

2.3 Das relationale Modell

- Datenstrukturierung in Tabellenform (Relationen)
- Beziehungen zwischen Tabellen durch entsprechende Spaltenwerte (Attributwerte)
- Manipulation der Tabellen durch Mengenoperationen (Relationenalgebra)

◆ **Beispiel :**

Auftr.-Nr.	Datum	Kd-Nr	Kd-Name
10228	1.1.92	204342	...
10229	1.3.92	252554	...
24224	7.5.92	321001	...
...

Artikel-Nr	Bezeichn.	Preis	...
4022	Blech	2043,42	...
8029	Welle	252,54	...
2243	Achse	3210,00	...
...

Auftr.-Nr.	Artikel-Nr.	Menge	Preis
10228	8029	200	...
10228	4022	54	...
10229

Zeile
Tupel

◆ **Probleme :** Wie werden Tabellen festgelegt ?
 Wie erreicht man Redundanzfreiheit ?
 Wie wird mit den Tabellen gearbeitet ?

◆ **Lösung :** Relationentheorie, Normalisierung, Relationenalgebra

Bild 2.3 -1
Datenbank- und Informationssysteme
Sep-02 (c) Prof. Dr. F. Laux

Die mit Abstand bekannteste Datenstrukturierung ist das weitverbreitete **relationale Modell**. Es führt alle Datenstrukturen des zu modellierenden Sachverhalts auf *Relationen* zurück. Dabei werden die Eigenschaften eines Objekts als Attribute von Relationen abgebildet. Normalisierte Relationen können anschaulich durch zweidimensionale *Tabellen* dargestellt werden. Die Beliebtheit dieses Modell für den kommerziellen Einsatz erklärt sich nicht nur aus der Tatsache, dass viele betriebliche Informationen ebenfalls tabellarisch dargestellt werden, sondern auch aus der formalen Schlichtheit des Modells und seiner Flexibilität.


Wenn wir das Beispiel von 2.3-1 betrachten, fällt auf, dass sich unser relationales Modell auf 2 betriebliche Begriffe (Objekte) bezieht: Artikel und Auftrag. Die Artikelrelation ist eine direkte Abbildung der betrieblichen Situation, während ein Auftrag durch 2 Relationen dargestellt wird. Wir können an diesem einfachen Beispiel schon erkennen, dass das Modell vor allem für Sachverhalte geeignet ist, die sich in natürlicher Weise tabellarisch darstellen lassen.

Die *Beziehungen* der Relationen zueinander werden über korrespondierende Attributwerte definiert. Diese Strukturinformationen müssen bei jeder Abfrage neu aufgebaut werden, was zu Kritik bezüglich der Effizienz von Datenrecherchen geführt hat. Soll z. B. ein Auftrag rekonstruiert werden, so müssen in unserem Beispiel Daten aus drei Relationen (wieder) zusammengefügt werden. Dies geschieht durch Mengenoperationen auf den Relationen (*Relationenalgebra*, Bild 2.3-7).

Wir nennen ein Element einer Relation *Tupel*. Jede *Zeile* einer Tabelle entspricht einem Tupel der Relation. Die möglichen Werte eines Attributs nennen wir seinen Wertebereich. In einer tabellarischen Darstellung entstammen alle Attributwerte (Spalte) dem Wertebereich des Attributs.

Komplexe Information kann auf verschiedene Arten strukturiert werden. Nur wenige Strukturen sind "gut". Die Relationentheorie liefert uns die Regeln dafür.

ÜA: In unserem Beispiel hätte der Auftrag auch durch eine Relation dargestellt werden können. Wie müsste die Relation dann aussehen? Welche Nachteile sehen Sie dabei?



2.3.1 Relationen

- ◆ **Relation :** Seien A_1, A_2, \dots, A_n Mengen (**Domänen** genannt), dann ist $R \subseteq A_1 * A_2 * \dots * A_n$ eine **n-stellige Relation** (Teilmenge des kartesischen Produkts über die A_i)
- ◆ **(n-)Tupel :** Ein Element $r = (r_1, r_2, \dots, r_n) \in R$ mit $r_i \in A_i$ heißt **n-Tupel** (oder einfacher **Tupel**) der Relation R
- ◆ **Attribut:** Die Spalten einer Relation nennen wir **Attribute**. Ein Element $r_i \in A_i$ eines n-Tupels heißt **Attributwert** (Datenelement) von $r \in R$

Bild 2.3 -2
Datenbank- und Informationssysteme
Sep-02 (c) Prof. Dr. F. Laux

Zur Präzisierung des Begriffs der Relation legen wir folgendes fest:

Eine *Domäne* (*Wertebereich*) A_i repräsentiert die Menge der zulässigen Werte eines Attributs. Die Werte können beliebige Daten darstellen.

Eine *Relation* R auf den Domänen A_1, A_2, \dots, A_n ist definiert als Teilmenge des kartesischen Produkts $\times A_i$ ($i=1 \dots n$).

Ein Element r von R heißt *Tupel* mit den Komponenten $r_i \in A_i$.


Vereinfacht kann man eine Relation als *zweidimensionale Tabelle* auffassen. Dabei entsprechen die *Spalten* den Domänen, die *Zeilen* den Tupeln und ein Tabellenwert repräsentiert einen Attributwert. Wird jede Spalte mit einem Bezeichner, dem Attributnamen, versehen, so wird die Spaltenreihenfolge bedeutungslos. Dieser Vergleich ist allerdings nicht ganz exakt, denn die Tupel einer Relation bilden eine Menge, können demzufolge keine Dubletten aufweisen, während eine Tabelle gleiche Zeilen enthalten kann.

- Beispiele: Domänen: N = Menge der natürlichen Zahlen
 Q = Menge der rationalen Zahlen
 string = Menge der Zeichenfolgen
 $[0:120]$ = Alle Zahlen im Intervall 0 bis 120
- Attribute: Name, Alter, Budget, Personal-Nr

Relation:

Mitarbeiter (Personal-Nr	Name	Alter)	
1234	Einstein	76	<- Tupel
0007	Bond	36	<- Tupel
0815	Nobody	22	<- Tupel

ÜA: Weisen Sie den Attributen *Name, Alter, Budget, Personal-Nr* Domänen zu.



2.3.1 Relationenschema

Relationenschema : Beschreibung einer Klasse von Relationen eines bestimmten Typs (Relationentyp)
 Der Ausdruck $RS (a_1 : A_1, a_2 : A_2, \dots a_n : A_n)$ heißt **Relationenschema**, wobei :

- RS = Name des Schemas
- a_i = Attributname
- $A_i = \text{dom}(a_i)$ Wertebereich (Domäne) des Attributs a_i

Beispiel : **Teilestamm (** Teilenummer : {10000 ... 99999},
 Bezeichnung : {Text | max. 40 Zeichen},
 Preis : {Währung DM},
 Herstellung : {Eigen-, Fremdfertigung},
 Mengeneinheit : {Stck, kg, l})

Bild 2.3 -3 Datenbank- und Informationssysteme Sep-02 (c) Prof. Dr. F. Laux

Im Hinblick auf die Datenmodellierung ist eine zeitinvariante Beschreibung von Relationen wichtig, denn wir wollen nicht die augenblickliche Situation sondern ihre statische Struktur ausdrücken.

Relationenschema :

Wir definieren ein *Relationenschema* als Beschreibung (Definition) einer Klasse gleichartiger Relationen. Seine Definition enthält den *Namen des Schema*, die *Attributnamen* und ihre *Domänen*. D. h. eine Relation ist eine mögliche Ausprägung (ein Exemplar) eines Relationenschemas.

Attribut : Jede Menge A_i des RS wird eindeutig durch einen Namen identifiziert. Dieser Name wird als **Attribut(name)** des RS bezeichnet

Domäne : Jede Menge A_i des RS wird als **Wertebereich (Domäne)** des Attributs a_i bezeichnet.

Bei den Grundbegriffen zu Kap. 2 haben wir zwischen Exemplar und seiner Beschreibung bzw. Struktur unterschieden. Relation und Relationenschema sind ein weiteres Beispiel dafür:

Bsp: Relationenschema: *Mitarbeiter* (Personal-Nr:N, Name:string, Alter:[0:120]);
 Attribute: Personal-Nr, Name, Alter;
 zugehörige Domänen: N, string, [0:120]

Eine Ausprägung dieses Relationenschemas finden Sie auf der vorherigen Seite.

Wie wir bereits in einem Beispiel gesehen haben, wird die Realität durch mehrere, voneinander abhängigen Relationen beschrieben. Wir definieren deshalb ein **relationales Datenbankschema** als ein Quadrupel (**{RS}**, **{FD}**, **{DI}**, **DML**), wobei **{RS}** eine Menge von *Relationenschemata*, **{FD}** eine Menge von *funktionalen Abhängigkeiten* zwischen Attributen einer Relation, **{DI}** eine Menge von *Abhängigkeiten zwischen Relationen* (**referentielle Integrität**) darstellt und **DML** eine relationale *Manipulationssprache* bezeichnet.

Bsp: Datenbankschema: **Projektverwaltung** (vgl. ERM auf Bild 2.1-5)

RS: Schema 1: *Mitarbeiter* (Pers-Nr:N, Name:string)
 Schema 2: *Projekt* (Proj-Name:string, Budget:DM)
 Schema 3: *arbeitet* (MA-Nr:N, Proj-Name:string, Stunden:decimal)
 (N = natürliche Zahlen, string = Zeichenfolgen, DM = Deutsche Mark
 (Währung), decimal = reelle Zahlen)

FD: Pers-Nr -> Name (es gibt eine Funktion f mit: Name = f(Pers-Nr))
 Proj-Name -> Budget
 MA-Nr, Proj-Name -> Stunden

DI: arbeitet.MA-Nr \in {Mitarbeiter.Pers-Nr} (Jeder Attributwert von MA-Nr in der Relation *arbeitet* muss auch in der Relation *Mitarbeiter* vorkommen)

Damit wird ausgesagt, daß nur Werte aus der Relationen *Mitarbeiter* zugelassen sind, d.h. der Wertebereich ist von Attributwerten einer anderen Relation abhängig und damit zeitlich veränderlich (*dynamisch*). Diesen Sachverhalt nennen wir **referentielle Integrität**.


arbeitet.Proj-Name \in {Projekt.Proj-Name} (beachten Sie, dass die Namensgleichheit zwar nicht zufällig ist, aber formal keine Bedeutung (Semantik) besitzt, deshalb ist diese Bedingung ebenso notwendig wie die vorangehende.

Die referentiellen Bedingungen werden oft schon bei der Definition eines Relationenschemas angegeben. Z. B. kann die Relation *arbeitet* geschrieben werden als:


arbeitet (MA-Nr:Mitarbeiter.Pers-Nr, Proj-Name:Projekt.Proj-Name, Stunden:decimal)

DML: Relationenalgebra (mit den Mengenoperationen *Projektion*, *Vereinigung*, *Selektion*, *kartesisches Produkt*, etc., siehe Bild 2.3-7)

ÜA: Erstellen Sie von den Relationen in Bild 2.3-1 ein vollständiges Datenbankschema.



2.3.2 Schlüsselattribute



Sei $RS(a_1:A_1, a_2:A_2, \dots, a_n:A_n)$ ein Relationenschema.
 Jede minimale Attributkombination $K \subseteq \{a_1, a_2, a_3, \dots, a_n\}$, die jedes Tupel r einer Relation über RS eindeutig identifizieren kann, heißt **Schlüsselkandidat**.

Schlüsselkandidat:

Identifikations-schlüssel: Eine fest ausgewählter Schlüsselkandidat $s \in K$ heißt **Identifikationsschlüssel (Identschlüssel, ID, primary key)**.
 [Er identifiziert eindeutig jedes Tupel einer Relation]

Schlüsselement: Jedes Attribut $s_i \in s$ heißt **Schlüsselement**.

Fremdschlüssel: Eine Attributkombination f des RS heißt **Fremdschlüssel (foreign key)**, wenn f auch Identschlüssel in einer anderen Relation ist.

Primärschlüssel: Ein Primärschlüssel ist ein **Datenfeld** einer Datei, das für die **primäre Dateioorganisation** verwendet wird. (Er lokalisiert genau einen Datensatz)

Sekundärschlüssel: Ein Sekundärschlüssel ist ein **Datenfeld** einer Datei, das als **Zugriffshilfe** (z.B Index) verwendet wird.

Bild 2.3 -5
Datenbank- und Informationssysteme
Sep-02 (c) Prof. Dr. F. Laux

Zur Illustration der Schlüsselbegriffe betrachten wir unser Beispiel auf der vorherigen Seite, das relationale Datenbankschema *Projektverwaltung*.


Das Attribut *Pers-Nr* ist der einzige Schlüsselkandidat der Relation *Mitarbeiter*, denn der Name ist nicht ausreichend, um ein Tupel zu identifizieren. Damit ist *Pers-Nr* auch der *Identifikationsschlüssel (ID)*.

Der ID der Relation *arbeitet* ist die Attributmenge {MA-Nr, Proj-Name} (warum?). Beide Attribute werden auch als *Schlüsselemente* bezeichnet, da sie zum ID gehören. Ausserdem sind sie auch *Fremdschlüssel* bezüglich den Relationen *Mitarbeiter* und *Projekt*.

Im Gegensatz zu den o.g. Schlüsseln sind *Primär- und Sekundärschlüssel* Begriffe der physischen Datenorganisation (Dateistruktur). In der Literatur wird dieser Unterschied gelegentlich verwechselt und der ID als Primärschlüssel bezeichnet. Dies hat vermutlich zwei Gründe:

- in der englischen Literatur heißt der ID *primary key*
- der Primärschlüssel wird häufig auch als ID verwendet. (dies erleichtert die technische Realisierung der Eindeutigkeitsforderung für den ID)

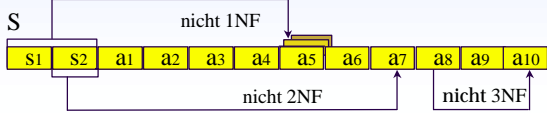
Beispiel: Wenn die Matrikelnr zugleich die Satznummer in einer Studentendatei ist, wird automatisch sichergestellt, daß die Matrikelnr eindeutig ist, denn es gibt für eine Nummer nur eine physische Adresse. In unserem Beispiel hat dies den Nachteil, dass die Datei viele Lücken (nicht verwendete Matrikelnummern) aufweist. Aus diesem Gesichtspunkt wäre eine andere Organisation der Daten und damit eine Unterscheidung zwischen Primärschlüssel und ID günstig.



2.3.3 Normalformen

funktionale Sei $R (a_1:A_1, a_2:A_2, \dots, a_n:A_n)$ ein RS und $\{a_j\} \subseteq \{a_i\}$, $A_i = \text{dom}(a_i)$ ($i, j := 1, 2, \dots, n$)

Abhängigkeit: Ein Attribut $a_k \in \{a_i\}$ ist von der Attributmenge $\{a_j\}$ **funktional abhängig**, wenn es eine skalare (einwertige) Funktion gibt mit: $y = f(x)$ für alle $y \in A_k$ und geeignetem $x \in X_{a_j}$
 Notation: $\{a_j\} \rightarrow a_k$



Sei S (z. Bsp. $S = \{s_1, s_2\}$) der Identischlüssel von RS.

1NF : Eine Relation R heißt **einfach**, bzw. in **erster Normalform**, wenn für alle Attribute a_i gilt: $S \rightarrow a_i$

2NF : Eine Relation R ist in **zweiter Normalform**, wenn sie in 1NF ist und keine funktionale Abhängigkeit von einer Teilmenge von S besteht. D.h. alle a_i sind **voll funktional** von S (Identischlüssel) abhängig.

3NF : Eine Relation R ist in **dritter Normalform**, wenn sie in 2NF ist und keine funktionale Abhängigkeit von Nichtschlüsselattributen $a_n (\in \{a_j\} \setminus S)$ vorliegt.
 Das bedeutet, es gibt **keine transitiven Abhängigkeiten** [$S \rightarrow a_n \rightarrow a_i$] vom Identischlüssel.

Bild 2.3 -6
Datenbank- und Informationssysteme
Sep-02 (c) Prof. Dr. F. Laux

Das Ziel der Normalisierung ist:

- Redundanzen zu verringern
- Anomalien zu vermeiden.

Betrachten wir das folgende Beispiel:

Relationenschema: Auftrag (Auftr.-Nr., Datum, Kd-Nr, Kd-Name, Kd-Adresse) (wir lassen die Domänen der Einfachheit halber weg und unterstreichen den ID)

Wenn ein Kunde mehrere Aufträge erteilt, muss jedesmal seine Nummer, Name und Adresse erfasst werden.

Diese Redundanz kann zu inkonsistenten Kundendaten führen: nehmen wir an, daß ein Kunde zwei Aufträge erteilt, dann müssen seine Daten 2 mal identisch erfasst werden. Ist dies nicht der Fall (z.B. können unterschiedliche Kundennummern oder Namen für einen Kunden erfasst werden) treten Widersprüche auf (z.B. ein Kunde wird als 2 verschiedene Kunden geführt oder ein Kunde ist unter verschiedenen Namen bekannt).

Wenn alle Aufträge eines Kunden gelöscht werden, verschwinden auch die Kundendaten (Löschanomalie). Es ist unmöglich, einen Kunden zu erfassen, ohne gleichzeitig auch einen Auftrag zu erzeugen.

2.3.4 Relationenalgebra

Mit den folgenden Operationen erzeugt eine Relationenmenge eine Algebra:

Projektion

Vereinigung

*natürlicher Verbund
(natural join)*

Selektion

Differenz

Umbenennung

Bild 2.3 -7
Datenbank- und Informationssysteme
Sep-02 (c) Prof. Dr. F. Laux

Die *Relationenalgebra* (RA) dient als Maßstab für *Datenmanipulationssprachen* (DML). Obwohl sie keine kommerzielle Bedeutung besitzt, werden alle realisierten relationalen Sprachen mit dieser Algebra verglichen. Die RA bietet mengenbasierte Grundoperationen auf Relationen: Vereinigung, Durchschnitt (Selektion), Projektion, Differenz und Verbundoperation.

Der *natürliche Verbund* (*natural join*) ist eine Teilmenge des kartesischen Produktes zweier Relationen wobei Tupel mit korrespondierenden Attributen verknüpft werden.

Beispiel für einen natürlichen Verbund:

Seien folgende Relationen gegeben:

Auftrag :

<u>(AuftrNr,</u>	Datum,	Artikel,	Kd_Nr)
1	12/10/95	Welle	101
2	14/10/95	Schraube	101
3	16/10/95	Stahl ST37	102

und

Kunde :

<u>(Kd_Nr,</u>	Name)
101	Miterrand
102	Chirac

Der *natürliche Verbund* über den Fremdschlüssel Kd_Nr lautet dann:

<u>(AuftrNr ,</u>	Datum,	Artikel,	Kd_Nr	Name)
1	12/10/95	Welle	101	Miterrand
2	14/10/95	Schraube	101	Miterrand
3	16/10/95	Stahl ST37	102	Chirac



2.3.5 Regeln für das Relationenmodell

- ◆ **Identschlüssel :** Zu jeder Relation muß ein Identschlüssel (**ID**) existieren
- ◆ **3NF :** Jede Relation muß in die dritte Normalform (3NF) gebracht werden
- ◆ **Fremdschlüssel :** Attribute mit dynamischem Wertebereich sind als Fremdschlüssel zu modellieren
- ◆ **keine Rekursionen :** Es darf keine Beziehungen von Nichtschlüsselattributen zw. Relationen geben
- ◆ **keine Attributstruktur :** Attribute dürfen keine Struktur aufweisen oder überlappen

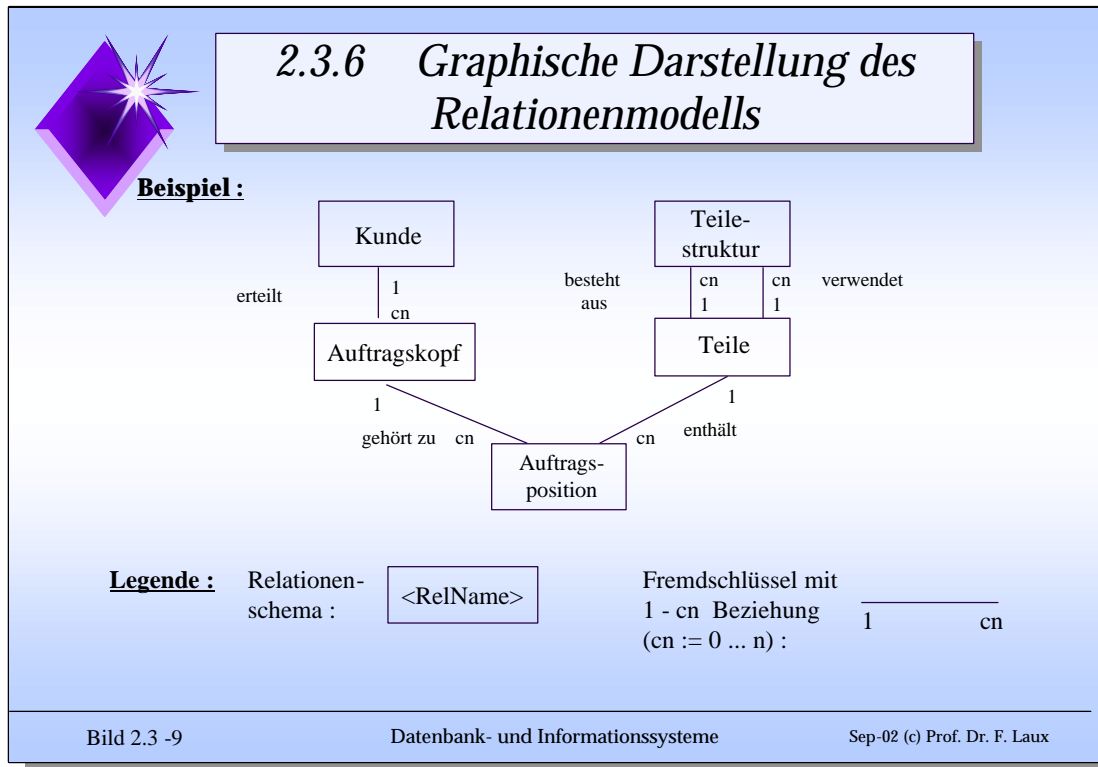
Bild 2.3 -8 Datenbank- und Informationssysteme Sep-02 (c) Prof. Dr. F. Laux

Die Datenmodellierung wird als eine intuitive, kreative Aufgabe angesehen, welche viel Erfahrung erfordert.

Trotzdem gibt es einige Regeln, die uns beim Modellierungsprozess unterstützen:

Der Ident-Schlüssel und andere funktionale Abhängigkeiten (FD) helfen uns bei der Normalisierung. Besonders zu beachten sind transitive Abhängigkeiten, welche uns zur 3. Normalform (3NF) und zu den Fremdschlüsseln führen. Im Sinne einer klaren, stabilen und widerspruchsfreien Datenstruktur sind Rekursionen von Nichtschlüsselementen zu vermeiden.

Wenn ein Attribut eine Unterstruktur aufweist, sollte es in seine atomaren Komponenten zerlegt werden, denn sonst bleibt die Semantik dieser Struktur dem Modell verborgen.



Eine graphische Darstellung hilft uns, die Datenstrukturen leichter zu erfassen und gibt uns einen Überblick über das Datenmodell.

Wir verwenden das **Bachman Diagramm**, um Relationen und ihre Beziehungen graphisch darzustellen:

Ein Rechteck steht für ein Relationenschema und eine Verbindungslinie zwischen zwei Schemata zeigt eine Fremdschlüsselbeziehung an. Wir verwenden Pfeile oder Kardinalitätssymbole *1 - cn* um die Mächtigkeit (Anzahl Verbindungen der Tupel untereinander) der Beziehung anzuzeigen.

Unser Beispiel zeigt die Struktur eines Auftragsverwaltungssystems.

Ein Kunde kann mehrere Aufträge erteilen. Jeder Auftrag enthält eine oder mehrere Positionen, wobei jede Position genau einen Artikel umfasst. Artikel - außer Einzelteile - besitzen eine baumartige (hierarchische) Komponentenstruktur. Umgekehrt kann ein Teil in anderen Artikeln enthalten sein (Verwendungsnachweis).

Wie man leicht sieht, wird eine viele-zu-viele (n:m) Beziehung zwischen Auftrag und Artikel durch zwei 1:cn Beziehungen (gehört zu, enthält) und eine Beziehungsrelation (Auftragsposition) hergestellt. Diese Hilfskonstruktion ist nötig, da wir im relationalen Modell nur 1:cn Beziehungen darstellen können.

ÜA: Welche Beziehungstypen kennen Sie? (Hinweis: siehe 1.4, 2.2.3)