



Arbeitsbereich NATS

Prof. Menzel

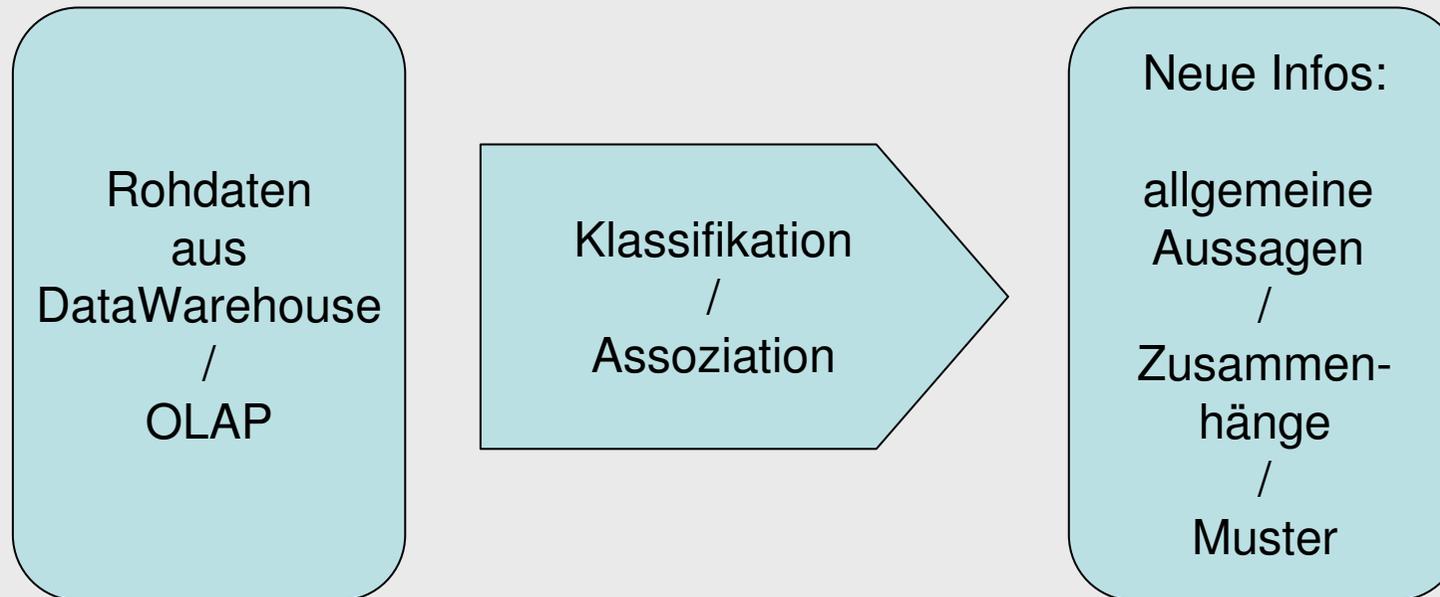
Seminar Data Mining

Datenschutzgerechtes Data Mining

Seminarvortrag von Simon Boese
Student der Wirtschaftsinformatik

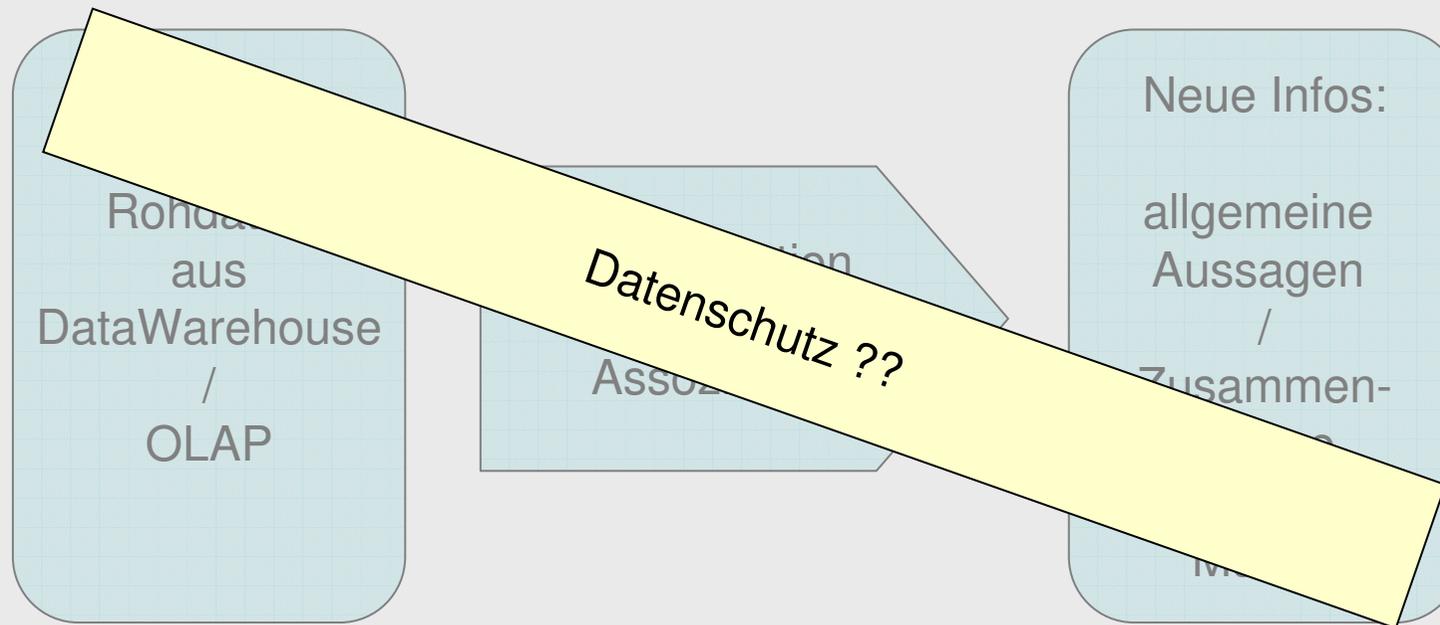


Wissensgewinnung





Wissensgewinnung





Arbeitsbereich NATS

Prof. Menzel

Seminar Data Mining

Agenda

- Ziele des Data Mining und Problemstellung
 - Privatsphäre
 - Was muss sich ändern?
 - Vorgehensweisen
 - Zusammenfassung
-
-



Was soll erreicht werden?

↳ Wissensgenerierung

-
- Individuelle Informationen bilden die Basis
 - ▶ Kundenverhaltensdaten etc.
 - Klassifikation, Regression, Assoziation, Clustering
 - Optimierung der
 - ▶ Vertriebsstruktur
 - ◊ Produktplatzierung
 - ◊ gezielte Werbung
 - ▶ SCM
 - ◊ Produktionsverfahren
 - ◊ Lieferantenbeziehungen
-



...was noch?

- Überwachung
- Sicherheit
- Prävention

**Wer darf
was
und
WOZU
?**



EU-Richtlinie

(95/46/EG)

Richtlinie über „die Verarbeitung personenbezogener Daten und den Schutz der Privatsphäre in der elektronischen Kommunikation“

- Jede Person hat das Recht, sich keiner (teil-) automatisierten Entscheidung zu unterwerfen
 - ▶ Entscheidungen aufgrund Data Mining Resultate

- Ausnahmen für Strafverfolgung und Geheimdienste



Data Mining Reporting Act₁

US-Senat 2003

Definition des Data Mining

Eine Anfrage, Suche oder Analyse auf einer oder mehrerer Datenbanken, wobei

- mind. eine DB von nicht-staatlichen Stellen erzeugt/erhoben oder kontrolliert wurde/wird ODER die ursprüngliche Erhebung durch staatliche Stelle keinen geheimdienstlichen oder strafverfolgungsrelevanten Hintergründe hatte.



Data Mining Reporting Act₂

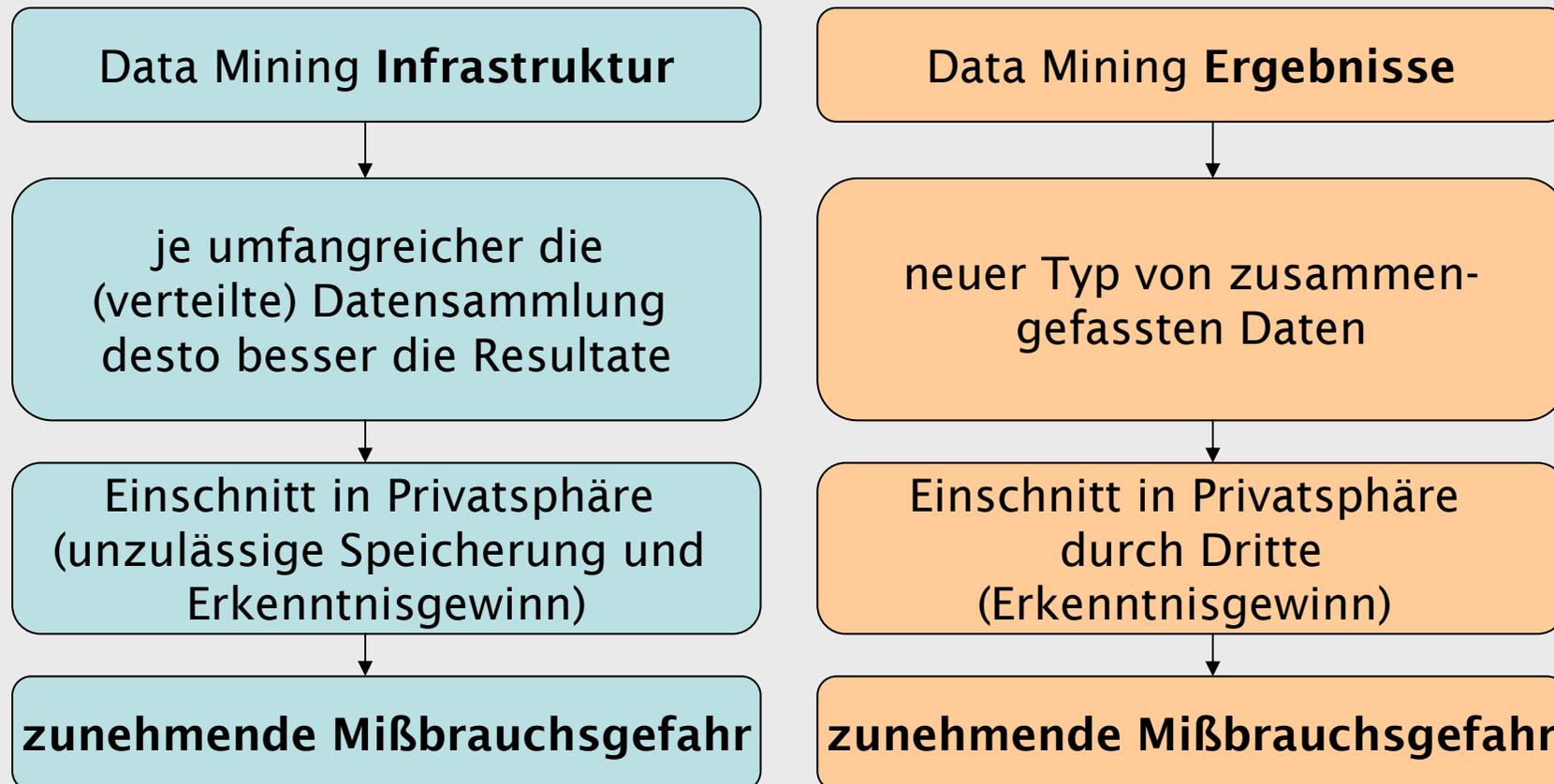
US-Senat 2003

Definition des Data Mining

Eine Anfrage, Suche oder Analyse auf einer oder mehrerer Datenbanken, wobei

- bei der Suche keine personenbezogenen Daten verwendet werden dürfen, um nach Informationen über diese bestimmte Person zu suchen,
- eine staatl. Behörde auch nach Mustern, die auf terroristische oder andere kriminelle Aktivitäten hinweisen, suchen darf.

Problemstellung





Arbeitsbereich NATS

Prof. Menzel

Seminar Data Mining

Agenda

- Ziele des Data Mining und Problemstellung

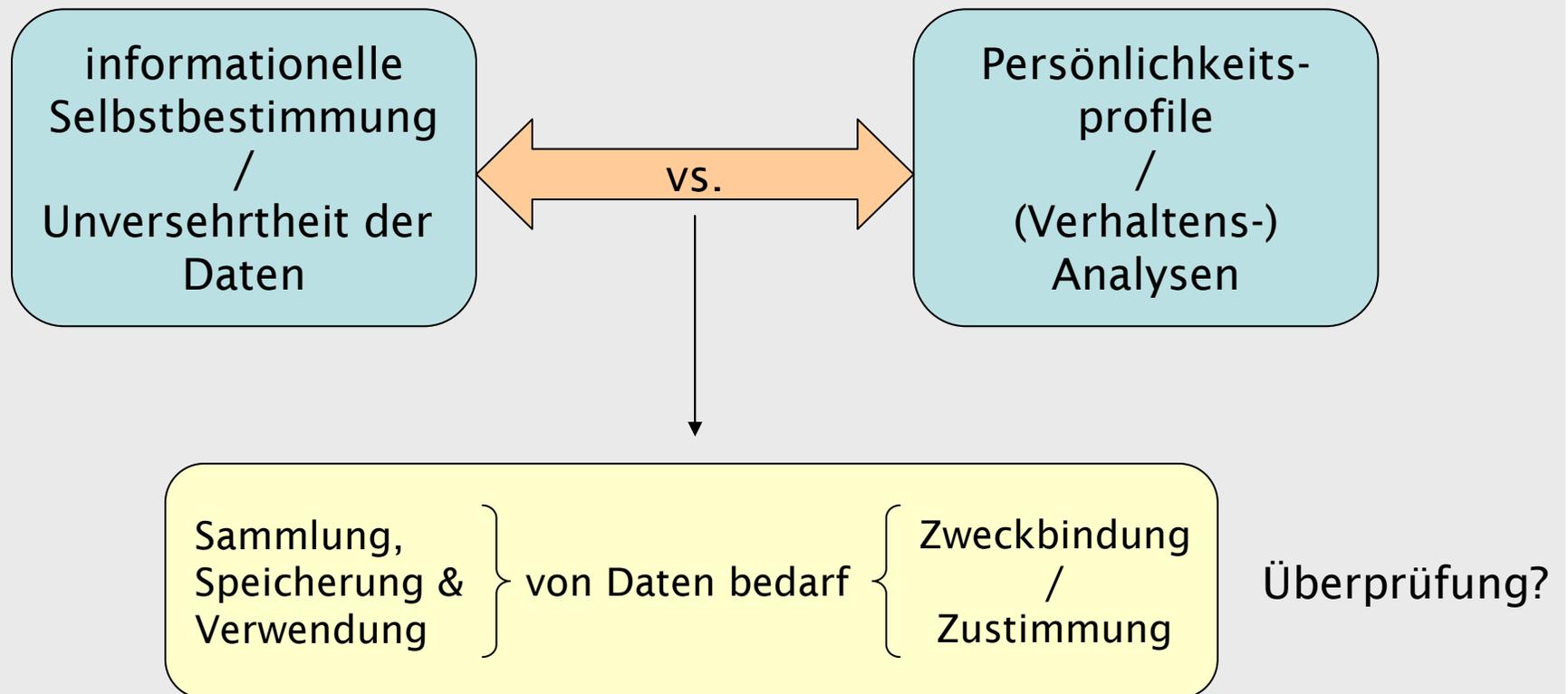
- Privatsphäre

- Was muss sich ändern?

- Vorgehensweisen

- Zusammenfassung

Privatsphäre₁





Privatsphäre₂

Wann ist die Privatsphäre verletzt?

- Solange keine konkrete Person identifiziert werden kann, stellen Data Mining Ergebnisse keinen Eingriff dar.
- Jede Wissenserweiterung (über die Basisdaten hinausgehend) ist bereits eine Störung der Privatsphäre.



Privatsphäre₃

Wissenserweiterung ist jedoch oberstes Ziel!

- Der Erkenntnisgewinn muss bewertet werden und den Möglichkeiten, aus den Daten Rückschlüsse auf Individuen zu ziehen, gegenüber gestellt werden:
 - ▶ Messung der Schwere des Eingriffs
 - ▶ Definition von Akzeptanzwerten



Erhebungsprozess

Die Privatsphäre muss während des gesamten Prozesses geschützt bleiben:

- **Datenerhebung** (Anonymisierung und Pseudonymisierung)
- **Teilnehmende Parteien** (Datensicherheit)
- **Veröffentlichung/Nutzung**



Arbeitsbereich NATS

Prof. Menzel

Seminar Data Mining

Agenda

- Ziele des Data Mining und Problemstellung
 - Privatsphäre
 - Was muss sich ändern?
 - Vorgehensweisen
 - Zusammenfassung
-
-



Anforderungen₁

Datenschutzgerechtes Data Mining muss zu-
sichern, dass die offen gelegten Daten...

- nicht zu einem Individuum zurückverfolgt werden können
- keine Verletzung der Privatsphäre darstellen



Anforderungen₂

Datenschutzgerechtes Data Mining muss zu-
sichern, dass die offen gelegten Daten...

- das Risiko einer Identifikation minimieren
 - auch nicht in Kombination mit anderen verfügbaren Daten zu einer Verletzung der Privatsphäre führen
-
-



Arbeitsbereich NATS

Prof. Menzel

Seminar Data Mining

Agenda

- Ziele des Data Mining und Problemstellung
 - Privatsphäre
 - Was muss sich ändern?
 - Vorgehensweisen
 - Zusammenfassung
-
-



k-Anonymität₁

-
- Grundgedanke: Bildung einer Gruppe von k Datensätzen
 - ▶ Maximal: Identifikation einer Gruppe
 - Quasi-Identifikatoren (QI): alle Informationen, die in Kombination mit regulär zugänglichen Daten die Identifikation eines Einzelnen ermöglichen
-



k-Anonymität₂

Definition:

Kein Datensatz darf in seinen QI einzigartig sein. Es müssen immer mind. k Datensätze die selben QI aufweisen



k-Anonymität₃

- Um die QI bereinigten Daten, machen die Identifikation eines Individuums unmöglich
 - Unsicherheit bei Informationen über eine konkrete Person
 - Rückschlüsse nur auf eine Gruppe
-
-



Verschleierung

- Anhand einer bestimmten Verteilung die Originaldaten verzerren.
 - Statistisch ungenau machen.
 - Verwendbarkeit der Ergebnisse?
 - ▶ Verzerrung mittelt sich möglicherweise heraus.
 - ▶ Rückrechnung schwierig, falls überhaupt die Methode der Verschleierung bekannt ist.
-



Schwellenwerte und Generalisierung₁

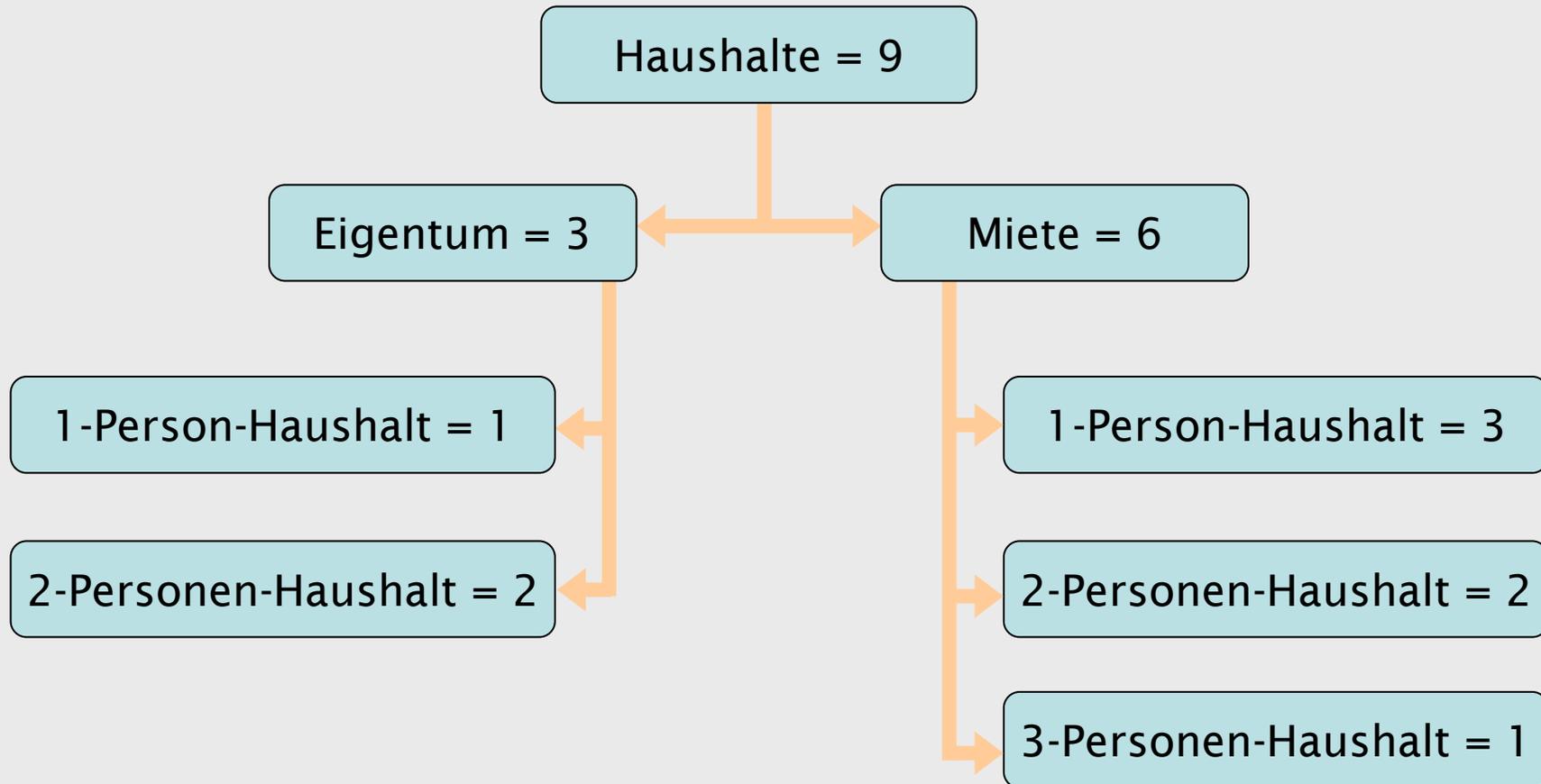
- Schwellenwerte:
 - ▶ Kriterien erst ab gewisser Anzahl an Treffern veröffentlichen

 - Generalisierung:
 - ▶ weniger Details, Zusammenfassung

 - „beschneiden“ des Entscheidungsbaumes

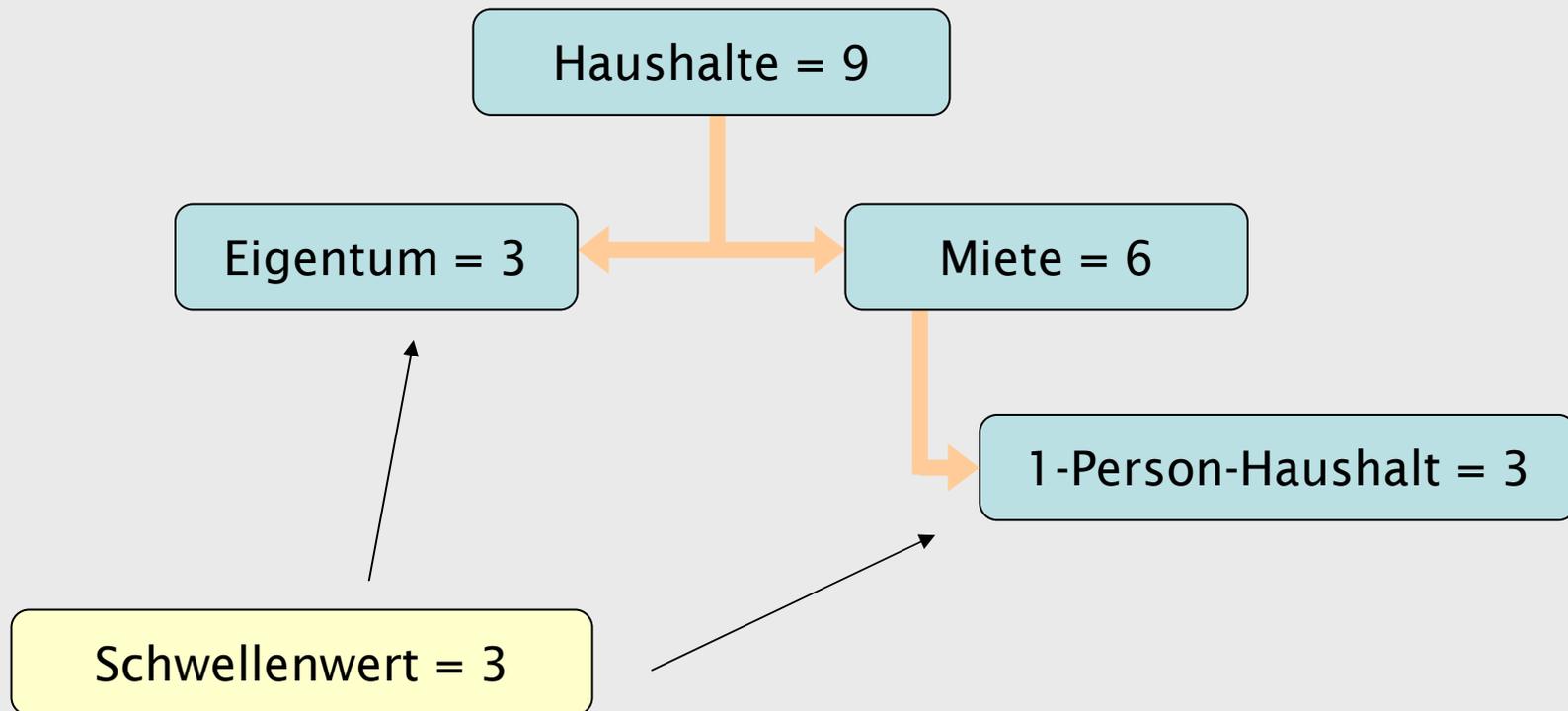
 - können auch bei der k-Anonymität verwendet werden
-

Schwellenwerte und Generalisierung₂ – Bsp.



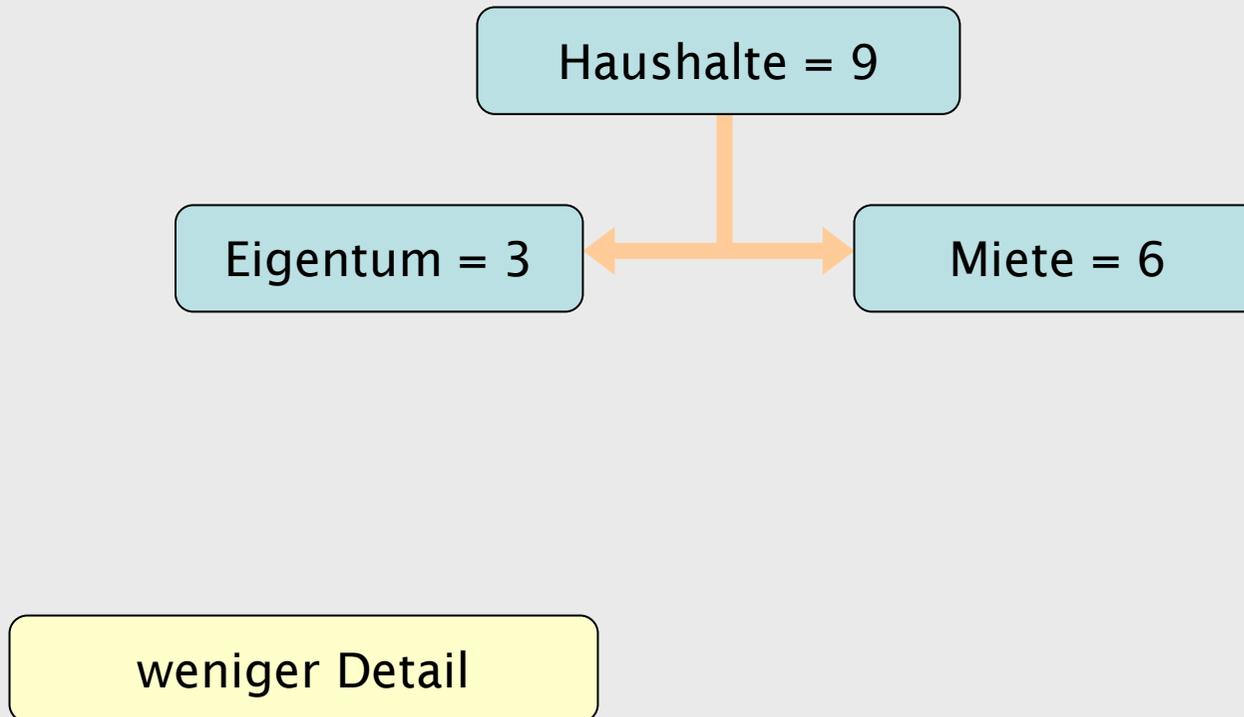


Beispiel: Schwellenwerte





Beispiel: Generalisierung





Sicherheit

-
- Bisheriger Fokus: Data Mining Ergebnisse
 - ▶ Verletzung der Privatsphäre einzelner Individuen bei Veröffentlichung

 - Neuer Fokus: Data Mining Prozesse
 - ▶ „collaboration meets competition“
 - Datensicherheit
 - Kryptographie



Verteilung der Daten₁

- Wann muss überhaupt datenschutzgerechtes Data Mining eingesetzt werden?
 - ▶ Wenn Daten zentralisiert bzw. verteilt vorliegen, aber nur von einer Stelle kontrolliert, besteht keine zusätzliche Gefahr einer Verletzung der Privatsphäre.
 - ↳ anderer Ansatz



Verteilung der Daten₂

- Wann muss überhaupt datenschutzgerechtes Data Mining eingesetzt werden?
 - ▶ Liegen Daten verteilt vor und werden von unterschiedlichen Stellen kontrolliert, könnten durch Kooperation dieser Stellen sensible Daten offen gelegt werden.



Verteilung der Daten₃

An k verschiedenen Orten P werden in einer Datensammlung D bestehend aus den Spaltenbezeichnungen I und den Entitäten E Informationen gespeichert:

$$D_k = (E_k, I_k)$$

- Horizontal/Homogen verteilte Daten
- Vertikal/Heterogen verteilte Daten

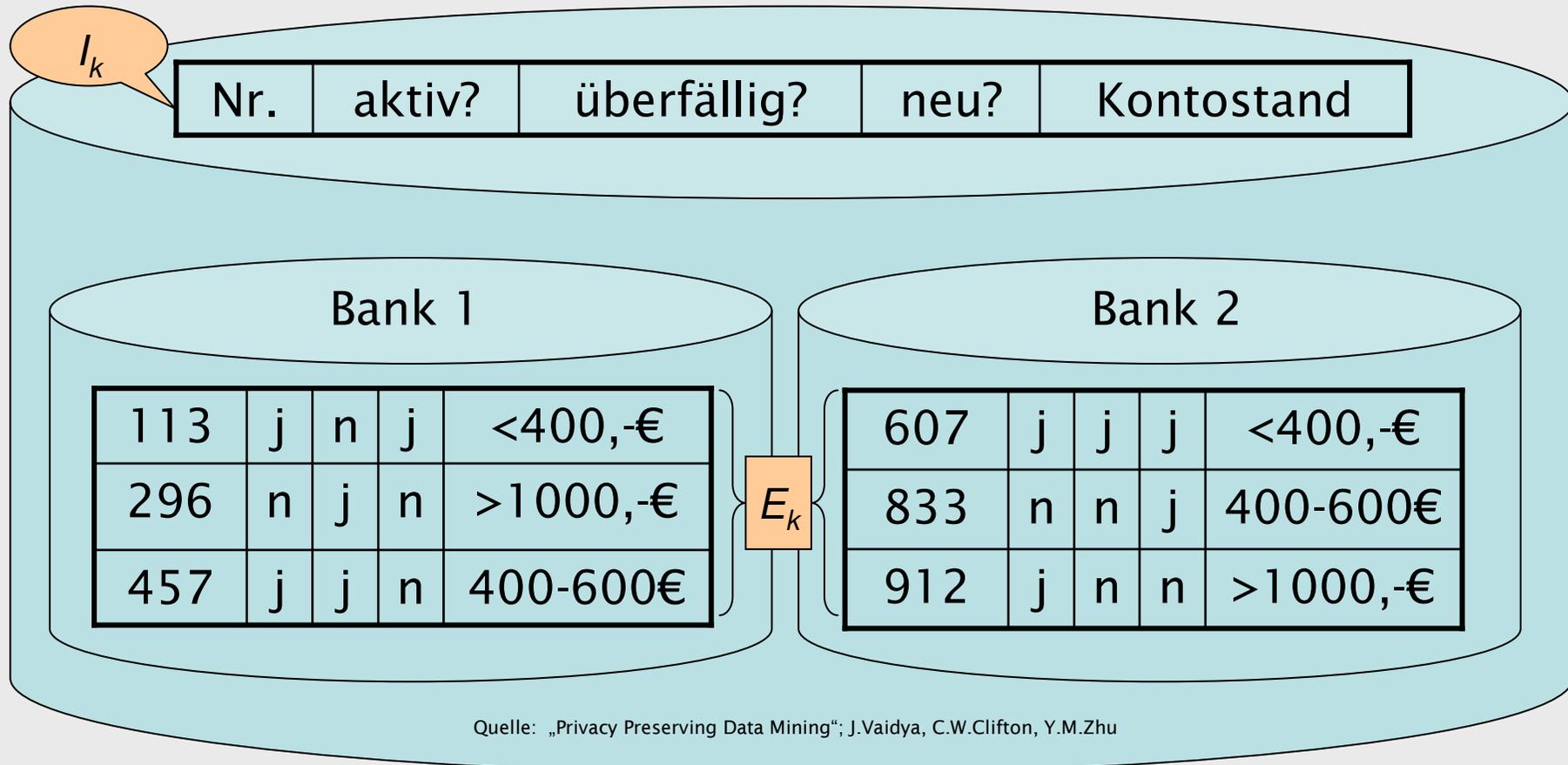
Horizontale Datenverteilung₁

- an verschiedenen Orten sind gleiche Daten über verschiedene Objekte vorhanden
- Datenmenge erhöhen, statistische Genauigkeit steigt

$$E_G = \bigcup_i E_i = E_1 \cup \dots \cup E_k$$

$$I_G = \bigcap_i I_i = I_1 \cap \dots \cap I_k$$

Horizontale Datenverteilung₂





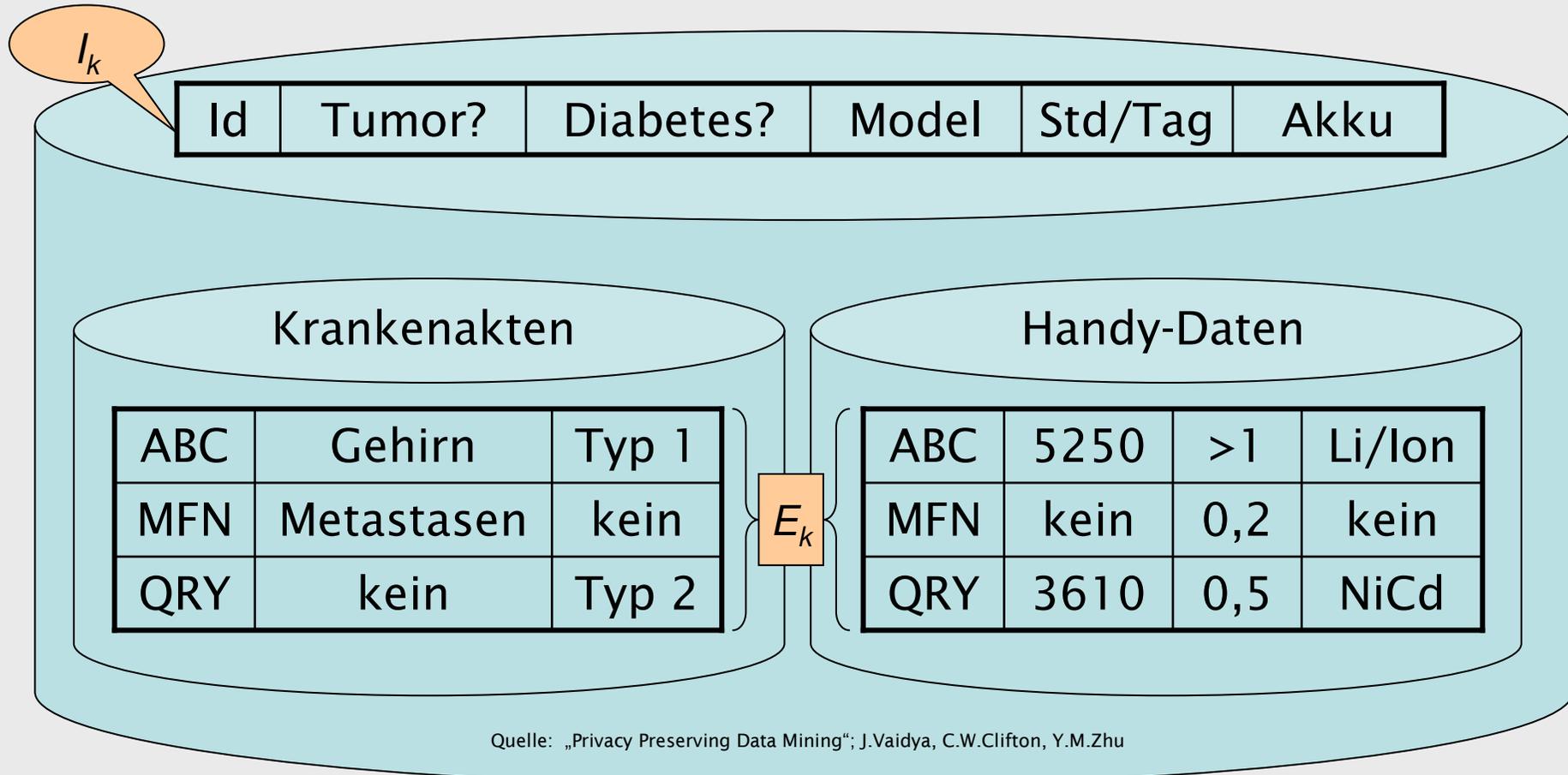
Vertikale Datenverteilung₁

- an verschiedenen Orten liegen unterschiedliche Daten über gleiche Objekte
- Datenumfang erhöhen, globale Zusammenhänge erkennen

$$E_G = \bigcap_i E_i = E_1 \cap \dots \cap E_k$$

$$I_G = \bigcup_i I_i = I_1 \cup \dots \cup I_k$$

Vertikale Datenverteilung₂





Secure Multi-party Computation (SMC)

- Generelles Problem bei der Berechnung von Funktionen, wenn die Input-Daten verteilt – vertikal/horizontal – vorliegen
 - Kommunikation via Protokoll
 - Ziel: Außer dem Ergebnis erfährt keine Partei etwas Neues
 - Erfolgt eine Berechnung sicher, dann war sie auch **privat!** (Goldreich et al. basierend auf Yaos Millionärs Protokoll)
-



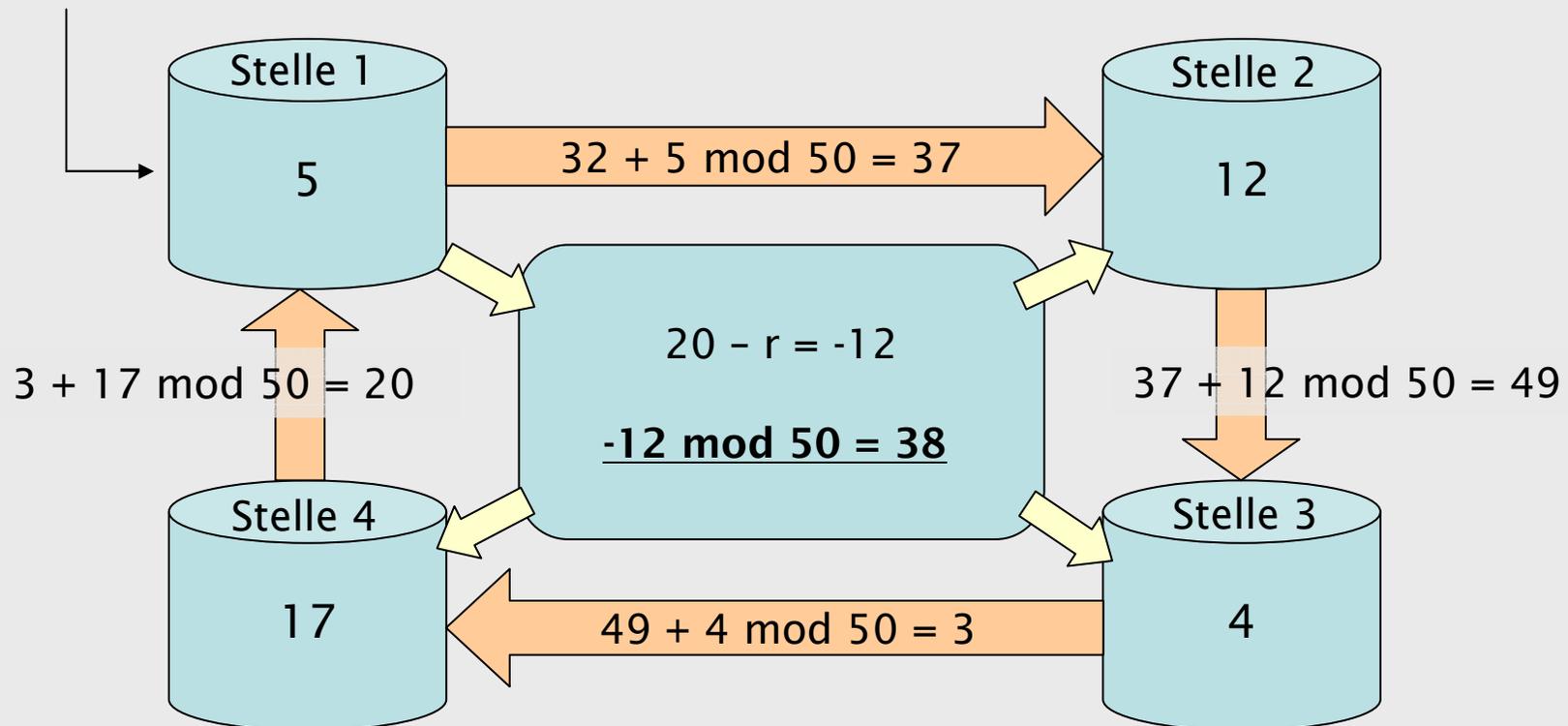
SMC – Beispiel₁

- Berechnung einer Summe
 - Jede (beteiligte) Partei hält einen Wert, der nicht bekannt gegeben werden soll
 - Annahme: der Wertebereich Ω der zu berechnenden Summe ist bekannt
 - gleichverteilte Zufallszahl $r \in \Omega$ als Verschleierung
-



SMC – Beispiel₂

$|\Omega| = 50$ $r = 32$



Quelle: „Privacy Preserving Data Mining“; J.Vaidya, C.W.Clifton, Y.M.Zhu



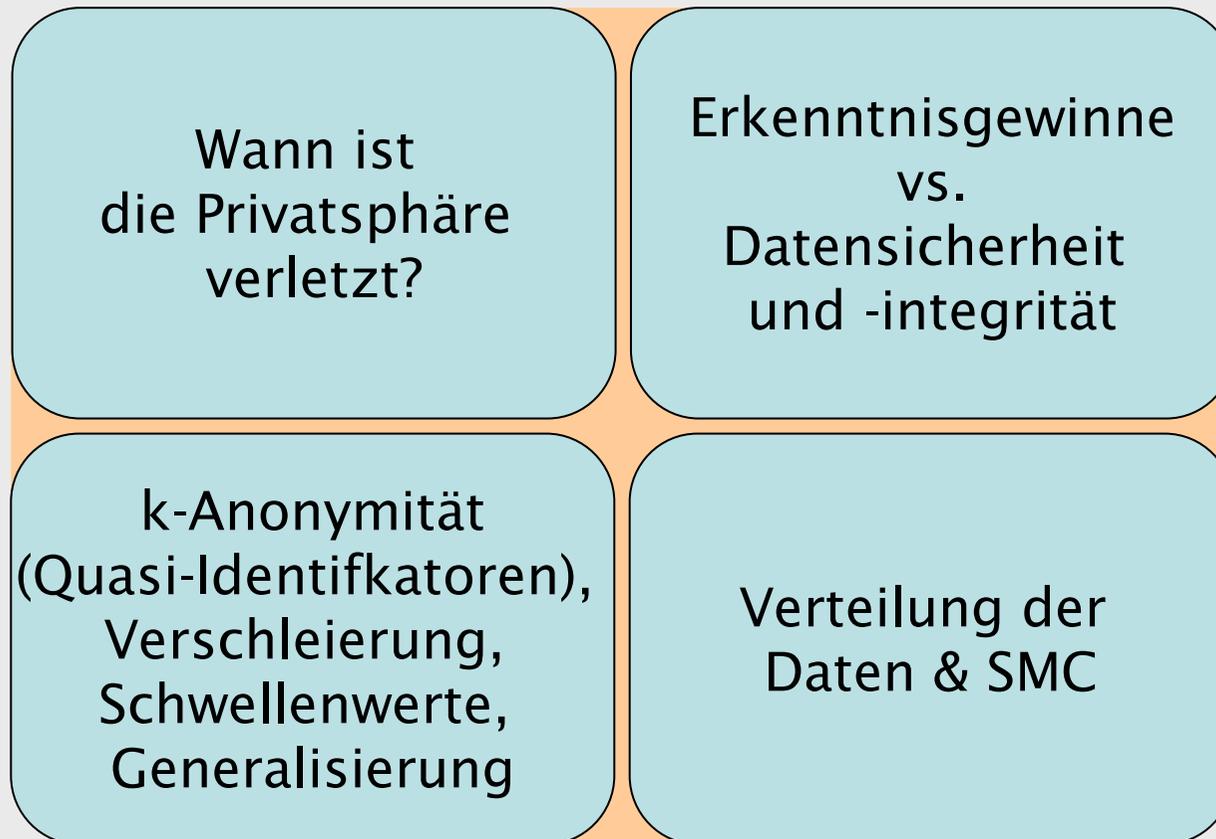
Arbeitsbereich NATS

Prof. Menzel

Seminar Data Mining

Agenda

-
- Ziele des Data Mining und Problemstellung
 - Privatsphäre
 - Was muss sich ändern?
 - Vorgehensweisen
- **Zusammenfassung**





Arbeitsbereich NATS

Prof. Menzel

Seminar Data Mining

Literatur

- J. Vaidya, Ch.W. Clifton, Y.M. Zhu;
„Privacy Preserving Data Mining“;
Springer Science+Business Media;
New York 2006



Arbeitsbereich NATS

Prof. Menzel

Seminar Data Mining

Ende

Fragen , Anregungen und/oder Kritik?

Vielen Dank für Ihre/Eure Aufmerksamkeit.