstitutionsfehler

Relative Anzahl der Fehler bei Entscheidung in β :

$$R(\beta) \stackrel{\text{def}}{=} \frac{\# \text{ falsch klassifiziert in } \beta}{\# \text{ alle Objekte}} = \frac{|\omega(\beta)| - \max_{\kappa} |\omega_{\kappa}(\beta)|}{|\omega(\beta_{\Delta})|}$$

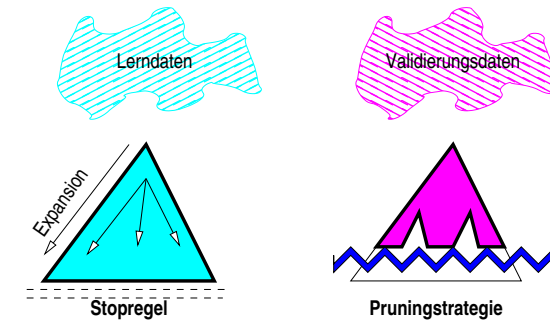
Kumulativer Resubstitutionsfehler

$\mathcal{B}_{\ell}^{\beta}$ = Menge aller Blattknoten in dem von β dominierten Teilbaum

$$R^*(\beta) \stackrel{\text{def}}{=} \sum_{\beta' \in \mathcal{B}_{\ell}^{\beta}} R(\beta')$$

Bemerkung Es gilt für alle $\beta \in \mathcal{B}$: $R^*(\beta) \leq R(\beta)$

Überanpassung an die Lernbeispiele



Fragmentierung der Lerndaten

- \mathcal{B} reinklassig $\rightsquigarrow \omega$ perfekt klassifiziert
- viele Lerndaten \rightsquigarrow großer Entscheidungsbaum
- Insignifikante Fragen in unteren Zweigen
- Unzuverlässige Entscheidung in den Blättern
- Stoppregeln sind „kurzsichtig“

Abhilfe

- „early stopping“
- Züchten & Zurückschneiden

Strafterme versus Kreuzvalidierung

Effizienz eines Teilbaums

Gut entmischende Teilbäume werden belohnt, aber zersplitterungsverdächtige Teilbäume werden bestraft!

$$\Delta_{\text{eff}}(\beta) \stackrel{\text{def}}{=} \frac{\text{Fehlerzuwachs in } \beta}{\# \text{ eingesparte Knoten}} = \frac{R(\beta) - R^*(\beta)}{|\mathcal{B}_{\ell}^{\beta}| - 1}$$

Kreuzvalidierungsfehler

Jedem Objekt $\mathbf{x} \in \tilde{\omega}$ wird durch einen Entscheidungsbaum ein Blattknoten $\beta(\mathbf{x})$ und damit auch eine Klassenmarkierung $\delta_{\ell}(\beta(\mathbf{x}))$ zugeordnet.

$$\tilde{\varepsilon}(\mathcal{B}) \stackrel{\text{def}}{=} \frac{\sum_{\kappa=1}^K |\{\mathbf{x} \in \tilde{\omega}_{\kappa} \mid \delta_{\ell}(\beta(\mathbf{x})) \neq \kappa\}|}{|\tilde{\omega}(\beta_{\Delta})|}$$

CART Pruning-Algorithmus

Breiman, Friedman, Olshen & Stone (1984)

1 ZÜCHTEN

Expandiere initialen Baum $\mathcal{B}^{(0)}$ mittels Lerndaten $\omega_1, \dots, \omega_K$ unter Einhaltung des „Reinheitsgebotes“.

2 SUKZESSIVES ZURÜCKSCHNEIDEN

Erzeuge eine Folge gestutzter Teilbäume von $\mathcal{B}^{(0)}$

- Setze $i \rightarrow 0$.
- Berechne alle Effizienzwerte $\Delta_{\text{eff}}(\beta)$, $\beta \in \mathcal{B}^{(i)}$.
- Wähle Knoten $\beta^* \in \mathcal{B}^{(i)}$ mit minimaler Effizienz.
- Kappe den Teilbaum unterhalb β^* .
- Setze $i \leftarrow i + 1$ und bezeichne gekürzten Baum als $\mathcal{B}^{(i)}$.
- Ist $\mathcal{B}^{(i)} \neq \{\beta_{\Delta}\}$, dann \rightsquigarrow b.

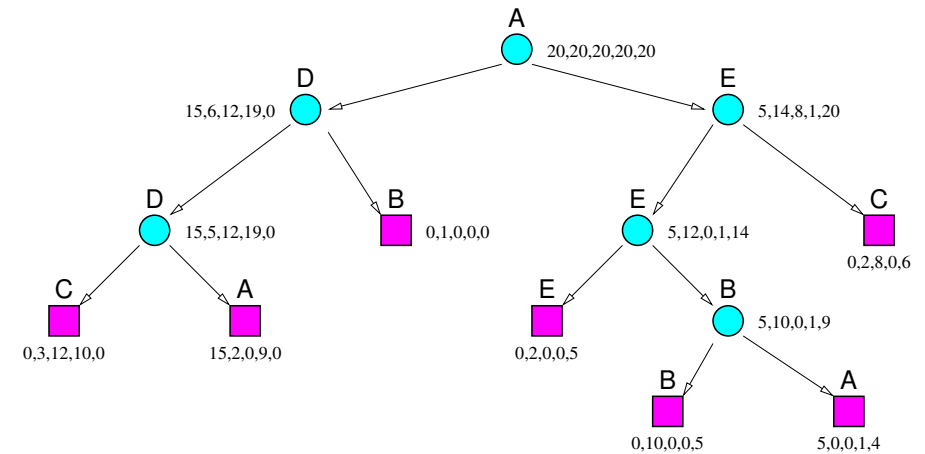
3 AUSWAHL NACH VALIDIERUNGSFEHLER

Wähle aus $\{\mathcal{B}^{(i)} \mid i = 0, 1, 2, \dots\}$ denjenigen Baum mit geringstem Fehler auf den Validierungsdaten $\tilde{\omega}_1, \dots, \tilde{\omega}_K$.

(sumfithog1A)

Beispiel — CART-Algorithmus

5 Klassen · 6 Fragen · 7 Blätter · 100 Objekte

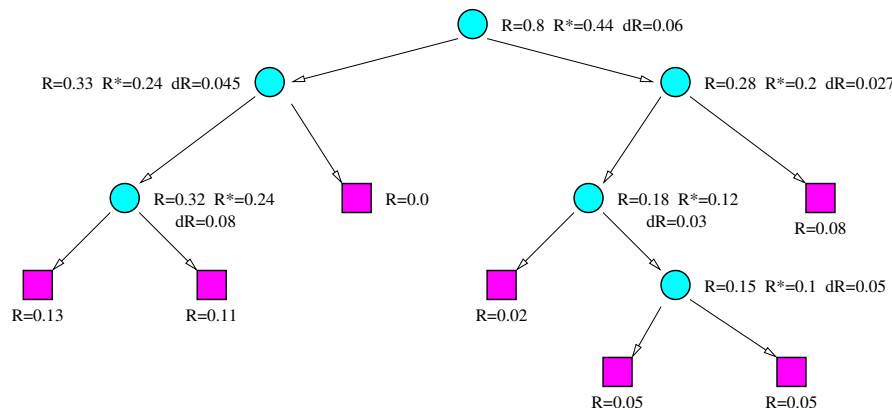


Klassenhäufigkeiten je Knoten

Bestklassenmerkung je Knoten

Beispiel — CART-Algorithmus

5 Klassen · 6 Fragen · 7 Blätter · 100 Objekte



Lokale Resubstitutionsfehlerraten

Kumulative Resubstitutionsfehlerraten

Effizienzen — nur innere Knoten werden gezählt

Kreuzvalidierendes Stutzen der Äste

„Frühe Validierung“ — schon zur Bewertung statt erst zur Auswahl

Lokale Fehlerrate

im Knoten β nach Einschleusen der Konterdaten ω :

$$\varepsilon(\beta) \stackrel{\text{def}}{=} 1 - \frac{|\omega_{\kappa}(\beta)|}{|\omega(\beta)|} \quad \text{mit } \kappa := \delta_{\ell}(\beta) \text{ oder } \kappa := \underset{\lambda}{\operatorname{argmax}} |\omega_{\lambda}(\beta)|$$

Kumulative Fehlerrate

nach Durchschleusen von ω bis zu den Blättern:

$$\varepsilon^*(\beta) \stackrel{\text{def}}{=} \sum_{\beta' \in \mathcal{B}_{\ell}^{\beta}} \frac{|\omega(\beta')|}{|\omega(\beta)|} \cdot \varepsilon(\beta')$$

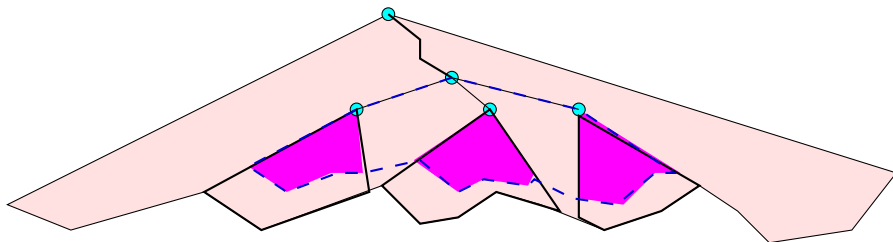
Die **Gesamtfehlerrate** ist $\varepsilon(\mathcal{B}) = \varepsilon^*(\beta_{\Delta})$

Minimale Fehlerrate

aller Teilbäume \mathcal{B}^{β} unterm Knoten β :

$$\varepsilon^{\forall}(\beta) \stackrel{\text{def}}{=} \min \{ \varepsilon(\mathcal{B}') \mid \mathcal{B}' \text{ Teilbaum von } \mathcal{B}^{\beta} \}$$

Induktive Bottom-Up Beschneidung



Lemma

Sei $(\mathcal{B}, Q, \delta_\ell)$ ein binärer Entscheidungsbaum über $\Omega = \mathbb{R}^D$. Die optimale Fehlerrate des Teilbaums \mathcal{B}^β berechnet sich nach folgender Rekursion:

$$\varepsilon^\forall(\beta) = \begin{cases} \varepsilon(\beta) & \beta \in \mathcal{B}_\ell \\ \min \left(\sum_{\beta' \prec \beta} p(\beta'|\beta) \cdot \varepsilon^\forall(\beta') \right) & \beta \notin \mathcal{B}_\ell \end{cases}$$

Gelfands IEP-Algorithmus

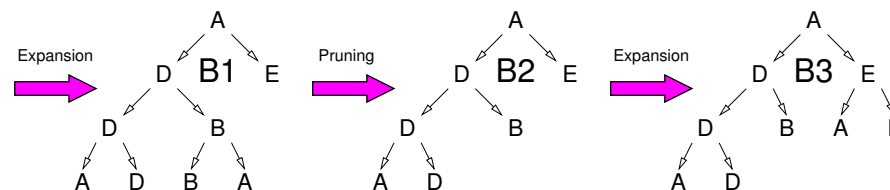
„Iterative Expansion-Pruning“

- 1 INITIALISIERUNG
Setze $i \leftarrow 0$ und $\mathcal{B}^{(0)} \leftarrow \{\beta_\Delta\}$.
- 2 ERSTES EXPANDIEREN
Expandiere $\mathcal{B}^{(i)}$ mit den Daten ω^a . $\rightsquigarrow \mathcal{B}^{(i+1)}$
- 3 ERSTES STUTZEN
Beschneide $\mathcal{B}^{(i+1)}$ mit den Daten ω^b . $\rightsquigarrow \mathcal{B}^{(i+2)}$
- 4 ZWEITES EXPANDIEREN
Expandiere $\mathcal{B}^{(i+2)}$ mit den Daten ω^b . $\rightsquigarrow \mathcal{B}^{(i+3)}$
- 5 ZWEITES STUTZEN
Beschneide $\mathcal{B}^{(i+3)}$ mit den Daten ω^a . $\rightsquigarrow \mathcal{B}^{(i+4)}$
- 6 TERMINIERUNG
Falls $\mathcal{B}^{(i+2)} \equiv \mathcal{B}^{(i+4)}$, dann \rightsquigarrow ENDE.
- 7 WIEDERHOLUNG
Setze $i \leftarrow i + 4$ und weiter bei \rightsquigarrow 2.

$$\left\{ \begin{array}{l} \omega_1^a \cup \omega_2^a \cup \dots \cup \omega_K^a \\ \omega_1^b \cup \omega_2^b \cup \dots \cup \omega_K^b \end{array} \right\} = \left\{ \begin{array}{l} \omega^a \\ \omega^b \end{array} \right\}$$

Disjunkt:

Wiederholtes Züchten und Beschneiden



Expansionsphase · top-down

1. Schleuse die Daten ω^a bis zu den Blattknoten von $\mathcal{B}^{(i)}$.
2. Bestimme die Mengen $\omega_\kappa^a(\beta)$ für alle κ, β .
3. Züchte für alle Blattknoten $\beta \in \mathcal{B}_\ell^{(i)}$ einen Teilbaum unter β mittels $\omega^a(\beta)$.

Pruningphase · bottom-up

1. Schleuse die Daten ω^a bis zu den Blattknoten von $\mathcal{B}^{(i)}$.
2. Markiere alle $\beta \in \mathcal{B}^{(i)}$ mit neuen Klassen $\delta_\ell(\beta)$.
3. Überprüfe alle β durch Vergleich von lokaler und minimaler RFR auf Eliminierbarkeit.

Die Auswahl der besten Frage

Monothetische Knoten \rightsquigarrow keine Attributkombinationen

Problem

Die Expansion eines jeden Knotens β im TDI-Algorithmus erfordert die $\Delta_Q(\beta)$ -Bewertung **jeder Frage** Q zu **jedem Attribut** \mathcal{X}_n !

Nominale Attribute

Wieviele Zwei- oder Mehrwege-Fragen sind zu testen?

- Wertverzweigung eine Frage/Attribut
- Attribut-Wert-Gleichung $|\mathcal{X}_n|$ Targets/Attribut
- Literalkomplex $2^{|\mathcal{X}_n|}/2$ Mengen/Attribut

Numerische und ordinale Attribute

Wieviele Schwellenwert-Fragen sind zu testen?

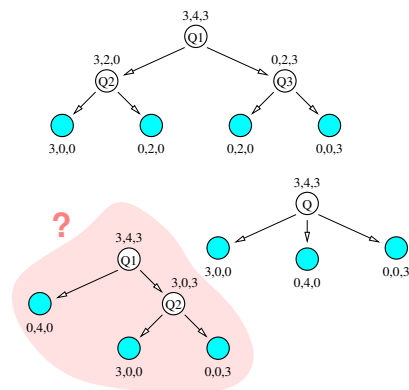
- Ordinale Attribute $|\mathcal{X}_n| - 1$ Schwellen/Attribut
- Numerische Attribute $|\omega(\beta)| - 1$ Schwellen/Attribut

Die Befragung nominaler Attribute

Symmetrische Verzweigung $x_n = ?$ versus asymmetrische Verzweigung $x_n \stackrel{?}{=} \xi_\ell$

Datenfragmentierung

Die minimal zersplitternde Folge **binärer** Fragen wird nicht automatisch gefunden.



Unbalancierte Auswahl

Die Maximierung der Entscheidungssicherheit bevorzugt systematisch Fragen mit **hohem Verzweigungsfaktor**.

Gain Ratio Impurity

Abhilfe schafft Normierung auf die maximale Entropie:

$$\Delta'_Q(\beta) \stackrel{\text{def}}{=} \frac{\mathfrak{I}(\beta) - \sum_{j=1}^L p_j \cdot \mathfrak{I}(\beta_j)}{\mathcal{H}(p_1, \dots, p_L)}$$

Literalkomplexe in Zweiklassen-Szenarien

Auswahl der besten Teilmenge

Aufgabenstellung

Finde zum Attribut \mathcal{X}_n in β diejenige Teilmengenanfrage

$$Q : x_n \mapsto \begin{cases} 1 & x_n \in U \\ 0 & x_n \notin U \end{cases}, \quad U \subset \mathcal{X}_n$$

mit der max. Reduktion $\Delta_Q(\beta)$ der Entscheidungsunsicherheit.

Premiumschlitten & Volumenmodelle

Objekte = Fahrzeuge · Klassen Ω_1 und Ω_2 · Attribut x_{19} (Hersteller)

\mathcal{X}_{19}	VW	Benz	Alfa	Dacia	BMW	Porsche
Ω_1	112	9	3	1	28	5
Ω_2	112	1	2	4	12	0
$\hat{P}(1 \xi)$	0.5	0.9	0.6	0.2	0.7	1.0

Aufsteigende $\hat{P}(1|\xi)$ -Sortierung → nur 5 mögliche **optimale** Fragen:

$$\begin{aligned} x_{19} &\in \{\text{Dacia, VW}\} & x_{19} &\in \{\text{Dacia, VW, Alfa, BMW}\} \\ x_{19} &\in \{\text{Dacia}\} & x_{19} &\in \{\text{Dacia, VW, Alfa}\} & x_{19} &\in \{\text{Dacia, VW, Alfa, BMW, Benz}\} \end{aligned}$$



Der Zwillingsatz („Twoing Theorem“)

Linearer Suchaufwand für entropiegesteuertes Zweiklassen-Lernen

Satz

Es sei $\mathcal{X}_n = \{\xi_1, \dots, \xi_L\}$ der (nominale) Wertebereich des n -ten Attributs, und es zerfalle die Lernstichprobe $\omega \subset \Omega$ in zwei Klassenbereiche ω_1, ω_2 . Mit den Bezeichnungen

$$\hat{P}(\kappa|\xi_\ell) \stackrel{\text{def}}{=} \frac{|\{\mathbf{x} \in \omega_\kappa \mid x_n = \xi_\ell\}|}{|\{\mathbf{x} \in \omega \mid x_n = \xi_\ell\}|}$$

für $\kappa = 1, 2$ und $\ell = 1, \dots, L$ seien infolge geeigneter Sortierung der ξ_ℓ die Häufigkeitsbeziehungen

$$\hat{P}(1|\xi_1) \leq \hat{P}(1|\xi_2) \leq \dots \leq \hat{P}(1|\xi_L)$$

gültig. Dann besitzt die Teilmengenanfrage mit der maximalen Homogenitätsreduktion in Bezug auf das Entropiemaß die Gestalt

$$x_n \in \{\xi_1, \dots, \xi_\ell\}$$

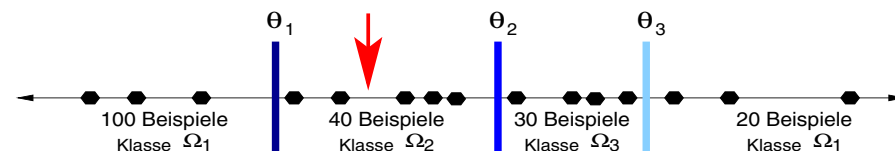
für ein geeignetes ℓ mit $1 < \ell < L$.

Schwellenwertfragen

Numerische und ordinale Attribute · zwei oder mehr Klassen

Reduzierter Suchaufwand für $\theta \in \mathcal{X}_n = \mathbb{R}$

- Nur $T_\beta = |\omega(\beta)|$ **Mittelpunktschwellen** zu prüfen.
- Nur **klassentrennende** Schwellen können $\Delta_Q(\beta)$ -maximal sein.
- Es gibt eine **Rekursionsformel** für $\Delta_{Q,n,\theta}(\beta)$.



Sortierung $O(T \log T)$

Aufsteigendes Sortieren der \mathcal{X}_n -Attributwerte in $\omega(\beta)$:

$$a_1 < a_2 < a_3 \dots < a_t < \dots < a_{T_\beta}$$

Mittelpunktschwellen

Suffizienter Satz von Schwellenwerten für $Q_{n,\theta}$:

$$\theta_t = \frac{a_{t+1} - a_t}{2}, \quad t = 1, 2, \dots, T_\beta - 1$$

Separierende Schwellenwerte

Definition

Eine Mittelpunktschwelle θ_t von $\{x_n \mid \mathbf{x} \in \omega(\beta)\}$ heißt **innere Schwelle** von \mathcal{X}_n in β , falls alle Objekte $\mathbf{x} \in \omega(\beta)$ mit $x_n = a_t$ oder $x_n = a_{t+1}$ zu einundderselben Klasse Ω_κ gehören.

Andernfalls heißt θ_n **separierende Schwelle** oder **Klassengrenze**.

Lemma (Fayyad & Irani, 1992)

Sind $[\theta_t]$ die Mittelpunktschwellen zur assoziierten Stichprobe $[\omega_\kappa(\beta)]$ von β zum Attribut \mathcal{X}_n , und gilt

$$\theta_{t^*} = \underset{\theta_t}{\operatorname{argmax}} \Delta_{\{x_n \leq \theta_t\}}(\beta)$$

für die entropiebezogene Entscheidungsunsicherheit, so ist θ_{t^*} notwendigerweise eine Klassengrenze.

Bemerkung

Je stärker sich die Objekte klassenweise auf der \mathcal{X}_n -Achse häufen, desto weniger Reduktionswerte müssen berechnet werden.

Attribute mit Fehlanzeigen

Imputation

Wenn $x_n = ?$, so setze einen Standardwert $\hat{\xi}$ ein.

- Wähle für $\hat{\xi}$ das globale Attributmittel μ_n .
- Wähle für $\hat{\xi}$ das **lokale** Attributmittel $\mu_n(\beta)$.

Überlagerung

Wenn $x_n = ?$, so folge in der Abrufphase parallel allen Verzweigungen.

- Während der Lernphase werden defiziente Objekte bei der $\Delta_Q(\beta)$ -Berechnung **ignoriert** oder **pejorisiert**.

Surrogate Split

Wenn $x_i = ?$, so beantworte in der Abrufphase die/eine Ersatzfrage.

- In der Lernphase merkt man/frau sich die besten Fragen zum zweitbesten Attribut (ggf. weitere Alternativen).

Inkrementelle $\Delta_Q(\beta)$ -Berechnung

Lemma

Es seien $\theta_1 < \theta_2$ zwei benachbarte Klassengrenzen für \mathcal{X}_n in $\omega(\beta)$, zwischen denen genau m Muster der Klasse Ω_κ liegen. Dann gilt die Rekursionsformel

$$\begin{aligned} \Delta_{\{x_n \leq \theta_2\}}(\beta) &= \Delta_{\{x_n \leq \theta_1\}}(\beta) \\ &+ \frac{h(\ell, r) - h(\ell + m, r + m) + h(\ell_\kappa + m, r_\kappa + m) - h(\ell_\kappa, r_\kappa)}{T} \end{aligned}$$

mit den Abkürzungen

$$h(p, q) = p \log_2 p - q \log_2 q$$

und den Zählwerten

$$\begin{aligned} \ell_\kappa &= |\{x_d < \theta_1 \mid \mathbf{x} \in \omega_\kappa\}| & \ell &= \sum_\kappa \ell_\kappa \\ r_\kappa &= |\{x_d > \theta_1 \mid \mathbf{x} \in \omega_\kappa\}| & r &= \sum_\kappa r_\kappa \end{aligned}$$

Polythetische Entscheidungsfragen

Über das Züchten „schiefer“ statt achsenparalleler Entscheidungsbäume

Attributübergreifende Dichotomien (linear)

$$a_0 + \sum_{i=1}^N a_i x_i \stackrel{?}{\leq} 0$$

Trennfunktionsparameter mit guter Klassenentmischung!

- **CART/LC**
Gradientenabstieg via $\Delta_a(\beta)$
- **SADT**
Simulated Annealing of Decision Trees
- **LMDT**
Linear Machine Decision Trees („ADALINE-Knoten“)
- **QUEST**
Multivariate Variante des QUEST-Algorithmus

QUEST-Algorithmus

Quick Unbiased Efficient Statistical Tree

(Algorithmus)

1 KLASSENBEDINGTE MITTELWERTE

Berechne β -lokale klassenbezogene Mittelwertvektoren μ_1, \dots, μ_K .

2 STATISTISCHER HYPOTHESENTEST

Fisher-Test für die Nullhypothesen

$$H_0 : \mu_1^{(n)} = \mu_2^{(n)} = \dots = \mu_K^{(n)}$$

3 ZENTREN CLUSTERN („2-means“)

Partitioniere für das Gewinnerattribut $n^* \in \{1 : N\}$ die K Mittelwerte $\mu_1^{(n^*)}, \mu_2^{(n^*)}, \dots, \mu_K^{(n^*)}$.

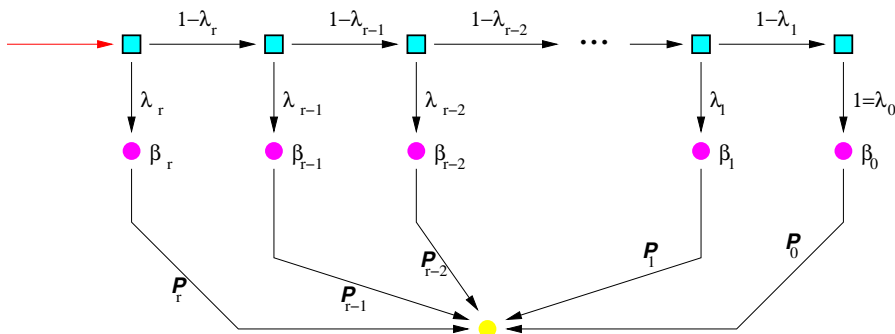
4 KONSTRUIERE TRENNFRAGE

Berechne NV-Dichteparameter für die beiden Cluster.

$$\mathcal{N}(x_{n^*} | \mu_1, \sigma_1^2) \stackrel{?}{\leq} \mathcal{N}(x_{n^*} | \mu_2, \sigma_2^2)$$

(sumfithogfA)

Lineare Interpolation von Klassenprädiktoren



Interpolationsformel für a posteriori-Klassenwahrscheinlichkeiten

Maximum-Likelihood-Koeffizienten nach EM-Algorithmus

$$\tilde{p}_\kappa(\beta_r) = \begin{cases} 1 \cdot \hat{p}_\kappa(\beta_\Delta) & r = 0 \\ \lambda_r \cdot \hat{p}_\kappa(\beta_r) + (1 - \lambda_r) \cdot \tilde{p}_\kappa(\beta_{r-1}) & r > 1 \end{cases}$$

Klassenprädiktoren

in den inneren und den Blattknoten des Entscheidungsbaumes

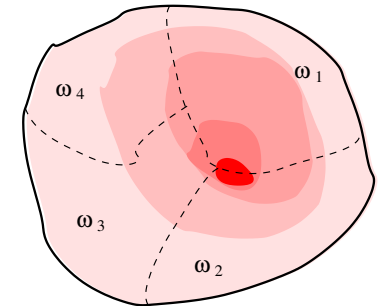
Schleusungspfad

Jedes Objekt $\mathbf{x} \in \mathcal{X}$ beschreibt einen Pfad

$$\beta_\Delta = \beta_0(\mathbf{x}) \prec \beta_1(\mathbf{x}) \prec \beta_2(\mathbf{x}) \prec \dots \prec \beta_{r-1}(\mathbf{x}) \prec \beta_r(\mathbf{x}) = \beta(\mathbf{x})$$

Lokale Prädiktoren

$\hat{p}_\kappa(\beta_0)$	$=$	$ \omega_\kappa $	$/$	$ \omega $
\mathbb{I}		\vee		\vee
$\hat{p}_\kappa(\beta_1)$	$=$	$ \omega_\kappa(\beta_1) $	$/$	$ \omega(\beta_1) $
\mathbb{I}		\vee		\vee
$\hat{p}_\kappa(\beta_2)$	$=$	$ \omega_\kappa(\beta_2) $	$/$	$ \omega(\beta_2) $
\mathbb{I}		\vee		\vee
\dots		\dots		\dots
\mathbb{I}		\vee		\vee
$\hat{p}_\kappa(\beta_r)$	$=$	$ \omega_\kappa(\beta_r) $	$/$	$ \omega(\beta_r) $



Leo Breimans Random Forests

Zweifache Ensembletechnik: Objekte (*bagging*) & Attribute

(Algorithmus)

Lernprobe ω , Wälder $M \in \mathbb{N}$, Auswahl $T_b \leq |\omega|$ und $N_b \ll N$.

1 LERNPHASE

Erzeuge Bäume $\mathcal{B}^{(m)}$, $m = 1, \dots, M$:

- Lernprobe $\omega^{(m)} \subset \omega$ via T_b -**Bootstrap** mit Ersetzen
- Zufallsbaum $\mathcal{B}^{(m)}$ via TDI-Algorithmus
- EINGESCHRÄNKTE LOKALE FRAGEAUSWAHL:
 $\mathcal{A}_\beta \subset \{\mathcal{X}_1, \dots, \mathcal{X}_N\}$, $|\mathcal{A}_\beta| = N_b$ via N_b -**Bootstrap** ohne Ers.
- Kein Zurückstutzen!

2 ABRUFPHASE

Mehrheitsentscheidung unter allen Bäumen des Waldes:

$$\kappa^*(\mathbf{x}) = \operatorname{argmax}_\kappa \{ |\mathcal{B}^{(m)} | \delta_\ell(\beta^{(m)}(\mathbf{x})) = \kappa \}$$

(sumfithogfA)

Bemerkung

Pro: Effizient, skalierbar (N , T), exzellentes Erkennungsverhalten.

Contra: Überanpassung, Reproduzierbarkeit, Präferenz stufenreicher Nominalattribute.

Prädiktion, Regression & Klassifikation

Konzeptlernen

Versionenräume

Naive Bayesregel

Multivariate lineare Regression

Logistische Regression

Ordinale Regression und Präferenzmodelle

Statistische Entscheidungsbäume

Zusammenfassung

Zusammenfassung (4)

1. Der **Konzeptraum** enthält die **zu lernenden**, der **Hypothesenraum** die **lernbaren** Teilmengen des Objektraums.
2. Der **Versionenraum** besteht aus allen **konsistenten** Hypothesen und ist als **Halbordnungsintervall** darstellbar.
3. Die Hypothesen des **Sterns** grenzen ein Positivbeispiel gegen alle Negativbeispiele ab.
4. **Lineare Diskriminanten** approximieren die **ideale Trennfunktion** im Quadratmittelsinn.
5. **Loglineare Diskriminanten** approximieren die **a posteriori** Klassenwahrscheinlichkeiten.
6. Beide Lernverfahren lassen sich **regularisieren** und **dualisieren**.
7. **Entscheidungsbäume** klassifizieren durch hierarchische Befragung **numerischer & diskreter** Attribute.
8. Sie werden durch ein **gieriges Top-Down-Verfahren** aus den Daten gelernt.
9. Für die lokale Suche nach der maximal **klassenentmischenden** Frage gibt es effiziente Verfahren.