

# Anonymisierungsverfahren

Seminar - Statistische Herausforderungen im Umgang  
mit fehlenden bzw. fehlerbehafteten Daten

**Ye Bin Park**

**Betreuer:** Prof. Dr. Thomas Augustin

**Institut für Statistik LMU**

**19.Dezember 2014**

# Gliederung

## 1. Einführung

### 1.1. Definitionen

### 1.2. Gründe & Ziele der Anonymisierung

### 1.3. Stufen der Anonymisierung

### 1.4. Möglichkeit der Identifizierung

## 2. Anonymisierungsverfahren

### 2.1. Verfahren zur Informationsreduktion

- merkmalssträgerbezogene Verfahren
- merkmalsbezogene Verfahren
- musprägungsbezogene Verfahren

### 2.2. Datenverändernde Verfahren

## 3. Auswahl der Verfahren

# Gliederung

## 1. Einführung

### 1.1. Definitionen

### 1.2. Gründe & Ziele der Anonymisierung

### 1.3. Stufen der Anonymisierung

### 1.4. Möglichkeit der Identifizierung

## 2. Anonymisierungsverfahren

### 2.1. Verfahren zur Informationsreduktion

- merkmalssträgerbezogene Verfahren
- merkmalsbezogene Verfahren
- musprägungsbezogene Verfahren

### 2.2. Datenverändernde Verfahren

## 3. Auswahl der Verfahren

# Definitionen

**Anonymität** ist der Zustand, wenn eine Person oder eine Gruppe nicht identifiziert werden kann.

- Ist gegeben, wenn die vorliegenden Daten nicht zur Gewinnung von Informationen über die einzelnen statistischen Objekte dienen können.

**Anonymisierung** ist das Verändern personenbezogener Daten derart, dass die Einzelangaben über persönliche oder sachliche Verhältnisse nicht mehr oder nur mit einem unverhältnismäßig großen Aufwand an Zeit, Kosten und Arbeitskraft einer bestimmten oder bestimmbaren natürlichen Person zugeordnet werden können .

- Bundesdatenschutzgesetz §3 Abs.6

# Definitionen

**Zusatzwissen** sind zusätzliche Informationen, die entweder vorliegen oder beschaffbar sind, zur Hilfe der De-Anonymisierung.

**Überschneidungsmerkmal** sind Merkmale, die im Zusatzwissen und im Datensatz vorkommen bzw. enthalten sind.

**De-Anonymisierung bzw. Re-Identifikation** ist die gezielte Aufhebung vorher durchgeführter Anonymisierung von Daten.

Mithilfe von Überschneidungsmerkmalen werden Zusatzwissen und Datensatz verknüpft.

**Mikrodaten** sind die Originaldaten statistischer Erhebungen, die sich auf **Individuen** beziehen.

**Makrodaten** sind aggregierte statistische Daten, die üblicherweise nur einer **Gruppe** von statistischen Objekten zugeordnet werden.

# Gründe & Ziele der Anonymisierung

## Gründe:

- Gesetzliche Verpflichtung
- Sicherung der Auskunftsbereitschaft der Befragten

## Ziele:

- Schutz vor der Re-Identifizierung der einzelnen Personen oder der einzelnen Unternehmen
- Ermöglichung der sinnvollen Ergebnissen der statistischen Analyse (möglichst zu den Originaldaten ähnlich)

# Stufen der Anonymisierung

- **Formale Anonymisierung**

Die direkten Identifikationsmerkmale werden vom Datensatz entfernt

- **Faktische Anonymisierung**

Die Zuordnung der Einheiten ist nur mit unverhältnismäßigem  
Zeit- und Arbeitsaufwand möglich

Daten werden faktisch anonym bezeichnet, wenn die Deanonymisierung zwar nicht gänzlich ausgeschlossen werden kann, die Angaben jedoch nur mit einem unverhältnismäßig hohen Aufwand an Zeit, Kosten, Arbeitskraft dem jeweiligen Merkmalsträger zugeordnet werden können.

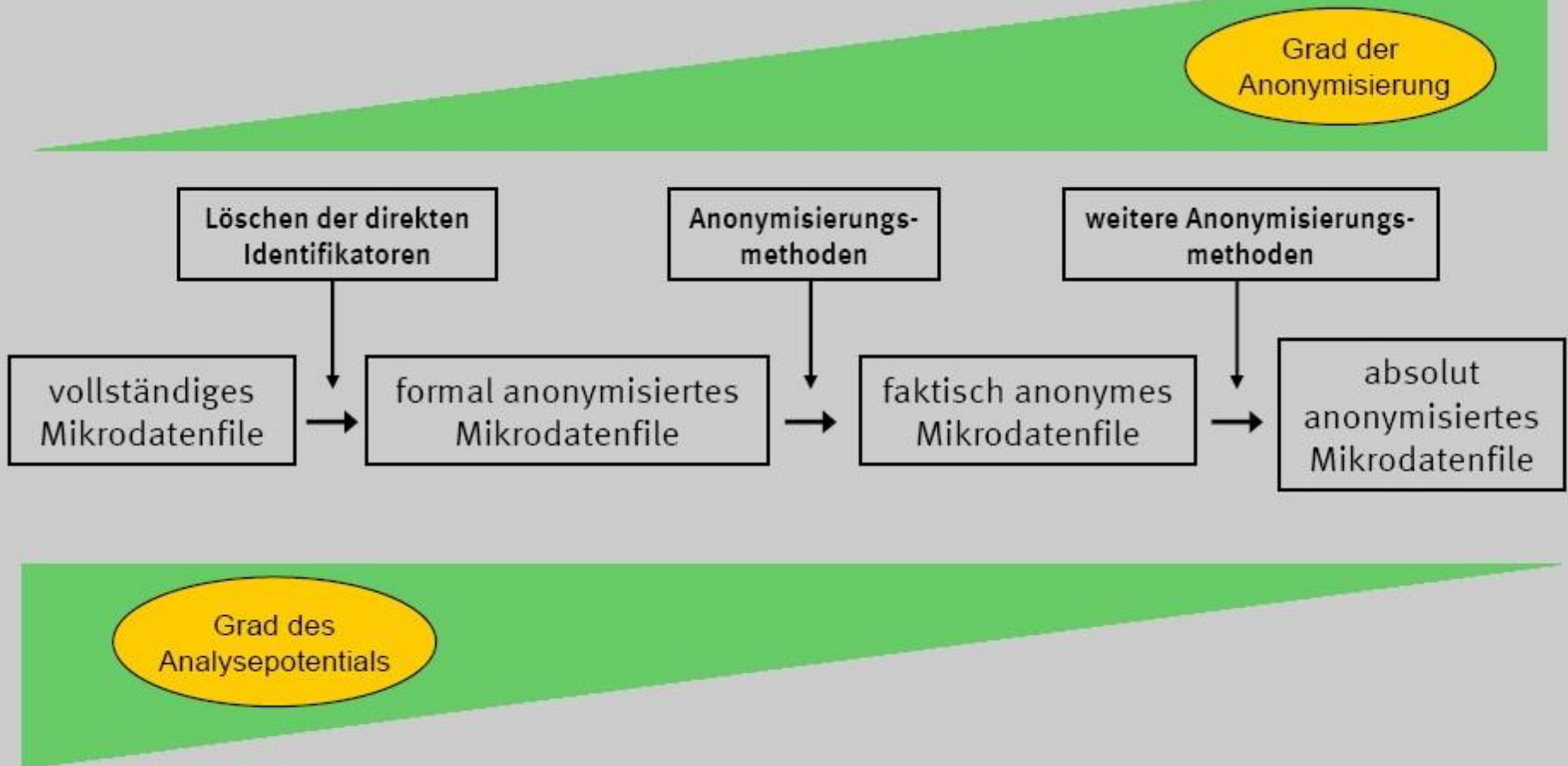
- **Absolute Anonymisierung**

Die Zuordnung der Einheiten ist ausgeschlossen

Absolut anonymisierte Daten werden durch Vergrößerung oder Entfernung einzelner Merkmale so weit verändert, dass eine Identifizierung der Auskunftgebenden unmöglich gemacht wird.

# Stufen der Anonymisierung

## Grad der Anonymisierung





## Zusammenhang zwischen dem Grad der Anonymisierung und dem Grad des Analysepotenzials

- Je stärker die Anonymisierung, desto geringer wird das statistische Analysepotenzial
- Der Informationsverlust beeinträchtigt die Aussagekraft und die Güte der Ergebnisse

Um möglichst ähnliche Auswertungsergebnisse zu erzeugen, sollte sehr wenig oder optimal gar nicht am Original-Datensatz geändert werden. Aber um den Informationsgewinn von Datenangreifern zu verhindern, ist eine sehr starke Veränderungen der Daten notwendig.

-> **KOMPROMISS**: Suche nach einer Lösung, die eine ausreichende Datensicherheit gewährleistet und gleichzeitig eine möglichst hohe Analysequalität erreicht

Die Daten sollen nur soweit verändert werden, wie es für die Erreichung der Anonymität erforderlich ist und dabei sind Verfahren anzuwenden, die die Analysepotenzial möglichst wenig beeinflussen.

# Möglichkeit der Identifizierung

**Die Möglichkeit der Identifizierung einer Person hängt von vielen Faktoren ab, wie z.B von:**

- dem möglichen Zusatzwissen
- der zur Verfügung stehenden technischen Möglichkeiten der Datenverarbeitung
- der zur Verfügung stehenden Zeit
- der zur Verfügung stehenden finanziellen Mitteln

**⇒ absolut anonymisierte Daten äußerst selten in der Praxis**

# Gliederung

## 1. Einführung

### 1.1. Definitionen

### 1.2. Gründe & Ziele der Anonymisierung

### 1.3. Stufen der Anonymisierung

### 1.4. Möglichkeit der Identifizierung

## 2. Anonymisierungsverfahren

### 2.1. Verfahren zur Informationsreduktion

- merkmalssträgerbezogene Verfahren
- merkmalsbezogene Verfahren
- musprägungsbezogene Verfahren

### 2.2. Datenverändernde Verfahren

## 3. Auswahl der Verfahren

# Anonymisierungsverfahren

Anonymisierungsverfahren können in **zwei** Gruppen eingeteilt werden:

- **Verfahren zur Informationsreduktion**
  - merkmalssträgerbezogene Verfahren
  - merkmalsbezogene Verfahren
  - musprägungsbezogene Verfahren
- **Datenverändernde Verfahren**
  - Kategoriale Variable
  - Metrische Variable

# Gliederung

## 1. Einführung

### 1.1. Definitionen

### 1.2. Gründe & Ziele der Anonymisierung

### 1.3. Stufen der Anonymisierung

### 1.4. Möglichkeit der Identifizierung

## 2. Anonymisierungsverfahren

### 2.1. Verfahren zur Informationsreduktion

- merkmalssträgerbezogene Verfahren
- merkmalsbezogene Verfahren
- musprägungsbezogene Verfahren

### 2.2. Datenverändernde Verfahren

## 3. Auswahl der Verfahren

- **Entfernen auffälliger Merkmalsträger**

Entfernen von Ausreißern, d.h. besonders auffällige und daher reidentifikationsgefährdete Merkmalsträger

- **Systematische Einschränkung der Grundgesamtheit**

Entfernen einer kompletten Teilgesamtheit, welche besonders hohen Reidentifikationsrisiko ausgesetzt ist

## **VORTEIL**

- die entfernten Merkmalsträger/ TG können nicht mehr reidentifiziert werden

## **NACHTEIL**

- die ausgeschlossene Beobachtungen können nicht mehr in die Analyse einbezogen werden

- Keine Reduzierung der Reidentifikationsgefahr der im Datenbestand verbliebenen Merkmalsträger

- Informationsunvollständigkeit

- **(Sub-)Stichprobenziehung**

Durch die Ziehung einer (Sub-)Stichprobe wird die Teilnahmewahrscheinlichkeit jedes Merkmalsträgers verringert

- Erzeugung von Unsicherheit, ob das gesuchte Objekt noch im Datenbestand ist oder nicht

### **VORTEIL**

- Schutz des gesamten Datenbestandes vor Reidentifikation
- Erzeugung von Risiko einer falschen Reidentifikation

### **NACHTEIL**

- Informationsunvollständigkeit, da eine Stichprobe immer nur eine Teilmenge der Population darstellt.

# Informationsreduktion Beispiel

\* Monat

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
Sendling	ledig	3600	700
Maxvorstadt	ledig	2900	350
Bogenhausen	verheiratet	5700	500
Schwabing-West	ledig	3420	442
Au-Haidhausen	verheiratet	3700	590
Altstadt-Lehel	verheiratet	3300	210

- **Gegeben:** 6 von 30 Arbeitnehmern eines Münchner Unternehmens
- **Sensible Information :** Einkommen
- **Bekannt:** Wohnort, Freizeitsausgaben



# Informationsreduktion

Beispiel - Merkmalsträgerbezogen

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
Sendling	ledig	3600	700
Maxvorstadt	ledig	2900	350
<del>Bogenhausen</del>	<del>verheiratet</del>	<del>5700</del>	<del>500</del>
Schwabing-West	ledig	3420	442
Au-Haidhausen	verheiratet	3700	590
Altstadt-Lehel	verheiratet	3300	210

➔ Entfernung auffälliger Merkmalsträger

# Informationsreduktion

## merkmalsbezogene Verfahren

- **Behandeln von einzelnen oder mehreren Merkmale**
- **Werden i.d.R auf Überschneidungsmerkmal angewendet** (um die Zuordnung zu verhindern) **oder auf besonders sensible Merkmale** (um die wahren und exakten Werte vor Enthüllungen zu schützen)

- **Beseitigung, Ersetzung oder Zusammenfassung von Merkmalen**

Die Merkmale werden vollständig beseitigt oder durch adäquate Linearkombinationen, Kennziffern oder Indizes ersetzt

### **VORTEIL**

- Entfernung der Überschneidungsmerkmale : Zuordnungswahrscheinlichkeit sinkt
- Entfernung der sensiblen Variablen: Anreiz, eine Reidentifikation vorzunehmen, sinkt

### **NACHTEIL**

- Erheblicher Informationsverlust

- **Vergrößerung von Merkmalsausprägungen**

- Gruppierung von metrischen Merkmalen (z.B. Bildung von Umsatzgrößenklassen)
- Vergrößerung durch Rundung der Werte metrischer Variablen (z.B. Rundung von Umsatzwerten auf ganze Tausenderbeträge)
- Zusammenfassung bereits existierender Kategorien (z.B. Zusammenfassung von zwei benachbarten Umsatzgrößenklassen)

## **VORTEIL**

- Erhöhung der Unsicherheit für Angreifer, da die Wahrscheinlichkeit von Falschzuordnung steigt
- Sinkung der Nutzen durch eine Enthüllung, da mit der Vergrößerung ein Informationsverlust verbunden ist

## **NACHTEIL**

- Erhebliche Verringerung des Informationsgehalts

# Informationsreduktion Beispiel - Merkmalsbezogen

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
Sendling	ledig	3600	700
Maxvorstadt	ledig	2900	350
Bogenhausen	verheiratet	5700	500
Schwabing-West	ledig	3420	442
Au-Haidhausen	verheiratet	3700	590
Altstadt-Lehel	verheiratet	3300	210

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
München-Süd	ledig	> 3500	700
München-West	ledig	0 - 3500	350
München-Ost	verheiratet	> 3500	500
München-West	ledig	0 - 3500	442
München-Ost	verheiratet	> 3500	590
München-Zentrum	verheiratet	0 - 3500	210

➔ Zusammenfassung von Kategorien, Gruppierung von metrischen Variablen

- **Local Suppression - Unterdrückung einzelner Werte**

Beobachtungen mit Ausprägungen oder Ausprägungskombinationen, die in der SP sehr selten oder einzigartig sind, werden unterdrückt

-> es entsteht „Missing Values“

## **VORTEIL**

- Vorher seltene oder einmalige Schlüsselkombinationen sind nicht mehr aufdeckbar

## **NACHTEIL**

- Erheblicher Informationsverlust

# Informationsreduktion

Beispiel - Ausprägungsbezogen

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
Sendling	ledig	3600	700
Maxvorstadt	ledig	2900	350
Bogenhausen	verheiratet	NA <del>5700</del>	500
Schwabing-West	ledig	3420	442
Au-Haidhausen	verheiratet	3700	590
Altstadt-Lehel	verheiratet	3300	210

➔ Entfernung von seltenem Wert

# Gliederung

## 1. Einführung

### 1.1. Definitionen

### 1.2. Gründe & Ziele der Anonymisierung

### 1.3. Stufen der Anonymisierung

### 1.4. Möglichkeit der Identifizierung

## 2. Anonymisierungsverfahren

### 2.1. Verfahren zur Informationsreduktion

- merkmalssträgerbezogene Verfahren
- merkmalsbezogene Verfahren
- musprägungsbezogene Verfahren

### 2.2. Datenverändernde Verfahren

## 3. Auswahl der Verfahren

# Datenverändernde Verfahren Swapping

- Basiert auf der Vertauschung von existierenden Merkmalsausprägungen zwischen verschiedenen Merkmalsträgern
- Bei mehreren Merkmalen wird die Vertauschung für jedes Merkmal getrennt vorgenommen
- **Einfaches Data-Swapping** (Kategoriale Variable)
  - Inhaltliche Zusammenfassung der Merkmalsträger, d.h. Gruppierung anhand ausgewählter kategorialer Merkmale
  - Die Merkmalswerte werden dann innerhalb der Gruppen für jedes Merkmal getrennt zufällig getauscht
- **Rank-Swapping** (Metrische Variable)
  - Sortierung der Merkmalswerte für jede einzelne Variable nach ihrer Größe
  - Definierung der Nachbarschaftsbereiche, auf die der Tausch beschränkt wird



# Datenverändernde Verfahren

Beispiel – Data-Swapping

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
Sendling	ledig	3600	700
Maxvorstadt	ledig	2900	350
Bogenhausen	verheiratet	5700	500
Schwabing-West	ledig	3420	442
Au-Haidhausen	verheiratet	3700	590
Altstadt-Lehel	verheiratet	3300	210

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
Sendling	<b>ledig</b>	<b>2900</b>	<b>442</b>
Maxvorstadt	<b>ledig</b>	<b>3420</b>	<b>700</b>
Bogenhausen	<b>verheiratet</b>	<b>3300</b>	<b>590</b>
Schwabing-West	<b>ledig</b>	<b>3600</b>	<b>350</b>
Au-Haidhausen	<b>verheiratet</b>	<b>5700</b>	<b>210</b>
Altstadt-Lehel	<b>verheiratet</b>	<b>3700</b>	<b>500</b>

➔ Gruppierung nach Familienstand

# Zusatzfolie

## Beispiel – Data-Swapping

- Die Arbeitnehmer werden anhand des ausgewählten Merkmals, hier **Familienstand**, gruppiert
- Innerhalb dieser Gruppe werden die Werte von den Merkmalen **Einkommen und Freizeitsausgaben** zufällig getauscht

# Datenverändernde Verfahren

Beispiel – Rank-Swapping

Wohnort	Familienstand	Einkommen*	Freizeits-Ausgaben*
Sendling	ledig	3600	700
Maxvorstadt	ledig	2900	350
Bogenhausen	verheiratet	5700	500
Schwabing-West	ledig	3420	442
Au-Haidhausen	verheiratet	3700	590
Altstadt-Lehel	verheiratet	3300	210

Wohnort	Familienstand	Einkommen sortiert*	Freizeits- Ausgaben*
Maxvorstadt	ledig	<b>2900</b>	350
Altstadt-Lehel	verheiratet	<b>3300</b>	210
Schwabing-West	ledig	<b>3420</b>	442
Sendling	ledig	<b>3600</b>	700
Au-Haidhausen	verheiratet	<b>3700</b>	590
Bogenhausen	verheiratet	<b>5700</b>	500

➔ Sortierung nach Einkommen

# Datenverändernde Verfahren Beispiel – Rank-Swapping

- Sortierung der Merkmalswerte **Einkommen** nach ihrer Größe
- Definierung der Nachbarschaftsbereiche : **3 Zeilen**

Wohnort	Familienstand	Einkommen sortiert*		1*		2*		3*		Ende	Freizeits-Ausgaben*
Maxvorstadt	ledig	<b>2900</b>	↑ ↓	3300	↑ ↓	3300	↑ ↓	3300	↑ ↓	3300	350
Altstadt-Lehel	verheiratet	<b>3300</b>		3420		3600		3600		3600	210
Schwabing-West	ledig	<b>3420</b>		2900		3420		3700		3700	442
Sendling	ledig	3600		3600		2900		3420		5700	700
Au-Haidhausen	verheiratet	3700		3700		3700		2900		3420	590
Bogenhausen	verheiratet	5700		5700		5700		5700		2900	500

- Dann Durchführung mit der Merkmalswerte **Freizeitsausgaben**

## Zwei Arten von Mikroaggregation (Metrische Variable)

- **Deterministische Mikroaggregation**
  - **Stochastische Mikroaggregation**
- Zusammenfassung der Objekte zu Gruppen
  - Ersetzung der Ursprungswerte jeweils durch das arithmetische Gruppenmittel
  - Gruppengröße mindestens drei Merkmalsträger
- ➔ das Reidentifikationsrisiko wird gesenkt,  
der Nutzen von Reidentifikationen reduziert

- **Gemeinsame Mikroaggregation**

- > **nach einer Variablen**

- Heraussuchen der dominierenden Variable
- Sortierung der Daten nach dieser Variable
- Zusammenfassung der **drei** benachbarte Merkmalsträger in einer Gruppe
- Ersetzung der Ursprungswerte durch den Durchschnitt der Werte

- > **nach einer Hilfsvariablen**

- Sortierung anhand von Hilfsvariablen (z.B. Die Hauptkomponente)

- > **nach allen metrischen Variablen**

- Gruppenbildung auf Basis der euklidischen Distanz

$$d(x_i, x_k) = \sqrt{\sum_{j=1}^p (x_{i,j} - x_{k,j})^2}$$

- Es werden die beiden Merkmalsträger herausgesucht, die den größten Abstand untereinander haben
- Danach werden diesen beiden jeweils die zwei dichtesten Merkmalsträger hinzu gruppiert

## Gemeinsame Mikroaggregation nach Variable Einkommen

Wohnort	Einkommen*	Freizeits-Ausgaben*
Maxvorstadt	2900	700
Altstadt-Lehel	3300	350
Schwabing-West	3420	500
Sendling	3600	442
Au-Haidhausen	3700	590
Bogenhausen	5700	210

Wohnort	Einkommen*	Freizeits-Ausgaben*
Maxvorstadt	3200	516,66
Altstadt-Lehel	3200	516,66
Schwabing-West	3200	516,66
Sendling	4333,33	423,33
Au-Haidhausen	4333,33	423,33
Bogenhausen	4333,33	423,33

$$\frac{1}{3} (2900 + 3300 + 3420) = 3200$$

$$\frac{1}{3} (700 + 350 + 500) = 516,66$$

$$\frac{1}{3} (3600 + 3700 + 5700) = 4333,33$$

$$\frac{1}{3} (442 + 590 + 210) = 423,33$$

# Datenverändernde Verfahren Deterministische Mikroaggregation

- Getrennte Mikroaggregation**

Durchführung der Mikroaggregation für jedes Merkmal einzeln

Wohnort	Einkommen*	Einkommen neu*
Maxvorstadt	2900	3200
Altstadt-Lehel	3300	3200
Schwabing-West	3420	3200
Sendling	3600	4333,33
Au-Haidhausen	3700	4333,33
Bogenhausen	5700	4333,33

Wohnort	Freizeits-Ausgaben*	Freizeits-Ausgaben neu*
Bogenhausen	210	333.33
Altstadt-Lehel	350	333.33
Sendling	442	333.33
Schwabing-West	500	596,66
Au-Haidhausen	590	596,66
Maxvorstadt	700	596,66

Wohnort	Einkommen*	Freizeits-Ausgaben*
Sendling	4333,33	333.33
Maxvorstadt	3200	596,66
Bogenhausen	4333,33	333.33
Schwabing-West	3200	596,66
Au-Haidhausen	4333,33	596,66
Altstadt-Lehel	3200	333.33



- **Zufällige Mikroaggregation**

- Zufällige Gruppenbildung von Merkmalsträgern
- Ersetzung der Ursprungswerte durch den Durchschnitt der Werte
- Gemeinsame (alle Variablen zusammen) oder getrennte (alle Variablen einzeln) Mikroaggregation

- **Bootstrap-Mikroaggregation**

- Für jeden Merkmalsträger: zufälliges Ziehen der zwei weiteren Merkmalsträgern
- Ziehung mit Zurücklegen
- Diese drei Merkmalsträger bilden eine Gruppe
- Ersetzung der Ursprungswerte durch den Durchschnitt der Werte

- Hinzufügen eines zufälligen Messfehlers zu den **metrischen Variabeln**
- Zufallszahlenaddierung oder –multiplizierung
- I.d.R. Normalverteilung

- **Additive stochastische Überlagerung (NV)**

$$X^a = X + W$$

- **Multiplikative stochastische Überlagerung (NV)**

$$X^a = X * W$$

## **Annahme:**

- *Nicht-Negativität der Elemente von  $W$*
- *$X$  unabhängig von  $W$*
- *$E(W) = 0$  (add.Ü)*
- *$E(W) = 1$  (multipl.Ü)*

$X$  := Originalwert

$W$  := Matrix aus Zufallszahlen

$X^a$  := überlagertes Wert

$*$  := elementweise Multiplikation

# Datenverändernde Verfahren stochastische Überlagerung

- $E(W) = 0$  (add.Ü)
- $E(W) = 1$  (multipl.Ü)

Erforderlich, damit  $E(X^a) = E(X)$  erfüllt wird.

## Additive stochastische Überlagerung mit $E(W) = 0$

$$\begin{aligned}E(X^a) &= E(X+W) \\ &= E(X) + E(W) \\ &= E(X)\end{aligned}$$

## Multiplikative stochastische Überlagerung $E(W) = 0$

$$\begin{aligned}E(X^a) &= E(X*W) \\ &= E(X) * E(W) \\ &= E(X)\end{aligned}$$

$X$  := Originalwert  
 $W$  := Matrix aus Zufallszahlen  
 $X^a$  := überlagerter Wert

## Problem bei der Normalverteilung

- Die größte Wahrscheinlichkeitsdichte um den Erwartungswert
- **Anwendung gestutzter NV**
  - Definierung der Verteilung der Zufallszahlen als gestutzte NV
  - Bereiche nahe dem EW und extrem außerhalb des EWs nicht zulässig
- **Anwendung Mischverteilungen aus mehreren NV**
  - Einzelne Elemente haben nicht den gesuchten EW
  - Mehrere NVs so kombiniert, dass die gewünschte Eigenschaft bzgl. des Ews erreicht wird

# Datenverändernde Verfahren Post-Randomisierung (PRAM)

- Ein Verfahren der Zufallsüberlagerung
  - Randomisierung **kategorialer Variablen** durch die Definition von Übergangswahrscheinlichkeiten
- 1) Festlegung der Übergangswahrscheinlichkeit**
  - 2) Transformierung der Merkmale mit der Übergangswahrscheinlichkeit**

## Beispiel: dichotome Variable – Geschlecht

Übergangswahrscheinlichkeit  $p = \begin{pmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{pmatrix}$

$p_{00}$  = Mann bleibt Mann

$p_{01}$  = Mann wird zu Frau

$p_{10}$  = Frau wird zu Mann

$p_{11}$  = Frau bleibt Frau

$$p_{jk} \equiv P(Y^a = j | Y = k)$$

Mit  $j, k \in \{0,1\}$

$$p_{j0} + p_{j1} = 1 \text{ für } j = 0,1$$

# Datenverändernde Verfahren SAFE-Verfahren

- Ein Verfahren der Mikroaggregation
- Veränderung von **kategorialen Variablen**
- Erzeugung eines Datenbestandes, in dem jeder Merkmalsträger bzgl. aller betrachteten Merkmale mit mindestens zwei weiteren Merkmalsträgern identisch ist

## **Grundlage der Veränderung**

- Minimierung der Abweichung in den Häufigkeitsverteilung

# Datenverändernde Verfahren simulationsverfahren

- Erzeugung synthetischer Merkmalsträger
- Anzahl der synth.Merkmalsträger nicht zwingend gleich der Anzahl der Merkmalsträger vom Originaldatensatz
- **Resampling**
  - Schätzung der mehrdimensionalen Kerndichte des gesamten Datenbestandes
  - Mit Hilfe dieser Dichte Erzeugung von synthetischen Merkmalsträger
- **Latin Hypercube Sampling (LHS)**
  - Ausgang: Anzahl  $n$  an gewünschten synthetischen Datensätze
  - Simulation der Merkmalswerte
  - Mit Hilfe der geglätteten empirischen Verteilungsfunktion / einer theoretischen Verteilungsfunktion werden für die einzelnen Variablen aus gleichverteilten Zufallswerten erzeugt
  - Umordnung der synthetischen Merkmalswerte durch Swapping-Verfahren



## Austausch von Angaben durch eingeschätzte Werte

- Besonders sensible Merkmalwerte werden durch geschätzte Werte ersetzt
- **Basis:** Regressionsmodell, wo alle vorhandenen Beobachtungen einbezogen werden

- **Einfache Imputation**

Einmalige Schätzung auf der Grundlage eines Modells

- **Multiple Imputation**

Durchführung der Schätzungen mit mehreren Modellen

-> Entstehung mehrerer anonymisierten Datensätze

# Gliederung

## 1. Einführung

### 1.1. Definitionen

### 1.2. Gründe & Ziele der Anonymisierung

### 1.3. Stufen der Anonymisierung

### 1.4. Möglichkeit der Identifizierung

## 2. Anonymisierungsverfahren

### 2.1. Verfahren zur Informationsreduktion

- merkmalssträgerbezogene Verfahren
- merkmalsbezogene Verfahren
- musprägungsbezogene Verfahren

### 2.2. Datenverändernde Verfahren

## 3. Auswahl der Verfahren

# Kriterien für die Auswahl der Verfahren

- **Leichte Handhabbarkeit des Verfahrens**

Eine leichte Handhabbarkeit ist notwendig, da die Verfahren später durch das Personal durchgeführt werden müssen.

- **Erfolgsaussichten des Verfahrens**

Es sollten die Verfahren, die erfolgversprechend sind, genutzt werden

-> Bei manchen Verfahren sind Nachteile bekannt, die durch andere Verfahren bereits gelöst sind. Andere Verfahren existieren erst in der Theorie, sind also noch nicht einsetzbar (einige Simulationsverfahren)

- **Repräsentative Vertretung der Verfahrensgruppen**

In der Verfahrensauswahl sollten möglichst alle Verfahrensgruppen vertreten sein, da jede Verfahrensgruppe einen anderen Ansatz der Anonymisierung repräsentiert

- **Methodenmix von Verfahren**

Da eine wirkungsvolle Anonymisierung oftmals durch Verwendung mehrerer Verfahren möglich ist, müssen alle Verfahren berücksichtigt werden, die zu einem solchen Methodenmix beitragen könnten

**Vielen Dank für Ihre/Eure  
Aufmerksamkeit!**

# Literaturverzeichnis

- Jörg Höhne (2010). *Statistik und Wissenschaft - Verfahren zur Anonymisierung von Einzeldaten*, Statistisches Bundesamt, Wiesbaden.
- Ronning G., Sturm R., Höhne J., Lenz R., Rosemann M., Scheffler M. und Vorgrimler D. (2005). *Statistik und Wissenschaft - Handbuch zur Anonymisierung wirtschaftsstatistischer Mikrodaten*, Statistisches Bundesamt, Wiesbaden.
- Augustin, T. and Wiencierz, A. (2012). *Wirtschafts- und Sozialstatistik Foliensatz 4.4. 06.12.2012.*  
[http://www.statistik.lmu.de/institut/ag/agmg/lehre/2011\\_WiSe/wiso/WiSo\\_folien\\_Kap\\_4.4\\_20120102.pdf](http://www.statistik.lmu.de/institut/ag/agmg/lehre/2011_WiSe/wiso/WiSo_folien_Kap_4.4_20120102.pdf)
- *Anonymität von Mikrodaten*. Statistische Ämter des Bundes und der Länder.  
<http://www.forschungsdatenzentrum.de/anonymisierung.asp>

# Abbildungsverzeichnis

## Folie 7

*Anonymisierung von Mikrodaten.* Forschungsdatenzentrum.

<http://www.empiwifo.uni-freiburg.de/lehre-teaching-1/Summer-term-10/Mat-Wirt-Sta/anonym>