

Informationsmanagement

Σ Alphabet = Menge von Symbolen/Zeichen

$s = s_1 \dots s_n \in \Sigma^*$ Zeichenkette

Problem-
Lösung

Daten - ausgetauschte Nachrichten

Interpretation
Zeichenketten nach Grammatik gebildet

Informationen + Interpretationsbezug

Wissen
Semantik für Empfänger

Wissen - Gesamtheit Kenntnisse, Fähigkeiten zur Problemlösung

+ Vernetzung zu bekannten Infos

Strukturierungsgrad

strukturierte Daten

gleichartige Struktur, Datenmodell folgend

z.B. Listen, Diagramme,

Tabelle, Java Objekte

semistrukturierte Daten

ohne festes Datenmodell

aber Daten implizieren Struktur oder erweiterbares Datenmodell

z.B. XML-Dokument

Object Exchange Model

„unstrukturierte Daten“

ohne formalisierte (inhaltliche) Struktur

z.B. natürlichsprachlicher Text

Bild, Musik (schwer für Computer)

(Text Mining, Natural Language Processing)

Inf. Management

Daten organisieren, strukturieren, speichern, abfragen, löschen, ändern, pflegen, interpretieren, vernetzen, analysieren, verifizieren

Inf. System Anforderungen

Informationsenthalt
nötige Datei ablesbar

Konsistenzherhaltung
nur vernünftige Daten speichern
möglichst Redundanzfrei
Platz sparen, Anomalien vermeiden

Entwurfsprozess

1. Anforderungsanalyse

Fachterminologie, Inf. bedarf

⇒ Informelle Dokumentation z.B. Interviews, Texte

4. Logischer Entwurf inkl. spezifisch, benötigt unabhängig

konzeptionelle Modelle → Konzepte eingesetztes Inf. System

⇒ Logische Datenmodelle z.B. relationales Modell

Wie?

2. Konzeptioneller Entwurf / design

Abstrakte Modellierung Domäne

Objekte, Beziehungen?

⇒ formale Beschreibung z.B. ER-Modell

Welche Daten?

UML-Strukturm.

5. Datendefinition

Deklaration / Programmierung Datenmodell in Inf. System

⇒ Definitionssprachen z.B. SQL

6. Physischer Entwurf

Zugriffs-, Speicherstrukturen auf disk

3. Verteilungsentwurf / distributed system design

Fragmentieren Daten,

Sync & Replica

7. Implementierung & Wartung

Interne Verarbeitung Anfragen
Anfrageoptimierung

Relationale DB's

Backend Unternehmensdaten
Web-Bereich

Android, Mac OS, SQLite, ...
SQL wichtiger Skill (noch Python, R)

(2) Konzeptionelle Datenmodellierung (Externe Ebene DB)

Anforderungen → Konzep. Datenmodell

Konzeptionelle Datenmodellierung

Objekte, Merkmale, Beziehungen, Regeln?

Entity-Relationship (ER) Modelle

statische Eigenschaften (keine Funktionen/ Methoden)

beliebt \Rightarrow „Laien“ verständlich
intuitiv

graphische Notation nach Chen

Entitätsktyp

Menge Entitäten gleichen Typs
Individualisierbare Objekte
Instanzen, E-Typ

Entitätsktyp

Rechteck

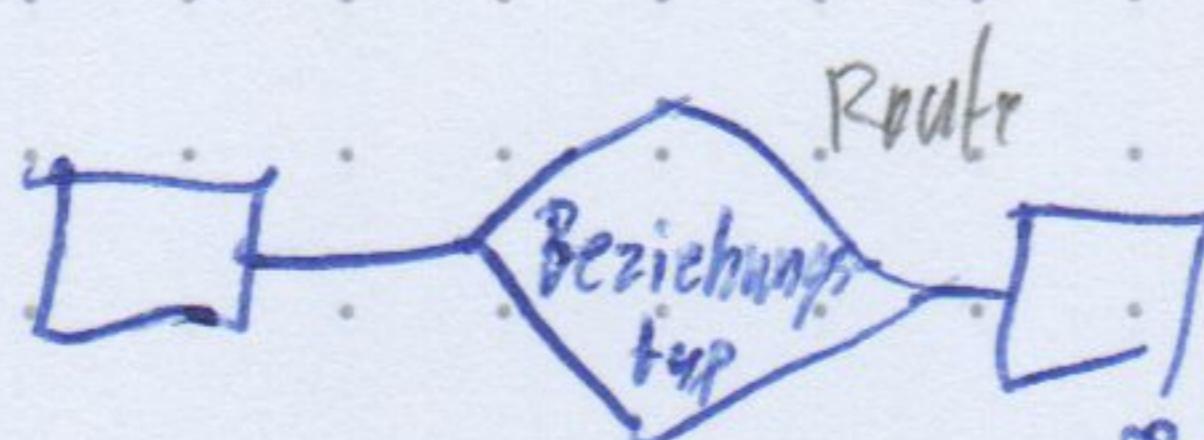
Entitätsktyp (Attributes, ...)

Beziehungsktyp

zwischen ≥ 2 Entitätsktypen

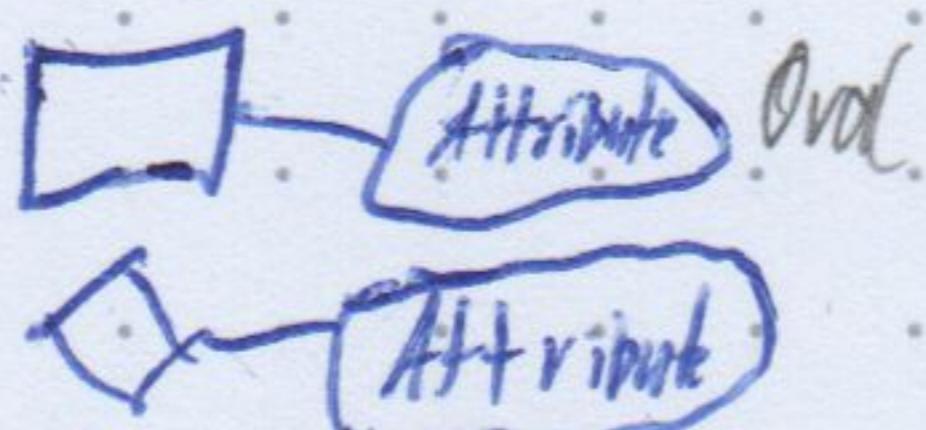
Relation $E_1 \times \dots \times E_n$

Beziehung kontakt zwischen Entitäten



Beziehungsktyp ($E_1, E_2, \dots; \text{Attributes}, \dots$)

Attribut / Merkmale



Schlüsselattribute

\subseteq Attribute Entitätsktyp.

identifiziert Entität eindeutig

Modellierungsentscheidung

ggf. künstlich z.B. ID

nu Entitätsktypen

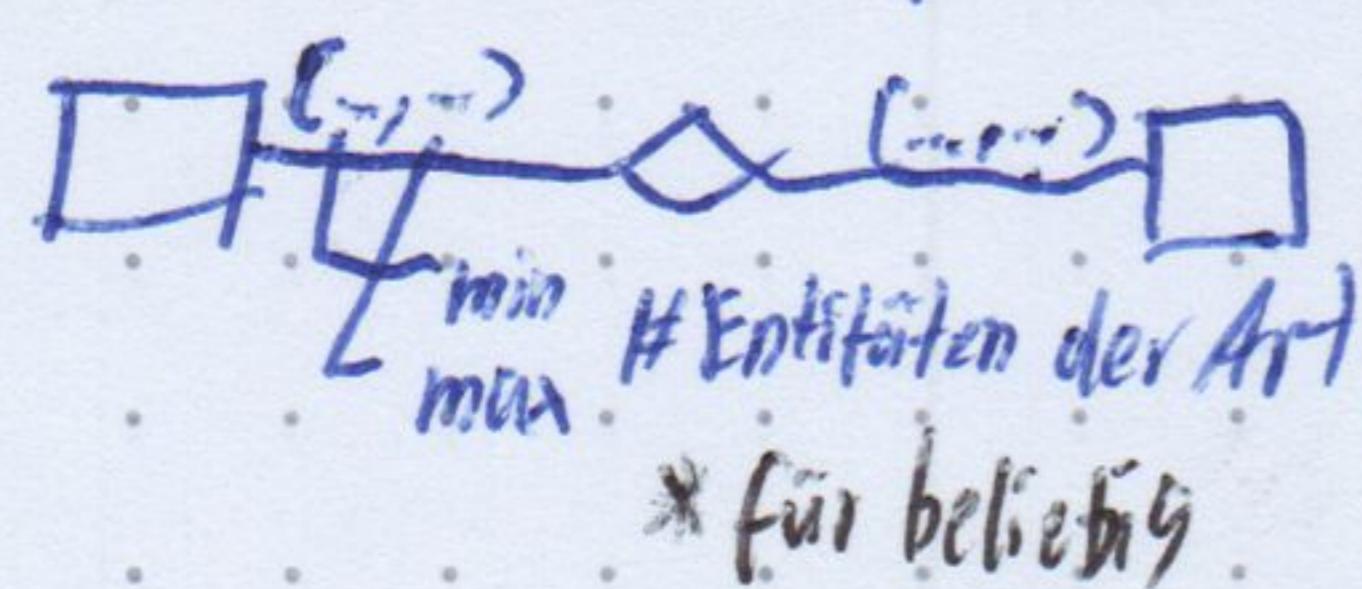
Unterstreichen

Attribut



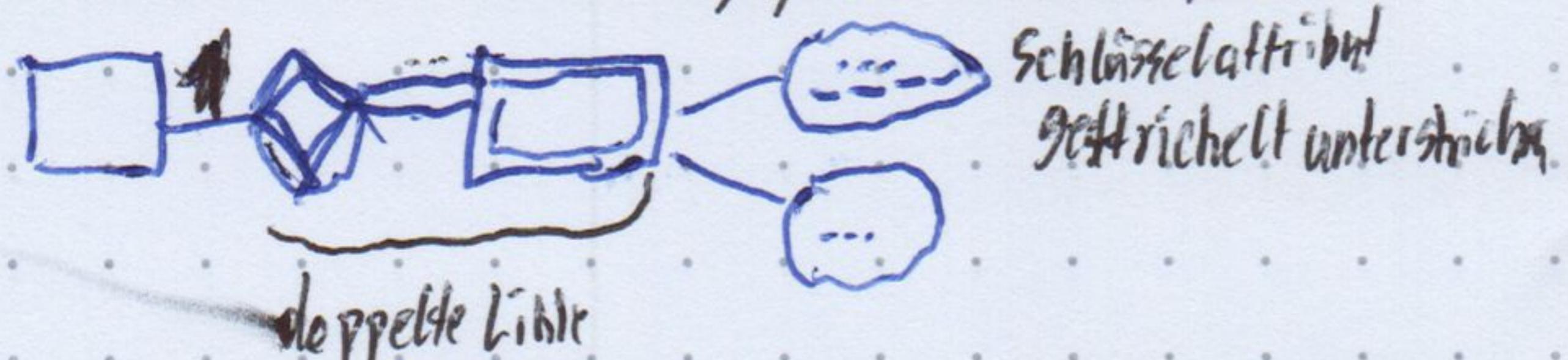
z.B. 1:1
1:N
N:1
N:M

\blacktriangleleft (min, max) - Notation präziser als Kardinalitäten



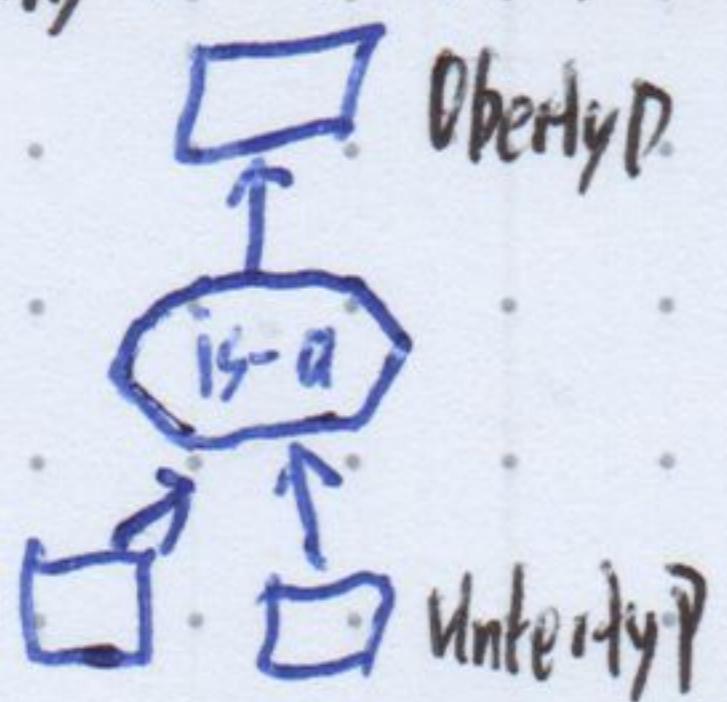
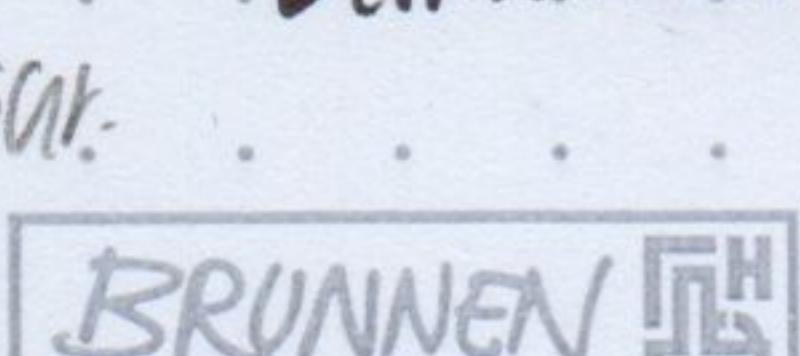
$\blacktriangleleft 1:N \Rightarrow (0,*) \blacktriangleright (1,1)$

Schwache (Existenzabhängige) Entitätsktypen



nicht
klassif.
relevant

Generalisierung / Vererbung der Attribute



Aggregation (Teil-Ganzes Beziehung)



Relationales Datenmodell beim logischen Entwurf

aus Relationen \rightarrow Tabelle

\hookrightarrow keine Reihenfolge Attribute
Tupel

$D = \{D_1, \dots, D_n\}$ Menge Datentypen/Domänen

$R = \{A_1, \dots, A_n\}$ Schema
 $\text{dom}: R \rightarrow D$

$r(R) \subseteq \text{dom}(A_1) \times \dots \times \text{dom}(A_n)$
Relation-Ausprägung von R

$t \in r$ Tupel

$t(A_i)$ Wert Attribut
 $t(\alpha)$ Teiltupel $\alpha \in R$

$R(A_1:MT, A_2:MT, \dots :MT)$
 \uparrow
PK

FKs implizit durch Namensgleichheit

null Wert möglich.

Schlüssel Attributmenge K $\left\{ \begin{array}{l} \text{einfach } |K|=1 \\ \text{zusammengesetzt } |K| > 1 \\ \text{trivial } K=R \end{array} \right.$

Superschlüssel

$\forall t_1, t_2 \in r: t_1(K) = t_2(K) \Rightarrow t_1 = t_2$

Schlüsselkandidat / Schlüssel

K minimal $\Leftrightarrow \forall K' \subseteq K: K'$ schlüsselleidentisch

$K_R = \{E_1, \dots, E_n\}$

Varianten

Primärschlüssel (PK)

1 gewählter Schlüsselkandidat

$PK_R = \{E_1\}$

Attribut

Fremdschlüssel

$K \subseteq \text{Schema } S$ bezüglich Schema R $\Leftrightarrow p(K) = \text{PK von } R$

Referenz

Referenzierte Wert erlaubt

Referenzielle Integrität

$\forall t, t_S \in S: t_E(K) = t_R(p(K))$

$FK_R = \{D \rightarrow S.D\}$

$A \rightarrow EA$

ER \rightarrow Relationales Modell

Regel 1

EntitätsTyp \rightarrow Schema

\hookrightarrow Schlüssel \Rightarrow PK falls minimal
sonst neu

Regel 4

Schwache Entitäten

Fremdschlüssel für starke
Teil PK

Regel 2

Beziehungstyp \rightarrow Schema

$R(E_1, \dots, E_m)$ $R(A_1, \dots, A_m)$

PKs aller Attribut

von E. mit Kardinalität
 $N_1 \neq$ falls keiner einer
mit 1

Fremdschlüssel E_j Eigene Attr

\hookrightarrow Benennung

f Name der referenzierten/falls eindeutig
künstlicher Name z.B. EntitätsnameAttrName
Rollenname falls vorhanden

Regel 3

Zusammenfassen Schemata mit gleichen PKs

Ausnahme: schlecht wenn Ergebnis dünn besetzt

Relationale Algebra Abfragesprache

RWK - Relationenwertebereichkalkül

RTK - Relationaltupelkalkül

RA - Relationenalgebra

) gleiche Aussagekraft

Keine Duplikate

Codd's Algebra

Selection

$\sigma_p(S)$

Konstantenselektion
Attributselektion
Logische Verknüpfung \vee, \wedge, \neg

$\{1=5, 3 \leq 3\}$

Attribut

Konstante

Attribut

Mengenoperatoren bei Vereinigungsvereträglichkeit
auf Tupelmengen
gleicher Grad = # Attribute
gleiche Domänen

Vereinigung

$R \cup S$

Durchschnitt

$R \cap S = R - (R - S)$

Differenz

$R - S$

Projektion

$\Pi_{\alpha}(S)$

$L \subseteq R$

Kartesisches Produkt

Erweiterungen

Join \bowtie natürlich $\Leftrightarrow \bowtie_p$ allgemein
'vereinigt gleichnamige' $\left[\begin{array}{l} R_1 = \dots \\ \vdots \\ R_n = \dots \\ \text{Attr. bleiben doppelt} \end{array} \right]$

Join

$R \bowtie S = \text{in beiden drin}$

Full Outer Join

$R \bowtie L = \text{fehlende Tupel Attr. Null}$

Left

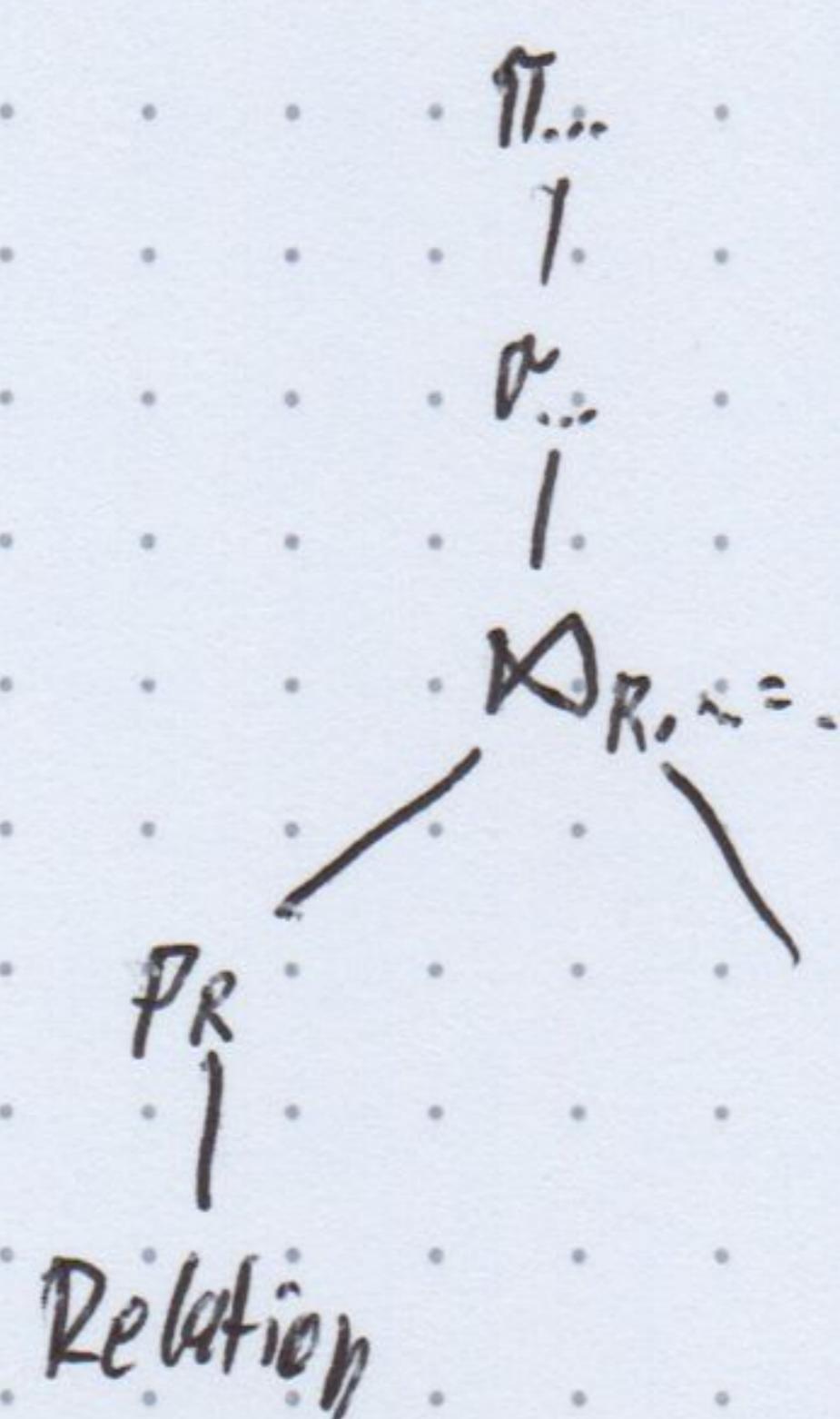
~~$R \bowtie L$~~ $R \bowtie L = \text{alle aus } R$

Right

$R \bowtie L = \text{alle aus } L$

Left / Right Semi-Join

$R \bowtie L = \Pi_R(R \bowtie L)$



Relationale Division

$R \div S = \{ t \in R-S \mid \forall s \in S : \exists r \in R : r \cdot s = t \}$

$= \Pi_{(R-S)}(R) - \Pi_{(R-S)}((\Pi_{(R-S)}(R) \times S) - R)$

gruppieren nach R-S alle R-S zurückgeben für die Einträge für R-S entfernen

Um benennung

Relation $P_{R1}(S) = V_1$

Attribute $P_{R1, A_1, \dots, A_n}(S)$

BRUNNEN

SQL - Structured English Query Language

Relationale Algebra
auf Multimengen
deklarative Anfragesprache für rel. DB

DQL - Data Query Language
DDL " Definition "
DML " Manipulation "
DCL " Control "

DQL

select [distinct] * | Attr₁, Attr₂, ...

from Table₁, ...

[where Prädikat]

↳ kombinieren AND, OR, NOT

↳ Vergleichop. =, <, >, ?, !=

↳ between ... and ...

case (when ... then ...) * [else ... end]

[natural join Table]^{*}

[cross join Table on Table.Attr = Table.Attr]^{*} Reihenfolge kann durch Klammern geändert werden

[inner], [left|right|full] outer join

Pattern Matching Strings ... like ' ... '

↳ beliebiges Zeichen

↳ % Zeichenkette auch €

[group by Attr₁, ...]

[having ...]

// Aggregationsfunktionen

count, sum, avg, max, min

])

// Relationale Algebra

Π select (* or having (X group by (or where (from))))

alle in select, having, in group by oder Aggregate

[order by Attr₁ [asc|desc], ...]

↳ ascending
default asc.

[limit n] ; MySQL, PostgreSQL; select top n : SQL Server; fetch first n rows only : Oracle

Umbenennung von Tabellen in from [as] newName
Attr. in select [as] newName

Unterabfragen (select...) in Äußerem Select, wenn nur ein Aggregat in select Unterabfrage alt Wert

↳ innere

from als Tabelle

where als Wert

als Tabelle

↳ EXISTS T - T nicht leer

s IN T

s OR ALL T

ANY

cast (value/attr as type)

Relationale Division select (R-S) from R group by (R-S) having count(*) = (select count(*) from S)

NULL-Werte

vergleichsoperator unterscheiden \Rightarrow min 1 Operand Null

zusätzliche Ops

dreiwertige Logik and, or, not \Rightarrow unterscheiden außer Kurzschlusslogik

... is NULL

Aggregatfunc. ignorieren NULL

... is not NULL

Gruppierung NULL eigene Gruppe

Recent Directions

Data Lake - Tabellen ohne Schema damit automated Join Detection

NL-TD-SQL Natural Language \xrightarrow{LLM} SQL start

Relationale Entwurfstheorie

Redundanzen

- Speicherplatzverschwendun

kann
 \Rightarrow potentiell Anomalien \Rightarrow Inkonsistenz
Update nicht alle Stellen eines Tupels benötigt
Einfüge benötigt erl. Dummy daten
Löschen betrifft z.B. Best. löschen Lieferanten mit,
wenn in keinen anderen Best.

$\beta \subseteq R$ von $\alpha \subseteq R$ funktional abhängig (FD-funktional abhängig) (α Determinante für β)
 $\Leftrightarrow \forall t_1, t_2 \in r(R): t_1(\alpha) = t_2(\alpha) \Rightarrow t_1(\beta) = t_2(\beta)$ $\alpha \rightarrow \beta$
für jeden Wert von α genau ein Wert von β

voll FD $\Leftrightarrow \alpha \rightarrow \beta \wedge \nexists \alpha' (\alpha \subset \alpha' \wedge \alpha' \rightarrow \beta)$
partiell " \exists

β transitiv abhängig von α
 $\Leftrightarrow \exists r: r$ voll FD von $\alpha \wedge \beta$ voll FD von r
 $\alpha \rightarrow r \rightarrow \beta$

Attributhülle $\alpha^+ =$ Menge von α funkt. abhängige Attribute
Algorithmus
 $\alpha^+ = \alpha$
while Änderungen
 foreach FDA-Y
 if $\beta \subset \alpha^+$ then $\alpha^+ = \alpha^+ \vee Y$

Schlüsselkandidaten aus FDs

1. Attr. nicht auf rechter Seite FDs
 \Rightarrow muss in Schlüsselk. vorkommen
2. Attributhüllen der Determinanten bestimmen
 \Rightarrow Liste Superschlüssel bestimmen und miteinander fechten
3. Minimieren Superschlüssel

Schlüsselmenge

$KEYS(R) =$ Menge Schlüsselkandidaten

alternative Schlüssel nicht PK Schlüsselkandidaten

Primattribut $\exists K \subseteq Keys(R): A \in K \rightarrow$ Teil eines Schlüssels

NPA - Nichtprimattribut -

NPA(R) - Menge dieser

Normalisierung Prozess schlechter Rel. Modell \rightarrow gutes durch Zerlegung

Anforderungen

1. redundanzfreie Relationen
2. große Relation verlustfrei ableitbar

Kriterien korrekte Zerlegung

1. Verlustlosigkeit

$$r = r_1 \bowtie r_2$$

2. Abhängigkeitsbewahrung

1 NF - Erste Normalform

\Leftrightarrow V Attribute Werte atomar

L nicht mehrwertig von Interpretation

L z.B. Listen abhängig

\Rightarrow teilen

VI

2 NF

\Leftrightarrow in 1 NF $\wedge \forall A \in \text{NPA}(R) : \forall K \in \text{Keys}(R) : A$ voll abhängig von K

\Rightarrow Intuitive Zerlegung

VII

3 NF

\Leftrightarrow in 2 NF $\wedge \forall A \in \text{NPA}(R) : \forall K \in \text{Keys}(R) : \exists Y \in \text{Schema}(R) : K \rightarrow Y \rightarrow A$ (transitive FD)

\Rightarrow Intuitive Zerlegung

Synthesealgorithmus R mit FD's \mapsto gutes 3 NF mit Anf.

1. Bestimme „kanonische Überdeckung“ aller FDs in R

Linksreduktion

Rechtsreduktion

Entfernung FDs der Form $x \rightarrow \emptyset$

Zusammenfassen gleicher linker Seiten

2. $\forall \alpha \rightarrow \beta : \text{Erstelle } R := \underbrace{\dots}_{R^*} \alpha \cup \beta$

3. Schlüssel

4. Entferne R wenn $\exists R' : R \subseteq R'$

Anfrageverarbeitung

Anfrage $\text{select } A_1 \dots \text{ from } R_1 \dots \text{ where } P$

J.A.-Compiler

Nicht-optimierter
Logischer Plan

R_1

J.A.-Optimierer

Optimierter
Physischer Plan

Aufführung

Optimierung

Regel-basiert Datenunabhängig
auf algebraischer Ebene

e.g.

1. Selektions-Aufspaltung

$$\text{Konjunktion } \sigma_{C_1, \dots, C_n}(R) = \sigma_{C_1}(\sigma_{C_2}(\dots(R)\dots))$$

weitere Äquivalenzen

$\wedge, \vee, \neg, \Delta$ kommutativ

$\wedge, \vee, \neg, \Delta$ einzeln assoziativ

σ -distributiv mit $\vee, \neg, -$

2. Selections-Pushdown

$$\text{Kommutativ } \sigma_{C_1}(\sigma_{C_2}(R)) = \sigma_{C_2}(\sigma_{C_1}(R))$$

$$\pi_{A_1, \dots, A_n}(\sigma_C(R)) = \sigma_C(\pi_{A_1, \dots, A_n}(R)) \text{ falls es sich nur auf } A_i \text{ bezieht}$$

$$\sigma_C(R \Delta S) = \sigma_C(R) \Delta S$$

3. Join Ersetzung

$$\sigma_C(R \times S) = R \Delta_C S$$

4. Projektions-Pushdown

$$\pi_{L_1}(R) = \pi_{L_1}(\pi_{L_2}(\dots(R)\dots)) \text{ falls } L_1 \subseteq L_2 \subseteq \dots$$

$$\pi(\sigma_C(R)) = \sigma_C(\pi(R)) \text{ falls } \sigma \text{ noch möglich}$$

$$\pi_A(R \Delta S) = \pi_A(R) \Delta \pi_A(S) \text{ falls } A = LNR$$

$$\pi(R \times S) = \pi(R) \vee \pi(S) \text{ } \pi \text{-distributiv mit } \vee$$

5. Zusammenfassen Operationsfolgen gleicher Op.

Kostenbasierte Optimierung Datenabhängig

Schätzung durch Kostenmodell

1. Auswahl Join-Reihenfolge

Schätzung Zwischenergebnis-Größen

BottomUp | Tabellen-Größen aus DB-Katalog

Schätzung \uparrow Annahme Gleichverteilte Attrib. Werte

$$\text{sel}(\sigma_p(R)) = |A_p(R)| / |R|$$

$$\text{sel}(\sigma_{A=\text{konstant}}(R)) = 1 / |A|$$

$$\text{sel}(R \Delta S) = |R \Delta S| / |R \times S|$$

$$\text{sel}(\sigma_{A,B}(R)) = 1 / \max(|A|, |B|)$$

2. Auswahl phy. Operatoren

1 über Selektivität der Operatoren ≥ 0.5

Stichprobenverfahren

Selektivität für Stichprobe berechnen

Histogramm für Attrib. Werte im Wertebereich

entl. Ausreißer extra; manuelle Änderungen

multidimensionale Hist. für zusammengesetzte o. Red.

$$P_1 \wedge P_2 \uparrow \approx \text{sel}(\sigma_{p_1}(R)) + \text{sel}(\sigma_{p_2}(R))$$

$$\text{or } P_1 + P_2 - (P_1 \cdot P_2)$$

$$\text{sel}(R \Delta \text{frenschlüssig } S) = |S| / |R \times S|$$

$$\text{sel}(R \Delta_{A=B} S) = 1 / \max(|A|, |B|)$$

Suchbaum

Binärbäum mit $n+1$ Blättern = Catalan-Zahl $C_n = (2n)! / ((n+1)! \cdot n!)$

Unterschiedliche

BRUNNEN

→ Klassischer Ansatz: Dynamische Programmierung

BottomUp in Phasen je Optimum für Pläne mit k Relationen

Pruning schlecht

9

Informations Management

Data: organise, structure, query
store, create, modify, delete
interpret, link, analyse, verify

Data Growth of unstructured Data

Science
Business
Industry
World Wide Web

containing implicit information
e.g. for Recommender System, best products for supermarket
predict subset of items R ⊂ interesting for user

Machine Learning - ML
tasks experience, (feedback)
⇒ perform well

Document Classification
document of topics \mapsto
 \Rightarrow $t \in T$ describing d

NLP - Natural Language Processing | Human Language Technology focus
Computational Linguistics tech

NLU - understanding transform language \rightarrow machine representation
G - generation

Machine Translation - MT
text in ~~lang~~ Lang A
 \Rightarrow text in B with same meaning

Information Retrieval - IR
find subset of document collection D
maximally relevant to query

Text Generators
prompt \Rightarrow continuing/fulfilling text

Ambiguity

Homograph \leftrightarrow Homophone
same writing voice / Aussprache
diff meaning
Adverb \leftrightarrow Adjective
Pragmatics
utterance

System Evaluation

intrinsic - in isolation, compare to predefined expectation (gold standard)

extrinsic - in use in larger system/experimental setup

Natural Language Text

Script Σ

Symbol/Character $s \in \Sigma$

String $s = s_1 \dots s_n \in \Sigma^*$

Message string following grammar

Data exchanged messages (not interpreted)

Language \leftrightarrow Script

{ vocabulary

} "alphabet"

grammar

alphabetic Latin, ...

Logographic / syllabic (Hanzi, ...)

consonant-based (Abjad) Arabic, ...

segment - (Abugida) Devanagari, ...

Writing direction

Text (quasi-) written coherent state of language

| connecting ideas
| \Rightarrow semantic meaning

Properties:

type/genre article, ...

register used lang. variety (dialect, formal, ...)

style how written (funny, polite, ...)

domain/topic

time of writing (lang. changes over time)

Electronic Text Formats

Plain text) mostly unstructured

HTML

) semi-structured

XML

e.g. paragraphs

.odt, .docx

) XML-based formats e.g. ZIPed XML files

.doc, .rtf

) binary content

.pdf

) mixed text, binary

Storing text in Relational DBs

with fixed or maximum length

(CHAR(n), VARCHAR(n) typically ≤ 4000)

arbitrary length

BLOB - Binary Large Object

e.g. for files e.g. images

CLOB - Character Large Object

e.g. text

typically

no sorting/grouping

limited comparison (=, <, >, BETWEEN, ...)

Encoding Issues

default of source files IDE

string types e.g. UTF-16 in Java

files OS

table

database connection

Variable length $\Rightarrow n^{th}$ char

Search

Sort - lang. specific

Character Encoding

functions

$$\text{enc}: \Sigma^* \rightarrow \{0,1\}^*$$

$$\text{dec}: \{0,1\}^* \rightarrow \Sigma^*$$

typically ~~Zentrale~~ character set

$$\Psi: \Sigma \rightarrow \{0,1\}^*$$

$$\text{enc}(c_1 \dots c_n) = \Psi(c_1) \dots \Psi(c_n)$$

Types

$$\text{Single Byte } \Psi: \Sigma \rightarrow \{0,1\}^8$$

e.g. ASCII, ISO 8859-1, CP-1252

$$\text{Multi Byte } \Psi: \Sigma \rightarrow \{0,1\}^{n \geq 8}$$

e.g. UCS-2, UTF-32

Variable Length depending on char and

e.g. UTF-8, UTF-16 predefined rules

Escape-code-based switching between
charsets,
e.g. ISO 2022

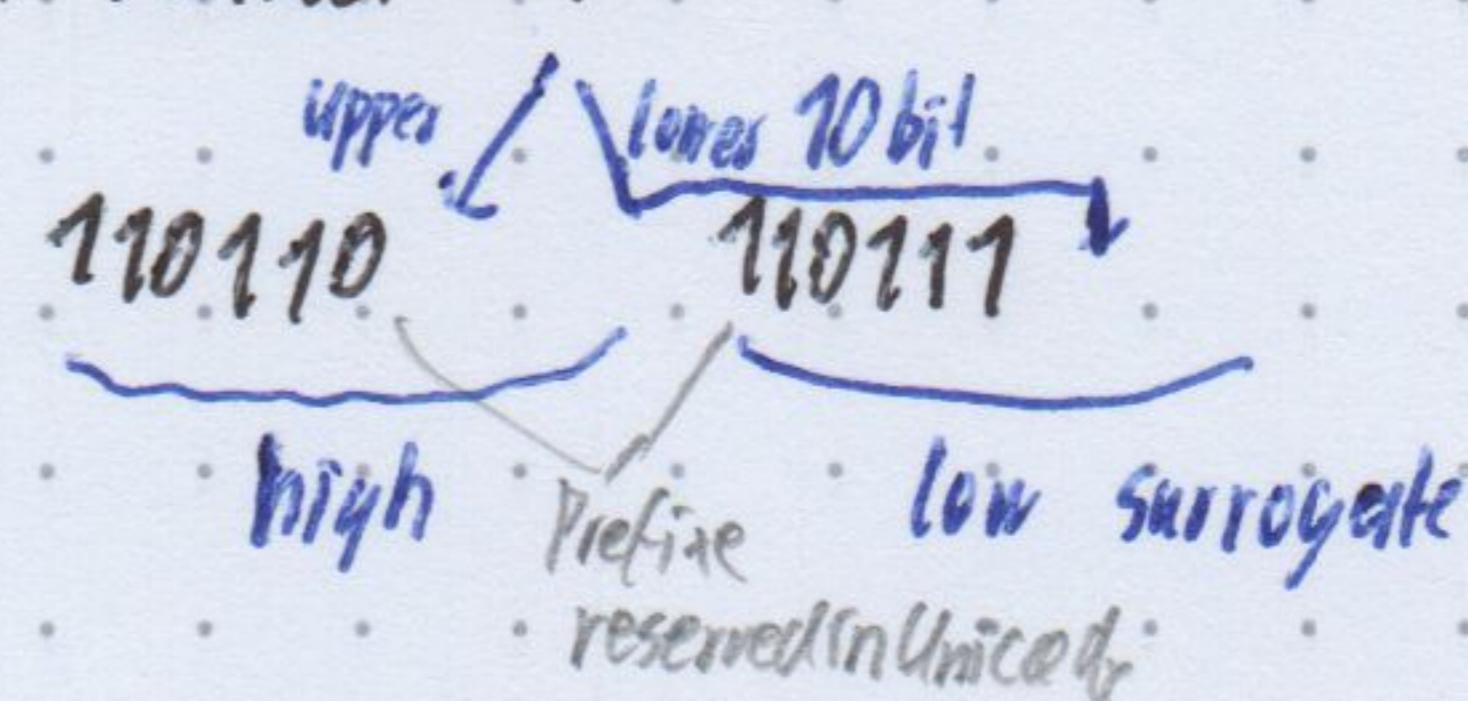
UTF-32 4 byte direct Unicode correspondence - slower, bigger than ISO 8859-1

UCS-2 2 byte only BMP deprecated

UTF-16 2-4 byte

| BMP \Rightarrow direct correspondence

| supplementary Plane \Rightarrow subtract 0x10000



UTF-8 1-4 byte, direct ASCII correspondence

leading 1s followed by 0 = H bytes (0 statt 10)

continuation bytes start with 10

(theoretically allows up to 6 bytes)

ASCII 7 significant bits, fahrende 0

control codes

Sonderzeichen

Zahlen ab 30₁₆

Buchstaben groß 41₁₆

Klein 61₁₆

Single
Byte

ISO 8859-x Single Byte

lower 128 symbols \Rightarrow ASCII

upper \Rightarrow language specific

e.g. x=1 \Rightarrow Western Europe

basically one per major lang.

Unicode idea: 1 charset for all i.e. logographic, fantasy lang.

currency symbols, emojis

4 Byte

21 significant bit

written U+ppxxxx

17 Planes 0x00 to 0x10

PP=0 - BMP - Basic Multilingual Plane

Latin, Arabic, most Chinese, ...

PP=1 - SMP - Supplementary Multilingual Plane

Historic lang., emojis, ...

2 - SIP - Supplementary Ideographic Plane

Extension for Chinese, ...

3-13 empty

14 - SSP - Supplementary Special Purpose

control chars

15-16 - PUA-A/B - Supplementary Private Use Area

private chars of certain font

Linguistic Preprocessing

Segmentation

Tokenization input stream → ordered seq. of tokens

segmented into

split at whitespaces?

cents. \$1.25

Ambiguities

Periods

| sentence termination → separate token

| abbreviations, ordinal numbers, urls

Comma, Whitespace

|

| in number

Single quote

| enclosing quote

| contractions, elisions

Dash

| part of token

| ranges

Colon

| sentence delimiter

| time expr.

↳ inflected word forms

| e.g. Chinese: no spaces

need context info

← Morphology - study of word forms / formation

Morphemes - smallest meaning-bearing units

| Free " / Stamm - useable in isolation

| open class words e.g. nouns, ocat =

| closed conjunctions

| Bound " / Affix - to words, often for inflection

| Suffix

| Prefix e.g. -s für Plural

| Infix fan + bloody + tastie e.g. Fugenlaute

| Circumfix get sagt

Word Formation / -bildung

Derivation

Stamm + derivation affix

Conversion / zero derivation

change part of speech of stamm without affix

Composition / compounding

linking Stämme

Decomposition / unbinding useful for "productive" lang.
e.g. German

Agglutinative lang. combine affixes
e.g. Turkish

Not phonological Normalization

Stemming - strip word endings

same family → similar stemmed repr.

Errors

Understemming related words → diff. stamm

Over diff. words → same stamm

Ambiguity

e.g. Homographs

Lemmatization - undo inflection / base form

typically requires part of speech

deals with irregular forms

Syntax - regularities & constraints: word & phrase structure

Part of Speech Tagging (POS Tagging)

Tagset of word classes

e.g.

Noun	Proposition of by, to
Verb	Preposition
Adjective	Determiner the, a, that
Adverb	

L = lexical class

⇒ valuable info for

word formation, lemmatization

possible neighbors, LLM

prenunciation, speech synthesis

semantic word sense disambiguation

shallow parsing

Ambiguities

Approaches

Rule-based

dictionary: word form → pos tag ~ maybe 90% correct

Rule unknown to most common (noun)

iterative (learned) ordered transformation rules

Probabilistic

estimate $P(\text{word}, \text{context}, \text{pos tag})$

⇒ most probable

Learned from manually labeled training data

Parsing: determine grammatical structure of sentence

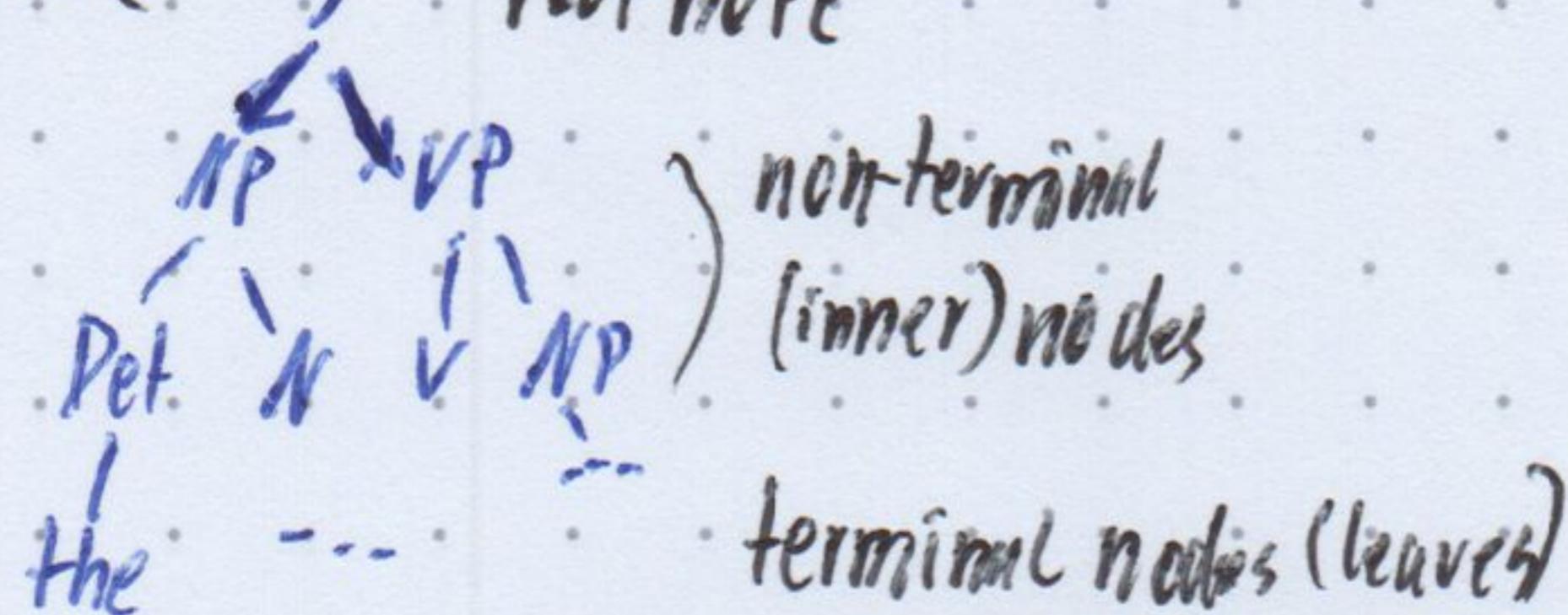
Phrase structure grammar

decompose sentence into constituents

L = wordgroup behaving as single unit mostly Phrases - word seq., yielding syntactic Unit
generally allow:
Substitution e.g. with it
Movement as group in sentence
Coordination group and on
Question yielding constituent

head / "key" word determines syntactic type
Grammatical modifier optional, "modify"
Pre modifier before head meaning
Post " after

Parse tree



Bracketed Notation [S [NP ...] ...]

Parenthesized notation (S (NP ...))

e.g.

NP - Noun Phrase the black cat

PP - Prepositional phrase in love

VP - Verb phrase eat cheese

AP - Adjectival phrase full of toys

AdvP - Adverbial phrase dearly

Syntactic Ambiguity ≥ 1 separated sync. structure

Attachment " - constituent can be added to tree at diff places

Coordination " - varying conjunction scope

Garden Path sentence ⇒ obvious parse doesn't work

e.g. The old man the boat

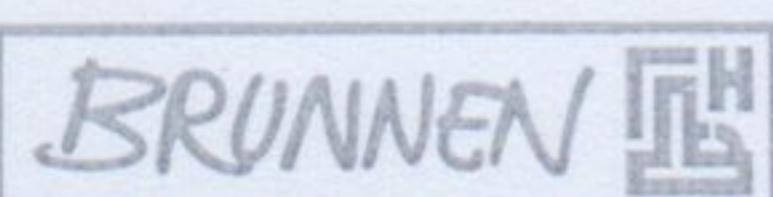
Parsing

sentence → in Lang? /

generateable by grammar? /

(un-)grammatical sentence?

Indicate by * prefix



e.g. CFG - context free grammar

$G = (T, N, S, R)$

terminals, non-t., startsymbol, Production Rules, $N \rightarrow (T \cup N)^*$

14

Semantics - meaning

Lexical " - " of lexical items (words)

Structural " - relationships meaning of lexical items in context

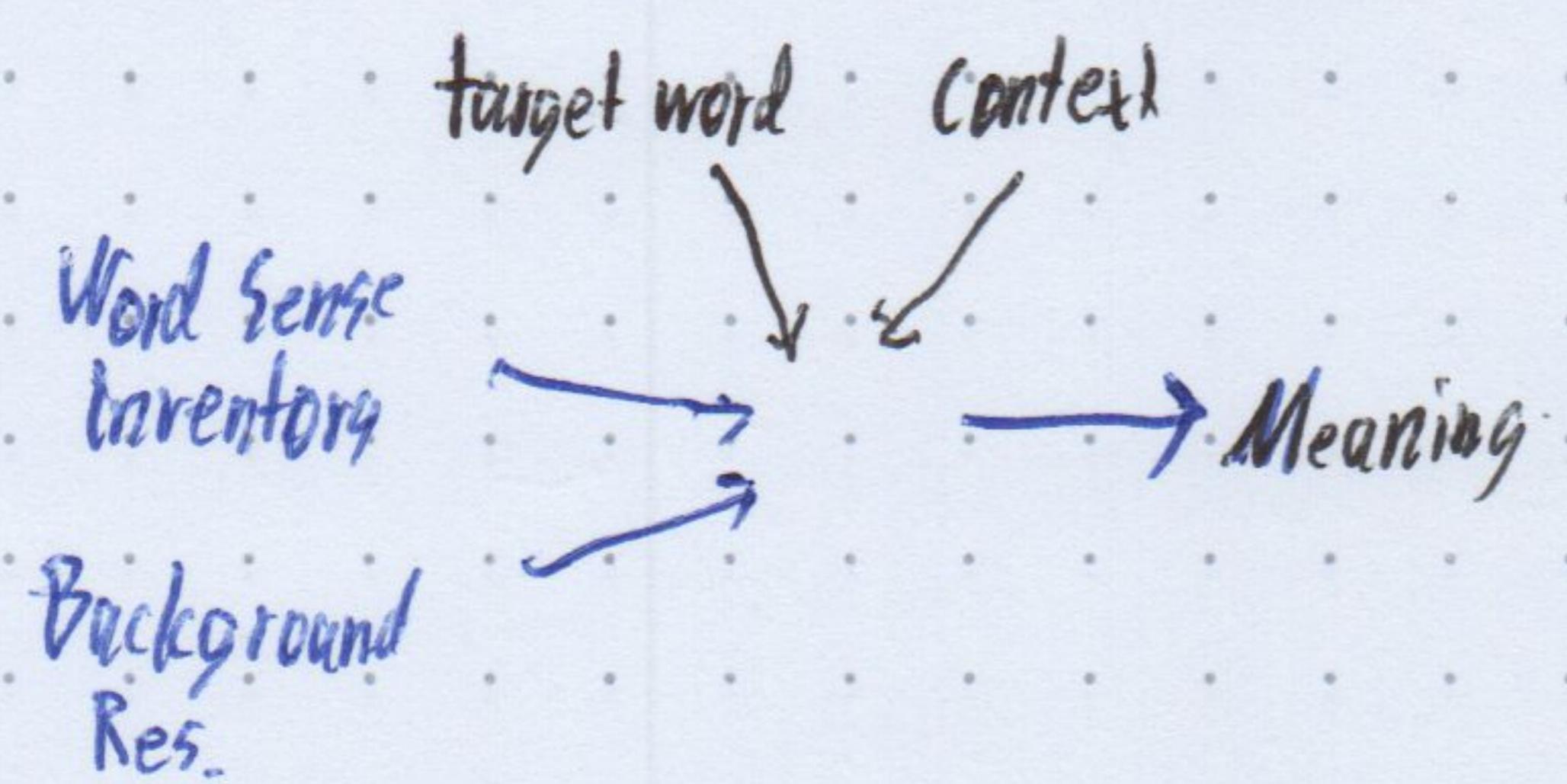
Entities

Affr.

Events

Lexical Ambiguity
syntactic

WSD - Word Sense Disambiguation



Pragmatics purpose/intention of utterance

Text Corpus - Text collection

real word
+ frequency analysis
+ basis for experimental res.

e.g. Brown Corpus - 1961 500 texts POS-tagged

Common Crawl - including copyrighted web

Treebanks - POS, parse tree

Parameters

Language

Monolingual

Multi-

Parallel - translated texts
ideally sentence aligned

Genre / Text type

Domain / Topic

Time of compilation

Size

Static \leftrightarrow Dynamic

Communication

Written Mix
Spoken

Annotation level

Word - POS, lemma, sense - expensive

Phrase entities, multi-word expr - time consuming

Sentence syntactic tree, boundaries

Discourse - co-referential chain
discourse segments

Storing Annotations

Inline - changes original

e.g. likes / POS-tag

Stand-off \rightarrow separate file

e.g. at position ... starts ... of length ...

IOB-Annotation Scheme for spans

| L-Begin
| - Outside
| - Inside & continue
e.g. B O B B I

Corpora in Relational Databases

Words (Word-ID, Word, Frequency)

Word-Index (Word-ID, Sent-ID, Position)

Sentence-Index (Sentence-ID, Sent-ID)

Sentences (Sent-ID, Sentence)

Sources (Sentence-ID, Source)

Bigrams (Word 1, Word 2, Frequency, Significance)

Co-Occurrence within sentence (Word 1, Word 2, Frequency, Significance)

Workflow

- 1) retrieve, store original documents
- 2) convert to plaintext
- 3) segment
- 4) annotations
- 5) format
- 6) analyse/use

Information Retrieval (IR) on unstructured

Information Need / intend
state requiring info
to solve problem

\Rightarrow retrieve relevant documents
RE - Retrieval Engine

\leftrightarrow DBMS on structured data

Query

Rep. info need so interpretable

Relevance of document

to info need

minimal: about topic
useful, interesting, satisfying
user depending
 \Rightarrow no clear cut-off

Core Challenges

results e.g. Millions

Relevance?

Intend? suboptimal query?

Lexical gap! synonyme phrasen

Ambiguity! Homographie

Activities

Indexing - document repr. enhancing search

Searching - interpret query, relevance score, ranking,
improve of user feedback

Boolean Retrieval

- document ~ set of words

Query word - document contains word f_{term}

Index word \mapsto Vektor of documents $\{ \begin{matrix} 1 & \text{in d} \\ 0 & \text{sonst} \end{matrix} \}$
in ~~Term~~ Term-Document Matrix

Δ large, many 0s

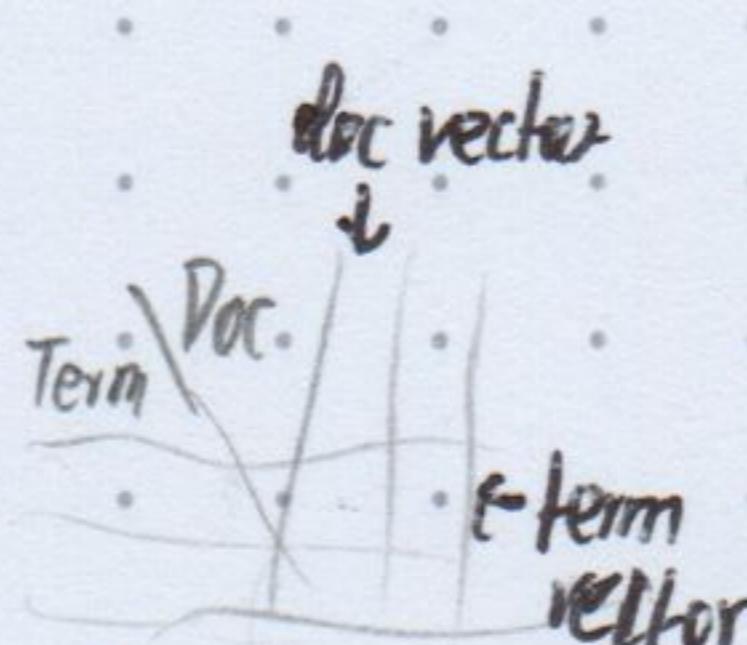
\Rightarrow Inverted Index

word \mapsto List of documents

Dictionary Postings

- user needs syntax

- "feast or famine" - often too few or many results



Search

bitweise Ops res 1 d relevant

0 net

Vector Space Model (VSM)

- document ~ Bag of Words

Multimenge

TF - Term Frequency - # occurrences in document $f_{tf}(t)$

(1) DF - (Inverse) Document Frequency - # docs containing word

query vector in doc vector space

euclidean distance? - depends on vec length

\Rightarrow cosine similarity $\frac{q \cdot r}{\|q\| \cdot \|r\|}$

Vector weights

$$w_t = \begin{cases} 0 & \text{if } f_{tf} = 0 \\ 1 & \text{Bindry} \\ \frac{f_{tf}(t)}{\log(f_{df}(t)) + 1} & \text{TF} \\ \text{mehrlich linear hilfreicher?} \end{cases}$$

$$\underbrace{\log(f_{df}(t)) + 1}_{\text{Normalized TF}}$$

$$\underbrace{1 / \log(f_{df}(t))}_{\text{TF-IDF}}$$

still considers useless determiners (a, the)

Ranked Retrieval

solves feast or famine
user decides # results

/ all \geq threshold

IR Evaluation

$$P = \text{Precision} = \frac{\text{retrieved \& relevant}}{\text{retrieved}} = \frac{\# TP}{\# TP + \# FP}$$

		Relevant	Irrelevant	Σ
True Positive	TP	TP	FP	TP + FP
False Negative	FN	FN	TP	FN + TP
True Negat.	TN	TN	FN	TN + FN

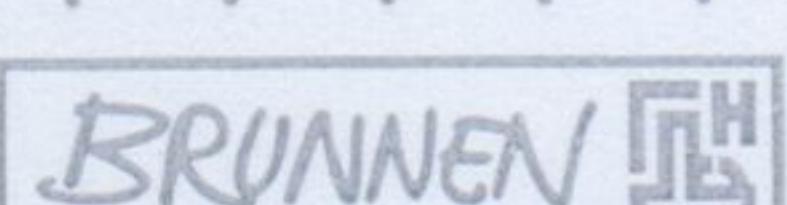
\oplus achieving one is easy

\Rightarrow F₁ score / measure

$$F_1 = \frac{2 \cdot P \cdot R}{P + R}$$

$$F_B = \frac{(1 + B^2) \cdot PR}{(B^2 P) + R}$$

harmonic mean



$$\text{Precision at Rank } P@R = \frac{\# \text{relevant docs before R}}{R}$$

IE - Information Extraction

Given: D document collection

⇒ structured knowledge

/set of structured facts

entities, events
relationships, facts, ...

often dependent
Predefined e.g. SQL
+ supported
query types
directly answer
for user

↔ IR

easier

Domain-independent

Query types usually unconstrained
faster

- less effective for user goal (the answer)

Entity Recognition

Challenges

Entity ↔ No.Entity
mobile phone ↔ Mobil

Coverage issue

complete list of possible Entities?

Variation, e.g. titles, abbreviations
Mr., Mr. Dr., Dr. First, last name
A, ...

Ambiguity

Darmstadt (German ↔ US city)

time dependency

Deutscher Bundeskanzler?

Multi-Word expressions boundaries?

"Biersch-Stiftung an der TU Darmstadt"

Metonymy figure of speech, Entity by associated Thing

Deutschland für → Fußball-Nationalmannschaft

List Lookup Approach

+ simple, fast

+ domain adaptable

- collect/maintain lists

- No Variations

Ambiguity

gazetteers - lists of locations

Rule-based

hand crafted set of regular expressions

+ fairly good performance

- Labour intensive

mid 1990s

⇒ Supervised Machine Learning from (hand-labeled) database

Corpus-based probabilistic models

probability language structures from large corpora

After IE \Rightarrow store in Knowledge Base - database for knowledge management
collect, organize, share, info
search utilize

Knowledge Graphs

fact \sim binary Relationship

e.g. SPO-triple (subject $\xrightarrow{\text{predicate}}$ object)

entities \sim nodes

Science Fiction
Genre

Star Trek

Automatic Classification

data points = instances annotated with class labels
features \sim vectorspace coordinates

\Rightarrow learn, generalise
instances \rightarrow class labels

binary \leftrightarrow multi-class
(decomposable in binary)

single \leftrightarrow multi-label per instance

sequence classification
instance sequence jointly classified

Split the data

Dev-Set

Training-Set

Dev-Test Set - analyse errors
optimize param

Test-Set - no optimization with it
of model

supervised - trained on labeled instances

Training

Neural network

