

MASCHINELLES LERNEN & DATAMINING

Vorlesung im Wintersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 25. August 2017

Grundbegriffe des Data Mining

Datensätze mit expliziter oder impliziter Objektcharakterisierung

Datensatz

Menge oder **Folge** von **Objekten** („Instanzen“) des Aufgabenbereichs Ω mit ihren Eigenschaften und/oder Beziehungen

Attribut

Objekteigenschaft $\hat{=}$ Element eines **Wertebereichs** \mathcal{X} („Skala“)

Beziehung

Relation $\mathcal{R} \subset \Omega \times \Omega$ zwischen Objekten oder ...

Abstand/Ähnlichkeit $d : \Omega \times \Omega \rightarrow \mathbb{R}$ zwischen Objekten

	o_1	o_2	o_3	o_4		o_1	o_2	o_3	o_4		\mathcal{X}_1	\mathcal{X}_2	\mathcal{X}_3	\mathcal{X}_4
o_1			∞		o_1	0	3	8	15	o_1	+	low	1.2	+4
o_2				∞	o_2	3	0	5	12	o_2	-	hi	0.5	-3
o_3	∞				o_3	8	5	0	7	o_3	+	hi	2.3	-3
o_4		∞			o_4	15	12	7	0	o_4	+	med	2.1	-7

Teil II

Datenaufbereitung

Datenmatrix $\hat{=}$ Objekte \times Attribute

Objekteigenschaften erster Ordnung

Definition

Sind $\mathcal{X}_1, \dots, \mathcal{X}_N$ die Attribute eines Objekts, so heißen die Elemente

$$\mathbf{x} = (x_1, \dots, x_N)^T \in \mathcal{X}_1 \times \dots \times \mathcal{X}_N =: \mathcal{X}$$

Variable

$\left\{ \begin{array}{l} \text{Typ} \\ \text{Name} \\ \text{Wert} \end{array} \right\}$

Datenvektoren des Objekts.

Die Menge \mathcal{X} heißt **Wertebereich** des Objekts.

Eine (Multi-)Menge $\{\mathbf{x}_1, \dots, \mathbf{x}_T\} \subset \mathcal{X}$ oder eine Folge $(\mathbf{x}_1, \dots, \mathbf{x}_T) \in \mathcal{X}^T$ bezeichnen wir als **Datenmatrix** oder ggf. als **Meßreihe**.

Datum

$\left\{ \begin{array}{l} \text{Attribut} \\ \text{Objekt} \\ \text{Eintrag} \end{array} \right\}$

Schreibweise

Datenvektoren

Meßwerte

Reelle Datenmatrix
$$\begin{pmatrix} x_{1,1} & \dots & x_{1,N} \\ \vdots & \ddots & \vdots \\ x_{T,1} & \dots & x_{T,N} \end{pmatrix}$$

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_t, \dots$$

$$\mathbf{x}_t = (x_{t,1}, \dots, x_{t,N})^T$$

$$\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_T)^T \in \mathbb{R}^{T \times N}$$

Werteskalen

Relationen und Distanzen

Skalenkonversion

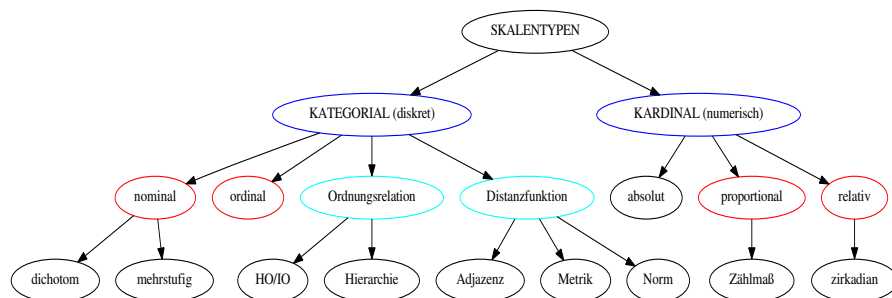
Detektion von Ausreißern

Imputation von Fehlanzeigen

Zusammenfassung

Skalentypen

Objektattribut $\hat{=}$ (Wertebereich, Operatorenmenge)



Diskrete Skala

Endlicher Wertebereich

$\mathcal{X} = \{\xi_1, \xi_2, \dots, \xi_K\}$

Numerische Skala

Kontinuierlicher Wertebereich

$\mathcal{X} \subseteq \mathbb{R}$

Attribute und ihre Skalentypen

Was bedeuten die Spalteneinträge einer Datenmatrix ?

Beispiel

	vorbestraft	Partei	Abinote	Geburt	Spenden
Angela	F	CDU	gut	1954	$345 \cdot 10^3$
Guido	F	FDP	ausreichend	1961	$137 \cdot 10^3$
Roland	T	CDU	gut	1958	$3.6 \cdot 10^6$
Gregor	F	PDS	sehr gut	1948	NA
Linus	F	Pirat	NA	1969	0
Bill	F	Rep	mangelhaft	1955	$-4.2 \cdot 10^9$
Roman	T	1933	...
⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮
⋮	⋮	⋮	⋮	⋮	⋮
	$\{T, F\}$	$\{\pi_1, \dots, \pi_9\}$	$\{\nu_1, \dots, \nu_5\}$	$\mathbb{Z} \subset \mathbb{R}$	\mathbb{R}

Typische Wertebereiche

Nominalskala

- Dichotomien $\{0, 1\}, \{T, F\}, \{+, -\}, \{m, f\} \dots$
- Zeichensätze $\{C, G, A, T\}$
- Farben („red“, green, blue)
- Gruppen & Prädikate

Ordinalskala

- Notenskala
„sehr gut“, „gut“, „befriedigend“, ...
- Unscharfe Prädikate
„kalt“, „kühl“, „lau“, „warm“, „heiß“
- Eingefrorene Quantitäten
„2-türig“, „4-türig“, „5-türig“

Intervallskala (relativ)

- Temperaturen $20^\circ\text{C}, 451^\circ\text{F}$
- Zeitangaben
1066, 2001/09/11, 469 v.Chr., ...

Verhältnisskala (absolut/proport.)

- absol. Temperatur 273°K
- Dauer 45 min, $13.7 \cdot 10^9$ Jahre
- Mengenangaben
 $C = 2.98, 17 \text{ cm}, 8 \mu\text{g}, \dots$

Binäre Skalenoperationen

Vergleichsoperationen $\mathcal{X} \times \mathcal{X} \rightarrow \{T, F\}$ · Rechenoperationen $\mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

Nominalskala

= Gleichheitstest

Alle Attributwerte ξ_ℓ sind gleichberechtigt.

Intervallskala (relativ)

– Differenzbildung

Unterschiede sind durch $x_1 - x_2$ quantifizierbar.

Ordinalskala

< Vergleichbarkeit

Abschnittsbildung nach Totalordnung:
 $\{\xi \mid \xi \preceq \xi_\ell\}$

Verhältnisskala (absolut/proport.)

÷ Quotientenbildung

Wohldefiniert:
„Nullpunkt“, „doppelt“, „Drittel“

Durchschnittswerte

Verallgemeinerung auf Metriken und normierte Vektorräume

Beispiel

Für die Wertemenge $\{1, 1, 1, 2, 2, 5, 9\}$ gilt:

$$\mu^{\text{mod}} = 1, \quad \mu^{\text{med}} = 2, \quad \mu^{\text{mean}} = 3, \quad \mu^{\text{geo}} = 2.0998$$

Definition

In einem metrischen Raum (\mathcal{X}, d) heißt der Wert

$$\mu^{\text{zen}} = \operatorname{argmin}_{z \in \mathcal{X}} \left(\sum_{t=1}^T d(z, x_t) \right)$$

das **Zentroid** der (Multi-)Menge $\{x_1, \dots, x_T\}$.

Lemma

- (1) Es ist $\mu^{\text{mean}}(\cdot)$ das Zentroid zur euklidischen Metrik $d(y, z) = (y - z)^2$.
- (2) Es ist $\mu^{\text{med}}(\cdot)$ das Zentroid zur Betragsmetrik $d(y, z) = |y - z|$.
- (3) Es ist $\mu^{\text{mod}}(\cdot)$ das Zentroid zur diskreten Metrik $d(y, z) = 1 - \delta_{y,z}$.

Durchschnittswerte

Wie berechnet man/frau einen für $(x_1, \dots, x_T) \in \mathcal{X}^T$ „(proto)typischen“ Wert ?

Nominalskala

Modus — der häufigste Wert:

$$\mu^{\text{mod}} = \xi_{\ell^*} \text{ mit}$$

$$\ell^* = \operatorname{argmax}_{\ell} N_{\ell}$$

mit den absoluten Häufigkeiten

$$N_{\ell} = \sum_{t=1}^T \delta_{x_t, \xi_{\ell}}$$

Ordinalskala

Median — der mittlere Wert:

$$\mu^{\text{med}} = \xi_{\ell^*} \text{ mit}$$

$$\sum_{k=1}^{\ell^*-1} N_k \leq \frac{T}{2} \leq \sum_{k=1}^{\ell^*} N_k$$

falls das Inventar \mathcal{X} geordnet ist:

$$\xi_1 < \xi_2 < \dots < \xi_{\ell} < \xi_{\ell+1} < \dots < \xi_L$$

Intervallskala (relativ)

Arithmetisches Mittel

$$\mu^{\text{mean}} = \frac{1}{T} \cdot \sum_{t=1}^T x_t = \frac{1}{T} \cdot \sum_{\ell=1}^L N_{\ell} \cdot \xi_{\ell}$$

Verhältnisskala (absolut/proport.)

Geometrisches Mittel

$$\mu^{\text{geo}} = \sqrt[T]{\prod_{t=1}^T x_t} = \sqrt[T]{\prod_{\ell=1}^L \xi_{\ell}^{N_{\ell}}}$$

Durchschnittswerte

Verallgemeinerung von (endlichen) Wertemengen auf diskrete Verteilungen

Definition

Es sei \mathbb{X} eine diskrete Zufallsvariable über dem Wertebereich $\mathcal{X} \subset \mathbb{R}$ mit der Wahrscheinlichkeitsfunktion $P(\cdot)$. Dann heißt

$$\mu(\mathbb{X}) = \mathcal{E}[\mathbb{X}] \stackrel{\text{def}}{=} \sum_{x \in \mathcal{X}} x \cdot P(\mathbb{X} = x)$$

der **Erwartungswert** von \mathbb{X} , es heißt

$$\mu^{\text{med}}(\mathbb{X}) = \xi \quad \text{mit} \quad \sum_{x < \xi} P(\mathbb{X} = x) \leq \frac{1}{2} \leq \sum_{x \leq \xi} P(\mathbb{X} = x)$$

der **Median** von \mathbb{X} , und es heißt

$$\mu^{\text{mod}}(\mathbb{X}) \stackrel{\text{def}}{=} \operatorname{argmax}_{x \in \mathcal{X}} P(\mathbb{X} = x)$$

der **Modus** von \mathbb{X} .

Durchschnittswerte

Verallgemeinerung von (endlichen) Wertemengen auf stetige Verteilungen

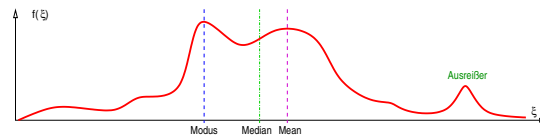
Definition

Für eine kontinuierliche Zufallsvariable über dem Wertebereich $\mathcal{X} = \mathbb{R}$ mit der Wahrscheinlichkeitsdichtefunktion $f_{\mathbb{X}}(\cdot)$ gilt entsprechend:

$$\begin{aligned}\mu(\mathbb{X}) &\stackrel{\text{def}}{=} \mathcal{E}[\mathbb{X}] = \int_{\mathbb{R}} x \cdot f_{\mathbb{X}}(x) dx \\ \mu^{\text{med}}(\mathbb{X}) &\stackrel{\text{def}}{=} \xi \quad \text{mit} \quad \int_{-\infty}^{\xi} f_{\mathbb{X}}(x) dx = \frac{1}{2} \\ \mu^{\text{mod}}(\mathbb{X}) &\stackrel{\text{def}}{=} \operatorname{argmax}_{x \in \mathbb{R}} f_{\mathbb{X}}(x)\end{aligned}$$

Bemerkung

Die Mediandefinition erfordert eine stetige und streng monotone Wahrscheinlichkeitsverteilungsfunktion.



Relationen auf diskreten Attributen

Spezialfall: Objekte besitzen genau ein Attribut \mathcal{X}

Adjazenz

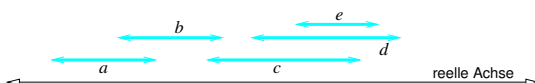
Die Matrix $\mathbf{A} \in \{0, 1\}^{L \times L}$ repräsentiert eine **(Objekt)nachbarschaft**.

- räumliche Nähe, Verwandtschaft, Interaktion ...
- „Elter-von“, Einflußnahme, ...

Präferenz

Die Relation $\mathcal{R} \subset \mathcal{X} \times \mathcal{X}$ repräsentiert eine (nicht notwendig totale) **Ordnung**.

- Halbordnung, Verband, Boolesche Algebra
- Turnier, (echte) Intervallordnung



Bemerkung

Zyklus: $c \preceq b \prec d \preceq c$
 \neg transitiv: $c \preceq b \preceq a$

Werteskalen

Relationen und Distanzen

Skalenkonversion

Detektion von Ausreißern

Imputation von Fehlanzeigen

Zusammenfassung

Abstände und Ähnlichkeiten

Diskrete metrische Attribute

Definition

Eine Abstandsfunktion $d : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ heißt **Metrik** auf \mathcal{X} , wenn $d(\cdot, \cdot)$ für alle $x, y, z \in \mathcal{X}$ die drei Eigenschaften

1. $d(x, y) = 0 \iff x = y$ Definitheit
2. $d(x, y) = d(y, x)$ Symmetrie
3. $d(x, z) \leq d(x, y) + d(y, z)$ Dreiecksungleichung

besitzt.

Bemerkungen

1. Jede Vektorraumnorm $\|\cdot\| : \mathcal{X} \rightarrow \mathbb{R}$ definiert eine Metrik $d(x, y) = \|x - y\|$.
2. Jedes innere Produkt $\langle \cdot, \cdot \rangle : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ definiert eine VR-Norm $\|x\| = \sqrt{\langle x, x \rangle}$.
3. Distanzen transformieren in Ähnlichkeiten $s(x, y) = \exp(-d(x, y)/2\sigma^2)$.
4. Ähnlichkeiten transformieren in Distanzen $d(x, y) = -2\sigma^2 \cdot \log(s(x, y))$.

Spezialfall Zeichenketten

Attribute mit einem diskreten Wertebereich $\mathcal{X} \subset \mathcal{A}^*$

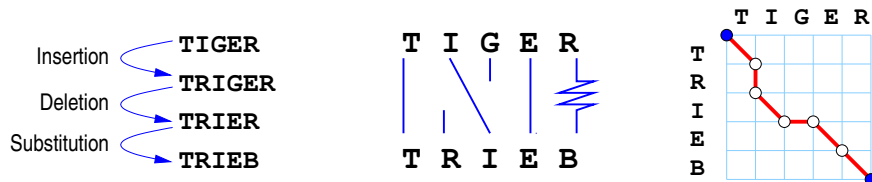
Elementare Operationen auf Zeichenketten

- Ersetzung eines Zeichens durch ein anderes
- Löschung eines Zeichens
- Einfügung eines Zeichens

substitution

deletion

insertion



Definition

Ist \mathcal{A} ein endliches Alphabet und sind v, w zwei Zeichenfolgen aus \mathcal{A}^* , so bezeichnet der **Levenshtein-Abstand** $d^{\text{lev}}(v, w)$ die minimale Anzahl von Elementaroperationen, mit denen v in w überführt werden kann.

Werteskalen

Relationen und Distanzen

Skalenkonversion

Detektion von Ausreißern

Imputation von Fehlanzeigen

Zusammenfassung

Spezialfall Zeichenketten

Zeichenkettenattribute sind metrisch und erlauben die Durchschnittsbildung

Lemma

Der Levenshtein-Abstand $d^{\text{lev}} : \mathcal{A}^* \times \mathcal{A}^* \rightarrow \mathbb{R}$ über dem Alphabet \mathcal{A} ist eine definite, symmetrische Distanzfunktion und erfüllt die Dreiecksungleichung — $(\mathcal{A}^*, d^{\text{lev}})$ ist folglich ein **metrischer Raum**.

Definition

Sei (\mathcal{X}, d) ein metrischer Raum und $(w_1, \dots, w_T) \in \mathcal{X}^T$ eine Auswahl (Multimenge) von Elementen. Der Wert

$$\mu^{\text{mid}} = \underset{z \in \{w_1, \dots, w_T\}}{\text{argmin}} \left(\sum_{t=1}^T d(z, w_t) \right)$$

heißt das **Medoid** der Menge bezüglich der Metrik $d(\cdot, \cdot)$.

Bemerkung

Das Medoid einer Wortmenge w_1, \dots, w_T mit maximaler Wortlänge N_{\max} läßt sich mit Aufwand $O(T^2 N_{\max}^2)$ berechnen.

Konversion der Attributskalen — wozu ?

Datensatz mit Attributen unterschiedlichen Skalentyps

Traditionelle Modellierungsverfahren erfordern einheitliche Skalen:

- **Numerische Skalen** $\mathcal{X} = \underbrace{\mathbb{R} \times \dots \times \mathbb{R}}_{N\text{-mal}} = \mathbb{R}^N$
Multivariate Normalverteilung

$$f(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{S}) = |\mathbf{S}|^{-1/2} \cdot \exp \left\{ -\frac{1}{2} \cdot (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

- **Diskrete Skalen** $\mathcal{X} = \mathcal{X}_1 \times \mathcal{X}_2 \times \dots \times \mathcal{X}_N$
N-dimensionale Wahrscheinlichkeitstabelle mit $\mathcal{X}_n = \{1, \dots, \ell_n\}$

$$P(\mathbf{x}) = p_{x_1, \dots, x_N} \quad \text{mit dem Tensor} \quad \mathbf{p} \in [0, 1]^{\ell_1 \times \ell_2 \times \dots \times \ell_N}$$

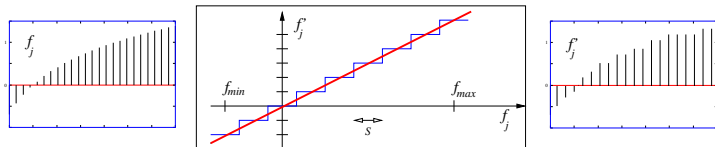
Option auf robusteres Datenmodell

Sind die Attributwerte wirklich normalverteilt ?

Kann ich mir eine Tabelle mit $\prod_n \ell_n$ Einträgen leisten ?

Diskretisierung numerischer Attribute (unüberwacht)

(kardinal → ordinal)



Äquidistante Intervalle

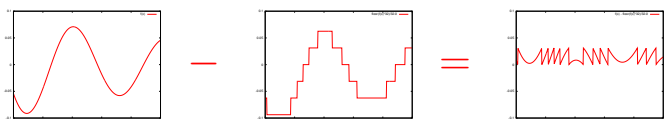
Mißachtet Datenverteilung ⇔ unglm. Zellenbesetzung & Übersteuern

Äquifrequente Intervalle

Histogrammegalisation ⇔ konstante Zellenbesetzung T/L

Nichtlinearer Skalarquantisierer

Minimiert den Störabstand (SNR): *mittlerer quadratischer Quantisierungsfehler*



Faustregel
 $L = \sqrt{T}$

Kardinalisierung nominaler Attribute

(nominal → numerisch)

Problem

Zahlreiche Methoden (*k*-nächste-Nachbarn, Bayesregel, Trennfunktionen) der Klassifikation und Vorhersage benötigen Objektstände oder numerische, besser noch gaußverteilte Objektattribute.

Nominale Entflechtung

Die nominale Skala mit Wertebereich

$\mathcal{X} = \{\xi_1, \dots, \xi_\ell\}$ wird auf einen Komplex

reellwertiger Attribute $\mathcal{X}_i = \{0, 1\}$, $i = 1, \dots, \ell$, abgebildet:

$$\phi(\xi_j) = (\underbrace{0, \dots, 0}_{(j-1)\text{-mal}}, \underbrace{1, 0, \dots, 0}_{(\ell-j)\text{-mal}}) \in \mathbb{R}^\ell$$

Für diese Darstellung gilt die **Äquidistanzeigenschaft**

$$d(\phi(\xi_i), \phi(\xi_j)) = \|\phi(\xi_i) - \phi(\xi_j)\| = \begin{cases} 0 & \xi_i = \xi_j \\ \sqrt{2} & \xi_i \neq \xi_j \end{cases}$$

Fall $\ell = 5$

\mathcal{X}	\mathbb{R}_1	\mathbb{R}_2	\mathbb{R}_3	\mathbb{R}_4	\mathbb{R}_5
ξ_1	1	0	0	0	0
ξ_2	0	1	0	0	0
ξ_3	0	0	1	0	0
ξ_4	0	0	0	1	0
ξ_5	0	0	0	0	1

Nominalisierung ordinaler Attribute

(ordinal → nominal)

Problem

Die Quantisierung numerischer Skalen liefert konstruktionsbedingt Werte einer **ordinalen** Skala.

Die immanente Reihenfolgeinformation wird aber von den einschlägigen Datenmodellen (W-Tabellen, lineare Modelle, Entscheidungsbäume) nicht genutzt.

Ordinale Entflechtung

Die ordinale Skala mit (sortiertem) Wertebereich

$\mathcal{X} = \{\xi_1, \dots, \xi_\ell\}$ wird auf einen Komplex **binärer** Attribute $\mathcal{X}_i = \{0, 1\}$, $i = 1, \dots, \ell - 1$, abgebildet:

$$\phi(\xi_j) = (\underbrace{0, \dots, 0}_{(j-1)\text{-mal}}, \underbrace{1, \dots, 1}_{(\ell-j)\text{-mal}}) \in \{0, 1\}^{\ell-1}$$

Für jede \mathcal{X} -Stufe $\xi - j$ gilt also:

$$\phi_i(\xi_j) = 1 \Leftrightarrow i \geq j.$$

Fall $\ell = 3$

\mathcal{X}	\mathcal{X}_1	\mathcal{X}_2
ξ_1	1	1
ξ_2	0	1
ξ_3	0	0

Fall $\ell = 5$

\mathcal{X}	\mathcal{X}_1	\mathcal{X}_2	\mathcal{X}_3	\mathcal{X}_4
ξ_1	1	1	1	1
ξ_2	0	1	1	1
ξ_3	0	0	1	1
ξ_4	0	0	0	1
ξ_5	0	0	0	0

Kontrastmatrizen

Auch im $\mathbb{R}^{\ell-1}$ ist genug Platz für ξ_1, \dots, ξ_ℓ

Ursprung & Einheiten

Einer-gegen-alle: treatment

ξ_1	0	0	0	0
ξ_2	1	0	0	0
ξ_3	0	1	0	0
ξ_4	0	0	1	0
ξ_5	0	0	0	1

Distanzen 0, $\sqrt{2}$, aber auch 1

Gestaffelt

Gegen-Anfangspartie: helmert

ξ_1	-1	-1	-1	-1
ξ_2	1	-1	-1	-1
ξ_3	0	2	-1	-1
ξ_4	0	0	3	-1
ξ_5	0	0	0	4

Distanzen 0 und viele andere ...

Spaltenmittelwertfrei

Einer-gegen-alle: sum

ξ_1	-1	-1	-1	-1
ξ_2	1	0	0	0
ξ_3	0	1	0	0
ξ_4	0	0	1	0
ξ_5	0	0	0	1

Distanzen 0, $\sqrt{2}$, aber auch $\sqrt{\ell+2}$

Äquidistant

Orthonormalpolynome: poly

ξ_1	$p_1(r_1)$	$p_1(r_2)$	$p_1(r_3)$	$p_1(r_4)$
ξ_2	$p_2(r_1)$	$p_2(r_2)$	$p_2(r_3)$	$p_2(r_4)$
ξ_3	$p_3(r_1)$	$p_3(r_2)$	$p_3(r_3)$	$p_3(r_4)$
ξ_4	$p_4(r_1)$	$p_4(r_2)$	$p_4(r_3)$	$p_4(r_4)$
ξ_5	$p_5(r_1)$	$p_5(r_2)$	$p_5(r_3)$	$p_5(r_4)$

$$\|u - v\|^2 = \|u\|^2 - 2\langle u, v \rangle + \|v\|^2 = 2$$

Redundante Kardinalisierung

Fehlererkennende und fehlerkorrigierende Codes

Indexcodierung

Holzhammermethode: $\phi(\xi_j) = j$ $\phi : \{\xi_1, \dots, \xi_\ell\} \rightarrow \mathbb{R}^1$

Dualcodierung

(Hamming)abstände „fehleranfällig“ $\phi : \{\xi_1, \dots, \xi_\ell\} \rightarrow \mathbb{R}^{\lceil \log_2 \ell \rceil}$

Vollständige Korrekturcodes

Erkennt und kompensiert Fehler in einer Komponente
(interessant ab $\ell = 4$) $\phi : \{\xi_1, \dots, \xi_\ell\} \rightarrow \mathbb{R}^L, L = 2^{\ell-1} - 1$

ξ_1 1 1 1 1 1 1 1
 ξ_2 0 0 0 0 1 1 1
 ξ_3 0 0 1 1 0 0 1
 ξ_4 0 1 0 1 0 1 0

ξ_1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
 ξ_2 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1
 ξ_3 0 0 0 0 1 1 1 1 0 0 0 0 1 1 1
 ξ_4 0 0 1 1 0 0 1 1 0 0 1 1 0 0 1
 ξ_5 0 1 0 1 0 1 0 1 0 1 0 1 0 1 0

Beinhaltet alle $\{0, 1\}^\ell$ -Spalten außer Komplementen und den uninformativen Attributen **0**, **1**.

Floyd-Warshall-Algorithmus

Schnelle Berechnung geodätischer Distanzen mittels dynamischer Programmierung

1 INITIALISIERUNG

Setze $D_{ij} = \begin{cases} 0 & i = j \\ 1 & \xi_i, \xi_j \text{ adjazent} \\ \infty & \text{sonst} \end{cases}$

2 REKURSION

Für alle $k, i, j \in \{1, \dots, L\}$:

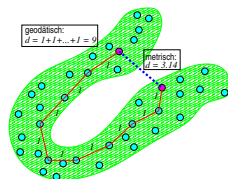
$$D_{ij} \leftarrow \min \{D_{ij}, D_{ik} + D_{kj}\}$$

3 TERMINIERUNG

Die Matrix **D** enthält alle minimalen Wegelängen zwischen Elementen ξ_i, ξ_j .

Wirkungsweise

Der FWA erzwingt in $O(L^3)$ Schritten die Gültigkeit der Dreiecksungleichung.



Bemerkung

Der Algorithmus ist auch anwendbar für gewichtete und nichtsymmetrische Adjazenzen.

Konversion von Distanzfunktionen

Nachbarschaft — Metrik — normierter Vektorraum

Metrik \Rightarrow symmetrische Nachbarschaft

Global operierende Schwellwertoperation ($0 < \delta_{\max} \in \mathbb{R}$)

$$\xi_i \propto \xi_j \Leftrightarrow d(\xi_i, \xi_j) \leq \delta_{\max}$$

Metrik \Rightarrow nichtsymmetrische Nachbarschaft

Lokale Umgebungsdefinition (k nächste Nachbarn, $k \in \mathbb{N}$)

$$\xi_i \propto \xi_j \Leftrightarrow \xi_j \in \mathcal{U}_{\mathcal{X}}^{(k)}(\xi_i)$$

Adjazenz \Rightarrow Metrik

Geodätische Abstände (minimale Pfadlängen im Adjazenzgraphen)

Metrik \Rightarrow (euklidischer) Vektorraum

Nicht jede metrische Distanz $D \in \mathbb{R}^{L \times L}$ ist im \mathbb{R}^{L-1} repräsentierbar.

Kardinalisierung von Präferenzrelationen

Schwache Ordnungsrelation (\mathcal{X}, \prec) \Rightarrow ein, zwei, mehrere relative Attribute

Intervallordnung

Repräsentation durch $\mathcal{X}_1 \times \mathcal{X}_2 = \mathbb{R}^2$ mit

$$a \prec b \Leftrightarrow a_2 < b_1$$

Inklusionsfreie Intervallordnung

Repräsentation durch $\mathcal{X}_1 = \mathbb{R}^1$ mit $\delta \in \mathbb{R}_+$ und

$$a \prec b \Leftrightarrow a_1 + \delta < b_1$$

Endliche Halbordnung

Repräsentation durch $\mathcal{X}_1 \times \dots \times \mathcal{X}_L = \mathbb{R}^L$ mit

$$a \prec b \Leftrightarrow \forall \ell = 1, \dots, L: a_\ell < b_\ell$$

Standardisierung numerischer Skalen

Vereinheitlichung von Wertebereichen u/o Dynamikeigenschaften

Min-Max-Normierung

$$f : \begin{cases} \mathbb{R} & \rightarrow [0, 1] \\ x & \mapsto \frac{x - x_{\min}}{x_{\max} - x_{\min}} \end{cases}, \quad f^{-1}(x) = (x_{\max} - x_{\min}) \cdot x + x_{\min}$$

Statistische Normierung

$$f : \begin{cases} \mathbb{R} & \rightarrow [\mu - C\sigma, \mu + C\sigma] \\ x & \mapsto \frac{x - \mu}{\sigma} \end{cases}, \quad f^{-1}(x) = \sigma \cdot x + \mu$$

Reziproke Transformation

$$f : \begin{cases} \mathbb{R} \setminus \{0\} & \rightarrow \mathbb{R} \setminus \{0\} \\ x & \mapsto 1/x \end{cases}, \quad f^{-1}(x) = 1/x$$

Werteskalen

Relationen und Distanzen

Skalenkonversion

Detektion von Ausreißern

Imputation von Fehlanzeigen

Zusammenfassung

Standardisierung numerischer Skalen

Vereinheitlichung von Wertebereichen u/o Dynamikeigenschaften

Wurzel-Transformation

$$f : \begin{cases} (C, \infty) & \rightarrow \mathbb{R}^+ \\ x & \mapsto \sqrt[B]{x - C} \end{cases}, \quad f^{-1}(x) = x^B + C$$

Logarithmus-Transformation

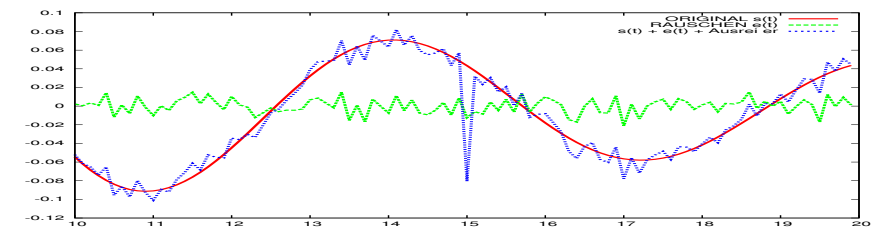
$$f : \begin{cases} (C, \infty) & \rightarrow \mathbb{R} \\ x & \mapsto \log_B(x - C) \end{cases}, \quad f^{-1}(x) = B^x + C$$

Fisher-Transformation

$$f : \begin{cases} (-1, +1) & \rightarrow \mathbb{R} \\ x & \mapsto \frac{1}{2} \log_e \frac{1+x}{1-x} \end{cases}, \quad f^{-1}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Meßfehler & Erhebungsfehler

Die „Rohdaten“ sind oft fehlerbehaftet, verrauscht, verzerrt



Zufällige Fehler

- Messungenauigkeit
- Übertragungsstrecke
- Modell *additives Rauschen*:
 $y_n = x_n + e_n, e_n \sim \mathcal{N}(0, \sigma^2)$

• **Ausreißer**

Systematische Fehler

- Kalibrierung
- Skalierung
- Trend, Drift, Saisoneffekt

• **Ausreißer**

Ausreißerdetektion

Was ist ein Ausreißer und wie erkenne ich ihn ?

Vertikale Detektion

Ein Wert x_{tj} fällt aus dem Rahmen seines **Attributs** \mathcal{X}_j .

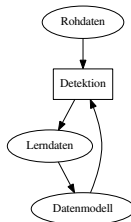
Kategoriale Attribute bieten *keine Handhabe* !

$$\mathcal{X}_j = \{m, f\}$$

Horizontale Detektion

Ein Wert x_{tj} fällt aus dem Rahmen seines **Objekts** \mathbf{o}_t .

Werden Objekte durch Ausreißer erst *interessant* ? $\mathbf{o}_t = („kath.“, „verh.“)$



Teufelskreis

vertikale **Detektion** \Leftarrow **Attributmodell**
 horizontale **Detektion** \Leftarrow **Objektmodell**

Hypothesentests für Ausreißer

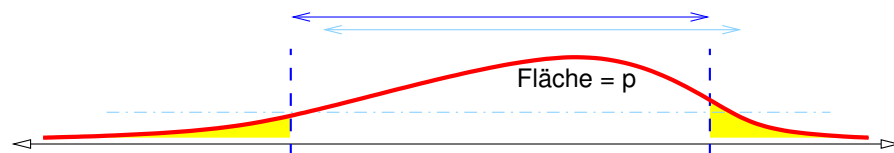
... bei bekannter unimodaler Verteilungsdichtefunktion

Definition (Quantilmethode)

Ein Wert $x_q \in \mathbb{R}$ heißt **q-Quantil** der Dichtefunktion $f_{\mathbb{X}}(\cdot)$ genau dann, wenn gilt:

$$F_{\mathbb{X}}(x_q) = P(\mathbb{X} \leq x_q) = q$$

Ein Wert $x \in \mathbb{R}$ heißt **Ausreißer** der Verteilung zum **Niveau** $p \in [0, 1]$, wenn er außerhalb des Akzeptanzintervalls $[x_{1/2-p/2}, x_{1/2+p/2}]$ liegt.



Bemerkungen

1. Für **symmetrische** Dichtefunktionen gilt für jedes $q \in [0, 1]$ die Identität $f_{\mathbb{X}}(x_q) = f_{\mathbb{X}}(x_{1-q})$. $[\mu - C\sigma, \mu + C\sigma]$
2. Für **multimodale** Dichtefunktionen ergibt das definierte Akzeptanzintervall keinen Sinn.

Hypothesentests für Ausreißer

Ausreißer $\hat{=}$ extrem unwahrscheinliche Attributwerte

Satz (Tschebyscheff)

Ist \mathbb{X} eine kontinuierliche Zufallsvariable mit dem Erwartungswert μ und der Varianz σ^2 , so gilt für jede Konstante $C > 0$ die Ungleichung:

$$P\left(\left|\frac{\mathbb{X} - \mu}{\sigma}\right| \geq C\right) \leq \frac{1}{C^2}$$

Beispiel

Zweiseitige Streuungswahrscheinlichkeiten m/o NV-Annahme:

	σ	2σ	3σ	4σ	5σ
Tschebyscheff	≤ 1	≤ 0.25	≤ 0.11	≤ 0.063	≤ 0.040
$\mathcal{N}(\mu, \sigma)$	$= 0.323$	$= 0.065$	$= 0.003$	$= 0.001$	≤ 0.0001

\Leftarrow ein „zahnloser“ Test ohne Kenntnis der Dichtefunktion !

Hypothesentests für Ausreißer

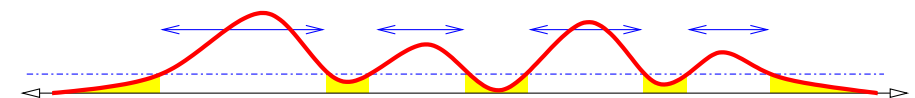
... bei bekannter multimodaler Verteilungsdichtefunktion

Definition (Bayesträgermethode)

Die Wertemenge $\mathcal{B}_c = \{x \mid f_{\mathbb{X}}(x) \geq c\}$ heißt **Bayesträger** der Verteilung $f_{\mathbb{X}}(\cdot)$ zum **Niveau** $p \in [0, 1]$, wenn gilt:

$$\int_{\mathcal{B}_c} f_{\mathbb{X}}(\xi) d\xi = p$$

Jeder Wert $x \in \mathbb{R}$ mit $f_{\mathbb{X}}(x) < c$ heißt **Ausreißer** der Verteilung zum **Niveau** p .



Bemerkungen

1. Für **symmetrisch-unimodale** Dichtefunktionen stimmen Bayesträger und Akzeptanzintervall überein.
2. Nicht verwechseln mit **Bayesintervall**, dem kürzesten Intervall mit Fläche p .

Faustregeln zur Ausreißerdetektion

Treffer als Fehlanzeige (NA=„not available“) markieren

Unimodal

- **Normalverteilung**

$$|x - \mu| > C \cdot \sigma$$

- **Gleichverteilung**

$$|x - \mu| > p\text{-Niveau}$$

- **Empirischer Trimm**

$$x \notin [x_{1/2-p/2}, x_{1/2+p/2}]$$

Multimodal

- **Tschebyscheff**

$$\frac{|x - \mu|}{\sigma} > \sqrt{\frac{1}{1-p}}$$

- **Gauß-Mischung**

$$(\forall \ell) |x - \mu_\ell| > C \cdot \sigma_\ell$$

- **Lonesome Cowperson**

$$|x - k\text{-NN}(x)| > d_{\max}$$

Werteskalen

Relationen und Distanzen

Skalenkonversion

Detektion von Ausreißern

Imputation von Fehlanzeigen

Zusammenfassung

Teufelskreis Parameterschätzung

Ausreißer verändern die genutzten Verteilungsparameter

Modellrechnung für die $C\sigma$ -Regel

- **Datensatz**

Eine $\mathcal{N}(\mu, \sigma)$ -verteilte Probe der Größe T
 zuzüglich M^+ Ausreißer der Gestalt $a^+ = \mu + c\sigma$
 zuzüglich M^- Ausreißer der Gestalt $a^- = \mu - c\sigma$
 ($T' = T + M^+ + M^-$)

- **Geschätzter Erwartungswert**

(Im Fall $M^+ = M^-$ gilt einfach $\hat{\mu} = \mu$.)

$$\hat{\mu} = \mu + \frac{M^+ - M^-}{T'} \cdot c\sigma$$

- **Geschätzte Varianz**

Gilt mit $M := M^+ / 2 = M^- / 2$ wegen
 $\frac{1}{T'} \sum_{x \in \omega'} x^2 = \mu^2 + \sigma^2 + \frac{M}{T+M} (c^2 - 1) \sigma^2$
 und der Abkürzung $r := M / (T+M)$.

$$\hat{\sigma} = \sigma \sqrt{1 + c^2 r - r}$$

Für eine nicht verschwindende Anzahl ($r \gg 0$) von markanten Ausreißern
 ($c \gg 1$) dominiert $c^2 r$ den Wurzelausdruck und die $C\sigma$ -Regel ist wegen
 $\hat{\sigma} \propto c$ entschärft!

Fehlanzeigen (a.k.a. „not available“)

Nicht zugängliche Attributwerte in der Datenmatrix

Fehlanzeige als Unfall

Sensorkomponente hat versagt
 Erhebungsprotokoll unvollständig
 Markierte Ausreißer

Fehlanzeige als Regelfall

Verzicht aus Kostengründen
 Nichthomogenes Warehousing
 Dünnbesetzung anwendungsbedingt
 z.B. Bewertungssysteme für Musik, Bücher,
 Restaurants, Webseiten, Bordellbetriebe, ...

Fehlanzeigenbehandlung

- **Objekt löschen**
 Können wir uns das leisten?
- **Eintrag markieren**
 und auf spezielle Weise weiterverarbeiten.
- **Imputieren**
 Leerstelle mit geeignetem Wert auffüllen.

Imputationstechniken

Welcherart Zusatzinformation wird zur Wertergänzung genutzt ?

$M_1 \ M_2 \ \dots \ M_n \ \dots \ M_N$
 $X_1 \ X_2 \ \dots \ X_n \ \dots \ X_N$

„Missing (Completely) At Random“

Kontextfrei (MCAR)

Attributstatistik (& Ausreißer)

- Ersetzen durch Datenmittel $\hat{\mu}$
- durch x_{\min} bzw. x_{\max}
- durch nächsten Nachbarn

$$x_n^* \stackrel{\text{def}}{=} \operatorname{argmin}_{\xi \in \bar{\omega}} d(x_n, \xi)$$

Regression (MAR)

Probabilistisches Datenmodell

$$x_n^* = \mathcal{E}[X_n \mid \dots, x_{n-1}, x_{n+1}, \dots]$$

Interpolation (MAR)

Zeile/Spalte $\hat{=}$ Meßreihe

- Linearer Ausgleich, z.B.

$$x_n^* \stackrel{\text{def}}{=} (x_{n-1} + x_{n+1}) / 2$$

- Polynome, Splines (nl.)
- Glättungsfilter

Matrixapproximation (MAR)

Lückenhafte (num.) Datenmatrix

$$X \stackrel{\text{NA}}{\approx} V^T D U$$

Regression für nominale Datensätze

Beispielszenarium: drei Attribute X_1, X_2, X_3 mit 2, 3 bzw. 4 Wertestufen

(Algorithmus)

1 ABSOLUTE HÄUFIGKEITEN

Erstelle Tabelle $f \in \mathbb{N}^{2 \cdot 3 \cdot 4}$ mit den 24 Auftretenszahlen f_{ijk} der Ereignisse $(x_1, x_2, x_3) = (\xi_i, \eta_j, \zeta_k)$.

2 EREIGNISWAHRSCHEINLICHKEITEN

Erstelle Tabelle der 24 ML-Schätzwerte $\hat{p}_{ijk} = f_{ijk} / T$.

3 BEDINGTE ATTRIBUTWAHRSCHEINLICHKEITEN

$$q_{i|jk}^{(1|23)} = \frac{\hat{p}_{ijk}}{\sum_{\ell} p_{\ell jk}}, \quad q_{j|ik}^{(2|13)} = \frac{\hat{p}_{ijk}}{\sum_{\ell} p_{i \ell k}}, \quad q_{k|ij}^{(3|12)} = \frac{\hat{p}_{ijk}}{\sum_{\ell} p_{ij \ell}}$$

4 IMPUTATION DES BEDINGTEN MODUS

$$\mu_{jk}^{(23)} = \operatorname{argmax}_i q_{i|jk}^{(1|23)}, \quad \mu_{ik}^{(13)} = \operatorname{argmax}_j q_{j|ik}^{(2|13)}, \quad \mu_{ij}^{(12)} = \operatorname{argmax}_k q_{k|ij}^{(3|12)}$$

(sumoflogIA)

Glättungsfilter für Meßreihenfehler

Imputation $\hat{=}$ kontextfrei Ersetzen + Filtern

Gleitender Mittelwert

der Ordnung $q = 2p + 1$, $p \in \mathbb{N}$:

$$\hat{x}_n = \frac{1}{q} \cdot \sum_{\ell=n-p}^{n+p} x_{\ell}$$

↘ Ausreißer, ↗ Phasentreue

Exponentialfilter

mit Abklingparameter $\alpha \in [0, 1]$:

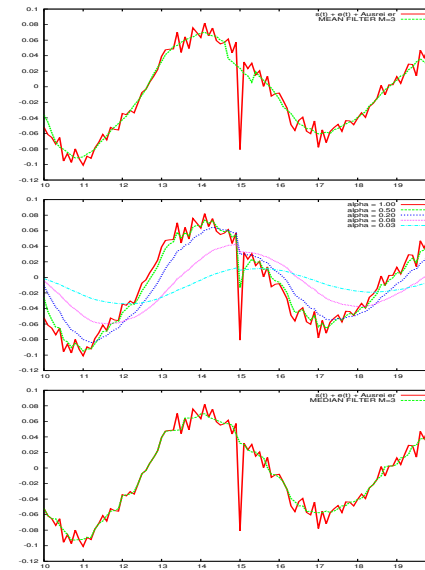
$$\hat{x}_n = \hat{x}_{n-1} + \alpha \cdot (x_n - \hat{x}_{n-1})$$

↗ Ausreißer,
 ↘ Phasentreue/Nivellierung

Gleitendes Medianfilter

der Ordnung $q = 2p + 1$, $p \in \mathbb{N}$:

$$\hat{x}_n = \mu^{\text{med}}(x_{n-p}, \dots, x_n, \dots, x_{n+p})$$



Werteskalen

Relationen und Distanzen

Skalenkonversion

Detektion von Ausreißern

Imputation von Fehlanzeigen

Zusammenfassung

Zusammenfassung (2)

1. Ein **Datensatz** besteht aus **Objekten**, die explizit durch eine Reihe von **Attributwerten** oder implizit durch Beziehungen wie **Abstand**, **Adjazenz** oder **Präferenz** charakterisiert sind.
2. Attribute besitzen eine **diskrete Skala** (**nominal** oder **ordinal**) oder eine **numerische Skala** (**relativ** oder **proportional**).
3. Die Skalen unterscheiden sich hinsichtlich ihres **Wertebereichs**, ihrer **Verknüpfungsoperationen** und ihrer **Durchschnittswertbildung**.
4. Auf **Zeichenketten** ist mit dem Levenshteinabstand eine **Metrik** und mit dem **Medoid** ein Durchschnitt definiert.
5. Skalen lassen sich nötigenfalls mittels **Quantisierung** (numerisch \rightsquigarrow ordinal), **Entflechtung** (ordinal \rightsquigarrow nominal) bzw. **Kontrastmatrizen** (nominal \rightsquigarrow numerisch) konvertieren.
6. Aus **Adjazenzen** leiten sich **geodätische Distanzen** her, aus **Präferenzen** ein oder zwei **Ordinalskalen**.
7. **Ausreißer** werden durch einen der attributbezogenen **Hypothesentests** detektiert.
8. Als **Ersatzwerte** für Ausreißer und andere **Fehlanzeigen** dienen Mittel- und Extremwerte; wenn möglich, imputieren wir durch **Interpolation** oder **Regression**.