

Markus Tausendpfund et al.

Quantitative Analyseverfahren

Eine Einführung

Fakultät für
**Kultur- und
Sozialwissen-
schaften**

Das Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere das Recht der Vervielfältigung und Verbreitung sowie der Übersetzung und des Nachdrucks, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Kein Teil des Werkes darf in irgendeiner Form (Druck, Fotokopie, Mikrofilm oder ein anderes Verfahren) ohne schriftliche Genehmigung der FernUniversität reproduziert oder unter Verwendung elektronischer Systeme verarbeitet, vervielfältigt oder verbreitet werden.

Der Inhalt dieses Studienbriefs wird gedruckt auf Recyclingpapier (80 g/m², weiß), hergestellt aus 100 % Altpapier.

Vorwort

In der quantitativen Sozialforschung wird zur Beschreibung von Daten und zur empirischen Überprüfung von Hypothesen auf statistische Verfahren zurückgegriffen. Wer eine (quantitative) Studie verstehen und kritisch bewerten möchte, der muss die grundlegenden Prinzipien, Anwendungsvoraussetzungen und auch Probleme der verwendeten statistischen Verfahren kennen. Für Sozialwissenschaftlerinnen und Sozialwissenschaftler sind deshalb elementare Kenntnisse dieser quantitativen Analyseverfahren unverzichtbar.

Für die Sozialwissenschaften stellt die Statistik eine zentrale Hilfswissenschaft dar. Während sich Statistiker – allgemeiner: Mathematikerinnen – häufig mit der Beweisführung und der Weiterentwicklung mathematischer Algorithmen beschäftigen, steht für Studierende der Politikwissenschaft, Verwaltungswissenschaft und Soziologie das Kennenlernen und die praktische Anwendung statistischer Verfahren im Vordergrund. Im Mittelpunkt des Kurses steht das Verständnis quantitativer Analyseverfahren, mit denen Sozialwissenschaftlerinnen und Sozialwissenschaftler bei der Auseinandersetzung mit quantitativen Studien konfrontiert werden.

Der vorliegende Kurs behandelt vier Themenbereiche: Univariate, bivariate und multivariate Datenanalyse sowie Grundlagen der Inferenzstatistik. Das Kapitel zur univariaten Datenanalyse behandelt die Häufigkeitsverteilung einzelner Merkmale. Dabei werden Lage- und Streuungsmaße sowie Formmaße vorgestellt. Die bivariate Datenanalyse untersucht Zusammenhänge zwischen zwei Merkmalen und Unterschiede zwischen zwei Merkmalen (Mittelwertvergleiche). Dabei werden Kreuztabellen sowie wichtige Zusammenhangsmaße behandelt. Bei der multivariaten Datenanalyse werden mit der linearen und logistischen Regression zwei zentrale Analyseverfahren der Sozialwissenschaften vorgestellt, die den Einfluss mehrerer unabhängiger Variablen auf eine abhängige Variable schätzen können. Aus zeitlichen, finanziellen und forschungspraktischen Gründen dominieren in den Sozialwissenschaften Stichproben. Deshalb behandelt der vierte Teil des Kurses die Grundlagen der Inferenzstatistik, die Instrumente zur Verfügung stellt, um zu entscheiden, ob und wie empirische Befunde aus Zufallsstichproben auf zugehörige Grundgesamtheiten übertragen werden dürfen.

Für das Sommersemester 2019 wurde der Kurs aktualisiert und erweitert. Neu sind die Kapitel „Grundlagen“ und „Statistisches Testen“. In der Moodle-Lernumgebung des Moduls M1 „Quantitative Methoden der Sozialwissenschaften“ findet sich eine Errata-Liste zu dem Kurs. Außerdem werden dort auch Videos und Übungsaufgaben veröffentlicht, die die Auseinandersetzung mit den Inhalten des Kurses fördern sollen.

Über Hinweise auf Fehler, Kommentare und Verbesserungsvorschläge freue ich mich. Sie erreichen mich unter der E-Mail-Adresse Markus.Tausendpfund@fernuni-hagen.de

Hagen, im Dezember 2018

Markus Tausendpfund

Inhaltsverzeichnis

Abbildungsverzeichnis.....	VII
Tabellenverzeichnis	VIII
1 Grundlagen.....	10
1.1 Einordnung im Forschungsprozess	10
1.2 Grundgesamtheit und Stichprobe	13
1.3 Klassifikationen von Variablen.....	15
2 Univariate Datenanalyse	18
2.1 Häufigkeitstabelle	18
2.2 Lagemaße.....	22
2.2.1 Modus.....	22
2.2.2 Median.....	23
2.2.3 Arithmetisches Mittel.....	25
2.3 Streuungsmaße	28
2.3.1 Varianz.....	28
2.3.2 Standardabweichung.....	32
2.4 Formmaße	32
2.4.1 Schiefe	33
2.4.2 Wölbung.....	36
2.5 Variablen standardisieren (z-Transformation).....	37
3 Bivariate Datenanalyse.....	39
3.1 Kreuztabellen	40
3.1.1 Relative Häufigkeiten und Prozentwerte	42
3.1.2 Gesamtprozente.....	43
3.1.3 Spaltenprozente	43
3.1.4 Zeilenprozente.....	44
3.1.5 Prozentsatzdifferenzen	45
3.1.6 Berechnung von Kreuztabellen	46
3.1.7 Hinweise und praktische Tipps.....	50
3.2 Cramér's V	51
3.2.1 Vorgehen bei der Berechnung von Cramér's V.....	53
3.2.2 Kontingenztabelle	54
3.2.3 Berechnung der Summenzahlen	55
3.2.4 Indifferenztabelle.....	55

3.2.5	Arbeitstabelle	57
3.2.6	Differenz von beobachteten und erwarteten Häufigkeiten	57
3.2.7	Quadrierung der Differenzen	58
3.2.8	Division der quadrierten Differenzen durch die erwarteten Häufigkeiten	59
3.2.9	Summe der Quotienten	59
3.2.10	Berechnung von Cramér's V	60
3.2.11	Interpretation von Cramér's V	60
3.3	Spearman's rho	61
3.3.1	Vorgehen bei der Berechnung von Spearman's rho	63
3.3.2	Bestimmung der Rangpositionen	64
3.3.3	Bestimmung der Differenzen zwischen den Rangpositionen	64
3.3.4	Bestimmung der quadrierten Differenzen	65
3.3.5	Bestimmung der Summe der quadrierten Differenzen	65
3.3.6	Berechnung von Spearman's rho	66
3.3.7	Interpretation von Spearman's rho	66
3.4	Pearson's r	67
3.4.1	Vorgehen bei der Berechnung von Pearson's r	69
3.4.2	Bestimmung der Produkte aus den Werten	70
3.4.3	Bestimmung der quadrierten Werte	71
3.4.4	Berechnung von Pearson's r	71
3.4.5	Interpretation von Pearson's r	72
4	Multivariate Datenanalyse	73
4.1	Einführung	73
4.2	Lineare Regression	75
4.2.1	Bivariate Regression	76
4.2.2	Multiple Regression	83
4.3	Logistische Regression	94
4.3.1	Bivariate Regression	95
4.3.2	Multiple Regression	98
5	Inferenzstatistik	104
5.1	Was ist das Problem?	104
5.2	Zentrale Konzepte der Inferenzstatistik	109
5.2.1	Zentraler Grenzwertsatz und Normalverteilung	109
5.2.2	Standardfehler	112

5.3	Schätzungsarten	118
5.3.1	Punktschätzung	118
5.3.2	Intervallschätzung	121
5.3.3	Berechnung der benötigten Fallzahl	129
5.3.4	Anwendungsprobleme in der Praxis	131
5.4	Statistisches Testen	133
5.4.1	Allgemeine Vorgehensweise bei einem Signifikanztest	135
5.4.2	Alpha- und Beta-Fehler	138
5.4.3	t-Test	139
6	Literatur	154

Abbildungsverzeichnis

Abbildung 1: Idealtypischer Ablauf eines quantitativen Forschungsprojekts.....	11
Abbildung 2: Grundgesamtheit und Stichprobe	14
Abbildung 3: Normalverteilung.....	33
Abbildung 4: Schiefe von Verteilungen.....	34
Abbildung 5: Empirische Verteilungen mit unterschiedlicher Schiefe	35
Abbildung 6: Wölbung.....	36
Abbildung 7: Streudiagramm mit Beispieldaten	78
Abbildung 8: Streudiagramm mit OLS-Regressionsgerade.....	80
Abbildung 9: Schematische Darstellung der vermuteten multivariaten Einflusstruktur	85
Abbildung 10: Bivariate logistische Regression (Beispieldaten).....	97
Abbildung 11: Grundgesamtheit und Stichprobe	104
Abbildung 12: Rückschluss von der Stichprobe auf die Grundgesamtheit.....	105
Abbildung 13: Wiederholte Ziehung von Zufallsstichproben.....	110
Abbildung 14: Normalverteilung.....	112
Abbildung 15: Abweichungen einzelner Stichprobenmittelwerte vom wahren Mittelwert	113
Abbildung 16: Stichprobenverteilungen bei unterschiedlicher Fallzahl	114
Abbildung 17: Ergebnisse des Politbarometers zu zwei Zeitpunkten (Angaben in Prozent)	122
Abbildung 18: 95-Prozent-Konfidenzintervall	123
Abbildung 19: 99-Prozent-Konfidenzintervall	124
Abbildung 20: Fiktive Befragung zur Wahlentscheidung von 1000 Personen (in Prozent)	126
Abbildung 21: 95-Prozent-Konfidenzintervalle (Stichprobengröße jeweils 1000 Personen)	128
Abbildung 22: Schätzen und Testen im Vergleich	134
Abbildung 23: t-Verteilung und Normalverteilung	141
Abbildung 24: Verschiedene t-Verteilungen.....	142
Abbildung 25: Varianten des t-Tests	142
Abbildung 26: Einseitiger und zweiseitiger t-Test.....	144

Tabellenverzeichnis

Tabelle 1: Zulässige Rechenoperationen in Abhängigkeit des Skalenniveaus	16
Tabelle 2: Interesse an Politik.....	18
Tabelle 3: Subjektive Schichtestufung	21
Tabelle 4: Lagemaße und Skalenniveau	22
Tabelle 5: Berechnung des Modus.....	22
Tabelle 6: Geschlecht	23
Tabelle 7: Berechnung des Medians (ungerade Fallzahl).....	24
Tabelle 8: Berechnung des Medians (gerade Fallzahl).....	24
Tabelle 9: Interesse an Politik.....	25
Tabelle 10: Berechnung des arithmetischen Mittels bei kleinen Fallzahlen	26
Tabelle 11: Lebenszufriedenheit	27
Tabelle 12: Mittelwerte und Ausreißer.....	27
Tabelle 13: Lebenszufriedenheit von zwei Gruppen	28
Tabelle 14: Arbeitstabelle für die Berechnung der Varianz (kleine Fallzahl)	30
Tabelle 15: Arbeitstabelle für die Berechnung der Varianz (große Fallzahl)	31
Tabelle 16: Variablen standardisieren	37
Tabelle 17: Wichtige Zusammenhangsmaße bei der bivariaten Datenanalyse.....	39
Tabelle 18: Kombinierte Häufigkeitsverteilung in einer Kreuztabelle	40
Tabelle 19: Beispiel mit zehn Murmeln (Urliste).....	41
Tabelle 20: Beispiel mit zehn Murmeln (Kreuztabelle)	42
Tabelle 21: Beispiel mit zehn Murmeln (Kreuztabelle mit Gesamtprozenten).....	43
Tabelle 22: Beispiel mit zehn Murmeln (Kreuztabelle mit Spaltenprozenten)	44
Tabelle 23: Beispiel mit zehn Murmeln (Kreuztabelle mit Zeilenprozenten).....	45
Tabelle 24: Interpretation von Prozentsatzdifferenzen	46
Tabelle 25: Geschlecht und Politikinteresse.....	46
Tabelle 26: Kreuztabelle mit absoluten Häufigkeiten (ohne Summenzahlen)	47
Tabelle 27: Kreuztabelle mit absoluten Häufigkeiten	48
Tabelle 28: Berechnung der Spaltenprozentage in einer Kreuztabelle (Beispiel 1).....	48
Tabelle 29: Berechnung der Spaltenprozentage in einer Kreuztabelle (Beispiel 2).....	49
Tabelle 30: Kreuztabelle mit absoluten Häufigkeiten und Spaltenprozenten	49
Tabelle 31: Kontingenztafel	52
Tabelle 32: Interpretation von Cramér's V	52
Tabelle 33: Geschlecht und politisches Interesse	54
Tabelle 34: Kontingenztafel	54
Tabelle 35: Berechnung der Summenzahlen	55
Tabelle 36: Indifferenztafel (Beispiel 1)	56
Tabelle 37: Indifferenztafel (Beispiel 2)	56
Tabelle 38: Vollständige Indifferenztafel.....	57
Tabelle 39: Arbeitstabelle	57
Tabelle 40: Berechnung der Differenz der beobachteten und erwarteten Häufigkeiten	58
Tabelle 41: Quadrierung der Differenz der beobachteten und erwarteten Häufigkeiten	58
Tabelle 42: Division der quadrierten Differenzen durch erwartete Häufigkeiten.....	59
Tabelle 43: Berechnung von Chi-Quadrat	59

Tabelle 44: Interpretation von Spearman's rho	62
Tabelle 45: Schulabschluss und politisches Interesse	63
Tabelle 46: Urliste.....	64
Tabelle 47: Rangpositionen bestimmen	64
Tabelle 48: Differenz der Rangpositionen	65
Tabelle 49: Rangpositionen quadriert	65
Tabelle 50: Summe der quadrierten Differenzen	65
Tabelle 51: Interpretation von Pearson's r	68
Tabelle 52: Urliste.....	69
Tabelle 53: Produkt berechnen	70
Tabelle 54: Berechnung der quadrierten Werte.....	71
Tabelle 55: Angaben zur Berechnung von Pearson's r	71
Tabelle 56: Unterschiedliche Bezeichnungen für Variablen der Regressionsanalyse	74
Tabelle 57: Fiktive Beispieldaten für bivariate Regression.....	77
Tabelle 58: Dummy-Kodierung für Familienstand.....	86
Tabelle 59: Bestimmungsfaktoren der Lebenszufriedenheit.....	88
Tabelle 60: Bestimmungsfaktoren der Lebenszufriedenheit (Regressionskoeffizienten)	92
Tabelle 61: Fiktive Beispieldaten für Wahlbeteiligung und Alter.....	96
Tabelle 62: Bestimmungsfaktoren der Wahlbeteiligung	100
Tabelle 63: Mittelwerte in Zufallsstichproben (Stichprobengröße jeweils 1000 Personen)	106
Tabelle 64: Mittelwerte von Zufallsstichproben (Stichprobengröße jeweils 1000 Personen).....	108
Tabelle 65: Vergleich zwischen Standardfehler und Standardabweichung	115
Tabelle 66: Mittelwerte von Zufallsstichproben	119
Tabelle 67: Erforderliche Stichprobengröße	130
Tabelle 68: Fehlerarten beim Hypothesentest.....	138
Tabelle 69: Lebenszufriedenheit von Frauen und Männern (fiktive Daten)	145
Tabelle 70: Kritische Werte der t-Verteilung.....	147
Tabelle 71: Lebenszufriedenheit von West- und Ostdeutschen (fiktive Daten).....	148
Tabelle 72: Zufriedenheit mit der Demokratie (fiktive Daten).....	150
Tabelle 73: Beispieldaten für die Berechnung eines t-Tests mit abhängigen Stichproben	151

1 Grundlagen

Vorschau



Markus Tausendpfund

Dieses Kapitel macht mit den Grundlagen der quantitativen Datenanalyse vertraut. Nach der Einordnung der Phase „Datenanalyse“ im Forschungsprozess werden die Begriffe „Grundgesamtheit“ und „Stichprobe“ behandelt. Bei einer empirischen Studie sind meist Aussagen über größere Gruppen angestrebt (z.B. die wahlberechtigte Bevölkerung in Deutschland). Allerdings liegen in den meisten Studien keine Informationen über alle Elemente dieser Gruppe vor, sondern nur über eine (zufällige) Auswahl dieser Gruppe. Die Gruppe, über die eine Aussage gemacht werden soll, wird als Grundgesamtheit bezeichnet. Die Gruppe, über die empirische Informationen vorliegen, wird als Stichprobe bezeichnet. Diese Begriffe werden knapp erläutert und die Voraussetzungen skizziert, um Befunde einer Stichprobe auf die zugehörige Grundgesamtheit übertragen zu können. Abschließend werden typische Klassifikationen von Variablen vorgestellt. Dabei liegt der Fokus auf dem Skalenniveau von Variablen, da das Skalenniveau eine wichtige Voraussetzung für die Anwendung bestimmter Analyseverfahren ist.

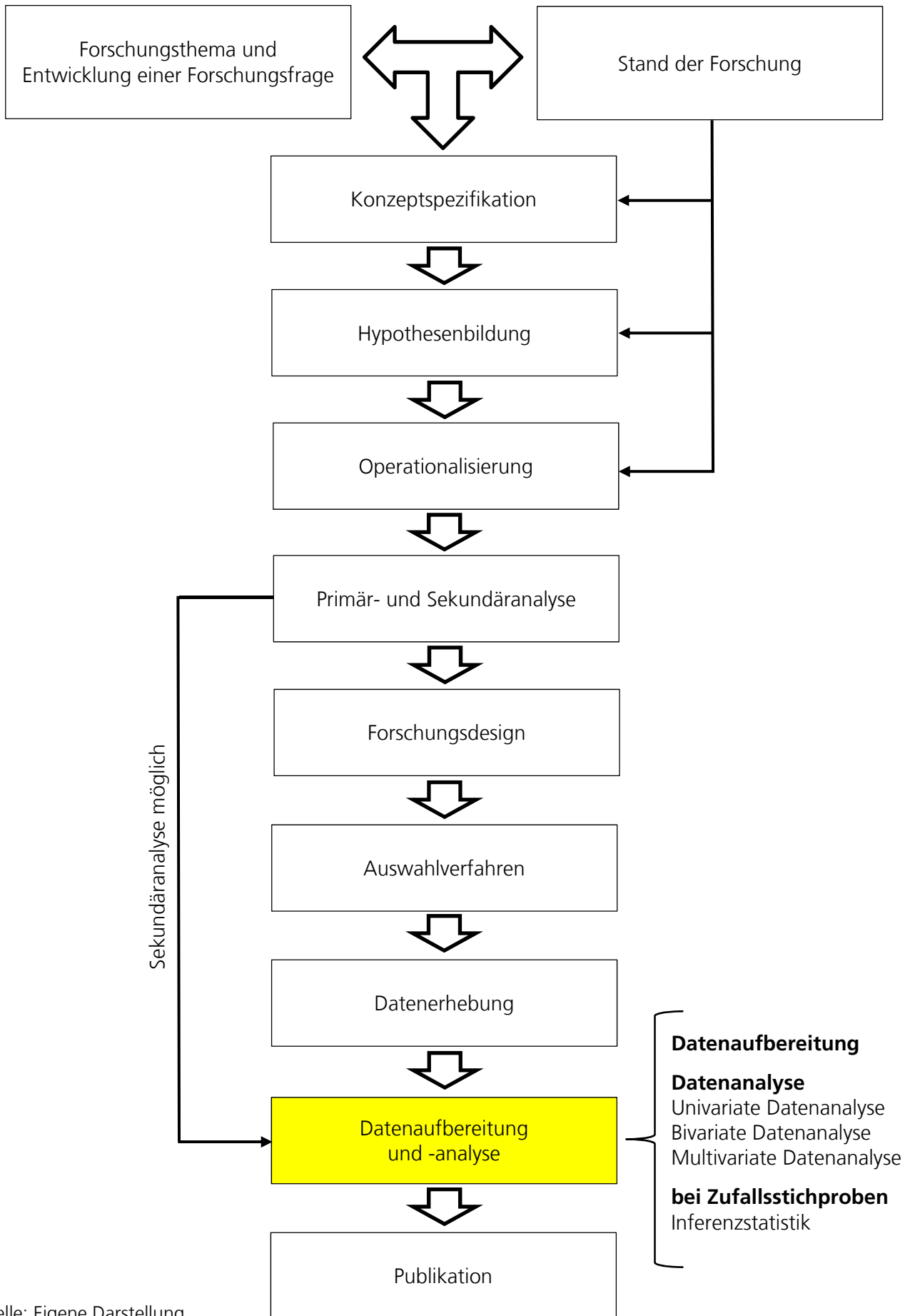
1.1 Einordnung im Forschungsprozess

Die quantitativen Analyseverfahren werden häufig mit dem quantitativen Forschungsprozess gleichgesetzt. Quantitativ arbeitende Sozialwissenschaftlerinnen nutzen statistische Analyseverfahren, um die theoretisch formulierten Hypothesen empirisch zu überprüfen. Sicherlich ist die Anwendung statistischer Analyseverfahren ein zentrales Merkmal des quantitativen Forschungsprozesses, aber die quantitative Datenanalyse sollte nicht isoliert betrachtet werden.

Vor der Datenanalyse bzw. Anwendung quantitativer Analyseverfahren müssen empirische Sozialforscher wichtige vorgelagerte Entscheidungen treffen, die unmittelbare Auswirkungen auf die empirischen Befunde haben. Wie Abbildung 1 zeigt, steht die Festlegung eines Forschungsthemas und die Entwicklung einer geeigneten Forschungsfrage am Beginn eines Forschungsprojekts. Auf dieser Grundlage werden die zentralen Konzepte identifiziert und theoretisch geklärt, ehe inhaltvolle Hypothesen formuliert und valide Operationalisierungen dieser Konzepte entwickelt werden. Diese Phasen in einem Forschungsprozess erfolgen in intensiver Auseinandersetzung mit dem existierenden Forschungsstand. Nur wer den Forschungsstand zu seinem Forschungsthema kennt, kann eine inhaltvolle Forschungsfrage entwickeln. Die Auseinandersetzung mit der Fachliteratur ist aber auch für die Konzeptspezifikation und die Entwicklung von Hypothesen erforderlich. Schließlich ist auch bei der „Übersetzung“ theoretischer Konzepte in empirische Indikatoren ein Überblick über existierende Operationalisierungen notwendig.

! Kein Analyseverfahren kann die intensive Auseinandersetzung mit dem existierenden Forschungsstand ersetzen. Ungeeignete Konzeptspezifikationen, schwammige Hypothesen oder auch ungültige Operationalisierungen führen zwangsläufig zu schlechten Daten und kein Analyseverfahren der Welt kann aus schlechten Daten valide empirische Befunde machen. Deshalb: Die Anwendung bzw. Durchführung quantitativer Analyseverfahren kann nur dann zu belastbaren empirischen Befunden führen, wenn die vorgelagerten Phasen erfolgreich bearbeitet wurden.

Abbildung 1: Idealtypischer Ablauf eines quantitativen Forschungsprojekts



Quelle: Eigene Darstellung

Wenn für die Bearbeitung einer Forschungsfrage und die Überprüfung der Hypothesen bereits geeignetes Datenmaterial existiert (z.B. ALLBUS), dann können die Phasen „Forschungsdesign“, „Auswahlverfahren“ und „Datenerhebung“ übersprungen werden. In einem solchen Fall führt die Sozialwissenschaftlerin eine Sekundäranalyse vor. Es werden existierende Daten genutzt, um die Forschungsfrage zu bearbeiten. Falls keine geeigneten Daten zur Verfügung stehen, dann bietet sich eine Primäranalyse an. Bei einer Primäranalyse werden neue Daten erhoben, um die Forschungsfrage zu beantworten.

Datenaufbereitung und -analyse

Die Phase „Datenaufbereitung und -analyse“ umfasst in der Regel mehrere Zwischenschritte (Stein 2014, S. 150; Tausendpfund 2018b, S. 50-51). Zunächst müssen die im Rahmen der Datenerhebung gesammelten empirischen Informationen systematisch in einen Datensatz aufgenommen werden (Kromrey et al. 2016, S. 217-218). Die Variablen müssen beschriftet und ein Codebuch muss angelegt werden (z.B. Lück und Baur 2011; Lück und Landrock 2014; Tausendpfund 2018b, S. 291-297). Bei der Arbeit mit qualitativ hochwertigen Sekundärdaten (z.B. ALLBUS) stehen meist „fertige“ Datensätze zur Verfügung. Insbesondere bei der eigenständigen Dateneingabe, aber auch bei der Arbeit mit Sekundärdaten, sind Fehlerkontrollen (z.B. Eingabefehler) und Plausibilitätstests erforderlich.

Vor der eigentlichen Datenanalyse müssen Variablen häufig verändert oder neu erstellt werden. Dieser Prozess wird häufig als Datenmodifikation oder Datentransformation bezeichnet (Fromm 2011; Kohler und Kreuter 2017, S. 91-130; Tausendpfund 2018a, S. 62-91). Dabei wird die Kodierung von Variablen angepasst, einzelne Subgruppen gebildet oder auf Basis der verfügbaren Informationen auch neue Variablen erstellt. Das Verändern und das Erstellen neuer Variablen dauert häufig länger als die eigentliche Datenanalyse. Eine sorgfältige Durchführung der einzelnen Schritte ist dabei eine Voraussetzung für die Gültigkeit der anschließenden Analysen.

Bei der anschließenden Datenanalyse lassen sich meist vier Schritte unterscheiden, die auch im Mittelpunkt dieses Kurses stehen: die univariate, die bivariate und die multivariate Datenanalyse sowie die Inferenzstatistik.

Univariate Datenanalyse

Die univariate Datenanalyse befasst sich mit einzelnen Variablen. In einem ersten Schritt werden die absoluten und relativen Häufigkeiten der einzelnen Ausprägungen einer Variable in Tabellen oder Grafiken dargestellt. In der quantitativen Sozialforschung sind wir allerdings in der Regel mit vielen Untersuchungsobjekten konfrontiert. Deshalb wird in einem zweiten Schritt die Informationsmenge von mehreren tausend Beobachtungen auf wenige Kennzahlen verdichtet. Dabei lassen sich Lage-, Streuungs- und Formmaße unterscheiden. Während Lagemaße über das Zentrum einer Verteilung informieren, beschreiben Streuungsmaße die Variation eines Merkmals in einer Verteilung. Mit Schiefe und Wölbung kann die Form einer Verteilung charakterisiert werden.

Bivariate Datenanalyse

Bei der bivariaten Datenanalyse werden immer genau zwei Variablen in Beziehung gesetzt (z.B. Bildung und Einkommen). Bivariate Analyseverfahren werden genutzt, um Zusammenhänge oder Unterschiede zwischen zwei Merkmalen zu untersuchen und Hypothesen empirisch zu überprüfen. Dafür nutzen wir Kreuztabellen und Zusammenhangsmaße. Kreuztabellen (engl. crosstabs) sind eine einfache und anschauliche Möglichkeit, die Beziehung von zwei Merkmalen in den Blick zu nehmen. Neben absoluten Häufigkeiten können auch die

Anteile der einzelnen Häufigkeiten (Anteile) berechnet werden. Die Stärke einer Beziehung zwischen zwei Merkmalen (z.B. Bildung und Einkommen) kann mit Zusammenhangsmaßen – sogenannten Koeffizienten – charakterisiert werden. Die bekanntesten Zusammenhangsmaße sind sicherlich Cramer's V, Spearman's rho und Pearson's r.

Mit bivariaten Analyseverfahren wird der Zusammenhang zwischen zwei Variablen untersucht. In der Realität können Merkmale wie Einkommen oder Wahlbeteiligung aber nicht durch eine Variable „erklärt“ werden, sondern in der Regel muss eine Vielzahl an Variablen gleichzeitig berücksichtigt werden. Mit der Regressionsanalyse steht ein sehr mächtiges Analyseinstrument zur Verfügung, um den Einfluss einer oder mehrerer unabhängiger Variablen auf eine abhängige Variable zu schätzen. Mit der linearen und logistischen Regression werden in diesem Kurs zwei multivariate Analyseverfahren vorgestellt, die in den Sozialwissenschaften häufig verwendet werden.

Multivariate Datenanalyse

Die univariate, bivariate und multivariate Datenanalyse haben das Ziel, die Verteilung von Variablen zu beschreiben und Zusammenhänge von zwei oder mehr Variablen zu untersuchen. Diese Datenanalyse basiert in der Regel auf Stichproben. Das heißt, es liegen nicht von allen Untersuchungsobjekten einer Grundgesamtheit empirische Informationen vor, sondern nur von einer (zufälligen) Auswahl. Die Inferenzstatistik beschäftigt sich mit der Frage, ob und wie Befunde von Zufallsstichproben auf zugehörige Grundgesamtheiten übertragen werden können (siehe auch Abschnitt 1.2).

Inferenzstatistik

Abbildung 1 soll verdeutlichen, dass die Datenanalyse bzw. die Anwendung statistischer Analyseverfahren immer nur eine Phase in einem sozialwissenschaftlichen Projekt darstellt. Empirische Befunde „sprechen“ niemals für sich selbst, sondern sind immer eingebunden in eine sozialwissenschaftliche Forschungsfrage. Ohne theoretische Vorarbeiten (z.B. Entwicklung von Hypothesen) kann eine quantitative Analyse nicht zielorientiert erfolgen. Deshalb muss eine quantitative Datenanalyse immer an den (theoretischen) Forschungsprozess zurückgekoppelt werden, um in Publikationen empirisch interessante und vor allem valide Schlussfolgerungen ziehen zu können.

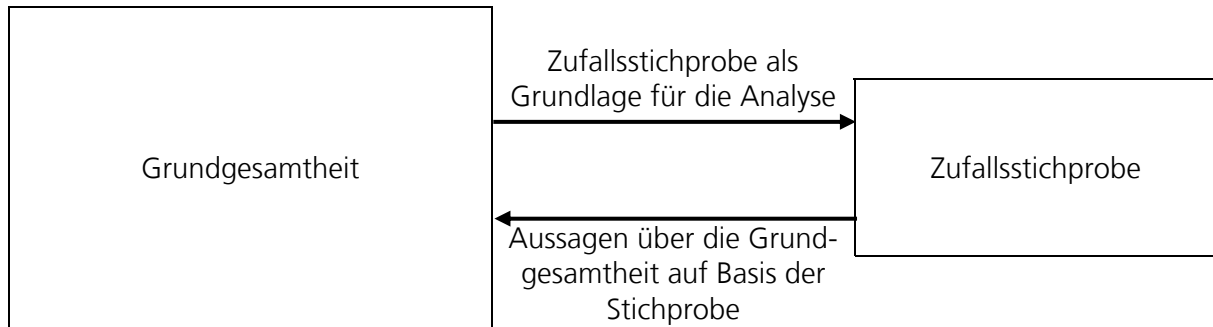
1.2 Grundgesamtheit und Stichprobe

Aus zeitlichen, finanziellen und forschungspraktischen Gründen dominieren in den Sozialwissenschaften Stichproben (Tausendpfund 2018b, S. 207-210). Bei empirischen Studien werden in der Regel nicht alle Elemente der Grundgesamtheit untersucht, sondern nur eine zufällige Auswahl dieser Elemente. Ein Beispiel: Bei der Bundestagswahl 2017 waren nach Angaben des Bundeswahlleiters 61.688.485 Personen wahlberechtigt. Bei einer Analyse des Wahlverhaltens bei der Bundestagswahl bilden diese Personen die Grundgesamtheit. Bei einer Vollerhebung würden empirische Informationen aller Untersuchungsobjekte der Grundgesamtheit erhoben. Die Kosten der Datenerhebung und die Dauer der Erhebung sprechen allerdings gegen eine solche Vollerhebung. Für die Analyse des Wahlverhaltens (z.B. im Rahmen der German Longitudinal Election Study) wird deshalb auch keine Vollerhebung angestrebt, sondern lediglich eine Zufallsstichprobe realisiert.

In Abbildung 2 wird der Zusammenhang zwischen Grundgesamtheit und Stichprobe illustriert. Im Rahmen eines Forschungsprojekts sollen Aussagen über die Grundgesamtheit gemacht werden.

In vielen Fällen ist allerdings eine Vollerhebung nicht möglich. Deshalb wird eine Zufallsstichprobe realisiert, die als Grundlage für empirische Analysen dient. Für die Berechnung einfacher Lage- und Streuungsmaße (univariate Datenanalyse), die Untersuchung von Zusammenhängen zwischen zwei Merkmalen (bivariate Datenanalyse) sowie die Schätzung von Regressionsmodellen (multivariate Datenanalyse) werden die Daten der Stichprobe genutzt.

Abbildung 2: Grundgesamtheit und Stichprobe



Quelle: Eigene Darstellung.

Bei Zufallsstichproben sind allerdings Stichprobenfehler unvermeidlich. Der Mittel- oder Anteilswert einer Stichprobe wird vom wahren Mittel- oder Anteilswert der Grundgesamtheit abweichen. Ein Beispiel: In der Stichprobe wird ein mittleres Alter von 45,2 Jahren ermittelt. Dieses mittlere Alter wird vom mittleren Alter in der Grundgesamtheit abweichen. Das mittlere Alter in der Grundgesamtheit ist möglicherweise 45,1 Jahre oder 45,3 Jahre, aber vermutlich nicht 45,2 Jahre. Diese Abweichung wird als Stichprobenfehler bezeichnet.

Die Inferenzstatistik stellt uns Instrumente bereit, um zu entscheiden, ob und wie empirische Befunde aus Zufallsstichproben auf zugehörige Grundgesamtheiten übertragen werden dürfen. Die Grundlagen der Inferenzstatistik und die Instrumente werden im Kapitel „Inferenzstatistik“ behandelt. Die grundsätzliche Frage, ob Stichprobenergebnisse auf die Grundgesamtheit übertragen werden dürfen, begegnet uns allerdings bereits bei der Darstellung der univariaten, bivariaten und multivariaten Datenanalyse. Wir werden entsprechende Fragen an den erforderlichen Stellen knapp beantworten und ggf. auf das ausführliche Kapitel am Ende dieses Kurses verweisen.

! An dieser Stelle möchten wir auf zwei häufige Fehler hinweisen, die wir im Zusammenhang mit Stichproben immer wieder beobachten. Erstens: Die Anwendung der Inferenzstatistik setzt eine Zufallsstichprobe voraus. Nur bei einer Zufallsstichprobe kann innerhalb statistischer Fehlergrenzen ein Befund auf die Grundgesamtheit übertragen werden. Zweitens: Bei einem sogenannten Signifikanztest (dabei handelt es sich um ein Instrument der Inferenzstatistik) wird geprüft, ob ein in der Stichprobe gefundener Zusammenhang (sehr) wahrscheinlich auch in der Grundgesamtheit existiert. Ein Befund wird als signifikant bezeichnet, wenn der Befund von der Stichprobe auf die Grundgesamtheit übertragen werden kann. Signifikant bedeutet aber nicht, dass es sich um einen wichtigen oder starken Zusammenhang von zwei Merkmalen (z.B. politisches Interesse und Wahlbeteiligung) handelt.

1.3 Klassifikationen von Variablen

Eine Variable ist ein sozialwissenschaftliches Merkmal mit mindestens zwei Ausprägungen. Das Geschlecht, der allgemeinbildende Schulabschluss oder auch das politische Interesse einer Person sind Beispiele für sozialwissenschaftliche Variablen. Sozialwissenschaftliche Merkmale bzw. Variablen können nach verschiedenen Kriterien klassifiziert werden. Wir unterscheiden vier Kriterien: Skalenniveau, diskrete und stetige Variablen, dichotome und polytome Variablen sowie manifeste und latente Variablen.

Eine wichtige Voraussetzung für die Anwendung bestimmter Analyseverfahren ist das Skalenniveau der Variable bzw. des Merkmals. In den Sozialwissenschaften werden meist die Skalenarten von Stevens (1946) verwendet, der vier Skalenniveaus unterscheidet: Nominal-, Ordinal-, Intervall- und Ratioskala. Intervall- und Ratioskalen werden auch metrische Skalen genannt (Tausendpfund 2018b, S. 119-124). Das jeweilige Skalenniveau bestimmt die zulässigen Rechenoperationen. Je höher das Skalenniveau, desto mehr Rechenoperationen sind möglich.

Verschiedene Skalenniveaus

Das nominale Skalenniveau ist das niedrigste Skalenniveau. Können die Ausprägungen eines Merkmals lediglich im Hinblick auf Gleichheit oder Ungleichheit verglichen werden, dann liegt ein nominales Skalenniveau vor (Gehring und Weins 2009, S. 43-47). Ein Beispiel für eine nominalskalierte Variable ist das Geschlecht. In vielen sozialwissenschaftlichen Datensätzen wird der Ausprägung „weiblich“ die Ziffer 1 und der Ausprägung „männlich“ die Ziffer 2 zugeordnet. Aber diese Zuordnung ist eine Konvention. Man könnte auch 1 für männlich und 2 für weiblich verwenden. Bei einer nominalskalierten Variable stellen die Ziffern lediglich eine Kennzeichnung dar, die nicht richtig oder falsch, sondern allenfalls mehr oder weniger sinnhaft ist. Die Möglichkeiten der quantitativen Datenanalyse bei nominalskalierten Variablen sind daher begrenzt.

Das ordinale Skalenniveau ist das nächsthöhere Skalenniveau. Bei einer ordinalskalierten Variable können die verschiedenen Ausprägungen einer Variable in eine Rangfolge gebracht werden. Beispiele für ordinalskalierte Variablen sind der Schulabschluss oder auch das politische Interesse. Die allgemeine Hochschulreife ist ein höherer Schulabschluss als die Mittlere Reife und die Mittlere Reife ist ein höherer Abschluss als ein Hauptschulabschluss. Ein „sehr starkes“ Interesse für Politik ist ein größeres Interesse als ein „mittleres“ Interesse für Politik. Bei einer ordinalskalierten Variable können zwar die einzelnen Ausprägungen in eine Rangfolge gebracht werden, aber die Abstände zwischen den Ausprägungen (z.B. Abstand zwischen „Hauptschulabschluss“ und „Mittlere Reife“ sowie zwischen „Mittlere Reife“ und „Allgemeine Hochschulreife“) sind nicht gleich. Über die Abstände zwischen den Ausprägungen von ordinalskalierten Variablen sind daher keine Aussagen möglich.

Variablen sind intervallskaliert, wenn deren Ausprägungen nicht nur in eine Rangfolge gebracht werden können, sondern auch die Abstände zwischen den Ausprägungen sinnvoll interpretiert werden können. Ein Beispiel ist die Temperaturmessung in Celsius. Der Abstand zwischen 15 und 20 Grad Celsius ist genau so groß wie der Abstand zwischen 20 und 25 Grad Celsius (jeweils fünf Grad Celsius). Intervallskalen besitzen allerdings keinen natürlichen Nullpunkt. Der Nullpunkt bei der Celsius-Skala wurde lediglich unter pragmatischen Gesichtspunkten gewählt; auch Tempera-

turen im negativen Bereich der Celsius-Skala sind immer noch eine „Temperatur“. Bei einer Intervallskala sind die Abstände zwischen den Merkmalsausprägungen interpretierbar, aber es können keine Verhältnisse berechnet werden.

Bei einer Ratioskala (auch Verhältnisskala genannt) existiert ein natürlicher (echter) Nullpunkt. Die Temperaturmessung in Kelvin erfolgt auf einer Ratioskala, da bei 0 Kelvin keine Temperatur (keine Bewegungsenergie) mehr feststellbar ist. Auch das Einkommen und das Alter sind Beispiele für ratioskalierte bzw. verhältnisskalierte Variablen. Dabei können nicht nur die Abstände zwischen zwei Ausprägungen, sondern auch die Verhältnisse von zwei Ausprägungen interpretiert werden. Ein Einkommen von 5000 Euro ist doppelt so hoch wie ein Einkommen von 2500 Euro. Eine 60jährige Person ist doppelt so alt wie eine 30jährige Person.

In Tabelle 1 sind die zulässigen Rechenoperationen in Abhängigkeit des Skalenniveaus dokumentiert. Wie Tabelle 1 zeigt, steigt mit dem Skalenniveau auch die Anzahl der Rechenoperationen. Bei einem nominalskalierten Merkmal können die Ausprägungen nur ausgezählt werden und bei einem ordinalskalierten Merkmal können die Ausprägungen in eine Reihenfolge gebracht werden. Bei intervallskalierten Variablen können Differenzen, bei ratioskalierten Variablen auch Verhältnisse gebildet werden.

Tabelle 1: Zulässige Rechenoperationen in Abhängigkeit des Skalenniveaus

	auszählen	ordnen	Differenz bilden	Quotienten bilden
Nominalskala	Ja	Nein	Nein	Nein
Ordinalskala	Ja	Ja	Nein	Nein
Intervallskala	Ja	Ja	Ja	Nein
Ratioskala	Ja	Ja	Ja	Ja

Quelle: Mittag (2017, S. 20)

Die Kenntnis des Skalenniveaus einer Variable ist eine wichtige Voraussetzung für die Wahl eines geeigneten Analyseverfahrens. Je höher das Skalenniveau einer Variable, desto mehr (und leistungsfähigere) Analyseverfahren stehen der Sozialwissenschaftlerin zur Verfügung. Die Kenntnis des Skalenniveaus einer Variable ist wichtig, um bei der Datenanalyse nur die zulässigen Analyseverfahren auszuwählen. Viele statistische Verfahren sind nur zulässig, wenn die Variable mindestens intervallskaliert ist.

Diskrete und stetige Variablen

Die Einteilung als diskrete oder stetige Variable basiert auf der Anzahl der möglichen Ausprägungen. Eine diskrete Variable ist eine Variable, die nur endlich viele Ausprägungen oder höchstens „abzählbar“ unendlich viele verschiedene Ausprägungen besitzt (Diaz-Bone 2018, S. 22; Mittag 2017, S. 18). Bei einer diskreten Variable sind keine Zwischenwerte zwischen zwei aufeinander folgenden Ausprägungen möglich. Beispiele für diskrete Variablen sind der Familienstand einer Person, die Anzahl der Fachsemester oder auch die Kinderzahl einer Familie. Bei diesen Variablen sind Zwischenwerte wie 5,6 Fachsemester oder 2,3 Kinder keine möglichen Ausprägungen. Eine stetige Variable ist dadurch gekennzeichnet, dass auch Zwischenwerte möglich sind. Typische Beispiele für stetige Variablen sind Zeit- und Größenangaben, aber auch monetäre Größen wie Einkommen oder Mietpreise. In der Praxis wird bei solchen Merkmalen

aber nur eine begrenzte Anzahl an Nachkommastellen erfasst, beispielsweise werden bei Größenangaben meist nur zwei Nachkommastellen angegeben. Grundsätzlich sind allerdings auch mehr Nachkommastellen möglich.

Eine diskrete Variable, die nur eine geringe Anzahl an Ausprägungen hat, wird als kategoriale Variable bezeichnet (Diaz-Bone 2018, S. 23). Hat eine kategoriale Variable nur zwei mögliche Ausprägungen, dann handelt es sich um eine dichotome Variable. Typische Beispiele für dichotome Variablen sind der Tabakkonsum oder auch die Wahlbeteiligung, bei denen nur die Ausprägungen „Ja“ und „Nein“ möglich sind. Eine diskrete Variable mit mehreren Ausprägungen wird als polytome Variable bezeichnet. Ein Beispiel für eine polytome Variable ist die Zugehörigkeit zu einer Religionsgemeinschaft mit den Ausprägungen „römisch-katholische Kirche“, „evangelische Kirche (ohne Freikirchen)“, „evangelische Freikirche“, „eine andere christliche Religionsgemeinschaft“, „eine andere, nicht-christliche Religionsgemeinschaft“ und „keiner Religionsgemeinschaft“.

Dichotome und polytome Variablen

Schließlich lassen sich auch manifeste und latente Variablen unterscheiden. Bei manifesten Variablen handelt es sich um Merkmale, die direkt beobachtbar sind. Eine manifeste Variable ist beispielsweise das Geschlecht oder die Haarfarbe einer Person. Dagegen handelt es sich bei latenten Variablen um Merkmale, die sich der direkten Beobachtung entziehen. Latente Variablen sind beispielsweise Intelligenz, Einstellungen wie die Zufriedenheit mit der Demokratie oder auch das soziale Vertrauen. Für eine empirische Untersuchung müssen latente Variablen erst „beobachtbar“ gemacht werden. Dieser Vorgang wird als Operationalisierung bezeichnet (Tausendpfund 2018b, S. 11).

Manifeste und latente Variablen