

Arbeitsbereich TEAS
Institut für Geographische Wissenschaften
Freie Universität Berlin

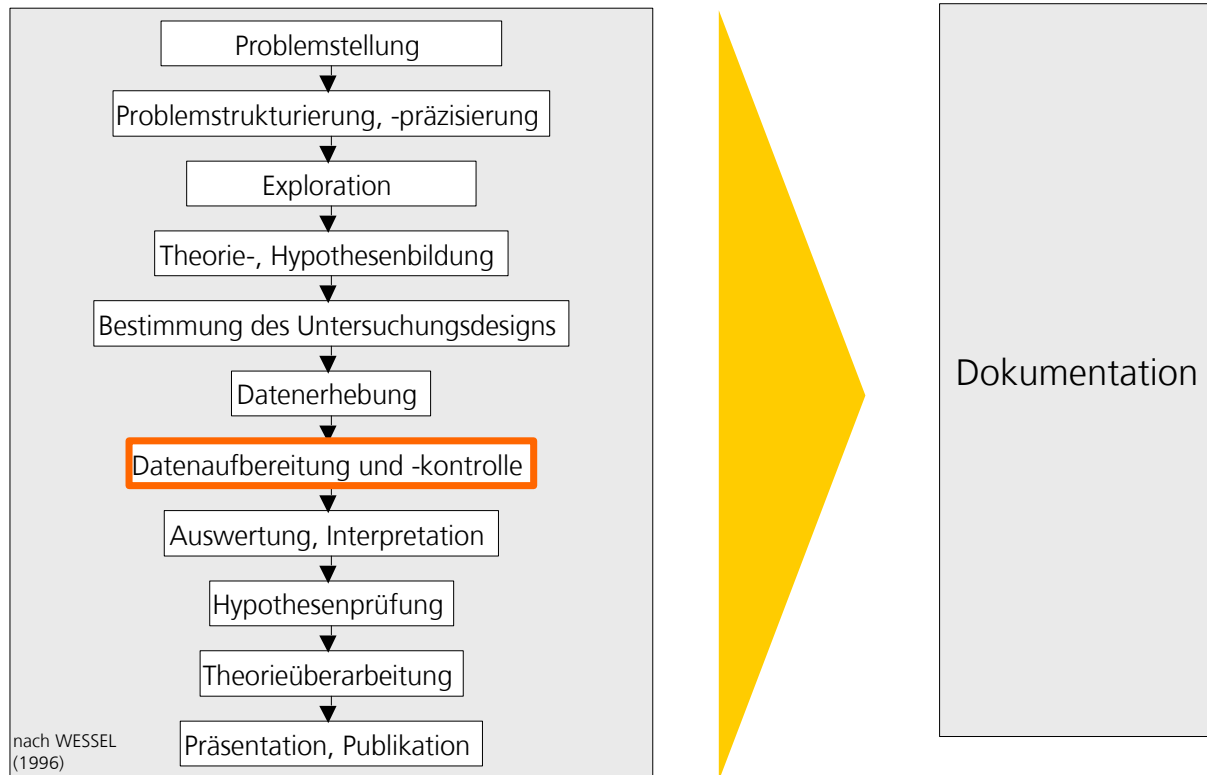
Aufbereitung empirischer Daten

Dozentin: Dipl.-Geogr. Angelika Schulz

INHALT

- Datenmatrix
- Definition der Variablen
- Dateneingabe (I-III)
- Fehlende Werte (I-III)
- Datensatzprüfung
- Datenschutz (I-II)
- Codeplanerstellung
- Primärauswertung

➔ Datenaufbereitung im wissenschaftlichen Forschungsprozess



➔ Erstellung einer Datenmatrix

Vorbereitung der informationstechnischen Verarbeitung der originalen Rohdaten

- Erstellung einer **Dummy-Matrix vor Beginn der Feldphase** zur Überprüfung der geplanten Vorgehensweise:
 - Umsetzung der Variablen
 - Erfassung aller Antwortmöglichkeiten mithilfe eines Codeplans

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	LAUFNR	ID	INT	METHOD	WO	TAG	DATUM	P	HH	P1	AGE	P1	SEX	P1	BER	FZG	F	FZG	M	FZG	P					
2	1																									
3	2																									
4	3																									
5	4																									
6	5																									
7	6																									
8	7																									

Untersuchungseinheiten, Fälle

Variablen

ebenfalls **vor** Beginn der Feldphase durchzuführen:

- **Pre-Test**
- **Probeauswertung**

➔ Definition der Variablen

- Ableitung der Variablen aus dem Erhebungsinstrument
- **eindeutige** Benennung nach einheitlichem Prinzip, z.B. ...
 - anhand der Fragennummer
 - Verwendung geeigneter von Abkürzungen
 - Kennzeichnung zusätzlich erstellter Variablen
- Festlegung eines **gültigen Wertebereichs** für jede Variable:
 - Berücksichtigung **aller** theoretisch möglichen Antworten
 - **eindeutige** Antwortkategorien ohne Überlappungsbereiche
- Berücksichtigung **softwareseitiger Einschränkungen**

Quelle: MiD 2002 – Mobilität in Deutschland 2002
(Haushaltsfragebogen und Datensatzbeschreibung)

H = Haushaltsfragebogen

1 = Haushaltsmitglied 1

2 = Haushaltsmitglied 2

Alter = Gegenstand der Variable

H1ALTER	Alter der Person 1 Druckformat: F8.2 Schreibformat: F8.2										
	<table><tr><th>Wert</th><th>Label</th></tr><tr><td>997,00</td><td>Verweigert</td></tr><tr><td>998,00</td><td>Weiß nicht</td></tr><tr><td>999,00</td><td>Keine Angabe</td></tr></table>	Wert	Label	997,00	Verweigert	998,00	Weiß nicht	999,00	Keine Angabe		
Wert	Label										
997,00	Verweigert										
998,00	Weiß nicht										
999,00	Keine Angabe										
H2ALTER	Alter der Person 2 Druckformat: F8.2 Schreibformat: F8.2 Missing-Bereich: -21,00										
	<table><tr><th>Wert</th><th>Label</th></tr><tr><td>-21,00 M</td><td>1-Personen Haushalt</td></tr><tr><td>997,00</td><td>Verweigert</td></tr><tr><td>998,00</td><td>Weiß nicht</td></tr><tr><td>999,00</td><td>Keine Angabe</td></tr></table>	Wert	Label	-21,00 M	1-Personen Haushalt	997,00	Verweigert	998,00	Weiß nicht	999,00	Keine Angabe
Wert	Label										
-21,00 M	1-Personen Haushalt										
997,00	Verweigert										
998,00	Weiß nicht										
999,00	Keine Angabe										

➔ Beispiel für Variablennamen

Ausschnitt aus einem Haushaltsfragebogen
(Angaben zu Haushaltsmitglied 1 und 2)

Ich selbst, Vorname: <input type="text"/>	Person 2, Vorname: <input type="text"/>
Ihr Geschlecht: <input type="checkbox"/> männlich <input type="checkbox"/> weiblich	Geschlecht: <input type="checkbox"/> männlich <input type="checkbox"/> weiblich
Ihr Alter: <input type="text"/> Jahre	Alter: <input type="text"/> Jahre

H = Haushaltsfragebogen

1 = Haushaltsmitglied 1

2 = Haushaltsmitglied 2

Alter = Gegenstand der Variable

H1ALTER

Alter der Person 1
Druckformat: F8.2
Schreibformat: F8.2

Wert	Label
997,00	Verweigert
998,00	Weiß nicht
999,00	Keine Angabe

H2ALTER

Alter der Person 2
Druckformat: F8.2
Schreibformat: F8.2
Missing-Bereich: -21,00

Wert	Label
-21,00 M	1-Personen Haushalt
997,00	Verweigert
998,00	Weiß nicht
999,00	Keine Angabe

Bsp. 1 Variablendefinition

Quelle: MiD 2002 – Mobilität in Deutschland 2002
(Haushaltsfragebogen und Datensatzbeschreibung)

➔ Beispiel für Variablennamen

Ausschnitt aus einem Personenfragebogen (Frage Nr. 2)

2. Besitzen Sie zurzeit ein verkehrstüchtiges Fahrrad?

☐ ja

☐ nein

Codeplan zu Frage 2 des Personenfragebogens

P = Personenfragebogen
02 = Frage Nr. 2

P02	Besitzen Sie zur Zeit ein verkehrstüchtiges Fahrrad?
	Druckformat: F8.2
	Schreibformat: F8.2
	Missing-Bereich: -200,00
Wert	Label
-200,00 M	Proxy ab 14 Jahre werden nicht befragt
1,00	Ja
2,00	Nein
7,00	Verweigert
8,00	Weiß nicht
9,00	Keine Angabe

Quelle: MiD 2002 – Mobilität in Deutschland 2002
(Personenfragebogen und Datensatzbeschreibung)

➔ Dateneingabe I

- **extrem fehleranfälliger Arbeitsschritt**
- **eindeutige** Kennzeichnung aller originalen Erhebungsunterlagen vor Beginn der Dateneingabe, z.B. durch laufende Nummerierung (eigene Variable im Datensatz)
 - erlaubt die Trennung anderer identifizierender Variablen vom Datensatz (z.B. Name, Adresse etc.)
 - erste – nicht hinreichende – Maßnahme des Datenschutzes (Anonymisierung)
- **verschiedene Verfahrensweisen** für Dateneingabe:
 - manuelle Eingabe (direkt vom Erhebungsinstrument oder anhand einer Datenmatrix)
 - maschinelles Einlesen (Voraussetzung: geeigneter Fragebogen)
 - computergestützte Datenerfassung

}

herkömmliche Erhebungsverfahren, PAPI

}

CATI, elektronische Fragebögen, online-Erhebung
- **verschiedene Auswirkungen** des gewählten Verfahrens auf ...
 - Bearbeitungsdauer
 - Kosten
 - Qualität der Daten

➔ Beispiel für Kennzeichnung des Fragebogens

Deckblatt eines Fragebogens

Bsp. Fragebogenkennzeichnung

The diagram shows a survey cover sheet titled "MOBILITÄT VON STUDENTEN IN BERLIN" with the subtitle "Haushalts- und Personenfragebogen". It includes a logo for "veritas iustitia libertas" and a list of instructions for coding. Two orange boxes highlight the ID marking fields: "Fragebogen-ID: ____ / ____" with a button "ID TEILN" and "Interviewer-ID: ____" with a button "ID INT". Arrows point from these boxes to labels on the right: "ID-Kennzeichnung des Fragebogens" and "ID-Kennzeichnung des Interviewers".

MOBILITÄT VON STUDENTEN IN BERLIN
Haushalts- und Personenfragebogen

Fragebogen-ID: ____ / ____ **ID TEILN**

Fragebogen-ID: ____ / ____ **ID TEILN**

Liebe StudenteInnen und Studienteilnehmer,

in diesem Teil der Befragung geht es um einige Angaben zu Ihrem Haushalt und zu Ihrer allgemeinen Verkehrsmittelnutzung. Im zweiten Teil folgen Fragen zu Ihren Stichtagen und den Trips, die Sie an diesen Tagen zurückgelegt haben.

Wir hoffen, dass Ihnen das Ausfüllen ein wenig Spaß macht und bedanken uns schon jetzt für Ihre Mitarbeit. Sie leisten damit einen wichtigen Beitrag zum Gelingen der Studie.

Ihr Projektteam der FU Berlin

Kodierungsplan und Hinweise:
(EXCEL-Datei: Eingabemaske HH-Pers Fragebogen.xls)

- Die Variablennamen stehen in den Kästen mit durchgehender Linie **F HH**
- Die jeweiligen einzutragenden Codes stehen in den gestrichelten Kästen **2** neben den auszufüllenden Feldern
- Bei nicht beantworteter Frage bitte 999 für „keine Antwort“ eintragen (Ist bei den jeweiligen Fragen als Möglichkeit mit angegeben)
- Die Felder mit Zahleneingabe (z.B. Alter) sind z.B. mit „1-x“ gekennzeichnet **1-x**

ID INT Interviewer-ID: ____

ID INT Interviewer-ID: ____

ID-Kennzeichnung des Fragebogens

ID-Kennzeichnung des Interviewers

➔ Dateneingabe II

- **einheitliche** und **eindeutige Codieranweisungen** erstellen:
 - Umgang mit nicht vorgesehenen oder uneindeutigen Antworten
 - insb. wenn mehrere Personen an der Dateneingabe beteiligt sind
- Die **Eingabe ungültiger Werte** lässt sich mit Hilfe entsprechender Software, z.B. Microsoft Access®, Microsoft Excel®, SPSS Data Entry®, von vornherein **vermeiden**:
 - fehlende Werte
 - Werte, die nicht codiert sind und somit außerhalb des Gültigkeitsbereiches liegen
- Das **Einlesen von Datendateien**, die im ASCII-Format vorliegen, erfordert einen **vollständigen** und **aktuellen Codeplan**, der alle Informationen über die Variablenformate enthält:
 - Bei der Festlegung der Variablenbreite ist besondere Sorgfalt hinsichtlich der Spaltentrennung erforderlich.
 - Bei numerischen Variablen ist auf korrekte Dezimaltrennung zu achten. Sonderformate wie Datum oder Währung sind zu berücksichtigen.
 - Nach dem Import sind **alle (!) Variablen** einzeln auf Richtigkeit zu **überprüfen** !

➔ Beispiel für eine Codieranleitung

Ausschnitt aus einer Codieranleitung – allgemeine Hinweise

Kodierungsplan und Hinweise:

(EXCEL-Datei: Eingabemaske HH-Pers Fragebogen.xls)

- Die Variablenamen stehen in den Kästen mit durchgehender Linie **F HH**
- Die jeweiligen einzutragenden Codes stehen in den gestrichelten Kästen **2** neben den auszufüllenden Feldern
- Bei nicht beantworteter Frage bitte 999 für „keine Antwort“ eintragen (Ist bei den jeweiligen Fragen als Möglichkeit mit angegeben)
- Die Felder mit Zahleneingabe (z.B. Alter) sind z.B. mit „1-x“ gekennzeichnet **1-x**

Ausschnitt aus einer Codieranleitung – Spezifikation der Codes

Bsp. 1 Codieranleitung

The diagram illustrates the mapping of variable names to specific code entries in a survey form. Orange boxes and arrows link the general instructions to the specific form fields:

- Benennung der Variablen:** Points to the variable names in the form, such as **F1 SEX**, **F2 SEX**, **F3 SEX**, **F4 SEX**, **F5 SEX**, and **F6 SEX**.
- Festlegung eines Wertebereiches:** Points to the code entries in the form, such as **2** for "männlich", **1** for "weiblich", and **999** for "keine Antwort".
- Codierung der Antwortkategorien (inkl. "999" = "keine Antwort"):** Points to the code entries in the form, such as **1-x** for "Alter" and **1**, **2**, **3**, **4**, and **999** for "Berufstätigkeit".

The form itself is a table with columns for different persons (Sie Selbst, Zweite Person, Dritte Person, Vierte Person, Fünfte Person, Sechste Person) and rows for different variables (Geschlecht, Alter, Berufstätigkeit). The code entries are specified in the form fields.

	F1 SEX	F2 SEX	F3 SEX	F4 SEX	F5 SEX	F6 SEX
Sie Selbst						
Zweite Person						
Dritte Person						
Vierte Person						
Fünfte Person						
Sechste Person						
Ihr Geschlecht	<input type="checkbox"/> männlich	<input type="checkbox"/> männlich	<input type="checkbox"/> männlich	<input type="checkbox"/> männlich	<input type="checkbox"/> männlich	<input type="checkbox"/> männlich
	<input type="checkbox"/> weiblich	<input type="checkbox"/> weiblich	<input type="checkbox"/> weiblich	<input type="checkbox"/> weiblich	<input type="checkbox"/> weiblich	<input type="checkbox"/> weiblich
Ihr Alter	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Ihre Jahre	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Ihre Jahre	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Ihre Jahre	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Ihre Jahre	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Ihre Jahre	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
Berufstätigkeit	<input type="checkbox"/> Vollzeit	<input type="checkbox"/> Vollzeit	<input type="checkbox"/> Vollzeit	<input type="checkbox"/> Vollzeit	<input type="checkbox"/> Vollzeit	<input type="checkbox"/> Vollzeit
	<input type="checkbox"/> Teilzeit	<input type="checkbox"/> Teilzeit	<input type="checkbox"/> Teilzeit	<input type="checkbox"/> Teilzeit	<input type="checkbox"/> Teilzeit	<input type="checkbox"/> Teilzeit
	<input type="checkbox"/> Azubi/Student(in)	<input type="checkbox"/> Azubi/Student(in)	<input type="checkbox"/> Azubi/Student(in)	<input type="checkbox"/> Azubi/Student(in)	<input type="checkbox"/> Azubi/Student(in)	<input type="checkbox"/> Azubi/Student(in)
	<input type="checkbox"/> nicht berufstätig	<input type="checkbox"/> nicht berufstätig	<input type="checkbox"/> nicht berufstätig	<input type="checkbox"/> nicht berufstätig	<input type="checkbox"/> nicht berufstätig	<input type="checkbox"/> nicht berufstätig

➔ Beispiel für eine Codieranleitung

Variable	Variablen-Name	CODES	SPALTE
LFN	Befragtennummer		1 - 4
V1	1. Geschlecht weiblich männlich	[1] [2]	5
V2	2. Alter in Jahren Weiß nicht Angabe verweigert	[98] [99]	6-7
V3	3. Familienstand ledig verheiratet geschieden getrennt lebend verwitwet sonstiges Weiß nicht Angabe verweigert	[1] [2] [3] [4] [5] [6] [98] [99]	8-9
V4	4. Beruf des Ehepartners (Coder: Eintrag klassifizieren:) Ungelernter Arbeiter Facharbeiter Meister Angestellter Leitender Angestellter Selbstkändiger Hausfrau	[1] [2] [3] [4] [5] [6] [7]	10-11

Ausschnitt aus einer Codieranleitung
– Spezifikation der Codes

Variablenname

Variablenlabel

Variablenausprägung (Antwortkategorien)

laufende Nummer der Spalten
in einer ASCII-Datei

Codierung der Antwortkategorien
(inkl. "999" = "keine Antwort")

Dateneingabe III – Datenblatt

Bei der Erstellung eines Datenblattes werden die Antworten zunächst manuell vom Fragebogen auf dieses übertragen; erst danach erfolgt die endgültige Eingabe in das gewünschte Datenverarbeitungsprogramm.

- Vorteile:
 - Identifizierung nicht vorhergesehener Antwortkategorien; Möglichkeit zur Ergänzung fehlender Codierungen
 - Konzentration auf jeweils eine einzige Aufgabe (a. "Entziffern" der Antworten und Codierung, b. Eintippen der codierten Antworten)
 - Gelegenheit zur ersten Konsistenzprüfung
- Nachteile:
 - zusätzlicher Bearbeitungsschritt erhöht die Wahrscheinlichkeit von Fehleingaben
 - erhöhter Zeitaufwand

➔ Fehlende Werte (missing values) I – Entstehung

- Fehlende Werte in Datensätzen ergeben sich **aus unterschiedlichen Gründen**:
 - Frage trifft für Probanden nicht zu
 - Filterfragen werden übersehen/ ignoriert oder Filterhinweise fehlen im Fragebogen
 - Probanden finden keine für sie zutreffende Antwortkategorie, sie weichen entweder auf "weiß nicht" aus oder antworten gar nicht.
 - Frage ist missverständlich gestellt
 - Eingabefehler
 - Beim Zusammenführen verschiedener Quellen gibt es keine "matches".
- Fehlende Werte können **vermieden** werden:
 - Fragebogen klar und übersichtlich gestalten (graphisch, sprachlich).
 - Fragen und Antwortkategorien eindeutig formulieren, ggf. mit Erläuterungen versehen.
 - Bei Verzweigungen (Filterfragen) eindeutige Sprungbefehle geben ("Sind sie nicht berufstätig, fahren Sie bitte fort mit Frage 10").

➔ Fehlende Werte (missing values) II - Umgang

Für den Umgang mit fehlenden Werten gibt es verschiedene Möglichkeiten:

- Datenzellen frei lassen (**blanks**)
 - Vorteil: Zahl der Fälle ist in jedem Fall korrekt zu bestimmen.
 - Nachteil: Beim Export in andere Datenformate werden leere Zellen manchmal mit "0" gefüllt. Das Ergebnis kann verfälscht werden (falsche Fallzahlen).
- Fehlende Werte mit **speziellen Codes** versehen, z.B. "999";
in keinem Fall Werte aus dem gültigen Wertebereich verwenden !
 - Vorteil: Fehlende Werte sind eindeutig zu identifizieren.
 - Nachteil: Werden die entsprechenden Codes (z.B. "999") nicht korrekt als "missing" definiert, werden sie bei Analysen als gültige Werte mit allen negativen Folgen mit berücksichtigt (falsche Fallzahlen, Mittelwerte etc.).
- **verschiedenen Ursachen** für fehlende Werte können erfasst werden
(z.B. "77" = "weiß nicht", "88" = "Antwort verweigert", "99" = "Frage trifft nicht zu (Filter)")
- variablenübergreifend **einheitliche Codierung** für fehlende Werte verwenden

➔ Fehlende Werte (missing values) III - Bereinigung

Die Bereinigung fehlender Werte ist immer **problematisch**, da der Originalzustand der Rohdaten nachträglich nicht mehr rekonstruiert werden kann.

Fehlende Werte können auf folgende Weise bereinigt werden:

- **intervallskalierte** Variablen: leere Zellen mit entsprechendem Mittelwert auffüllen
- **kategoriale** Variablen: leere Zellen anhand der Verteilung der Werte in den anderen Zellen füllen
- Fälle mit vielen fehlenden Werten können komplett entfernt werden ("listwise deletion").
- Alle durchgeführten **Bereinigungsverfahren** sind in jedem Fall **ausführlich** zu **dokumentieren** !

➔ Datensatzprüfung

Im Anschluss an die Dateneingabe können Datensätze **formal** und **inhaltlich** geprüft werden.

- **Aufwand** und **Nutzen** abwägen (Zeit, Kosten, Datenqualität)
- **Formale Prüfung** (fehlende Werte, nicht codierte Werte, ungültige Werte):
 - Erstellung von Häufigkeitstabellen (sehr einfach); wichtig: Filterkriterien berücksichtigen
 - gezielte Suchabfragen, z.B. nach Werten außerhalb des gültigen Wertebereichs (eher umständlich)
 - Ausdruck der Rohdaten als Liste, dann Sichtprüfung (nur bei überschaubaren Datensätzen mit vergleichsweise geringer Fallzahl sinnvoll)
- **Inhaltliche Prüfung**:
 - Kreuztabellierung von Variablen zur Identifizierung inkompatibler Wertebereiche (z.B. "Alter" = "< 18" und "Pkw-Führerscheinbesitz" = "ja")
 - ausführliche Einzelprüfung von 5-10% aller Datensätze hinsichtlich Plausibilität und Konsistenz zur Identifizierung von Fehlermustern

➔ Datenschutz I

- Bei der Verarbeitung personenbezogener Daten sind die **regional unterschiedlichen Bestimmungen der Datenschutzgesetze** zu berücksichtigen, insbesondere dann, wenn Dritte Zugang zu den erhobenen Daten haben sollen.
- **Personenbezogene Daten:** Einzelangaben über persönliche oder sachliche Verhältnisse einer bestimmten oder bestimmbaren natürlichen Person (in nicht aggregierter Form).
- **Anonymisierte** und **aggregierte** Daten fallen nicht mehr unter die Bestimmungen der Datenschutzgesetze:
 - **Anonymisieren** 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 (faktische Anonymisierung). Ein Restrisiko der Deanonymisierung kann bestehen bleiben.
- Zu berücksichtigen sind ggf. vorhandene **Verknüpfungsverbote**.
- Für die Nutzung durch Dritte (insbesondere im Rahmen von Sekundäranalysen) werden häufig **"public use files"** oder **"scientific use files"** erstellt, welche mindestens faktisch anonymisiert sind.

➔ Datenschutz II - Anonymisierung

- Zu unterscheiden sind Variablen, ...
 - die die **direkte Identifizierung** einer Person, eines Haushalts oder eines Unternehmens ermöglichen (z.B. Namen, Adressen, Telefonnummern, Sozialversicherungsnummern, in Einzelfällen auch PLZ-Gebiete),
 - die eine **indirekte Identifizierung** ermöglichen (z.B. räumliche Einheiten wie amtl. Gemeindeschlüssel, Zensusgebiete oder Wahlbezirke, exakte Daten wie Geburtsdatum, Führerscheinwerb, Jahr der Unternehmensgründung oder Insolvenz).
- Folgende Verfahren können für die Behandlung identifizierender Variablen angewendet werden:
 - Variable vollständig entfernen
 - Zusammenfassung einzelner Werte einer Variable zu größeren Wertebereichen (Recodierung), die jeweils unteren/ oberen Wertebereiche kappen (z.B. bei genauen Angaben zum Einkommen "Jahreseinkommen in Euro" = " >50000").
 - Zusammenfassung von zwei einzelnen Variablen zu einer neuen Variable
 - Einstreuung von Zufallsfehlern
 - Zerlegung des Datensatzes in Teildatensätze für einzelne Merkmalsbereiche
 - Überführung identifizierender Elemente mittels eines Schlüssels in ein Pseudonym (Pseudoanonymisierung)

Erstellung eines Codeplans

- Der **endgültige Datensatz** ist mittels eines **begleitenden Codeplans** zu dokumentieren.
Dieser stellt die minimale Form der Dokumentation dar und ist insbesondere bei Dateien im ASCII-Format erforderlich.
- Bestandteile des Codeplans sind folgende Angaben:
 - Anzahl der Zeilen pro Fall
 - Anzahl der Variablen pro Fall
 - Anzahl der Fälle im Datensatz
 - jeweilige Position der Variablen im Datensatz (Startspalte, Endspalte)
 - Breite der Variablen (Anzahl der Zeichen)
 - Namen und Label aller Variablen
 - Schlüssel für gültige Antwortkategorien inkl. fehlender Werte
 - Quellenangabe für verwendete Codierungsschemata (z.B. amtl. Gemeindekennzahlen, Klassifikation der Wirtschaftszweige, Güterverzeichnis für Produktionsstatistiken)

Dokumentation der Primärauswertung

- Primärauswertung als minimale Form der Datenaufbereitung in Form von einfachen Tabellen oder Kreuztabellen der einzelnen Variablen (Publikation i. d. R. in Tabellenbänden)
- Die folgenden Bestandteile einer Tabelle sind **zwingend** erforderlich:
 - Überschrift/ Titel
 - zugrunde liegende Grundgesamtheit
 - bei Prozentwerten: Angabe der zugrunde liegenden Fallzahl (n)
 - Hinweise auf geltende Filterbedingungen
 - ggf. Begriffsdefinitionen
 - Quelle der Daten