

Abhängigkeiten und Normalformen

- 6.1 Funktionale Abhängigkeiten
- 6.2 The Chase — Die Verfolgungsjagd
- 6.3 Normalformen — Informationstheoretischer Ansatz

Motivation (1/2)

Ziel beim Datenbank-Entwurf:

Ein DB-Schema entwickeln, so dass Informationen zum gewünschten Anwendungsbereich "sinnvoll" gespeichert werden können.
Insbesondere: Wenn möglich, Redundanzen und Inkonsistenzen vermeiden.

Beispiel: Relation *Warenlager*[Bauteil-Nr, Lager-Nr, Menge, Ort]

Warenlager:

Bauteil-Nr	Lager-Nr	Menge	Ort
2411	2	200	Riedberg
2412	3	300	Bornheim
3001	1	100	Hanau
2415	2	100	Riedberg

Unschön: Redundanz — der Ort von Lager 2 ist mehrfach gespeichert.

Dadurch können Inkonsistenzen auftreten:

Update-Anomalien = Inkonsistenzen, die durch Aktualisierung der DB auftreten können:

- ▶ **Änderungs-Anomalie:** den Ort in Zeile 1 durch "Westend" ersetzen
↪ Adresse von Lager 2 nicht mehr eindeutig
- ▶ **Lösch-Anomalie:** Löschen von Zeile 2
↪ Information über die Adresse von Lager 3 geht verloren
- ▶ **Einfüge-Anomalie:** Die Adresse eines neuen Lagers kann erst dann eingefügt werden, wenn mindestens ein Bauteil dort gelagert wird.

Zur Erinnerung — Benannte Perspektive

In Kapitel 1 hatten wir festgelegt:

- ▶ Eine abzählbar unendliche Menge **att** von **Attribut-Namen**.
Diese Menge ist **geordnet** via \leq_{att} .
- ▶ Eine abzählbar unendliche Menge **relname** von **Relations-Namen**.
Die Mengen **att**, **dom**, **relname** seien disjunkt.
- ▶ Eine Funktion **sort**: **relname** $\rightarrow \mathcal{P}^{\text{fin}}(\text{att})$, die jedem Relations-Namen eine endliche Menge von Attributen zuordnet ... und zwar so, dass f.a. $U \subseteq^{\text{fin}} \text{att}$ gilt: es gibt unendlich viele $R \in \text{relname}$ mit $\text{sort}(R) = U$.
- ▶ Ein **Relationsschema** ist einfach ein Relations-Name R .
- ▶ Manchmal schreiben kurz $R[U]$ für $\text{sort}(R) = U$.
- ▶ Ein **R-Tupel** ist eine Funktion $t: \text{sort}(R) \rightarrow \text{dom}$.
- ▶ Eine **R-Relation** ist eine endliche Menge von R-Tupeln.
- ▶ $\text{inst}(R)$ bezeichnet die Menge aller Relationen über R .
- ▶ $\text{inst}(U)$ bezeichnet die Menge aller Relationen über einem Relationsschema der Sorte U (für $U \subseteq^{\text{fin}} \text{att}$)

Motivation (2/2)

Zur Vermeidung dieser Update-Anomalien:

Informationen auf 2 Relationen *Adressen*[Lager-Nr, Ort] und *Lagerung*[Bauteil-Nr, Lager-Nr, Menge] aufteilen

Adressen:

Lager-Nr	Ort
2	Riedberg
3	Bornheim
1	Hanau

Abhängigkeit:
Lager-Nr \rightarrow Ort

Lagerung:

Bauteil-Nr	Lager-Nr	Menge
2411	2	200
2412	3	300
3001	1	100
2415	2	100

Abhängigkeit:
Bauteil-Nr, Lager-Nr \rightarrow Menge

Zur Anfrage-Optimierung:

Die "optimierte Anfrage" muss nur auf solchen Datenbanken äquivalent zur Original-Anfrage sein, die die obigen Abhängigkeiten erfüllen.

↪ die minimale, solchermaßen äquivalente Anfrage ist evtl. noch kleiner als die, die durch die Tableau-Minimierung aus Kapitel 5 gefunden wird.

Abhängigkeiten und Normalformen

6.1 Funktionale Abhängigkeiten

6.2 The Chase — Die Verfolgungsjagd

6.3 Normalformen — Informationstheoretischer Ansatz

Funktionale Abhängigkeiten

Definition 6.1

Sei U eine endliche Menge von Attribut-Namen.

- (a) Eine **funktionale Abhängigkeit** (kurz: **FD**) über U ist ein Ausdruck der Form $X \rightarrow Y$, wobei $X, Y \subseteq U$. (“FD” steht für “functional dependency”)
- (b) Eine Relation $I \in \text{inst}(U)$ **erfüllt** die FD $X \rightarrow Y$ (kurz: $I \models X \rightarrow Y$), falls für alle Tupel t und s aus I gilt:

$$\pi_X(t) = \pi_X(s) \implies \pi_Y(t) = \pi_Y(s)$$

(D.h.: Wenn t und s in sämtlichen Spalten aus X übereinstimmen, dann stimmen sie auch in jeder Spalte aus Y überein.)

- (c) Ist Σ eine Menge von FDs über U , so gilt

$$I \models \Sigma \iff \text{für alle } f \in \Sigma \text{ gilt } I \models f$$

- (d) Eine **Schlüsselbedingung** ist eine FD der Form $X \rightarrow U$.

Notation

Attribut-Namen : A, B, C, A_1, A_2, \dots

Attribut-Mengen : $X, Y, Z, X_1, X_2, \dots, U$

Relationen : I, J

$\{A, B, C\}$: ABC

$X \cup Y$: XY

Funktionale Abhängigkeiten und verlustfreie Joins

Proposition 6.2

Sei $X \rightarrow Y$ eine FD über U und sei $Z := U \setminus (X \cup Y)$.

Für jede Relation $I \in \text{inst}(U)$ gilt:

Falls $I \models X \rightarrow Y$, so ist $I = \pi_{XY}(I) \bowtie \pi_{XZ}(I)$.

Beweis: Übung.

Folgerung: Die in Relation I gespeicherte Information kann “verlustfrei” auf zwei Relationen aufgeteilt werden (eine mit den Spalten XY und eine mit den Spalten XZ), aus denen die Original-Relation rekonstruiert werden kann. (Stichwort: “lossless join”)

Beispiel für einen “verlustreichen Join”: Siehe Tafel.

Abhängigkeiten und Normalformen

6.1 Funktionale Abhängigkeiten

6.2 The Chase — Die Verfolgungsjagd

6.3 Normalformen — Informationstheoretischer Ansatz

Äquivalenz bzw. Query Containment bzgl. Σ

Definition 6.4

Sei R ein Relationenschema, sei $\mathcal{F} \subseteq \text{inst}(R)$, und seien Q_1 und Q_2 zwei Anfragen an Relationen aus $\text{inst}(R)$.

- ▶ $Q_1 \subseteq_{\mathcal{F}} Q_2$: $\iff \llbracket Q_1 \rrbracket(I) \subseteq \llbracket Q_2 \rrbracket(I)$, für alle $I \in \mathcal{F}$.
- ▶ $Q_1 \equiv_{\mathcal{F}} Q_2$: $\iff \llbracket Q_1 \rrbracket(I) = \llbracket Q_2 \rrbracket(I)$, für alle $I \in \mathcal{F}$.

Sei Σ eine Menge von FDs über R .

- ▶ $\text{sat}(R, \Sigma) := \text{sat}(\Sigma) := \{I \in \text{inst}(R) : I \models \Sigma\}$
- ▶ $Q_1 \subseteq_{\Sigma} Q_2$: $\iff Q_1 \subseteq_{\text{sat}(\Sigma)} Q_2$
- ▶ $Q_1 \equiv_{\Sigma} Q_2$: $\iff Q_1 \equiv_{\text{sat}(\Sigma)} Q_2$

Beispiel zu “The Chase — Die Verfolgungsjagd”

Beispiel 6.3

Tableau-Anfrage $Q = (T, t)$ mit

$$T = \begin{array}{|c|c|c|c|} \hline A & B & C & D \\ \hline w & x & y & z' \\ \hline w' & x & y' & z \\ \hline \end{array} \quad \text{und} \quad t = (w, x, y, z)$$

Klar: Die Anfrage $Q := (T, t)$ ist **minimal** im Sinne von Kapitel 5.1.

Situation jetzt:

- ▶ Gegeben sei die FD-Menge $\Sigma := \{B \rightarrow D\}$
- ▶ Q soll nur auf solchen DBs ausgewertet werden, die Σ erfüllen
- ▶ Ziel: Vereinfache (minimiere) Q .

Details: siehe Tafel.

Vereinbarungen für den Rest von Kapitel 6.2

Der Einfachheit halber betrachten wir im Folgenden

- ▶ ein festes Relationenschema R
- ▶ Mengen Σ von FDs über R , in denen oBdA jede FD von der Form $X \rightarrow A$ mit $X \subseteq \text{sort}(R)$ und $A \in \text{sort}(R)$ ist
- ▶ eine feste lineare Ordnung $<$ auf der Variablenmenge **var**
- ▶ nur Tableau-Anfragen $Q = (T, t)$ über R , in denen **keine Konstanten vorkommen**

Bemerkung: Die Ergebnisse aus Kapitel 6.2 können leicht verallgemeinert werden auf Anfragen mit Konstanten und auf Anfragen über einem Datenbankschema.

Regel für die Verfolgungsjagd

Definition 6.5

FD-Regel:

Sei $f := (X \rightarrow A)$ eine FD über R , sei (T, t) eine Tableau-Anfrage über R .
Seien u und v Zeilen von T mit $\pi_X(u) = \pi_X(v)$ und $u(A) \neq v(A)$.
Sei $\{x, y\} := \{u(A), v(A)\} \subseteq \mathbf{var}$ und sei $x < y$.

Anwenden der FD f auf u, v in (T, t) liefert die Tableau-Anfrage $(\zeta(T), \zeta(t))$, wobei ζ die Substitution mit $\zeta(y) := x$ und $\zeta(z) := z$ für alle $z \in \mathbf{Var}((T, t)) \setminus \{y\}$.

Verfolgungssequenzen

Definition 6.7

(a) Eine **Verfolgungssequenz** für (T, t) mittels Σ ist eine Folge

$$(T_0, t_0), (T_1, t_1), (T_2, t_2), \dots$$

für die gilt:

- ▶ $(T_0, t_0) = (T, t)$ und
- ▶ für jedes $i \geq 0$ entsteht (T_{i+1}, t_{i+1}) durch 1-maliges Anwenden der FD-Regel mit einer FD aus Σ auf (T_i, t_i) .

(b) Die Verfolgungssequenz ist **terminiert**, falls sie endlich ist und auf ihr letztes Element, (T_m, t_m) , keine FD-Regel mit einer FD aus Σ mehr angewendet werden kann.
 (T_m, t_m) heißt dann das **Resultat** der Sequenz.

Anwenden der FD-Regel erhält Äquivalenz bzgl. Σ

Proposition 6.6

Sei Σ eine Menge von FDs über R ,
sei $f := (X \rightarrow A) \in \Sigma$,
sei $Q := (T, t)$ eine Tableau-Anfrage über R und
sei $Q' := (T', t')$ eine Tableau-Anfrage, die durch 1-maliges Anwenden der FD-Regel mit einer FD $f \in \Sigma$ aus Q entsteht.

Dann gilt: $Q \equiv_{\Sigma} Q'$.

Beweis: Siehe Tafel.

Notation: $T \models \Sigma$

Definition 6.8

(a) Ein Tableau T über R **erfüllt** die FD $X \rightarrow A$ (kurz: $T \models X \rightarrow A$), falls für alle Zeilen u und v von T gilt:

$$\pi_X(u) = \pi_X(v) \implies u(A) = v(A)$$

(D.h.: Wenn u und v in sämtlichen Spalten aus X übereinstimmen, dann stimmen sie auch in der Spalte A überein.)

(b) Ist Σ eine Menge von FDs über R , so gilt

$$T \models \Sigma \quad : \iff \quad T \models f, \text{ für alle } f \in \Sigma$$

Eigenschaften des Resultats einer Verfolgungssequenz

Lemma 6.9

Sei (T', t') das Resultat einer terminierten Verfolgungssequenz für (T, t) mittels Σ .
Dann gilt: $(T', t') \equiv_{\Sigma} (T, t)$ und $T' \models \Sigma$.

Beweis: Übung.

Klar: Jede Verfolgungssequenz ist endlich und kann zu einer terminierten Sequenz vervollständigt werden.

Bemerkenswert: Man kann zeigen, dass alle terminierten Verfolgungssequenzen für (T, t) mittels Σ dasselbe Resultat liefern. Dies wird auch

Church-Rosser-Eigenschaft genannt.

Berechnung von $\text{chase}(T, t, \Sigma)$

Korollar 6.12

Es gibt einen **Polynomialzeit-Algorithmus**, der bei Eingabe einer Tableau-Anfrage (T, t) über R und einer Menge Σ von FDs über R die Tableau-Anfrage $\text{chase}(T, t, \Sigma)$ berechnet.

Beweis:

Algorithmus:

(1) Wiederhole so lange, bis keine FD-Regel bzgl. Σ mehr auf (T, t) anwendbar ist:

(1.1) (T', t') sei das Resultat der Anwendung einer FD-Regel bzgl. Σ auf (T, t) .

(1.2) Setze $(T, t) := (T', t')$.

(2) Gib (T, t) aus.

Korrektheit folgt direkt aus der Church-Rosser-Eigenschaft.

Polynomielle Laufzeit, da jede FD $f \in \Sigma$ auf jedes Paar u, v von Zeilen von T höchstens 1-mal angewendet werden kann und da jede einzelne Anwendung der FD-Regel nur polynomiell viel Zeit benötigt.

Church-Rosser-Eigenschaft der Verfolgungsjagd

Theorem 6.10

Sei (T, t) eine Tableau-Anfrage über R und sei Σ eine Menge von FDs über R .
Dann gilt: Alle terminierten Verfolgungssequenzen für (T, t) mittels Σ liefern dasselbe Resultat.

Hier ohne Beweis.

Ein Beweis findet sich am Ende von Kapitel 8.4 des Buchs [AHV].

Definition 6.11

Ist (T, t) eine Tableau-Anfrage über R und Σ eine Menge von FDs über R , so bezeichnet $\text{chase}(T, t, \Sigma)$ das Resultat einer (bzw. sämtlicher) terminierter Verfolgungssequenzen für (T, t) mittels Σ .

Bemerkung: Von Lemma 6.9 wissen wir, dass $\text{chase}(T, t, \Sigma) \equiv_{\Sigma} (T, t)$ und $\text{chase}(T, t, \Sigma) \models \Sigma$.

Äquivalenz von Anfragen bzgl. Σ

Theorem 6.13

Seien $Q_1 := (T_1, t_1)$ und $Q_2 := (T_2, t_2)$ Tableau-Anfragen über R und sei Σ eine Menge von FDs über R . Dann gilt:

(a) $Q_1 \subseteq_{\Sigma} Q_2 \iff \text{chase}(T_1, t_1, \Sigma) \subseteq \text{chase}(T_2, t_2, \Sigma)$ und

(b) $Q_1 \equiv_{\Sigma} Q_2 \iff \text{chase}(T_1, t_1, \Sigma) \equiv \text{chase}(T_2, t_2, \Sigma)$ und

Beweis: Siehe Tafel.

Bemerkung: Aus Theorem 6.13, Korollar 6.12 und Korollar 5.6 folgt insbesondere, dass " $Q_1 \subseteq_{\Sigma} Q_2$ " bzw. " $Q_1 \equiv_{\Sigma} Q_2$ " entscheidbar ist und zur Komplexitätsklasse NP gehört.

Anfrage-Minimierung bzgl. Σ

Vorgehensweise:

- ▶ **Eingabe:** Tableau-Anfrage $Q = (T, t)$ über R und Menge Σ von FDs über R
- ▶ **Schritt 1:** Berechne $Q' := (T', t') := \text{chase}(T, t, \Sigma)$. Klar: $Q' \equiv_{\Sigma} Q$
- ▶ **Schritt 2:** Nutze den Algorithmus aus Korollar 5.9(a), um eine **minimale** Tableau-Anfrage $Q'' := (T'', t'')$ mit $Q'' \equiv Q'$ zu berechnen
Insbesondere gilt: $Q'' \equiv_{\Sigma} Q' \equiv_{\Sigma} Q$

Notation: Ist $Q = (T, t)$ eine Tableau-Anfrage, so schreibe $\text{min}(Q) := \text{min}(T, t)$, um die gemäß Korollar 5.9(a), minimale zu Q äquivalente Tableau-Anfrage zu bezeichnen.

Lemma 6.14

Sei $Q = (T, t)$ eine Tableau-Anfrage über R und sei Σ eine Menge von FDs über R .
 Dann gilt: $|\text{min}(\text{chase}(T, t, \Sigma))| \leq |\text{min}(T, t)|$.

Hier ohne Beweis.

Eine Beweisskizze findet sich auf Seite 178 des Buchs [AHV].

Ziel der Normalisierung

Gegeben:

- ▶ eine endliche Menge U von Attribut-Namen
- ▶ eine Menge Σ von Abhängigkeiten zwischen den Attributen aus U

Ziel:

Finde ein Datenbankschema mit den Attributen aus U , das

- ▶ möglichst **Redundanz-frei** ist \rightsquigarrow spart Speicherplatz
 \rightsquigarrow verhindert Änderungs-Anomalien
- ▶ alle gewünschten Daten speichern kann
 "lossless", d.h. "**informationsverlustfrei**"
- ▶ die Abhängigkeiten aus Σ "respektiert"
 "dependency preserving", d.h. "**abhängigkeitstreu**"

Abhängigkeiten und Normalformen

6.1 Funktionale Abhängigkeiten

6.2 The Chase — Die Verfolgungsjagd

6.3 Normalformen — Informationstheoretischer Ansatz

Motivation (1/2)

Ziel beim Datenbank-Entwurf:

Ein DB-Schema entwickeln, so dass Informationen zum gewünschten Anwendungsbereich "sinnvoll" gespeichert werden können.
 Insbesondere: Wenn möglich, Redundanzen und Inkonsistenzen vermeiden.

Beispiel: Relation *Warenlager*[Bauteil-Nr, Lager-Nr, Menge, Ort]

Warenlager:

Bauteil-Nr	Lager-Nr	Menge	Ort
2411	2	200	Riedberg
2412	3	300	Bornheim
3001	1	100	Hanau
2415	2	100	Riedberg

Unschön: Redundanz — der Ort von Lager 2 ist mehrfach gespeichert.

Dadurch können Inkonsistenzen auftreten:

Update-Anomalien = Inkonsistenzen, die durch Aktualisierung der DB auftreten können:

- ▶ **Änderungs-Anomalie:** den Ort in Zeile 1 durch "Westend" ersetzen
 \rightsquigarrow Adresse von Lager 2 nicht mehr eindeutig
- ▶ **Lösch-Anomalie:** Löschen von Zeile 2
 \rightsquigarrow Information über die Adresse von Lager 3 geht verloren
- ▶ **Einfüge-Anomalie:** Die Adresse eines neuen Lagers kann erst dann eingefügt werden, wenn mindestens ein Bauteil dort gelagert wird.

Motivation (2/2)

Zur Vermeidung dieser Update-Anomalien:

Informationen auf 2 Relationen *Adressen*[Lager-Nr,Ort] und *Lagerung*[Bauteil-Nr,Lager-Nr,Menge] aufteilen

Adressen:

Lager-Nr	Ort
2	Riedberg
3	Bornheim
1	Hanau

Abhängigkeit:
Lager-Nr \rightarrow Ort

Lagerung:

Bauteil-Nr	Lager-Nr	Menge
2411	2	200
2412	3	300
3001	1	100
2415	2	100

Abhängigkeit:
Bauteil-Nr, Lager-Nr \rightarrow Menge

Zur Anfrage-Optimierung:

Die "optimierte Anfrage" muss nur auf solchen Datenbanken äquivalent zur Original-Anfrage sein, die die obigen Abhängigkeiten erfüllen.

\rightsquigarrow die minimale, solchermaßen äquivalente Anfrage ist evtl. noch kleiner als die, die durch die Tableau-Minimierung aus Kapitel 5 gefunden wird.

Test auf Implikation

Satz 6.15

Es gibt einen Algorithmus, der bei Eingabe einer endlichen Attributmenge U , einer Menge Σ von FDs über U und einer FD $X \rightarrow Y$ über U entscheidet, ob $\Sigma \models_U X \rightarrow Y$.

Beweis: siehe Tafel.

Notation

► **Relationsschema:**

bisher: Relations-Name R (dem eine Attributmenge U zugeordnet ist);
Schreibweise: $R[U]$

jetzt: $(R[U], \Sigma)$, wobei R ein Relations-Name der Sorte U und Σ eine Menge von FDs über U .

► **DB-Schema:**

bisher: Menge \mathbf{R} von Relations-Namen

jetzt: Menge $\{ (R_1[U_1], \Sigma_1), \dots, (R_n[U_n], \Sigma_n) \}$ von Relationsschemata

► **Implikation bzw. Äquivalenz von FD-Mengen:**

Σ und Γ seien zwei Mengen von FDs über U , f sei eine FD über U .

► $\Sigma \models_U f$: \iff für alle $I \in inst(U)$ gilt: falls $I \models \Sigma$, so $I \models f$.

► $\Sigma \models_U \Gamma$: \iff für alle $I \in inst(U)$ gilt: falls $I \models \Sigma$, so $I \models \Gamma$.

► $\Sigma \equiv_U \Gamma$: \iff $\Sigma \models_U \Gamma$ und $\Gamma \models_U \Sigma$.

Zerlegungen

Definition 6.16

(a) Eine **Zerlegung** (engl: **decomposition**) eines Relationsschemas $(R[U], \Sigma)$ ist ein DB-Schema $\{ (R_1[U_1], \Sigma_1), \dots, (R_n[U_n], \Sigma_n) \}$, so dass

$$U_1 \cup \dots \cup U_n = U \quad \text{und} \quad \Sigma \models_U \Sigma_1 \cup \dots \cup \Sigma_n.$$

(b) Eine Zerlegung heißt **abhängigkeitstreu** (engl: **dependency preserving**), falls $\Sigma \equiv_U \Sigma_1 \cup \dots \cup \Sigma_n$.

(c) Eine Zerlegung heißt **informationsverlustfrei** (engl: **lossless**), falls für alle Relationen $I \in inst(R)$ mit $I \models \Sigma$ gilt: $I = \pi_{U_1}(I) \bowtie \dots \bowtie \pi_{U_n}(I)$.

Zur Erinnerung: $I \subseteq \pi_{U_1}(I) \bowtie \dots \bowtie \pi_{U_n}(I)$ gilt sowieso.

Um Informationsverlustfreiheit nachzuweisen, genügt es also, nachzuprüfen, ob für alle $I \in inst(R)$ gilt: $\pi_{U_1}(I) \bowtie \dots \bowtie \pi_{U_n}(I) \subseteq I$.

Boyce-Codd Normalform (BCNF)

Intuition: "Do not represent the same fact twice" ... Elimination von Redundanzen

Definition 6.17

- (a) Ein Relationenschema $(R[U], \Sigma)$ ist in **Boyce-Codd Normalform (BCNF)**, falls für alle FDs $X \rightarrow A$ mit $X \subseteq U$ und $A \in U \setminus X$ gilt:

Falls $\Sigma \models X \rightarrow A$, so $\Sigma \models X \rightarrow U$ d.h. X ist ein "Superschlüssel"

- (b) Ein DB-Schema $\{ (R_1[U_1], \Sigma_1), \dots, (R_n[U_n], \Sigma_n) \}$ ist in BCNF, falls jedes seiner Relationenschemata (R_i, Σ_i) in BCNF ist.

Algorithmus zur BCNF-Dekomposition

Eingabe: Relationenschema $(R[U], \Sigma)$

Ausgabe: DB-Schema $D = \{ (R_1, \Sigma_1), \dots, (R_n, \Sigma_n) \}$ in BCNF

- (1) $D := \{ (R[U], \Sigma) \}$
- (2) Solange D nicht in BCNF ist, wiederhole:
 - (2.1) Wähle ein Rel.schema $(S[V], \Gamma) \in D$, das nicht in BCNF ist
 - (2.2) Wähle Mengen $\emptyset \neq X, Y, Z \subseteq V$ mit $V = X \dot{\cup} Y \dot{\cup} Z$ so dass $\Gamma \models X \rightarrow Y$ und, für alle $A \in Z$, $\Gamma \not\models X \rightarrow A$.
 - (2.3) Ersetze $(S[V], \Gamma)$ in D durch $(S_1[XY], \pi_{XY}(\Gamma))$ und $(S_2[XZ], \pi_{XZ}(\Gamma))$
 - (2.4) Für alle $(S'[V'], \Gamma') \in D$ und $(S''[V''], \Gamma'') \in D$ mit $V' \subseteq V''$, entferne $(S'[V'], \Gamma')$ aus D .

Satz 6.18

Der obige Algorithmus zerlegt ein Relationenschema **informationsverlustfrei** in ein Datenbankschema, das in BCNF ist.

Beweis:

Der Algorithmus terminiert mit D in BCNF: klar.

D ist eine **Zerlegung** von (R, Σ) : klar.

Informationsverlustfreiheit der Zerlegung folgt direkt aus Proposition 6.2.

Bemerkung: Exponentielle Laufzeit wegen Berechnung von $\pi_V(\Gamma)$ und Test auf BCNF-Eigenschaft.

Beispiel

- ▶ $R[U] = \text{Warenlager}[\text{Bauteil-Nr}, \text{Lager-Nr}, \text{Menge}, \text{Ort}]$
- ▶ $\Sigma = \{ \text{Lager-Nr} \rightarrow \text{Ort}, \text{Bauteil-Nr}, \text{Lager-Nr} \rightarrow \text{Menge} \}$
- ▶ $(R[U], \Sigma)$ ist **nicht in BCNF**,
denn $\Sigma \not\models \text{Lager-Nr} \rightarrow \text{Bauteil-Nr}, \text{Lager-Nr}, \text{Menge}, \text{Ort}$

- ▶ $R_1[U_1] = \text{Adressen}[\text{Lager-Nr}, \text{Ort}]$ und $\Sigma_1 = \{ \text{Lager-Nr} \rightarrow \text{Ort} \}$
- ▶ $R_2[U_2] = \text{Lagerung}[\text{Bauteil-Nr}, \text{Lager-Nr}, \text{Menge}]$ und $\Sigma_2 = \{ \text{Bauteil-Nr}, \text{Lager-Nr} \rightarrow \text{Menge} \}$
- ▶ das DB-Schema $\{ (R_1[U_1], \Sigma_1), (R_2[U_2], \Sigma_2) \}$ ist **in BCNF**

Notation: Für eine Menge Γ von FDs und eine Attributmeng $V \subseteq U$ setzen wir

$$\pi_V(\Gamma) := \{ X \rightarrow A : XA \subseteq V \text{ und } \Gamma \models X \rightarrow A \}.$$

Beachte: $\pi_V(\Gamma)$ kann mit Hilfe des Algorithmus aus Satz 6.15 berechnet werden.

BCNF vs. Abhängigkeitstreue

Beispiel 6.19

- ▶ $R[U] = \text{Adresse}[\text{Straße}, \text{Ort}, \text{PLZ}]$
 $\Sigma = \{ \text{PLZ} \rightarrow \text{Ort}, \text{Ort}, \text{Straße} \rightarrow \text{PLZ} \}$
- ▶ Algorithmus zur BCNF-Dekomposition liefert $D = \{ (S[\text{PLZ}, \text{Ort}], \{ \text{PLZ} \rightarrow \text{Ort} \}), (T[\text{PLZ}, \text{Straße}], \emptyset) \}$
- ▶ Laut Satz 6.18 ist D eine **informationsverlustfreie Zerlegung in BCNF**.
- ▶ D ist **nicht abhängigkeittreu** (weil die FD $\text{Ort}, \text{Straße} \rightarrow \text{PLZ}$ verloren geht)
- ▶ Es gilt sogar: **Es gibt keine BCNF-Zerlegung von**

$$\left(\text{Adresse}[\text{Straße}, \text{Ort}, \text{PLZ}], \{ \text{PLZ} \rightarrow \text{Ort}; \text{Ort}, \text{Straße} \rightarrow \text{PLZ} \} \right),$$

die abhängigkeittreu ist.

Informationsverlustfreiheit und Abhängigkeitstreue

Frage:

Gibt es an Stelle von BCNF eine andere Variante von “Normalform”, die informationsverlustfrei und abhängigkeitstreu ist?

Antwort:

Ja, die so genannte *dritte Normalform 3NF*

Im Vergleich zu BCNF ist die **3NF** weniger restriktiv, **eliminiert Redundanzen also weniger gründlich als BCNF**.

Dafür sind in 3NF aber informationsverlustfreie und abhängigkeitstreu Zerlegungen möglich.

In der Praxis ist 3NF meist interessanter als BCNF.

Gütekriterien für Datenbank-Entwürfe

Was ist ein gutes Datenbankschema?

- ▶ eins in BCNF oder 3NF (oder 2NF, 4NF, ...)

Warum sind BCNF oder 3NF “gut”?

- ▶ Vermeidung von Update-Anomalien (bei BCNF “gründlicher” als bei 3NF)
- ▶ Existenz von Algorithmen, die ein gegebenes DB-Schema in ein “gutes” DB-Schema überführen:
 - ▶ informationsverlustfrei
 - ▶ bei 3NF auch: **abhängigkeitstreu**

Dritte Normalform (3NF)

Sei $(R[U], \Sigma)$ ein Relationsschema.

- ▶ Ein **Superschlüssel** ist eine Menge $X \subseteq U$ mit $\Sigma \models X \rightarrow U$.
- ▶ Ein **Schlüssel** ist ein (bzgl \subseteq) minimaler Superschlüssel.
- ▶ Ein **Schlüsselattribut** ist ein Attribut $A \in U$, das zu mindestens einem Schlüssel gehört.

Definition 6.20

- (a) Ein Relationsschema $(R[U], \Sigma)$ ist in **dritter Normalform (3NF)**, falls für alle FDs $X \rightarrow A$ mit $X \subseteq U$ und $A \in U \setminus X$ gilt:

Falls $\Sigma \models X \rightarrow A$, so $\Sigma \models X \rightarrow U$ oder A ist ein Schlüsselattribut.

- (b) Ein DB-Schema $\{ (R_1[U_1], \Sigma_1), \dots, (R_n[U_n], \Sigma_n) \}$ ist in 3NF, falls jedes seiner Relationsschemata (R_i, Σ_i) in 3NF ist.

Beispiel: Das Relationsschema

$(\text{Adresse} [\text{Straße, Ort, PLZ}], \{ \text{PLZ} \rightarrow \text{Ort}, \text{Ort, Straße} \rightarrow \text{PLZ} \})$ ist in 3NF.

Jetzt: Weitere Rechtfertigung der Güte von BCNF mittels Informationstheorie

Vorteil:

- ▶ Dieser Ansatz benutzt ausschließlich Charakteristika der **Daten** (bzw. des Schemas (R, Σ)) und hängt nicht von Begriffen wie “Update-Anomalien” und “Informationsverlustfreiheit” ab.
- ▶ Der Ansatz funktioniert auch für andere Datenmodelle als das relationale Modell; insbesondere auch für **XML**.
Dort: noch einigermaßen unklar, was genau “Update-Anomalien” bzw. Update-Sprachen und was “Informationsverlustfreiheit” sein soll.

Literatur:

Marcelo Arenas und Leonid Libkin.

An Information-Theoretic Approach to Normal Forms for Relational and XML Data. Journal of the ACM, Volume 52, No. 2, 2005, pages 246–283.

(Diese Arbeit hat bei der Konferenz PODS 2003 den “Best Paper Award” gewonnen.)

“Crashkurs” in Informationstheorie

Informationstheorie — Motivation (2/2)

1. **Möglichkeit:** Jeden Buchstaben als Bitstring der Länge $\lceil \lg(26) \rceil = 5$ kodieren.
Originaltext der Länge n \rightsquigarrow Bitstring der Länge $5n$

Beobachtung: In Texten deutscher Sprache kommen manche Buchstaben häufiger vor als andere. Zum Beispiel:

- ▶ häufigster Buchstabe: **e** mit Wahrscheinlichkeit $P(e) = 0,175$
 (d.h.: jeder 6-te Buchstabe ist ein e)
- ▶ zweithäufigster Buchstabe: **n** mit Wahrscheinlichkeit $P(n) = 0,1$
 (d.h.: jeder 10-te Buchstabe ist ein n)
- ▶ seltenster Buchstabe: **q** mit Wahrscheinlichkeit $P(q) = 0,01$
 (d.h.: nur jeder 100-te Buchstabe ist ein q)

Naheliegende Idee:

Kodiere Buchstaben, die oft vorkommen durch *kurze* Bitstrings und Buchstaben, die selten vorkommen, durch längere Bitstrings.

Bei optimaler Kodierung dann:

Originaltext der Länge n \rightsquigarrow Bitstring der Länge $4n$

D.h.: im Schnitt nur 4 Bits pro Original-Buchstabe
 (zum Vergleich: 5 Bits bei herkömmlicher Kodierung)

Informationstheorie — Motivation (1/2)

- ▶ 1948 von **Claude Elwood Shannon** begründet
- ▶ **Fragestellung:**
 - ▶ Was ist Information?
 - ▶ Wie kann man den **Informationsgehalt** verschiedener Informationsquellen miteinander vergleichen?
- ▶ **Szenario: (Beispiel)**
 - ▶ eine Informationsquelle erzeugt Texte in deutscher Sprache
 - ▶ ein “Sender” liest die Texte und will sie an einen “Empfänger” schicken
 - ▶ dazu werden die Texte als Bitstrings repräsentiert
 - ▶ **Ziel: Bitstring soll möglichst kurz sein**
- ▶ Annahme: Originaltext mit Buchstaben aus $\{a, \dots, z\}$: 26 Stück

Informationsgehalt und Entropie — Anschaulich

Begriffe:

- ▶ **Entropie** von Texten deutscher Sprache ≈ 4
 \approx “durchschnittliche Länge des Bitstrings pro Buchstaben des Originaltextes”
 (bei optimaler Kodierung)
- ▶ **Informationsgehalt** eines Buchstabens $\alpha \in \{a, \dots, z\}$
 \approx “Länge des Bitstrings, der den Buchstaben α kodiert”
 (bei optimaler Kodierung)

Intuition: Lesen des Buchstabens “y” gibt mir mehr Information als Lesen des Buchstabens “e”, da es weniger Worte gibt, in denen “y” vorkommt als Worte, in denen “e” vorkommt.

Informationsgehalt und Entropie — Präzise

Definition 6.21

(a) Ein **Zufallsraum** $\mathcal{E} = (\{e_1, \dots, e_m\}, P_{\mathcal{E}})$ besteht aus einer Menge $\{e_1, \dots, e_m\}$ von Ereignissen, denen Wahrscheinlichkeiten $p_1 = P_{\mathcal{E}}(e_1), \dots, p_m = P_{\mathcal{E}}(e_m)$ zugeordnet sind, so dass $0 \leq p_i \leq 1$, für alle $i \in \{1, \dots, m\}$, und $\sum_{i=1}^m p_i = 1$.

(b) Der **Informationsgehalt von Ereignis e_i** ist $\lg\left(\frac{1}{p_i}\right)$ (\lg ist Logarithmus zur Basis 2)

$$\lg\left(\frac{1}{p_i}\right) \quad \left(= -\lg(p_i) \right)$$

(c) Die **Entropie von \mathcal{E}** , kurz: $H(\mathcal{E})$, ist

$$H(\mathcal{E}) := H(p_1, \dots, p_m) := \sum_{i=1}^m p_i \cdot \lg\left(\frac{1}{p_i}\right) \quad \left(= -\sum_{i=1}^m p_i \cdot \lg(p_i) \right)$$

Konvention für $p_i = 0$: $0 \cdot \lg\left(\frac{1}{0}\right) := -0 \cdot \lg(0) := 0$

Eigenschaften von Entropie und Informationsgehalt

Eigenschaften des Informationsgehalts:

► Ereignis, das mit **Wahrscheinlichkeit 1** eintrifft: **Informationsgehalt 0**
(denn: $\lg\left(\frac{1}{1}\right) = \lg(1) = 0$)

► Informationsgehalte zweier voneinander unabhängiger Ereignisse addieren sich:

$$\lg\left(\frac{1}{p_i \cdot p_j}\right) = \lg\left(\frac{1}{p_i}\right) + \lg\left(\frac{1}{p_j}\right)$$

Eigenschaften der Entropie:

► $H(p_1, \dots, p_m)$ ist **maximal**, falls $p_1 = \dots = p_m = \frac{1}{m}$.

$$\text{Es gilt: } H\left(\frac{1}{m}, \dots, \frac{1}{m}\right) = \sum_{i=1}^m \frac{1}{m} \lg(m) = \lg(m)$$

► $H(p_1, \dots, p_m)$ ist **minimal**, falls es ein $i \in \{1, \dots, m\}$ gibt mit $p_i = 1$ und $p_j = 0$ f.a. $j \neq i$.

$$\text{Es gilt: } H(1, 0, \dots, 0) = 1 \cdot \lg\left(\frac{1}{1}\right) + 0 + \dots + 0 = 0$$

► Somit gilt stets: $0 \leq H(p_1, \dots, p_m) \leq \lg(m)$

Lemma 6.22

Für jeden Zufallsraum $\mathcal{E} = (\{e_1, \dots, e_m\}, P_{\mathcal{E}})$ mit $(p_1, \dots, p_m) := (P_{\mathcal{E}}(e_1), \dots, P_{\mathcal{E}}(e_m))$ gilt: $H(p_1, \dots, p_m) = 0 \iff$ Es gibt ein $i \in \{1, \dots, m\}$ mit $p_i = 1$ und $p_j = 0$ f.a. $j \neq i$.

Beweis: Siehe Tafel.

Bemerkungen zu Entropie und Informationsgehalt

► Man kann beweisen, dass

$$H(\mathcal{E}) \leq \left(\text{durchschnittliche Länge des Bitstrings, der ein Ereignis aus } \{e_1, \dots, e_m\} \text{ kodiert (bei optimaler Kodierung)} \right) \leq H(\mathcal{E}) + 1$$

► Wegen $H(\mathcal{E}) = \sum_{i=1}^m p_i \cdot \lg\left(\frac{1}{p_i}\right)$ gilt somit:

$$\left(\text{Informationsgehalt von Ereignis } e_i \right) \stackrel{\text{Def}}{=} \lg\left(\frac{1}{p_i}\right) \approx \left(\text{Länge des Bitstrings, der } e_i \text{ kodiert (bei optimaler Kodierung)} \right)$$

Informationstheorie und Normalformen

Notation

Der Einfachheit halber sei jetzt $\mathbf{dom} := \mathbb{N}_{\geq 1} = \{1, 2, 3, \dots\}$

Definition 6.23

Sei R ein Relations-Name der Sorte $\mathit{sort}(R)$.

(a) Für eine Relation $I \in \mathit{inst}(R)$ ist

$$\mathit{Pos}(I) := \{ (R, t, A) : t \in I, A \in \mathit{sort}(R) \}$$

die Menge der **Positionen von I** .

(Klar: $|\mathit{Pos}(I)| = ||I|| \stackrel{\text{Def}}{=} |I| \cdot |\mathit{sort}(R)|$)

(b) Für ein Relationsschema (R, Σ) setzen wir

$$\mathit{inst}(R, \Sigma) := \{ I \in \mathit{inst}(R) : I \models \Sigma \}$$

und für ein Element $k \in \mathbf{dom} = \mathbb{N}_{\geq 1}$ setzen wir

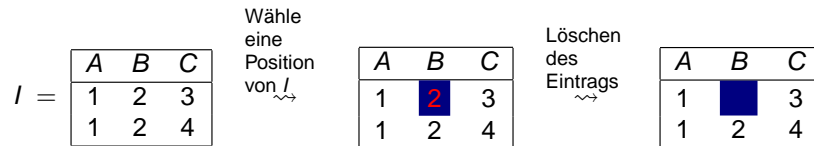
$$\mathit{inst}_k(R, \Sigma) := \{ I \in \mathit{inst}(R, \Sigma) : \mathit{adom}(I) \subseteq \{1, \dots, k\} \}$$

Entropie zum Messen von Redundanz — Bsp (2/2)

Beispiel: Relationsschema $R[ABC]$ mit $\Sigma = \{A \rightarrow B\}$

(Klar: Ist nicht in BCNF, da A kein Superschlüssel)

Instanz $I \in \mathit{inst}(R, \Sigma)$:



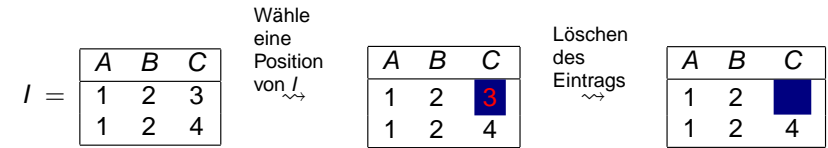
- ▶ Wie eben: Wähle eine endliche Teilmenge von \mathbf{dom} , die echt größer ist als $\mathit{adom}(I)$, etwa: $\{1, \dots, 6\}$
- ▶ Jetzt aber: Wähle eine andere Position von I
- ▶ Wahrscheinlichkeitsverteilung für's Einsetzen einzelner Elemente an der gelöschten Position:
 $P(2) = 1$ (denn: $I \models \Sigma$ und $(A \rightarrow B) \in \Sigma$) und
 f.a. $i \in \{1, \dots, 6\} \setminus \{2\}$ gilt: $P(i) = 0$.
- ▶ Entropie = $\lg(1) = 0$ = durchschnittlicher Informationsgehalt eines Eintrags an der ausgewählten Position

Entropie zum Messen von Redundanz — Bsp (1/2)

Beispiel: Relationsschema $R[ABC]$ mit $\Sigma = \{A \rightarrow B\}$

(Klar: Ist nicht in BCNF, da A kein Superschlüssel)

Instanz $I \in \mathit{inst}(R, \Sigma)$:



- ▶ Wähle eine endliche Teilmenge von \mathbf{dom} , die echt größer ist als $\mathit{adom}(I)$, etwa: $\{1, \dots, 6\}$
- ▶ Wahrscheinlichkeitsverteilung für's Einsetzen einzelner Elemente an der gelöschten Position:
 $P(4) = 0$ (denn: I hat 2 Zeilen!) und f.a. $i \in \{1, \dots, 6\} \setminus \{4\}$ gilt: $P(i) = \frac{1}{5}$.
- ▶ Entropie = $\lg(5) \approx 2,322$ = durchschnittlicher Informationsgehalt eines Eintrags an der ausgewählten Position

Entropie zum Messen von Redundanz

Grundidee:

Messen, wieviel Information man gewinnt, wenn der Eintrag an einer bestimmten Position verloren geht und dann an dieser Position wieder ein Eintrag eingesetzt wird (entweder der Originaleintrag oder ein anderer, der die FDs aus Σ nicht verletzt und die Zeilenanzahl von I nicht verändert).

Entropie zum Messen von Redundanz — Präzise

Definition 6.24

Sei (R, Σ) ein Relationsschema, $k \in \mathbf{dom} = \mathbb{N}_{\geq 1}$, $I \in \mathit{inst}_k(R, \Sigma)$ und sei $\mathit{pos} \in \mathit{Pos}(I)$ eine Position von I .

- (a) Für jedes $a \in \mathbf{dom}$ sei $I_{\mathit{pos} \leftarrow a}$ die Relation, die aus I entsteht, indem an Position pos der Wert a eingesetzt wird.

D.h.: Ist $\mathit{pos} = (R, t, A)$ mit $t \in I$ und $A \in \mathit{sort}(R)$, so ist $I_{\mathit{pos} \leftarrow a} := (I - \{t\}) \cup \{t_{\mathit{pos} \leftarrow a}\}$, wobei für alle $B \in \mathit{sort}(R)$ gilt:

$$t_{\mathit{pos} \leftarrow a}(B) := \begin{cases} a & \text{falls } B = A \\ t(B) & \text{sonst} \end{cases}$$

- (b) Der Zufallsraum $\mathcal{E}_{\Sigma}^k(I, \mathit{pos}) = (\{1, \dots, k+1\}, P)$ besteht aus den Ereignissen $\{1, \dots, k+1\}$, denen die folgenden Wahrscheinlichkeiten $p_a := P(a)$, für alle $a \in \{1, \dots, k+1\}$ zugeordnet werden:

$$p_a := P(a) := \begin{cases} 0, & \text{falls } (I_{\mathit{pos} \leftarrow a} \not\models \Sigma \text{ oder } |I_{\mathit{pos} \leftarrow a}| \neq |I|) \\ \frac{1}{|\{b \in \{1, \dots, k+1\} : I_{\mathit{pos} \leftarrow b} \models \Sigma \text{ und } |I_{\mathit{pos} \leftarrow b}| = |I|\}|} & , \text{sonst} \end{cases}$$

Insbesondere: Die Entropie $H(\mathcal{E}_{\Sigma}^k(I, \mathit{pos})) = \sum_{a \in \{1, \dots, k+1\}} p_a \cdot \lg\left(\frac{1}{p_a}\right)$ gibt in etwa den mittleren Informationsgehalt von Position pos von I an.

Bemerkungen

- ▶ Theorem 6.25 liefert eine elegante Charakterisierung der BCNF.
- ▶ **ABER:** Als "gutes" Maß für die Redundanz an Position pos ist $H(\mathcal{E}_{\Sigma}^k(I, \mathit{pos}))$ zu grobkörnig, da es z.B. keinen Unterschied zwischen den beiden folgenden Relationen macht:

$R[ABC]$ mit $\Sigma := \{A \rightarrow B\}$

$$I = \begin{array}{|c|c|c|} \hline A & B & C \\ \hline 1 & 2 & 3 \\ \hline 1 & \blacksquare & 4 \\ \hline \end{array}$$

Entropie $H(\mathcal{E}_{\Sigma}^k(I, \mathit{pos})) = 0$

$$J = \begin{array}{|c|c|c|} \hline A & B & C \\ \hline 1 & 2 & 3 \\ \hline 1 & 2 & 4 \\ \hline 1 & \blacksquare & 5 \\ \hline \end{array}$$

Entropie $H(\mathcal{E}_{\Sigma}^k(J, \mathit{pos})) = 0$

Somit: In beiden Fällen ist die Entropie gleich. Aber intuitiv ist die "Redundanz" bei J größer als bei I .

Informationstheoretische Charakterisierung von BCNF

Theorem 6.25 (Arenas, Libkin, 2003)

Sei (R, Σ) ein Relationsschema (insbes. sei Σ eine Menge von FDs über $\mathit{sort}(R)$). Dann gilt:

(R, Σ) ist in BCNF

$$\iff \text{Für alle } k \in \mathbf{dom} = \mathbb{N}_{\geq 1}, I \in \mathit{inst}_k(R, \Sigma) \text{ und } \mathit{pos} \in \mathit{Pos}(I) \text{ gilt: } H(\mathcal{E}_{\Sigma}^k(I, \mathit{pos})) > 0$$

Beweis: Siehe Tafel.

Bemerkung: Anschaulich besagt dieses Theorem, dass (R, Σ) genau dann in BCNF ist, wenn es keine Redundanz gibt, d.h. keine Position den mittleren Informationsgehalt 0 hat.

Ausblick

Marcelo Arenas und Leonid Libkin entwickelten ein komplizierteres Maß $\mathit{RIC}_I(\mathit{pos} | \Sigma)$, den so genannten relativen Informationsgehalt von I an Position pos unter den Abhängigkeiten Σ , so dass gilt:

- ▶ $0 \leq \mathit{RIC}_I(\mathit{pos} | \Sigma) \leq 1$
- ▶ je näher bei 1 desto weniger Redundanz, d.h. desto größer der Informationsgehalt
- ▶ für die Beispiele I und J von der vorherigen Folie gilt: $\mathit{RIC}_I(\mathit{pos} | \Sigma) = 0,875 > 0,781 = \mathit{RIC}_J(\mathit{pos} | \Sigma)$.
- ▶ Ein Relationsschema (R, Σ) heißt gut-entworfen (engl.: well-designed), falls für alle $I \in \mathit{inst}(R, \Sigma)$ und alle Positionen $\mathit{pos} \in \mathit{Pos}(I)$ gilt: $\mathit{RIC}_I(\mathit{pos} | \Sigma) = 1$ (d.h.: jede Instanz hat in jeder Position den maximalen Informationsgehalt)
- ▶ (R, Σ) ist gut-entworfen $\iff (R, \Sigma)$ ist in BCNF
- ▶ Beim Durchführen des Algorithmus zur BCNF-Dekomposition sinkt der Informationsgehalt an keiner Position.
- ▶ Ähnliche Resultate auch für andere Normalformen (3NF, 4NF, ...) sowie für XML.