

Studienbrief



Psychologie

Empirisches Forschungsprojekt

Grundlagen der Datenanalyse mit SPSS (III)

Dr. rer. nat. Antonia Zapf
Yvonne Ziert

3 EMF

Verfasserin

Prof. Dr. rer. nat. Antonia Zapf

Studium der Statistik mit Anwendungsgebiet medizinische Biometrie. Promotion an der Universität Göttingen (Abteilung für Medizinische Statistik). Tätigkeit als wissenschaftliche Mitarbeiterin am Institut für Biometrie an der Medizinischen Hochschule Hannover und der Universität Göttingen. Abteilungsleiterin Biometrie am UKE Hamburg. Arbeitsschwerpunkte: klinische Studien und statistische Methoden für Diagnosestudien.

Im Rahmen ihrer bisherigen wissenschaftlichen Tätigkeit hat sie verschiedene Lehrveranstaltungen im Bereich der medizinischen Biometrie durchgeführt.

Dipl. Sozialwiss. Yvonne Ziert, MPH

Studium der Sozialwissenschaften an der Leibniz Universität Hannover und der John Moores University in Liverpool. Master of Public Health an der MHH. Tätigkeit als wissenschaftliche Mitarbeiterin am Institut für Biometrie der MHH.

Der Studienbrief wurde von Frau Ziert durchgesehen und aktualisiert.

Lektorat

Wissenschaftliche Mitarbeiterinnen und Mitarbeiter der Hamburger Fern-Hochschule

Satz/Repro

Haussatz

Redaktionsschluss

Dezember 2019

1. Auflage 2019

© HFH · Hamburger Fern-Hochschule, Alter Teichweg 19, 22081 Hamburg

Das Werk ist urheberrechtlich geschützt. Die dadurch begründeten Rechte, insbesondere das Recht der Vervielfältigung und der Verbreitung sowie der Übersetzung und des Nachdrucks, bleiben, auch bei nur auszugsweiser Verwertung, vorbehalten. Kein Teil des Werkes darf in irgendeiner Form ohne schriftliche Genehmigung der Hamburger Fern-Hochschule reproduziert oder unter Verwendung elektronischer Systeme verarbeitet, vervielfältigt oder verbreitet werden.

Gedruckt auf 100% chlorfrei gebleichtem Papier.

Inhaltsverzeichnis

Abkürzungsverzeichnis	4
Einleitung	5
1 Theoretische Grundlagen der multivariaten Statistik	7
1.1 Ziele der multivariaten Statistik.....	7
1.2 Definition des statistischen Modells	8
1.3 Variablenselektion in multivariaten Analysen.....	10
Übungsaufgaben.....	11
2 Varianzanalyse.....	12
2.1 Voraussetzungen für die Varianzanalyse.....	13
2.2 Varianzanalyse (ANOVA).....	13
2.2.1 Einfaktorielle ANOVA.....	13
2.2.2 Mehrfaktorielle ANOVA.....	18
2.3 Kovarianzanalyse (ANCOVA).....	21
Übungsaufgaben.....	26
3 Regressionsanalyse	27
3.1 Lineare Regression	27
3.2 Logistische Regression	35
Übungsaufgaben.....	40
Zusammenfassung	41
Glossar	42
Lösungen zu den Übungsaufgaben	44
Literaturverzeichnis	51

Abkürzungsverzeichnis

BMI	Body Mass Index
m	männlich
OR	Odds Ratio
SQE	Sum of squares explained
SQR	Sum of squares residual
SQT	Sum of squares total
vs.	versus (bei Vergleichen)
w	weiblich

Einleitung

Mit dem Wissen, das mit den ersten Studienbriefen dieses Moduls erworben wurde, können Daten in SPSS eingegeben bzw. eingelesen und bei Bedarf Modifikationen vorgenommen werden. Weiterhin können deskriptive Analysen durchgeführt werden, die eine wichtige Grundlage für alle konfirmatorischen Analysen sind. Zudem wurden bivariate Analysen vorgestellt. Dabei wurde das Vorliegen eines Zusammenhangs zwischen zwei Variablen u. a. mit Hilfe von Korrelationsmaßen sowie des Chi-Quadrat- und t-Tests überprüft.

Die in der Psychologie betrachteten Fragestellungen und Zusammenhänge sind jedoch oftmals noch komplexer – ebenso wie die Daten, die einem zur Verfügung stehen. Bei der Evaluation von Interventionsmaßnahmen etwa ist es denkbar, dass beispielsweise die Patientenzufriedenheit nicht nur vom Therapieerfolg abhängt, sondern auch noch von weiteren Faktoren beeinflusst wird (z. B. persönliche Merkmale der Patienten). Statistische Methoden, die Zusammenhänge mehrerer Variablen gemeinsam überprüfen, werden multivariate Analyseverfahren genannt.

In diesem Studienbrief soll die Durchführung multivariater Analysen mit dem Statistikprogramm SPSS demonstriert werden. Die entsprechenden theoretischen Grundlagen wurden bereits in den Modulen zu Forschungsmethodik und Statistik vorgestellt und werden hier an den entsprechenden Stellen noch einmal kurz rekapituliert. Im vorliegenden Studienbrief wird die Anwendung der gängigsten multivariaten Verfahren erläutert. Andere Methoden oder mögliche Erweiterungen der Verfahren werden in kurzen Exkursen abgehandelt. Für die Anwendung weiterer multivariater Verfahren in SPSS sei hier auf die Literatur (z. B. Fromm, 2012, oder Rudolf & Müller 2012) verwiesen.

Als Beispieldaten werden wieder die Ergebnisse der fiktiven Studie zur Patientenzufriedenheit im Klinikverbund „Nordsterne“ verwendet (Datensatz *Zentren_gesamt.sav*). Diese Daten werden auch für die Übungsaufgaben verwendet.

Im *ersten Kapitel* werden die theoretischen Grundlagen der multivariaten Analyseverfahren dargestellt. Dazu gehören die Beschreibung der Ziele der multivariaten Verfahren und die Erläuterung der verschiedenen Komponenten, aus denen sich das statistische Modell zusammensetzt. Weiterhin werden Methoden der Variablenselektion bei multivariaten Analysen präsentiert, da solche Analysen schnell unübersichtlich werden können.

Im *zweiten Kapitel* wird die Varianzanalyse (ANOVA) vorgestellt und deren Anwendung in SPSS demonstriert. Neben der einfaktoriellen ANOVA soll zusätzlich die Anwendung der mehrfaktoriellen ANOVA und die Kovarianzanalyse (ANCOVA) demonstriert werden. In einem Exkurs wird auf den Effekt von Messwiederholungen (repeated measures ANOVA) eingegangen und erläutert, wie mehrere Zielgrößen mit Hilfe einer multivariaten Varianzanalyse (MANOVA) berücksichtigt werden können.

Das *dritte Kapitel* beschäftigt sich mit der Regressionsanalyse. Die häufigsten Varianten der Regressionsanalyse sind die lineare Regression für metrische Zielgrößen und die binäre logistische Regression für dichotome Zielgrößen. Die Durchführung dieser Regressionsanalysen mit dem Statistikprogramm SPSS wird ausführlich behandelt.

Anwendung multivariater Verfahren

Beispieldatensatz

Kapitel 1

Kapitel 2

Kapitel 3

Für alle vorgestellten statistischen Verfahren wird erklärt, wie mit Hilfe von SPSS überprüft werden kann, ob die (aus statistischer Sicht) nötigen Voraussetzungen erfüllt sind. Weiterhin wird Schritt für Schritt erläutert wie

- die Analysen durchgeführt werden,
- der SPSS-Output (Ergebnisfenster, auch Viewer genannt) aufgebaut ist und
- die Ergebnisse dargestellt und interpretiert werden können.

In diesem Studienbrief wird auf die SPSS-Syntax nicht genauer eingegangen. Es wird jedoch für alle in diesem Studienbrief durchgeführten Analysen die Syntax in kommentierter Form im WebCampus zur Verfügung gestellt.

Alle im Studienbrief verwendeten Abbildungen (*Screenshots*) wurden mit dem Programm *IBM SPSS Statistics 25* erzeugt.

Studienziele

Nach der Bearbeitung des vierten Studienbriefs sind die Studierenden in der Lage:

- ⇒ die Ziele der multivariaten Statistik zu erläutern,
- ⇒ die Komponenten eines statistischen Modells zu erklären und am Beispiel anzuwenden,
- ⇒ verschiedene Varianzanalysen (inklusive der Überprüfung der Voraussetzungen) durchzuführen und die Ergebnisse zu interpretieren,
- ⇒ lineare und logistische Regressionsanalysen mit Hilfe von SPSS durchzuführen und die Ergebnisse zu interpretieren,
- ⇒ die Vorgehensweise bei der multivariaten Datenanalyse wissenschaftlich angemessen darzustellen.

1 Theoretische Grundlagen der multivariaten Statistik

Wie bereits in der Einleitung erklärt wurde, konzentriert sich dieser Studienbrief auf die wichtigsten multivariaten Analyseverfahren: Varianzanalyse und Regressionsanalyse. In diesem Studienbrief soll es vorrangig darum gehen, wie diese Verfahren mit Hilfe von SPSS umgesetzt werden können. Der theoretische Hintergrund wird dabei jeweils kurz wiederholt.

In Kapitel 1 geht es zunächst um die theoretischen Grundlagen, die nicht verfahrensspezifisch sind. Abschnitt 1.1 behandelt zum Einstieg die grundsätzlichen Ziele der multivariaten Statistik.

Sowohl bei der Varianz- als auch bei der Regressionsanalyse gibt es verschiedene Varianten. Welche Variante sich für die Analyse eignet, entscheidet sich auf Basis des zugrundeliegenden statistischen Modells. Die verschiedenen Komponenten einer statistischen Modellierung werden in Abschnitt 1.2 vorgestellt. Die verschiedenen Möglichkeiten der Variablenselektion in multivariaten Analysen werden in Abschnitt 1.3 thematisiert.

1.1 Ziele der multivariaten Statistik

Der Vorteil multivariater Verfahren besteht darin, dass der Zusammenhang mehrerer Variablen gemeinsam untersucht werden kann. Das ist insbesondere dann von Interesse, wenn Einflussgrößen stark korreliert sind (Multikollinearität, für Details vgl. Fromm, 2012). So kann es sein, dass in bivariaten Analysen Zusammenhänge beobachtet werden, die in Wirklichkeit durch den Einfluss einer dritten Variablen zustande kommen (sog. Scheinkorrelation).

Vorteil

Ein oft verwendetes Beispiel für eine **Scheinkorrelation** ist die Beobachtung, dass in Ländern mit weniger Störchen auch weniger Kinder geboren werden. Das könnte zu dem Schluss führen, dass die Störche die Kinder bringen. Der wahre Grund für diesen beobachteten Zusammenhang ist natürlich ein anderer: In Ländern mit einem höherem Industrialisierungsgrad nimmt die Anzahl der Störche ab, da ihr Lebensraum zerstört wird. Andererseits ist bekannt, dass die Geburtenrate in Ländern mit zunehmender Industrialisierung sinkt. Untersucht man den Einfluss der Anzahl der Störche und des Industrialisierungsgrads auf die Anzahl der Geburten einzeln, kann man für beide Einflussgrößen einen Zusammenhang zeigen. Untersucht man aber den Einfluss der Storchenzahl und des Industrialisierungsgrads gemeinsam, bleibt der Effekt der Industrialisierung bestehen, während der der Storchenzahl verschwindet.

Beispiel

In diesem Beispiel sind die tatsächlichen Zusammenhänge mit Hilfe des gesunden Menschenverstands schnell aufgedeckt. In realen Datensätzen ist es oftmals wesentlich schwieriger, die Wahrheit herauszufinden. Denn wenn ein Datensatz Dutzende oder sogar hunderte möglicher Einflussgrößen beinhaltet und dabei wenig Vorwissen über Zusammenhänge zwischen den Variablen existiert, kann man sich in einer multivariaten Analyse schnell verlieren.

Im nächsten Abschnitt wird erläutert, welche Komponenten es gibt, die in der Analyse u. U. berücksichtigt werden müssen, und aus denen sich das statistische Modell ergibt.

1.2 Definition des statistischen Modells

Komponenten eines statistischen Modells

Ein statistisches Modell ergibt sich aus den verschiedenen Komponenten, die berücksichtigt werden sollen. Die wichtigsten Komponenten sind die **Zielgröße** (auch Zielvariable genannt) und die **Einflussgröße(n)** (auch Einflussvariable(n) genannt).

Ziel

Mit der statistischen Analyse soll der Effekt untersucht werden, den die Einflussgröße auf die Zielgröße hat.

Beispiel

Wenn man den **Einfluss des Raucherstatus auf die Zufriedenheit mit einer medizinischen Rehabilitationsmaßnahme** (im folgenden Reha-Maßnahme genannt) untersuchen will, stellt der Raucherstatus die Einflussgröße dar und die Zufriedenheit die Zielgröße.

Es ist darauf zu achten, dass die Einflussgröße zeitlich vor oder zugleich mit der Zielgröße erhoben wurde. Es würde z. B. keinen Sinn machen, zu untersuchen, ob eine Gewichtsreduktion während der Reha-Maßnahme einen Effekt auf den Raucherstatus zu Beginn der Reha hat.

Durch die Skalierung der Einfluss- und Zielgröße wird das in Frage kommende Analyseverfahren bestimmt (siehe **Tabelle 1.1**).

Tabelle 1.1: Skalierung der Variablen und Wahl des geeigneten Analyseverfahrens

Zielgröße	Einflussgröße	Analyseverfahren
metrisch	metrisch	Lineare Regressionsanalyse
metrisch	kategorial	Varianzanalyse Lineare Regressionsanalyse
kategorial	metrisch oder kategorial	Logistische Regressionsanalyse
binär (kategorial)	metrisch oder kategorial	Binäre logistische Regressionsanalyse

Kategoriale Einflussgrößen werden auch **Einflussfaktoren** genannt.

Anzahl von Einfluss- und Zielgrößen

Ein weiterer Faktor, der das statistische Modell definiert, ist die Anzahl von Einfluss- und Zielgrößen.

- Soll lediglich eine Einflussgröße berücksichtigt werden, spricht man von **einfaktorieller ANOVA bzw. Einfachregression**.
- Analysen mit mehreren Einflussgrößen werden als **mehrfaktorielle ANOVA bzw. multiple Regression** bezeichnet.
- Sollen mehrere Zielgrößen gemeinsam modelliert werden, spricht man von **multivariater ANOVA bzw. multivariater Regression**.
- Wenn die Zielvariable mehrmals erhoben wurde (sog. **Messwiederholungen**), spricht man von **repeated measures ANOVA bzw. generalisierten Regressionsmodellen**.

Die multivariate ANOVA und die repeated measures ANOVA werden in einem Exkurs (siehe Abschnitt 2.3) grob skizziert.

Wenn es Variablen gibt, die einen nicht primär interessieren, von denen aber angenommen wird, dass sie ebenfalls einen Einfluss auf die Zielgröße haben, können diese als **Kovariablen** im Regressionsmodell bzw. in der Kovarianzanalyse berücksichtigt werden. Dadurch wird der Effekt der eigentlich interessierenden

Einflussgröße um den Effekt der Kovariablen bereinigt. Man spricht in diesem Zusammenhang auch davon, dass für die Kovariable kontrolliert wurde.

Ein weiterer möglicher Störeffekt, der insbesondere mit Varianzanalysen aufgedeckt werden kann, ist die sog. **Interaktion** (Wechselwirkung) zwischen zwei Einflussfaktoren. Wenn eine Interaktion vorliegt, müssen die beiden interagierenden Faktoren getrennt ausgewertet werden.

Eine **Gewichtsreduktion** wirkt sich **bei unter- und übergewichtigen Menschen** unterschiedlich auf deren Gesundheitszustand aus. Während eine Gewichtsreduktion bei untergewichtigen Menschen zu einem schlechteren Gesundheitszustand führt, hat sie bei Übergewichtigen eine Verbesserung des Gesundheitszustands zur Folge. Entsprechend müssen Untergewichtige getrennt von Übergewichtigen analysiert werden, weil die Effekte in der mehrfaktoriellen ANOVA ansonsten nicht sinnvoll interpretiert werden können.

Wechselwirkung zwischen Einflussfaktoren

Beispiel

Weitere Erläuterungen zum Thema Interaktion finden sich z. B. bei Hartung et al. (2005) oder bei Janssen und Laatz (2013).

Literaturhinweis

Je nach Datenlage und Fragestellung ergibt sich dann das entsprechende statistische Modell mit der dazugehörigen Modellgleichung. Es sei darauf hingewiesen, dass in diesem Studienbrief nicht alle existierenden statistischen Modelle vorgestellt werden können. Wir konzentrieren uns hier auf die Varianz- und die Kovarianzanalyse sowie auf die lineare und binäre logistische Regressionsanalyse als gängigste Verfahren. Als weiterführende Literatur wird auf Hartung et al. (2005) oder Schumacher und Schulgen (2008) für die Theorie und Fromm (2012) und Rudolf und Müller (2012) für die Anwendung in SPSS verwiesen.

Herleitung von statistischen Modellen für zwei Fragestellungen aus der fiktiven Studie zur Patientenzufriedenheit in den Reha-Zentren des Klinikverbunds „Nordsterne“:

- 1.) Es soll untersucht werden, welche Parameter mit der Entscheidung der Ärzte und Ärztinnen zusammenhängen, ob Angehörige in die Therapie einbezogen werden. Die Entscheidung über die Einbeziehung von Angehörigen wird zu Beginn der Reha-Maßnahme getroffen. Es bestehen keine Vermutungen darüber, was die ärztliche Entscheidung beeinflusst. Es sollen folgende Einflussgrößen berücksichtigt werden: Geschlecht, Alter, Raucherstatus vor der Reha und BMI vor der Reha. Beim Raucherstatus wird hier nur zwischen Raucher (inkl. Gelegenheitsraucher) und Nichtraucher unterschieden. Die Zielgröße (Einbeziehung von Angehörigen) ist binär. Die Einflussgrößen sind teils metrisch, teils kategorial. Damit ergibt sich die binäre logistische Regression als statistisches Modell.
- 2.) Es soll der Einfluss des Raucherstatus und des Geschlechts gemeinsam auf die Zufriedenheit mit der Reha-Maßnahme untersucht werden. Dabei besteht die Vermutung, dass der Raucherstatus sich bei Männern und Frauen unterschiedlich auf die Zufriedenheit auswirkt. Weiterhin wird davon ausgegangen, dass die Gewichtsabnahme einen Einfluss auf die Zufriedenheit hat – dieser Zusammenhang ist allerdings nicht primär von Interesse und wird daher hier auch nicht weiter untersucht.

Da die Zufriedenheit als Zielgröße metrisch und der Raucherstatus und das Geschlecht als Einflussgrößen (Einflussfaktoren) kategorial skaliert sind, bietet sich eine Varianzanalyse an. Die Gewichtsreduktion sollte als Kovariable berücksichtigt werden. Es ergibt sich als statistisches Modell also eine mehrfaktorielle Kovarianzanalyse mit Zufriedenheit als Zielgröße, Raucherstatus und Geschlecht als Einflussgrößen und Gewichtsreduktion als Kovariable. Zusätzlich wird die Interaktion zwischen Raucherstatus und Geschlecht in das Modell mit aufgenommen.

Beispiel

1.3 Variablenselektion in multivariaten Analysen

Es wurde bereits erwähnt, dass multivariate Analysen schnell unübersichtlich werden können. Wenn es, wie im ersten Beispiel im vorherigen Abschnitt, keine konkreten Vermutungen über mögliche Zusammenhänge gibt und viele potenzielle Einflussgrößen existieren, braucht man Strategien, um die Anzahl der Variablen zu reduzieren.

Prozeduren der Variablenselektion

In SPSS gibt es verschiedene Prozeduren zur Variablenselektion. Die drei wichtigsten sind:

- Vorwärtss Selektion
- Rückwärtss Selektion
- schrittweise Selektion

Bei der **Vorwärtss Selektion** wird das Modell Stück für Stück aufgebaut, wohingegen bei der **Rückwärtss Selektion** das komplette Modell (mit allen Einflussgrößen) Stück für Stück reduziert wird. Die **schrittweise Selektion** ist die Kombination der Vorwärts- und der Rückwärtss Selektion. Alle drei Selektionsverfahren haben jedoch ihre Schwächen. Eine ausführliche Diskussion der Verfahren ist z. B. bei Rudolf und Müller (2012) oder bei Gaus und Muche (2014) zu finden.

Zweistufiges Verfahren der Variablenselektion

Für den Fall vieler Einflussgrößen wird an dieser Stelle darüber hinaus ein zweistufiges Verfahren empfohlen. Denn wenn die Anzahl der Einflussgrößen im Verhältnis zur Anzahl an Beobachtungen (z. B. Patienten) zu groß ist, sind die Modellparameter nicht schätzbar – das Modell „bricht zusammen“.

- Die *erste Stufe* beinhaltet die Überprüfung, ob die verschiedenen Variablen einzeln einen Einfluss auf die Zielgröße zu haben scheinen. Dies geschieht mit Hilfe der **einfaktoriellen ANOVA oder Einfachregression**. Dabei ist man jedoch nicht so streng, dass man Signifikanz ($p < 0,05$) fordert. Es genügt, wenn sich eine Tendenz zur Signifikanz (z. B. $p < 0,2$ oder $p < 0,3$) zeigt. Dadurch verringert man das Risiko, dass Variablen, die eigentlich wichtig sind, deren Effekt aber nicht groß genug für eine Signifikanz ist, verloren gehen. Durch diese Vorauswahl kann man im Allgemeinen die Anzahl der möglichen Einflussgrößen auf die Zielgröße ganz erheblich reduzieren.
- In der *zweiten Stufe* werden die ausgewählten Variablen gemeinsam in ein multivariates Modell (**mehrfaktorielle ANOVA oder multiples Regressionsmodell**) aufgenommen. Insbesondere bei der Varianzanalyse sollten zusätzlich zu den Einflussgrößen für alle kategorialen Variablen (Einflussfaktoren) auch die paarweisen Interaktionen in das Modell mit aufgenommen werden. Vergleicht man nun die Schätzer der Einflussgrößen aus dem einfaktoriellen und dem mehrfaktoriellen Modell, bekommt man bereits einen Eindruck von den Zusammenhängen zwischen den Einflussgrößen. Bleibt der Schätzer einer Einflussvariablen weitgehend unverändert, spricht das dafür, dass die Korrelation dieser Variable und der anderen Einflussgrößen gering ist. Um das mehrfaktorielle Modell nun so zu reduzieren, dass am Ende nur die relevanten Einflussgrößen übrigbleiben, kann man dann die oben erwähnten Verfahren, die SPSS zur Verfügung stellt, anwenden. Da man ja bereits das komplette Modell aufgestellt hat, bietet sich hier insbesondere die Rückwärtss Selektion an.

Für das Beispiel 1.) in Abschnitt 1.2 zur Herleitung von statistischen Modellen (Einfluss verschiedener Variablen auf die ärztliche Entscheidung, Angehörige in die Behandlung einzubeziehen) soll eine binäre logistische Regressionsanalyse durchgeführt werden. Die Zielgröße ist die **Einbeziehung von Angehörigen in die Behandlung**. Die möglichen Einflussgrößen sind: Geschlecht, Alter, Raucherstatus vor der Reha und BMI vor der Reha.

Die binäre logistische Einfachregression führt dann in der ersten Stufe zu folgenden Ergebnissen (siehe **Tabelle 1.2**):

Beispiel

Tabelle 1.2: Ergebnisse der binären logistischen Einfachregression bezüglich der Einbeziehung von Angehörigen

Einflussgröße	Odds Ratio	p-Wert	Aufnahme in das mehrfaktorielle Modell
Geschlecht (w vs. m)	1,565	0,353	Nein
Alter	1,093	0,195	Ja
Raucherstatus vor Reha (Nichtraucher vs. Raucher)	0,327	0,020	Ja
BMI vor Reha	1,001	0,995	Nein

Das komplette mehrfaktorielle Regressionsmodell mit allen Variablen, die in der Einfachregression tendenziell signifikant waren, führt zu folgenden Ergebnissen (siehe **Tabelle 1.3**):

Tabelle 1.3: Ergebnisse der multiplen binären logistischen Regression bezüglich der Einbeziehung von Angehörigen

Einflussgröße	Odds Ratio	p-Wert
Alter	1,084	0,256
Raucherstatus vor Reha (Nichtraucher vs. Raucher)	0,338	0,025

Mit der Rückwärtsselektion erhält man als endgültiges Ergebnis (siehe Syntax), dass lediglich der Raucherstatus vor der Reha-Maßnahme einen signifikanten Einfluss auf die Entscheidung der Ärzte und Ärztinnen hat, ob Angehörige in die Therapie mit einbezogen werden ($p\text{-Wert}=0,020$, $\text{Odds Ratio}=0,327$, 95 %-Konfidenzintervall $\text{OR}=0,127-0,841$). Dieses Ergebnis ist auch nachvollziehbar, da die Unterstützung der Angehörigen sehr wichtig ist, wenn das Rauchen aufgegeben werden soll.

Übungsaufgaben

- 1.1) Interpretieren Sie das Ergebnis bezüglich der Rolle des Raucherstatus vor der Reha-Maßnahme auf die ärztliche Entscheidung hat, ob Angehörige in die Therapie einbezogen werden: $p\text{-Wert}=0,020$, $\text{Odds Ratio}=0,327$, 95 %-Konfidenzintervall $\text{OR}=0,127-0,841$).
- 1.2) Es soll der Einfluss bestimmter Variablen, die den Verlauf der Reha betreffen, auf die Zufriedenheit mit der Reha-Maßnahme untersucht werden. In Frage kommen die durchschnittliche Zeit, die pro Tag Sport getrieben wird, die Veränderung des BMI und die Anzahl der wahrgenommenen Freizeitangebote. Welches Verfahren eignet sich für die Analyse und welche Komponenten beinhaltet das statistische Modell?

2 Varianzanalyse

Ziel Mit Hilfe der Varianzanalyse als Verallgemeinerung des t-Tests können die Mittelwerte mehrerer Gruppen verglichen werden. Das Prinzip der **Varianzanalyse** ist die sog. **Streuungszerlegung**. Das heißt, dass die gesamte Streuung (*sum of squares total*, SQT) aufgeteilt werden kann auf die Streuung innerhalb der Gruppen (*sum of squares residual*, SQR) und die Streuung zwischen den Gruppen (*sum of squares explained*, SQE). Die Teststatistik des F-Tests ist dann genau der Quotient der beiden Anteile ($F = \text{SQE} / \text{SQR}$). Je größer SQE im Verhältnis zu SQR ist, desto größer wird entsprechend die Teststatistik und desto kleiner wird der p-Wert, d.h., desto eher kann die Nullhypothese abgelehnt werden (Hartung et al., 2005; Janssen und Laatz, 2013).

Formen der Varianzanalyse

Es gibt verschiedene Formen der Varianzanalyse, die sich aus dem zugrunde liegenden statistischen Modell (siehe Abschnitt 1.2) ergeben. In Anlehnung an Fromm (2012) sind die verschiedenen Formen in **Tabelle 2.1** dargestellt. In diesem Kapitel wird die Anwendung der **ein- und mehrfaktoriellen Varianzanalyse** (ANOVA, Abschnitt 2.2) und der **Varianzanalyse mit Kovariablen** (ANCOVA, Abschnitt 2.3) in SPSS demonstriert. Die Durchführung einer **Varianzanalyse mit Messwiederholungen** (repeated measures ANOVA) und einer **multivariaten ANOVA** (MANOVA) in SPSS wird in einem Exkurs in Abschnitt 2.3 grob skizziert.

Hinweis

An dieser Stelle sei darauf hingewiesen, dass mit **multivariaten Analysen** üblicherweise (und dies gilt auch für diesen Studienbrief) Analysen gemeint sind, die **mehrere Variablen** berücksichtigen, egal ob es sich dabei um Einfluss- oder Zielvariablen handelt.

In SPSS dagegen bezeichnet die **multivariate ANOVA** ausschließlich eine ANOVA, bei der der Einfluss einer oder mehrerer Variablen auf **mehrere Zielgrößen** gemeinsam analysiert wird, wohingegen eine **multifaktorielle ANOVA** in SPSS eine ANOVA ist, bei der der Zusammenhang **mehrerer Einflussgrößen** und einer Zielgröße untersucht wird.

Tabelle 2.1: Verschiedene Formen der Varianzanalyse mit zugehörigen SPSS-Prozeduren

Bezeichnung des Verfahrens	Zahl der Zielvariablen	Zahl der Einflussvariablen	SPSS-Prozeduren
Einfaktorielle Varianzanalyse	1	1	Einfaktorielle ANOVA, Univariat
Mehrfaktorielle Varianzanalyse	1	>1	Univariat
Kovarianzanalyse	1	>1 (davon mind. 1 metrisch)	Univariat
Varianzanalyse mit Messwiederholungen	1 (mit Messwiederholungen)	≥ 1	Messwiederholung
Multivariate Varianzanalyse	>1	≥ 1	Multivariat

In einem ersten Schritt muss überprüft werden, ob die Voraussetzungen für eine Varianzanalyse erfüllt sind. Anschließend werden die Analysen durchgeführt, und die Ergebnisse zusammengefasst und interpretiert.

2.1 Voraussetzungen für die Varianzanalyse

Eine Varianzanalyse kann durchgeführt werden, wenn bestimmte Voraussetzungen erfüllt sind:

1. Die **Zielgröße** muss **metrisch oder quasi-metrisch** skaliert sein. Als quasi-metrisch wird eine ordinale Variable mit hinreichend vielen Kategorien (ca. ≥ 10) und gleichen Abständen zwischen den Ausprägungen bezeichnet.
2. In der **Gesamtpopulation** muss die **Zielgröße normalverteilt** sein. Ob diese Annahme erfüllt ist, kann mit Hilfe von Kenngrößen und/oder Grafiken beurteilt werden. Eine schiefe Verteilung kann u. U. durch eine geeignete Transformation (z. B. Wurzel- oder Logarithmus-Transformation) in eine Normalverteilung überführt werden (für Details siehe Hartung et al., 2005).
3. Es muss **mindestens eine kategoriale** (d. h. nominale oder ordinale) **Einflussgröße** geben (auch Einflussfaktor genannt). Metrische Einflussgrößen können durch sinnvolle Gruppierung auf eine kategoriale Skala übertragen werden. Mit einem Einflussfaktor ergibt sich die einfaktorielle ANOVA, mit mehreren die mehrfaktorielle ANOVA.
4. Die Varianzen in den Untergruppen, die sich durch die kategoriale/n Einflussgröße/n ergeben, müssen vergleichbar sein (bezeichnet als **Varianzhomogenität**). Ein $p\text{-Wert} \leq 0,2$ des Levene-Tests würde stark für eine Varianzheterogenität sprechen. Bei einem $p\text{-Wert} > 0,2$ können wir umgekehrt davon ausgehen, dass keine nennenswerte Varianzheterogenität vorliegt.
5. Bei Durchführung einer **ANCOVA** muss zusätzlich zu dem/den Einflussfaktor/en **mindestens eine metrische Kovariable** vorliegen.
6. Eine **repeated measures ANOVA** kommt zur Anwendung, wenn die **Zielgröße mindestens zweimal pro Person gemessen** wird.
7. Liegen **mehrere Zielgrößen** vor, die gemeinsam berücksichtigt werden sollen, führt das zur **MANOVA**.

Bei den verschiedenen SPSS-Prozeduren zur Varianzanalyse gibt es teilweise die Möglichkeit, Kriterien überprüfen zu lassen. Wie das geht, wird an den entsprechenden Stellen im Studienbrief erläutert.

2.2 Varianzanalyse (ANOVA)

2.2.1 Einfaktorielle ANOVA

Der einfachste Fall einer Varianzanalyse ist die **einfaktorielle ANOVA**. Auch wenn die einfaktorielle ANOVA im strengen Sinne kein multivariates (sondern ein bivariates) Verfahren ist, wird die Anwendung in SPSS im Folgenden demonstriert, da sie die Grundlage für alle weiteren Varianzanalysen darstellt. In Studienbrief 2 wurde die einfaktorielle ANOVA in einem Exkurs bereits grob erläutert.

Im Folgenden wird gezeigt, wie überprüft werden kann, ob die Voraussetzungen für die einfaktorielle ANOVA erfüllt sind und wie die Analyse durchgeführt wird. Zusätzlich werden die Inhalte der SPSS-Ausgabe und ihre Interpretation erläutert.

Einfaktorielle ANOVA

Beispiel

Das Ziel der Reha-Studie war, herauszufinden, welche Faktoren einen Einfluss auf die Patientenzufriedenheit haben. An dieser Stelle soll deshalb exemplarisch überprüft werden, ob es einen **Zusammenhang zwischen der Zufriedenheit mit der Reha-Maßnahme und dem Raucherstatus** (zu Beginn der Reha) der Patienten und Patientinnen gibt.

Da für die Analyse eine einfache ANOVA verwendet wird, müssen lediglich Voraussetzungen 1 bis 4 erfüllt sein.

- Die **Zielgröße** ist der Zufriedenheitsscore, die Einflussgröße ist der Raucherstatus zu Beginn der Reha.
- Der Zufriedenheitsscore ist eigentlich ein **ordinales Merkmal**. Da es sich dabei aber um einen Summenscore handelt, der 32 Kategorien (32=Maximalpunktzahl) mit gleichen Abständen aufweist, kann die Variable als **quasi-metrisch** bezeichnet werden.
- Der Raucherstatus als **Einflussgröße** wird in drei Kategorien erhoben:
 - Nichtraucher
 - Gelegenheitsraucher
 - Raucher

Beim Raucherstatus handelt es sich also um eine kategoriale Variable, d.h. einen Einflussfaktor. Voraussetzungen 1 und 3 sind also erfüllt. Voraussetzungen 2 und 4 können im Rahmen der Analyse überprüft werden.

Aus dem Menü ist die folgende Prozedur zu wählen: *Analysieren → Mittelwerte vergleichen → Einfaktorielle Varianzanalyse*. In dem sich öffnenden Fenster (siehe **Abb. 2.1**) zieht man die Zielgröße (hier Variable `v2_v9_Score`) in das Feld *Abhängige Variablen*, und den Einflussfaktor (hier `v14_1_rauch_vor`) in das Feld *Faktor*.

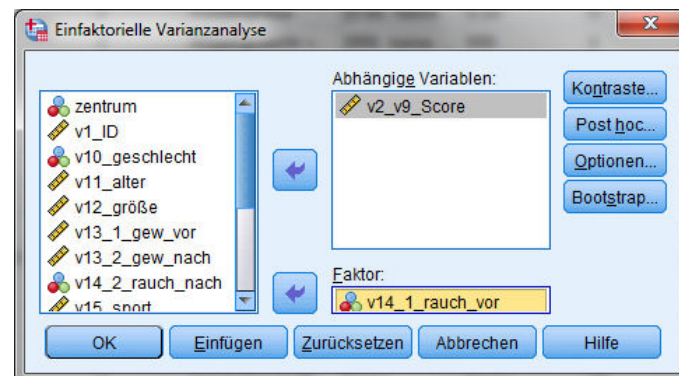


Abb. 2.1: Hauptfenster für die einfaktorielle ANOVA – Auswahl der Variablen

Würde man jetzt die Auswahl mit *OK* bestätigen, würde eine einfaktorielle ANOVA durchgeführt werden. Da wir jedoch auch das Vorliegen von Normalverteilung und Varianzhomogenität überprüfen wollen, wählen wir in dem Feld *Optionen* *Deskriptive Statistik* zur Überprüfung der Normalverteilung und *Test auf Homogenität der Varianzen* zur Überprüfung der Voraussetzung 4. Zusätzlich wählen wir noch *Diagramm der Mittelwerte* aus, um auch eine grafische Darstellung der Ergebnisse zu erhalten (siehe **Abb. 2.2**). Die Auswahl bestätigen wir mit *Weiter*.

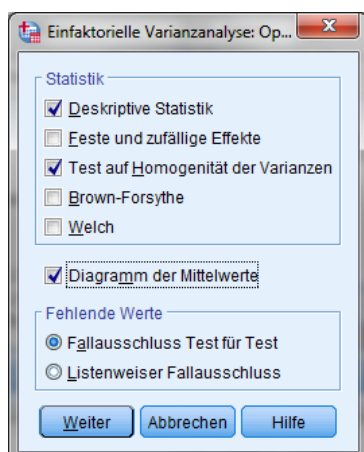


Abb. 2.2: Festlegung der Optionen bei einer einfaktoriellen ANOVA

Mit der ANOVA erhält man nur die Information, ob sich mindestens zwei Gruppen unterscheiden. Ist das der Fall, möchte man im Regelfall aber auch wissen, welche Gruppen sich genau unterscheiden (für weitere Informationen zum Thema post-hoc-Tests siehe Studienbrief 5 des Moduls Statistik II oder Janssen und Laatz, 2013). Dafür klicken wir auf das Feld *Post Hoc*. Dadurch öffnet sich ein Fenster, in dem man die gewünschte Paarvergleich-Methode auswählen kann (siehe Abb. 2.3). Da hier alle Paarvergleiche durchgeführt werden sollen und die Bonferroni-Korrektur bereits bekannt ist, wählen wir *Bonferroni*. Das *Signifikanzniveau*, das den α -Fehler bezeichnet, belassen wir bei 0,05.

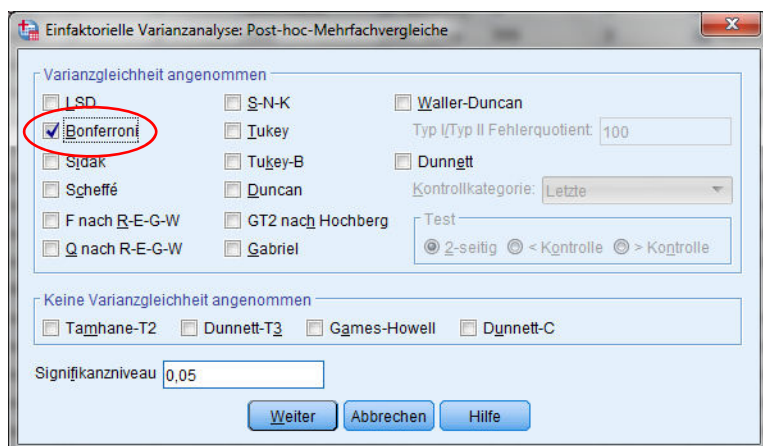


Abb. 2.3: Bestimmung der Post-Hoc-Mehrfachvergleiche in der einfaktoriellen Varianzanalyse

Wenn wir nun die Auswahl mit *Weiter* bestätigen und im Hauptfenster auf *OK* klicken, erhalten wir im Ausgabefenster die Ergebnisse in drei Abschnitten:

- Univariat
- Post-Hoc-Tests
- Mittelwert-Diagramme

**Ausgabe der Ergebnisse
in drei Abschnitten**

Im Abschnitt *Univariat* ist als erstes die Tabelle „ONEWAY deskriptive Statistiken“ angegeben (siehe Abb. 2.4). Da Voraussetzung 2 nur die Normalverteilung in der Gesamtpopulation verlangt, genügt es, die Ergebnisse in der Zeile *Gesamt* zu betrachten. Liegt der *Mittelwert* ungefähr in der Mitte zwischen *Minimum* und *Maximum*, kann i. A. von einer Normalverteilung ausgegangen werden. Da das hier der Fall ist (19,49 liegt ziemlich exakt in der Mitte zwischen 10 und 29), betrachten wir

also Voraussetzung 2 als erfüllt. Möchte man eine zuverlässigere Aussage treffen, sollte man ein Histogramm zeichnen (zur Durchführung siehe Studienbrief 2).

Die zweite Tabelle in der Ausgabe („Test der Homogenität der Varianzen“) beinhaltet die Ergebnisse des Levene-Tests (siehe **Abb. 2.4**). Dabei werden unterschiedliche Levene-Statistiken ausgegeben, von denen für uns die erste („Basiert auf dem Mittelwert“) von Bedeutung ist. Dabei interessiert uns v. a. die *Signifikanz* – in diesem Feld steht der p-Wert. Da der p-Wert hier bei 0,256 liegt, können wir also auch Voraussetzung 4 als erfüllt ansehen, d. h. es kann Varianzhomogenität angenommen werden.

ONEWAY deskriptive Statistiken

ZUF 8 Score								
	N	Mittelwert	Std.-Abweichung	Std.-Fehler	95%-Konfidenzintervall für den Mittelwert		Minimum	Maximum
					Untergrenze	Obergrenze		
Nichtraucher	64	17,70	3,959	,495	16,71	18,69	10	27
Gelegenheitsraucher	9	21,89	2,421	,807	20,03	23,75	19	26
Raucher	17	24,94	2,989	,725	23,40	26,48	21	29
Gesamt	90	19,49	4,667	,492	18,51	20,47	10	29

Test der Homogenität der Varianzen

		Levene-Statistik	df1	df2	Signifikanz
ZUF 8 Score	Basiert auf dem Mittelwert	1,384	2	87	,256
	Basiert auf dem Median	1,317	2	87	,273
	Basierend auf dem Median und mit angepassten df	1,317	2	77,427	,274
	Basiert auf dem getrimmten Mittel	1,377	2	87	,258

Abb. 2.4: Deskriptive Statistiken im Rahmen der einfaktoriellen ANOVA und Ergebnis des Levene-Tests

Damit kommen wir zu der Haupttabelle „Einfaktorielle ANOVA“, den Ergebnissen der Varianzanalyse (siehe **Abb. 2.5**). Hier ist in der zweiten Spalte die Quadratsumme angegeben: *Zwischen den Gruppen*=SQE, *Innerhalb der Gruppen*=SQR, und *Gesamt* =SQT. Die F-Teststatistik findet sich im Feld mit der Bezeichnung *F*, und der zugehörige p-Wert im Feld *Signifikanz*. Die Teststatistik ist mit 28,132 sehr groß, entsprechend ist der p-Wert sehr klein (0,000)

Einfaktorielle ANOVA

ZUF 8 Score					
	Quadratsumme	df	Mittel der Quadrate	F	Signifikanz
Zwischen den Gruppen	761,299	2	380,650	28,132	,000
Innerhalb der Gruppen	1177,189	87	13,531		
Gesamt	1938,489	89			

Abb. 2.5: Hauptergebnisse der einfaktoriellen ANOVA

Interpretation der Ergebnisse

Ist der p-Wert $< 0,001$, gibt SPSS ,000 aus. Wir empfehlen in diesen Fällen stattdessen in Berichten, Publikationen o. ä. die Schreibweise „ $< 0,001$ “ zu verwenden. Da hier der p-Wert $< 0,001$ ist, kann also von einem signifikanten Zusammenhang zwischen dem Raucherstatus vor der Reha und der Zufriedenheit mit der Reha-Maßnahme ausgegangen werden. Das heißt, dass sich mindestens eine der Raucherstatus-Gruppen von den beiden anderen hinsichtlich der Zufriedenheit signifikant unterscheidet.

Bei den *Post-Hoc-Tests* werden je nach gewählter Methode die Ergebnisse der Paarvergleiche ausgegeben (siehe **Abb. 2.6**). Damit können wir nun herausfinden, welche Gruppen genau sich unterscheiden. Dabei sind die Mittlere Differenz (I-J) mit zugehörigem 95 %-Konfidenzintervall und der p-Wert (Signifikanz) von besonderem Interesse.

Mehrfachvergleiche

Abhängige Variable: ZUF 8 Score

Bonferroni

(I) Raucherstatus vor Reha	(J) Raucherstatus nach Reha	Mittlere Differenz (I-J)	Std.-Fehler	Signifikanz	95%-Konfidenzintervall	
					Untergrenze	Obergrenze
Nichtraucher	Gelegenheitsraucher	-4,186*	1,310	,006	-7,38	-,99
	Raucher	-7,238*	1,004	,000	-9,69	-4,79
Gelegenheitsraucher	Nichtraucher	4,186*	1,310	,006	,99	7,38
	Raucher	-3,052	1,516	,142	-6,75	,65
Raucher	Nichtraucher	7,238*	1,004	,000	4,79	9,69
	Gelegenheitsraucher	3,052	1,516	,142	-,65	6,75

*. Die Differenz der Mittelwerte ist auf dem Niveau 0.05 signifikant.

Abb. 2.6: Post-hoc-Vergleiche bei der einfaktoriellen ANOVA

Da wir hier die **Bonferroni-Korrektur** gewählt haben, wurden sechs Paarvergleiche durchgeführt. Dabei sind allerdings jeweils zwei gleich, nur die Gruppen sind vertauscht (z. B. Nichtraucher vs. Raucher und Raucher vs. Nichtraucher). Entsprechend sind die Schätzer bis auf die unterschiedlichen Vorzeichen gleich, und die p-Werte sind dieselben. Hier zeigt sich also, dass sich die Zufriedenheit der Nichtraucher und der Gelegenheitsraucher, und die der Nichtraucher und der Raucher signifikant voneinander unterscheidet. Dagegen sieht man beim Vergleich der Gelegenheitsraucher und der Raucher hinsichtlich der Zufriedenheit lediglich eine Tendenz zur Signifikanz ($p = 0,142$).

Das *Mittelwert-Diagramm* stellt die Mittelwerte in den Untergruppen grafisch dar. Möchte man zusätzlich die Standardabweichung, Konfidenzintervalle o. ä. dargestellt haben, sollte man die Grafik direkt über das Menü *Diagramme* erstellen.

Beispiel

Die Ergebnisse des Beispiels lassen sich folgendermaßen zusammenfassen:

Die Voraussetzungen 1 bis 4 für die Durchführung einer einfaktoriellen ANOVA sind erfüllt. Es zeigt sich deutlich ein signifikanter globaler Unterschied zwischen den drei Raucherstatus-Gruppen hinsichtlich der Zufriedenheit (siehe **Abb. 2.7**; $p < 0,001$, siehe **Abb. 2.5**). Bei den Paarvergleichen zeigt sich, dass sich Nichtraucher (NR) von den Rauchern (R) und den Gelegenheitsrauchern (GR) unterscheiden (p-Werte: NR vs. R $< 0,001$, NR vs. GR 0,006, GR vs. R 0,142, siehe **Abb. 2.6**).

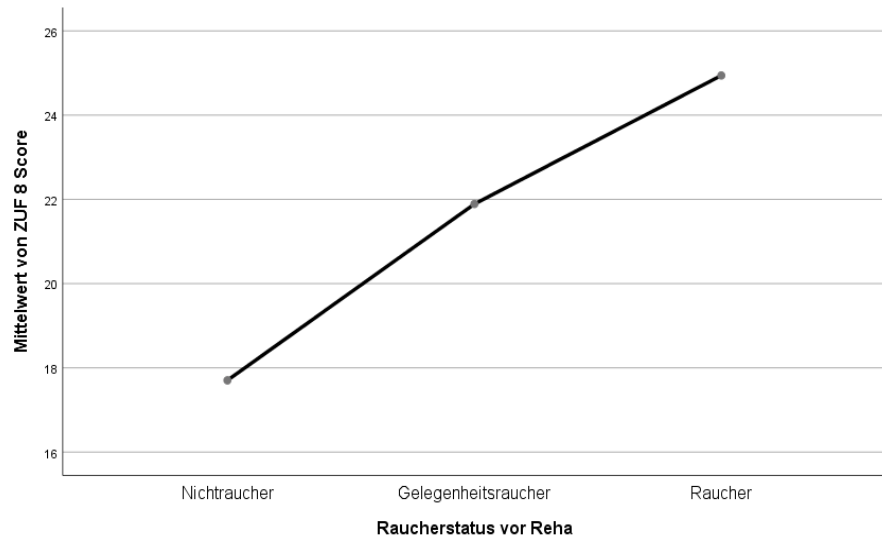


Abb. 2.7: Grafische Darstellung der mittleren Zufriedenheit in den drei Raucherstatusgruppen

2.2.2 Mehrfaktorielle ANOVA

Mehrfaktorielle ANOVA

Als Beispiel für eine **mehrfaktorielle ANOVA** wird folgendes Beispiel verwendet.

Beispiel

Gemeinsame Untersuchung des Einflusses von Geschlecht und Rauchen auf die Zufriedenheit mit der Reha-Maßnahme

Rauchen wird hier in der binären Variante verwendet: Nichtraucher vs. Raucher, wobei die Gelegenheitsraucher zu den Rauchern gezählt werden (siehe Beispiel 1.) in Abschnitt 1.2). Damit ergibt sich eine mehrfaktorielle ANOVA mit zwei Einflussfaktoren und einem Interaktionsterm.

Die Voraussetzungen, die für die mehrfaktorielle ANOVA erfüllt sein müssen, sind dieselben wie bei der einfaktoriellen ANOVA. Dass der Zufriedenheitsscore in der Gesamtpopulation normalverteilt ist, wurde im vorherigen Beispiel schon gezeigt (Voraussetzung 2). Die Voraussetzungen 1 (metrische Zielgröße) und 3 (kategoriale Einflussgrößen) sind ebenfalls erfüllt. Das Vorliegen von Varianzhomogenität (Voraussetzung 4) kann im Rahmen der Analyse überprüft werden.

Für die Durchführung einer mehrfaktoriellen ANOVA wählt man *Analysieren* → *Allgemeines lineares Modell* → *Univariat*. In dem Fenster, das sich öffnet (siehe **Abb. 2.8**), zieht man die Zielgröße ZUF 8 Score (*v2_v9_Score*) in das Feld *Abhängige Variable* und die Einflussfaktoren Geschlecht (*v10_geschlecht*) und rauchen_vor in das Feld *Feste Faktoren*.

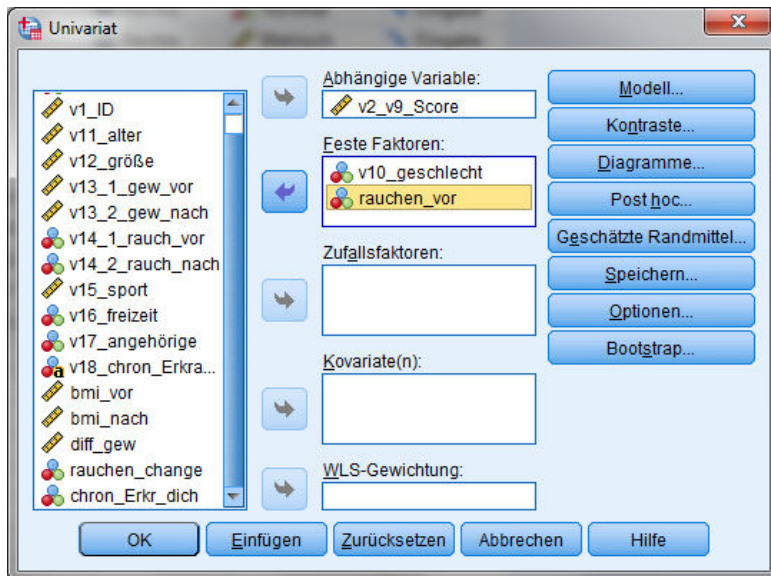


Abb. 2.8: Hauptfenster der allgemeinen Varianzanalyse zur Auswahl der Variablen

Für die Berücksichtigung der Interaktion wird das Feld *Modell* ausgewählt (siehe Abb. 2.9). Wenn hier *Gesättigtes Modell* ausgewählt wird, werden automatisch alle Haupteffekte und alle möglichen Interaktionen im Modell berücksichtigt. Wir entscheiden uns für diese Option.

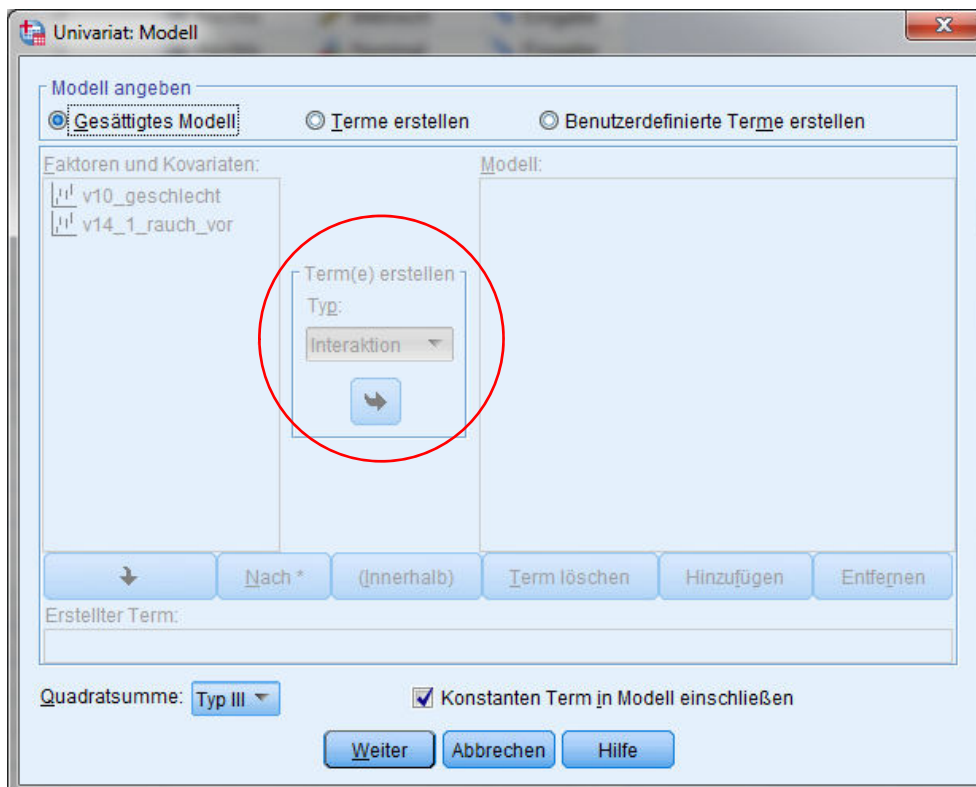


Abb. 2.9: Anpassung des Modells bei der allgemeinen ANOVA

Exkurs

Auswahl der Haupteffekte und Interaktionen

Wählt man das Feld *Terme erstellen*, kann man selbst die Haupteffekte und Interaktionen auswählen, die im Modell berücksichtigt werden sollen. Um einen Interaktionsterm in das Modell aufzunehmen, müssen in dem Feld *Faktoren und Kovariaten* die entsprechenden Variablen markiert werden. Anschließend wird bei *Term(e) erstellen* als *Typ Interaktion* ausgewählt und dann auf den Pfeil geklickt. Dadurch wird der Interaktionsterm in die Liste *Modell* hinzugefügt.

Da die Einflussgrößen hier jeweils binär sind, werden keine Post-Hoc-Analysen benötigt, sie könnten aber über das Feld *Post-Hoc-Mehrfachvergleiche für beobachteten Mittelwert* ausgewählt werden. Über das Feld *Diagramme* können Grafiken erstellt werden, mit deren Hilfe eine mögliche Interaktion optisch überprüft werden kann (siehe **Abb. 2.10**). Dafür zieht man eine Variable (hier *v10_geschlecht*) in das Feld *Horizontale Achse*, und die andere (hier *rauchen_vor*) in das Feld *Separate Linien*. Es ist wichtig, die Auswahl anschließend erst durch *Hinzufügen* zu bestätigen. Dann ist bei *Diagrammtyp* die Option *Liniendiagramm* zu aktivieren und anschließend alles durch *Weiter* zu bestätigen, da nur dann die Diagramme erstellt werden.

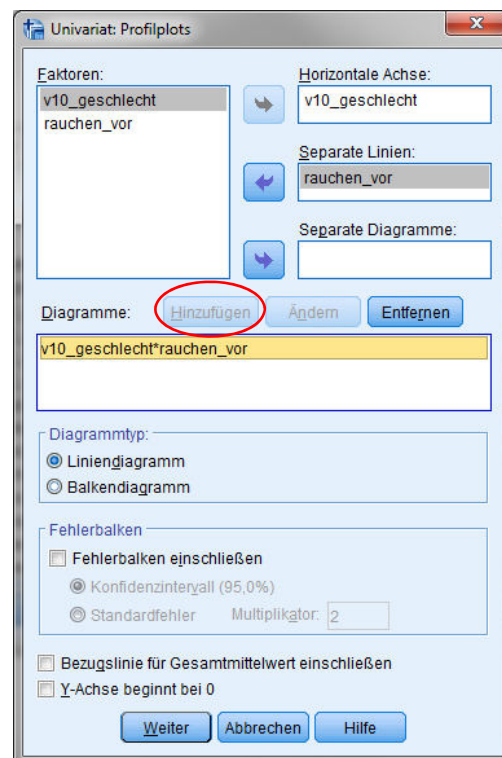


Abb. 2.10: Erstellung von benutzerdefinierten Profildiagrammen

Über das Feld *Optionen* kann die Durchführung eines Levene-Homogenitätstests veranlasst werden (siehe **Abb. 2.11**).

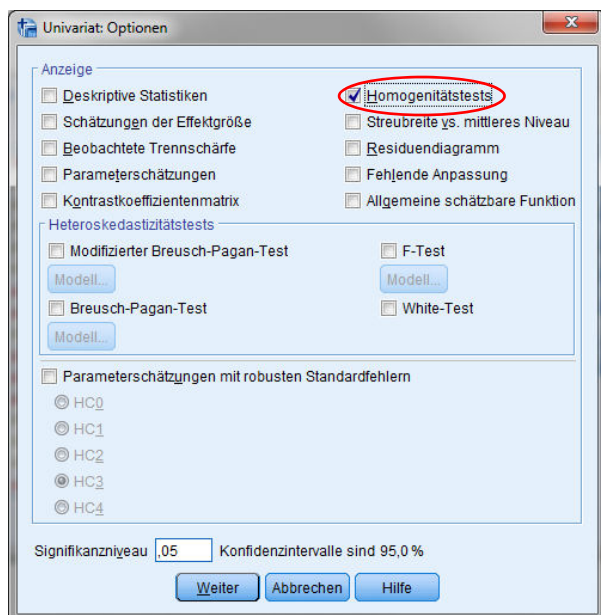


Abb. 2.11: Auswahl des Homogenitätstests im Rahmen einer allgemeinen ANOVA

Ergebnisse der mehrfaktoriellen ANOVA

Der Levene-Test auf Gleichheit der Fehlervarianzen liefert einen p-Wert von 0,816. Die Voraussetzung 4 kann daher also als erfüllt betrachtet werden. Die mehrfaktorielle ANOVA zeigt, dass keine Interaktion vorliegt (p-Wert=Sig.=0,974). Das zeigt sich auch in dem sog. Profildigramm (Abb. 2.12) – die Linien für die beiden Raucherstatusgruppen sind parallel.

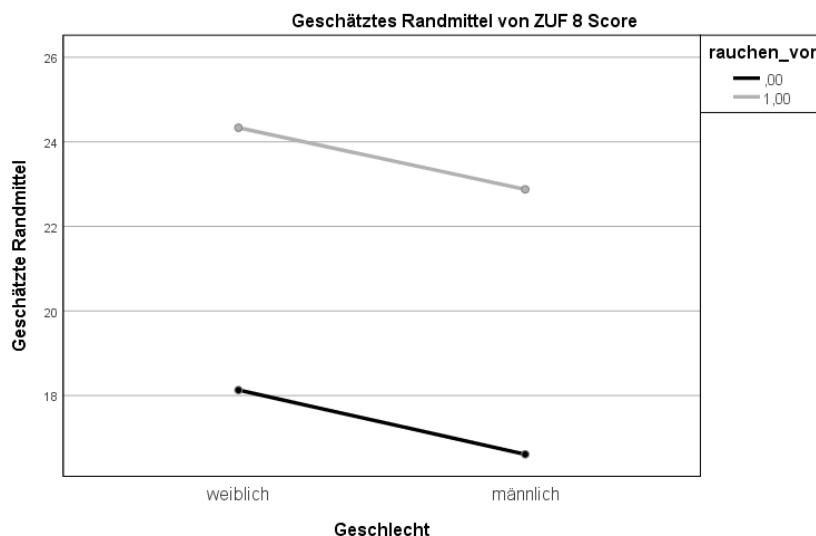


Abb. 2.12: Profildigramm zur Darstellung einer möglichen Interaktion

Der Effekt des Raucherstatus vor der Reha ist stark signifikant ($p < 0,001$), dagegen ist der Geschlechtseffekt nur tendenziell signifikant ($p = 0,119$).

Beispiel

2.3 Kovarianzanalyse (ANCOVA)

Wenn zusätzlich zu den interessierenden Faktoren eine oder mehrere Kovariablen berücksichtigt werden sollen (siehe Voraussetzung 5 in Abschnitt 2.1), empfiehlt sich die **Kovarianzanalyse** (ANCOVA) als Auswertungsmethode.

Definition

Kovariablen sind definiert als metrische Variablen, bei denen ein Einfluss auf die Zielgröße vermutet wird, der aber nicht primär von Interesse ist.

Durch die Berücksichtigung von Kovariablen kann die Reststreuung (auch Residuen genannt) reduziert werden. Das wiederum führt zu höherer Power, d. h. zu einer höheren Wahrscheinlichkeit, einen vorhandenen Effekt auch zu entdecken.

Beispiel

Gewichtsabnahme als Kovariable

Als Beispiel soll wieder die bereits in Studienbrief 4 des Moduls „Empirische Methoden I“ untersuchte Fragestellung dienen: Es soll der Einfluss des Raucherstatus vor der Reha sowie des Geschlechts auf die Zufriedenheit mit der Reha-Maßnahme um den vermuteten Effekt der Gewichtsveränderung während der Reha (Kovariable) bereinigt werden.

Das Vorgehen bei der Durchführung einer ANCOVA in SPSS ist fast identisch mit dem bei der Durchführung einer ANOVA. Die Voraussetzungen 1 bis 3 sind erfüllt (siehe Abschnitt 2.2.2).

Durchführung einer ANCOVA

Für die Durchführung einer ANCOVA wählt man: *Analysieren* → *Allgemeines lineares Modell* → *Univariat*. Die Kovariable (*diff_gew*) wird ausgewählt und in das Feld *Kovariate(n)* gezogen (siehe **Abb. 2.13**).

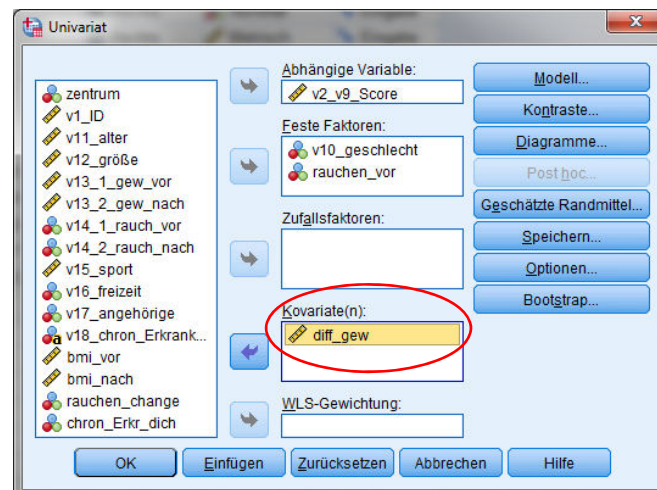


Abb. 2.13: Hauptfenster bei der ANCOVA zur Auswahl der Variablen

Wird in dem Feld *Modell* die Option *Gesättigtes Modell* ausgewählt, werden in dem statistischen Modell die Haupteffekte von Raucherstatus, Geschlecht und Gewichtsveränderung und die Interaktion zwischen Raucherstatus und Geschlecht berücksichtigt. Post-hoc-Vergleiche können bei einer ANCOVA nicht durchgeführt werden (das Feld *Post Hoc* ist inaktiv). Das „Profildigramm“ und Tests auf Varianzhomogenität können aber wie gehabt über *Diagramme* und *Optionen* ausgewählt werden.

Beispiel

Ergebnisinterpretation ANCOVA

In der Tabelle „Tests der Zwischensubjekteffekte“ (siehe **Abb. 2.14**) gibt der Eintrag in der Spalte *Quadratsumme vom Typ III* in der Zeile *Fehler* die Reststreuung (SQR) wieder. Bei der mehrfaktoriellen ANOVA aus dem vorherigen Beispiel lag die SQR bei 1190. Man sieht, dass sich die Reststreuung hier deutlich reduziert hat (auf 434).

Tests der Zwischensubjekteffekte

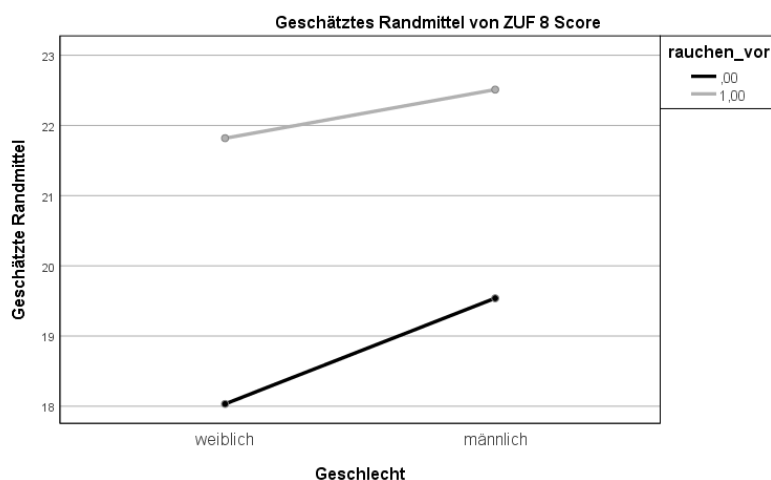
Abhängige Variable: ZUF 8 Score

Quelle	Quadratsumme vom Typ III	df	Mittel der Quadrate	F	Sig.
Korrigiertes Modell	1504,726 ^a	4	376,181	73,716	,000
Konstanter Term	3167,904	1	3167,904	620,781	,000
diff_gew	756,607	1	756,607	148,264	,000
v10_geschlecht	16,479	1	16,479	3,229	,076
rauchen_vor	151,790	1	151,790	29,745	,000
v10_geschlecht * rauchen_vor	2,552	1	2,552	,500	,481
Fehler	433,763	85	5,103		
Gesamt	36122,000	90			
Korrigierte Gesamtvariation	1938,489	89			

a. R-Quadrat = ,776 (korrigiertes R-Quadrat = ,766)

Abb. 2.14: Hauptergebnisse der ANCOVA

Die Gewichtsveränderung als Kovariable hat einen deutlichen Einfluss auf die Zufriedenheit ($p < 0,001$). Der starke Effekt des Raucherstatus bleibt erhalten ($p < 0,001$), der Effekt des Geschlechts tritt etwas deutlicher zu Tage ($p = 0,076$). Der p-Wert des Interaktionsterms ist durch die zusätzliche Berücksichtigung der Gewichtsveränderung kleiner geworden. Das spiegelt sich auch im Profildigramm wider (siehe **Abb. 2.15**); die Profillinien der beiden Raucherstatusgruppen sind nicht mehr parallel, sondern laufen leicht aufeinander zu. Die Wechselwirkung ist allerdings nicht so stark, dass eine separate Analyse durchgeführt werden muss.



Die Kovariaten im Modell werden anhand der folgenden Werte berechnet: Veränderung Körpergewicht (vor-nach) = 4,6000

Abb. 2.15: Profildigramm mit leichter Interaktion zwischen Raucherstatus und Geschlecht (aufeinander zulaufende Geraden)

Exkurs

Repeated measures ANOVA

Wie bereits in den Abschnitten 1.2 und 2.1 erwähnt, kann es sein, dass die Zielgröße mehrfach erhoben wird (Messwiederholungen) oder mehrere (verschiedene) Zielgrößen erhoben werden. Eine ausführliche Beschreibung zur Durchführung dieser Analysen in SPSS ist z. B. bei Bühl (2012) zu finden. Hier sei lediglich erwähnt, wo die Analysen in SPSS zu finden sind und was im Hauptfenster in welches Feld einzufügen ist. Die entsprechenden Verfahren der Varianzanalyse heißen **repeated measures ANOVA** und **MANOVA**.

Manchmal wird die Zielgröße nicht nur einmal pro Beobachtungseinheit (z. B. Person oder Klinik) gemessen, sondern mehrmals (sog. **Messwiederholungen**). Entweder erfolgt die Messwiederholung direkt hintereinander, um Messungenauigkeiten auszugleichen, oder mit zeitlichem Abstand, um eine Entwicklung über die Zeit abzubilden. Dies ist ganz typisch bei der Untersuchung der Effektivität von Trainings- oder Therapiemaßnahmen.

Über den Pfad *Analysieren* → *Allgemeines lineares Modell* → *Messwiederholung* gelangt man zum ersten Fenster (siehe **Abb. 2.16**).

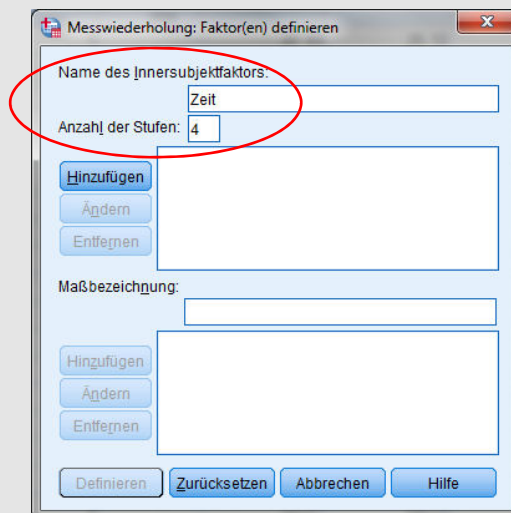


Abb. 2.16: Festlegung der Anzahl der Messwiederholungen und Benennung des Messwiederholungsfaktors

In dem Feld *Name des Innersubjektfaktors* kann man dem Messwiederholungsfaktor einen selbstgewählten Namen geben (hier beispielhaft: Zeit), in das Feld *Anzahl der Stufen* wird eingetragen, wie viele Messwiederholungen es gibt (hier beispielhaft: 4). Die stellt die Faktorstufen dar, ganz analog zu einem Zwischensubjektfaktor. Anschließend klickt man auf *Hinzufügen* und dann auf *Definieren*. In dem folgenden Fenster (siehe **Abb. 2.17**) fügt man bei *Innersubjektfaktoren* die Variablen ein, die die Zielgrößen zu den verschiedenen Zeitpunkten beinhalten. In das Feld *Zwischensubjektfaktoren* werden die Einflussfaktoren eingefügt und in das Feld *Kovariaten* die metrischen Kovariablen.

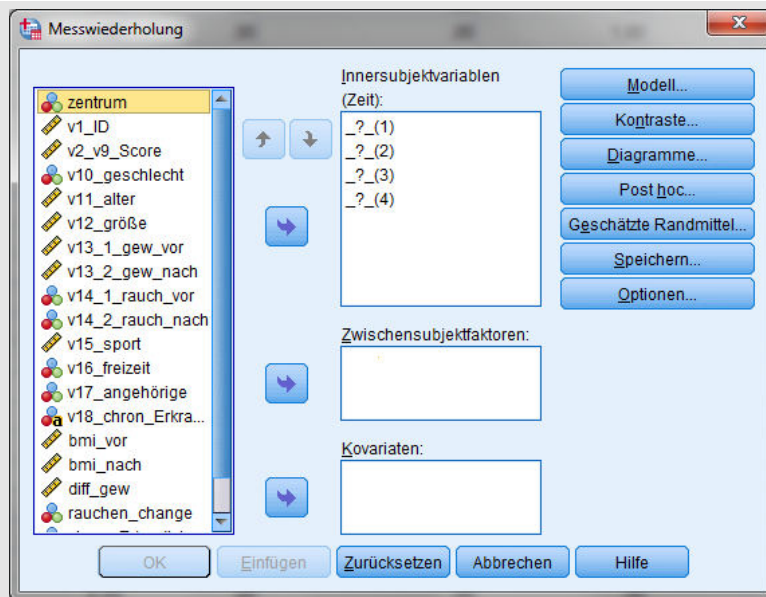


Abb. 2.17: Definition der Ziel- und Einflussgrößen bei der repeated measures ANOVA

Multivariate ANOVA (MANOVA)

Zur Durchführung einer MANOVA wählt man *Analysieren* → *Allgemeines Lineares Modell* → *Multivariat*. Das Feld, das sich öffnet (siehe Abb. 2.18) ähnelt sehr dem von *Univariat*. Hier können allerdings mehrere Zielvariablen in das Feld *Abhängige Variablen* eingetragen werden.

Exkurs

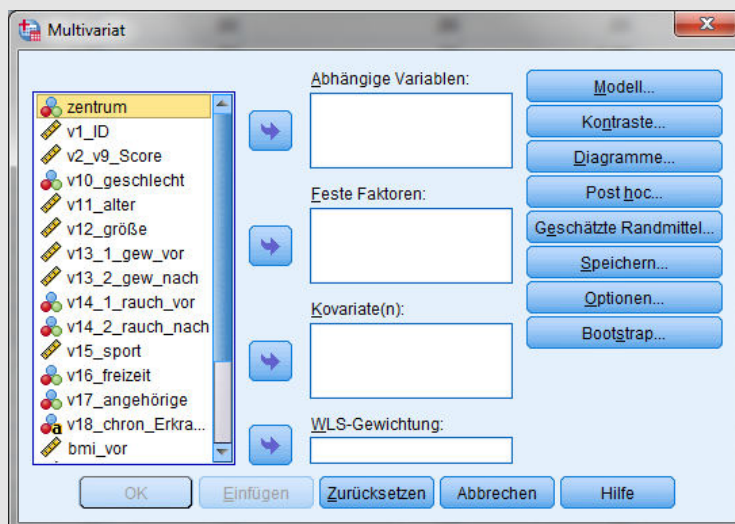


Abb. 2.18: Definition der Ziel- und Einflussgrößen bei einer MANOVA

Übungsaufgaben

- 2.1) Es soll untersucht werden, ob das Geschlecht und die Tatsache, dass Angehörige in die Behandlung einbezogen werden, einen Einfluss darauf haben, wie viele Minuten pro Tag im Durchschnitt Sport getrieben wird. Eine Wechselwirkung zwischen dem Geschlecht und der Einbeziehung Angehöriger kann nicht ausgeschlossen werden.

Wählen Sie die passende Form der Varianzanalyse. Überprüfen Sie, ob die entsprechenden Voraussetzungen erfüllt sind und führen Sie die Analyse durch. Interpretieren Sie abschließend die Ergebnisse.

- 2.2) Es wird vermutet, dass der BMI bei Antritt der Reha einen Effekt darauf hat, wie viel Sport die Patientinnen und Patienten treiben. Modifizieren Sie entsprechend die Analyse aus Aufgabe 2.1) und interpretieren Sie die neuen Ergebnisse.

3 Regressionsanalyse

Mit der Varianzanalyse können viele Fragestellungen bearbeitet werden. Voraussetzung ist allerdings, dass die Zielgröße metrisch ist, und dass es mindestens eine kategoriale Einflussgröße (Einflussfaktor) gibt. Die Varianzanalyse hat viele Gemeinsamkeiten mit der Regressionsanalyse. Die Regressionsanalyse ist jedoch auch für kategoriale Zielgrößen anwendbar, und es muss auch keine kategoriale Einflussgröße vorliegen. Es gibt verschiedene Formen der Regressionsanalyse, abhängig davon, wie die Zielgröße skaliert ist. Für metrische Zielgrößen wird die lineare Regressionsanalyse verwendet, für kategoriale Zielgrößen die logistische Regressionsanalyse. Bei der logistischen Regression werden wir uns hier auf die binäre logistische Regression für binäre Zielgrößen beschränken.

Das Prinzip der Regressionsanalyse ist, dass eine mathematische Gleichung aufgestellt wird, die den Zusammenhang zwischen zwei oder mehreren Variablen möglichst optimal widerspiegelt. Die Regressionskoeffizienten in dieser sog. Regressionsgleichung quantifizieren den Zusammenhang zwischen der Zielgröße und den Einflussgrößen. Wird nur eine Einflussgröße berücksichtigt, spricht man von der **Einfachregression**, bei mehreren Einflussgrößen von der **multiplen Regression**.

3.1 Lineare Regression

Die **lineare Regressionsanalyse** ist anwendbar, wenn die Zielgröße metrisch ist, und ein linearer Zusammenhang zwischen Ziel- und Einflussgröße besteht. Der lineare Zusammenhang lässt sich am leichtesten mit Hilfe eines Streudiagramms nachweisen. Im nächsten Beispiel wird demonstriert, wie in das Streudiagramm die Regressionsgerade eingezeichnet werden kann.

Für die Beobachtungseinheit i wird der lineare Zusammenhang zwischen der Zielgröße Y_i und den Einflussgrößen x_{i1} bis x_{ip} mit folgender Regressionsgleichung dargestellt:

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i, \quad i = 1, \dots, n$$

Die **Regressionskoeffizienten** β_0 bis β_p werden aus den Daten mit der **Methode der Kleinsten Quadrate** geschätzt. Der geschätzte Steigungsparameter $\hat{\beta}_j$ (mit $j = 1, \dots, p$) ist gut interpretierbar: wenn x_j um eine Einheit steigt, steigt Y um den Faktor $\hat{\beta}_j$. Dabei gibt das Vorzeichen von $\hat{\beta}_j$ an, ob der Zusammenhang zwischen Einfluss- und Zielgröße positiv oder negativ ist. Die Interpretation ist allerdings auf den Wertebereich der Einflussgrößen beschränkt, der durch die Daten abgedeckt ist.

Die **Fehlerterme** ε_i sind unbeobachtbare Zufallsvariablen, und werden auch **Residuen** genannt. Für die Residuen besteht die Annahme, dass sie unabhängig und identisch normalverteilt sind mit $E(\varepsilon_i) = 0$ und $Var(\varepsilon_i) = \sigma^2$. Dabei beinhaltet die Annahme über die Varianz der Residuen die Annahme von gleichen Varianzen (Varianzhomogenität). Ob diese Annahmen erfüllt sind, lässt sich durch grafische Darstellungen der normierten Residuen beurteilen. Die Anpassungsgüte des Regressionsmodells lässt sich mit Hilfe des Bestimmtheitsmaßes R^2 abschätzen. Für genauere Informationen zum Thema multiple lineare Regression sei hier für die Theorie beispielsweise auf Fahrmeir et al. (2012) oder Hartung et al. (2005) verwiesen, für die Anwendung in SPSS auf Fromm (2012) oder Duller (2013).

Voraussetzungen

Regressionsgleichung

Residuen

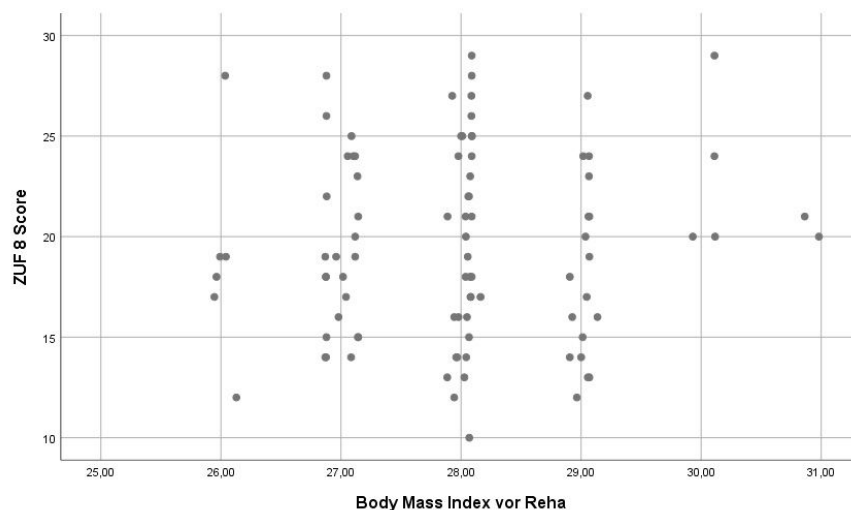
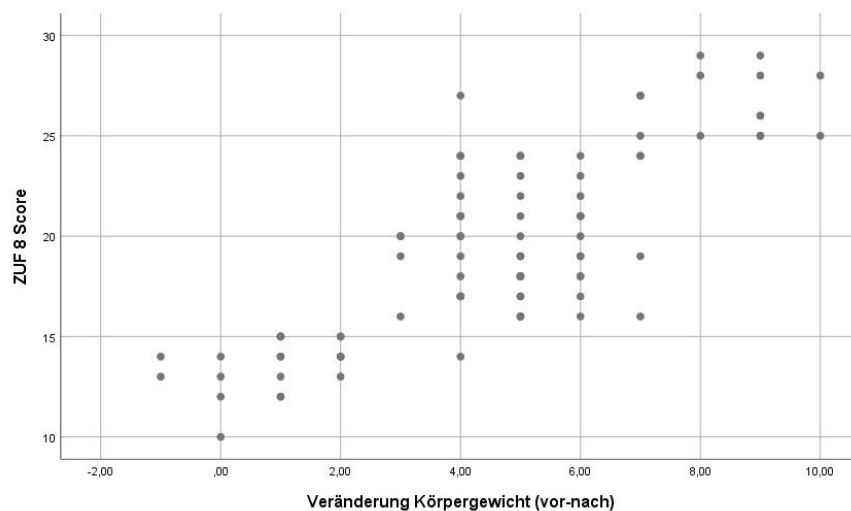
Bestimmtheitsmaß

Beispiel

Einflussfaktoren auf die Zufriedenheit mit der Reha (lineare Einfachregression)

Wir erweitern die Fragestellung, bei der der Einfluss der Gewichtsveränderung während der Reha auf die Zufriedenheit mit der Reha-Maßnahme untersucht wurde. Hierbei handelt es sich um eine lineare Einfachregression. Zusätzlich soll nun der Einfluss von Alter, BMI zu Beginn der Reha und die durchschnittliche tägliche sportliche Betätigung berücksichtigt werden.

In einem ersten Schritt muss überprüft werden, ob ein linearer Zusammenhang zwischen den Einflussgrößen und der Zielgröße angenommen werden kann. Dafür zeichnen wir paarweise Streudiagramme *Diagramme* → *Diagrammerstellung* → *Streu-/Punktdiagramm*. Für alle vier Einflussgrößen kann ein linearer Zusammenhang mit der Zielgröße angenommen werden (siehe **Abb. 3.1**; teilweise mit einer vermutlich horizontalen Regressionsgerade), auch wenn die Streuung teilweise sehr groß ist (z. B. beim BMI vor der Reha sowie dem Alter).



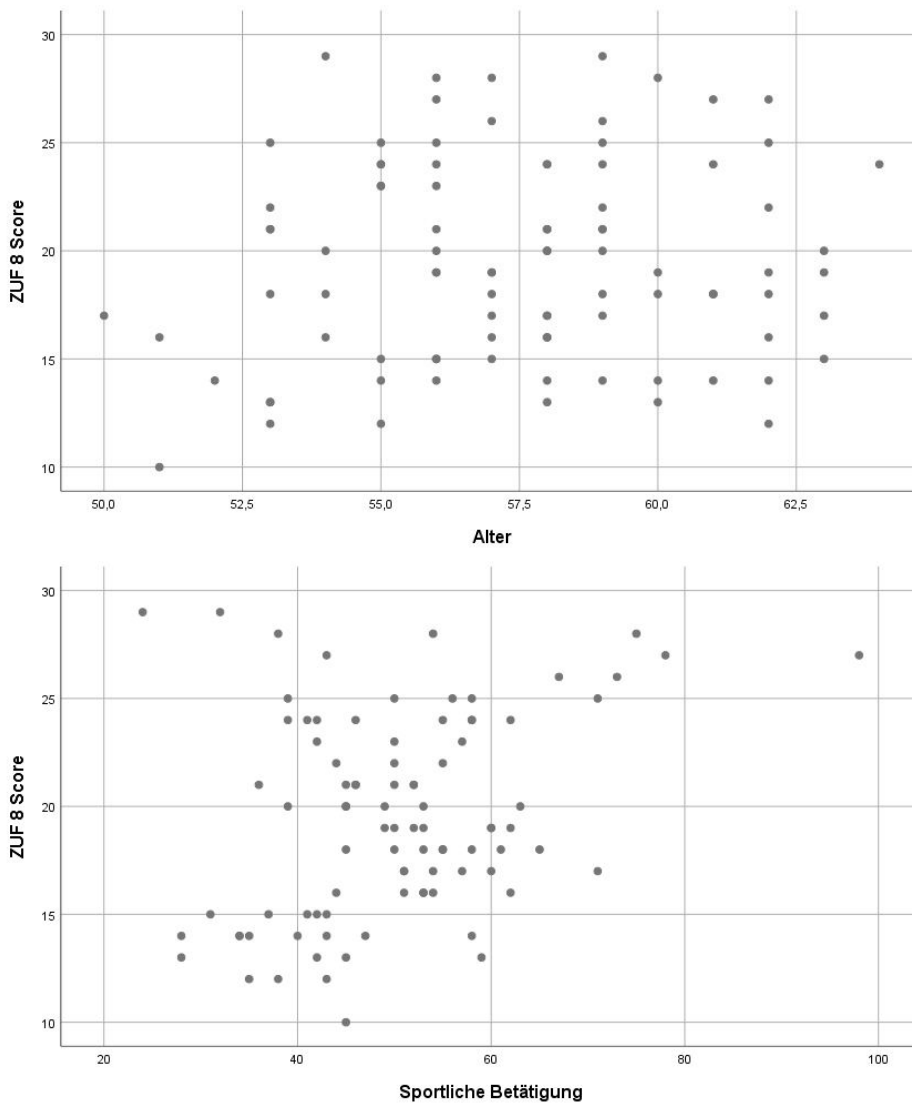


Abb. 3.1: Paarweise Streudiagramme zur Überprüfung der Linearitätsannahme vor der Durchführung einer linearen Regressionsanalyse

Für die Einflussgröße der Gewichtsveränderung soll an dieser Stelle demonstriert werden, wie in das Streudiagramm eine lineare Regressionsgerade eingefügt werden kann. Dafür klickt man in der Ausgabe das Streudiagramm doppelt an. Dadurch gelangt man zum *Diagramm-Editor* (siehe **Abb. 3.2**).

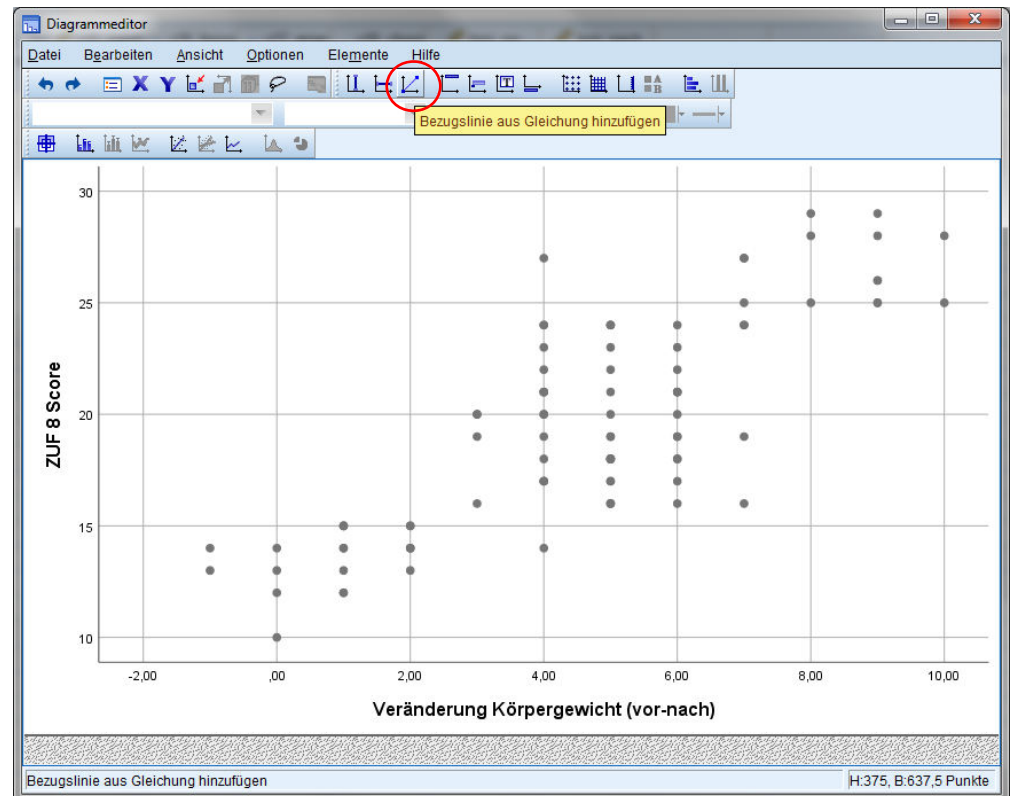


Abb. 3.2: Diagramm-Editor zum Bearbeiten einer Grafik (hier zum Hinzufügen der Regressionsgeraden)

Klickt man dort auf *Bezugslinie aus Gleichung hinzufügen*, wird automatisch die Regressionsgleichung in das Streudiagramm eingefügt, und es erscheint ein Fenster zu den *Eigenschaften* (siehe Abb. 3.3). Dort kann man bei Bedarf die Regressionsgleichung oder die Darstellungsoptionen ändern.

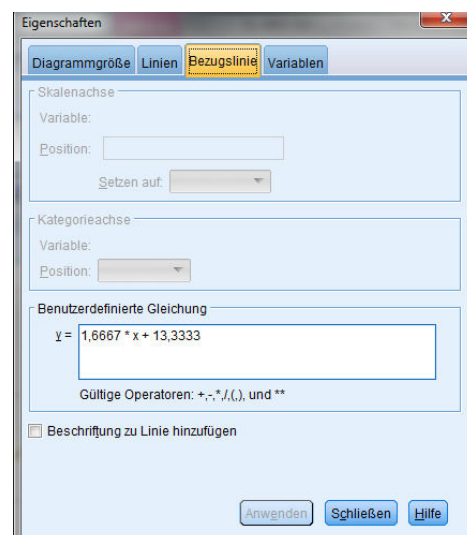


Abb. 3.3: Fenster zur möglichen Bearbeitung der Bezugslinie

Multiple lineare Regression

Um eine multiple lineare Regressionsanalyse in SPSS durchzuführen wählt man den Pfad *Analysieren → Regression → Linear*. In dem Fenster, das sich daraufhin öffnet (siehe Abb. 3.4), wird die Zielgröße (hier ZUF 8 Score – v2_v9_Score) in das Feld *Abhängige Variable* gezogen und alle Einflussgrößen (hier Veränderung Körper-

gewicht – *diff_gew*, Alter – *v11_alter*, Body Mass Index vor Reha – *bmi_vor* und Sportliche Betätigung – *v15_sport*) in das Feld *Unabhängige*. Bei *Methode* kann ein Verfahren zur Variablenselektion gewählt werden. Da wir aber vorerst keine Variablenselektion durchführen wollen, belassen wir es bei *Einschluss* (= keine Selektion).

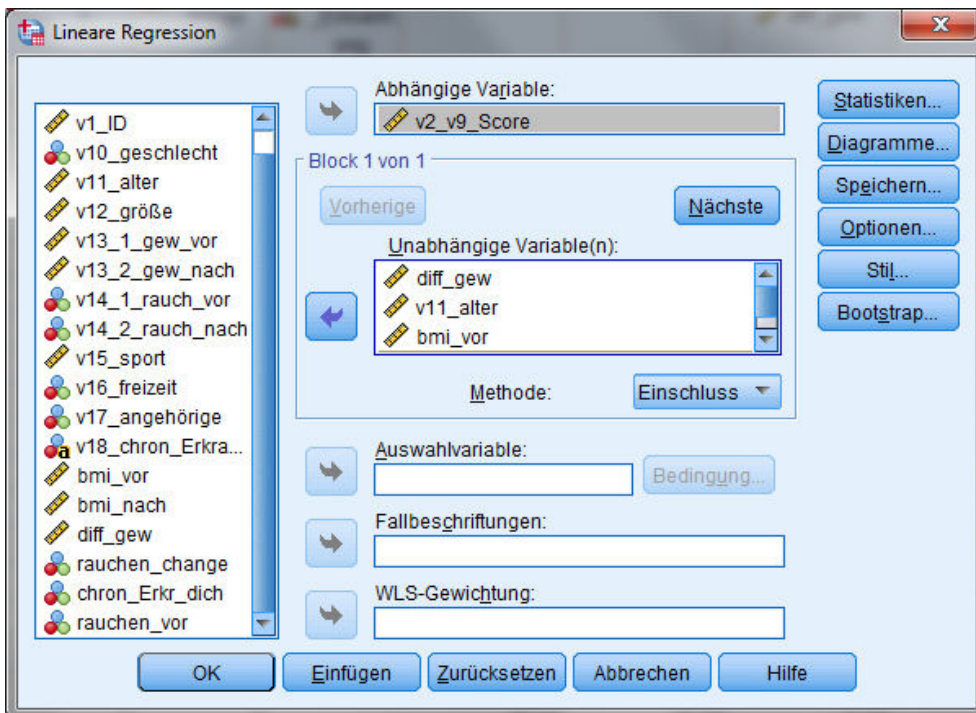


Abb. 3.4: Hauptfenster zur Auswahl der Variablen für eine lineare Regression

Wenn das Feld *Statistiken* angeklickt wird, öffnet sich ein Fenster, in dem man nützliche Zusatzausgaben auswählen kann (siehe Abb. 3.5). So sollten in der Ausgabe auf jeden Fall der *Schätzer* und die dazugehörigen 95%-*Konfidenzintervalle* enthalten sein. Außerdem sollte man sich die *Anpassungsgüte des Modells*, d.h. das r^2 , mit ausgeben lassen.

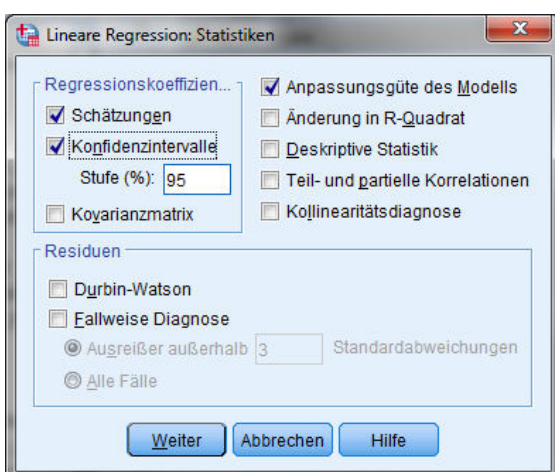


Abb. 3.5: Fenster zur Auswahl von zusätzlichen Statistiken im Rahmen der linearen Regression

Um zu überprüfen, ob die Residuen normalverteilt sind (mit Erwartungswert 0) und ob Varianzhomogenität vorliegt, klickt man das Feld *Diagramme* an (siehe Abb. 3.6). Klickt man das Feld *Normalverteilungsdiagramm* an, wird ein sog. Probability-

P-P-Diagramm

Probability-Diagramm (P-P-Diagramm) gezeichnet. Damit kann überprüft werden, ob die Residuen normalverteilt sind. Wenn man die standardisierten Residualwerte (*ZRESID) auf der x-Achse gegen die standardisierten Vorhersagewerte (*ZPRED) auf der y-Achse abzeichnet, erhält man ein Streudiagramm, das einem anzeigt, ob Varianzhomogenität angenommen werden kann.

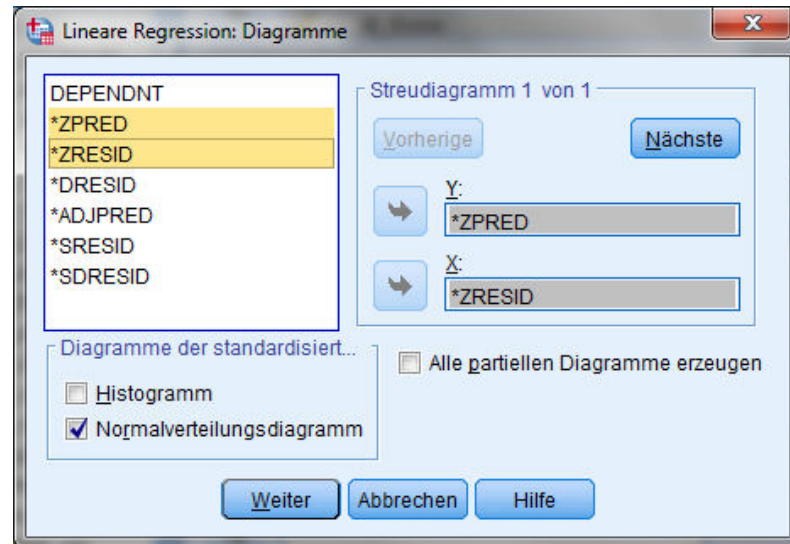


Abb. 3.6: Fenster zur Erstellung von Residualdiagrammen im Rahmen einer linearen Regression

Klickt man nun auf *Weiter* und anschließend im Hauptfenster auf *OK*, wird die multiple lineare Regression durchgeführt.

Beispiel

Überprüfung der Voraussetzungen für eine multiple lineare Regressionsanalyse

Die Punkte im P-P-Diagramm (siehe Abb. 3.7) liegen ohne große Abweichungen auf der Diagonalen. Das spricht dafür, dass die Residuen normalverteilt sind. Die Punktwolke in Abbildung 3.8 ist um die Nulllinie zentriert und zeigt keine zunehmende Streuung für größere Werte. Das heißt, es kann Varianzhomogenität angenommen werden. Gemeinsam mit den linearen Zusammenhängen, die sich in den Streudiagrammen gezeigt haben, sind nun also alle Bedingungen für die multiple lineare Regressionsanalyse erfüllt.

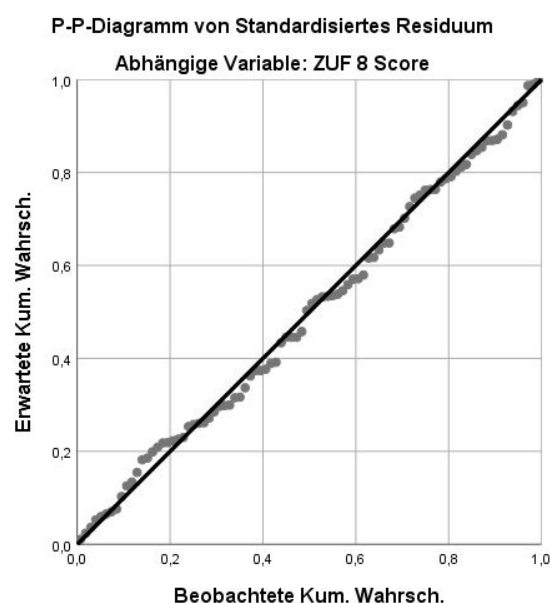


Abb. 3.7: P-P-Diagramm zur Überprüfung hinsichtlich der Normalverteilung der Residuen

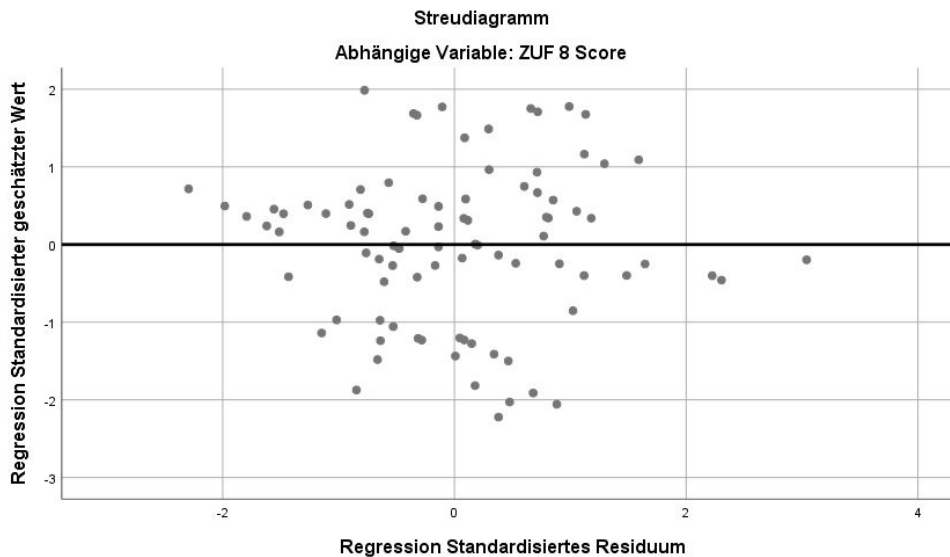


Abb. 3.8: Darstellung der standardisierten Residuen vs. der standardisierten Vorhersagewerte zur Überprüfung der Varianzhomogenität

In der Tabelle „Modellzusammenfassung“ ist ein *R-Quadrat* von 0,676 angegeben, und damit also eine ausreichende Modellgüte. In der Tabelle „Koeffizienten“ (siehe **Abb. 3.9**) sind die Hauptergebnisse der Regressionsanalyse abzulesen.

Koeffizienten^a

Modell		Nicht standardisierte Koeffizienten		Standardisierte Koeffizienten	T	Sig.	95,0% Konfidenzintervalle für B	
		Regressionskoeffizient B	Std.-Fehler	Beta			Untergrenze	Obergrenze
1	(Konstante)	-13,118	10,064		-1,303	,196	-33,128	6,892
	Veränderung Körpergewicht (vor-nach)	1,496	,124	,822	12,104	,000	1,250	1,741
	Alter	,057	,091	,039	,625	,534	-,125	,239
	Body Mass Index vor Reha	,799	,278	,182	2,874	,005	,246	1,351
	Sportliche Betätigung	,002	,027	,005	,079	,937	-,051	,055

a. Abhängige Variable: ZUF 8 Score

Abb. 3.9: Geschätzte Regressionskoeffizienten aus der multiplen linearen Regressionsanalyse

In der Spalte „Sig.“ sind die p-Werte aufgeführt. Es zeigt sich, dass das Alter und die durchschnittliche tägliche Zeit sportlicher Betätigung keinen Einfluss auf die Zufriedenheit haben ($p=0,534$ und $p=0,937$). Da es das Ziel sein sollte, ein statistisches Modell aufzustellen, dass nur relevante Variablen beinhaltet, führen wir nun eine **Rückwärtsselektion** durch. Dafür wählen wir im Hauptfenster der linearen Regressionsanalyse bei *Methode* als Verfahren die *Rückwärtsselektion*.

In der Ausgabe ist nun als erstes die Tabelle „Aufgenommene/Entfernte Variablen“ interessant (siehe **Abb. 3.10**).

Aufgenommene/Entfernte Variablen^a

Modell	Aufgenommene Variablen	Entfernte Variablen	Methode
1	Sportliche Betätigung, Alter, Body Mass Index vor Reha, Veränderung Körpergewicht (vor-nach) ^b	.	Einschluß
2	.	Sportliche Betätigung	Rückwärts (Kriterium: Wahrscheinlichkeit von F-Wert für Ausschluß $\geq ,100$).
3	.	Alter	Rückwärts (Kriterium: Wahrscheinlichkeit von F-Wert für Ausschluß $\geq ,100$).

a. Abhängige Variable: ZUF 8 Score

b. Alle gewünschten Variablen wurden eingegeben.

Abb. 3.10: Zusammenfassung der aufgenommenen und entfernten Variablen im Rahmen des Selektionsverfahrens

An der ersten Spalte (*Modell*) kann man ablesen, wie viele Modelle gerechnet wurden. In der zweiten Spalte (*Aufgenommene Variablen*) sind in der ersten Zeile alle Variablen des kompletten Modells aufgeführt. In der dritten Spalte (*Entfernte Variablen*) kann man sehen, wann welche Variablen aus dem Modell genommen wurden. Hier wurde erst „sportliche Betätigung“ aus dem Modell entfernt und im nächsten Schritt das „Alter“. Das dritte Modell stellt dementsprechend das endgültige Modell dar. Daher schauen wir uns die weiteren Ergebnisse nur für Modell 3 an. Anhand der ersten Spalte der entsprechenden Tabelle kann man die Ergebnisse zuordnen (siehe **Abb. 3.11**).

Das *R-Quadrat* ist vom kompletten zum endgültigen Modell praktisch unverändert geblieben (0,676 vs. 0,674). In der Tabelle „Koeffizienten“ sind für das dritte Modell, das nur noch die Einflussgrößen „Veränderung Körpergewicht (vor – nach)“ und „BMI vor Reha“ beinhaltet, u. a. die Koeffizientenschätzer (*Regressionskoeffizient* *B*), die zugehörigen 95,0%-Konfidenzintervalle für *B* und die p-Werte (*Sig.*) angegeben (siehe **Abb. 3.11**). „Veränderung Körpergewicht (vor – nach)“ und der „BMI vor Reha“ haben einen deutlich signifikanten Einfluss auf die Zufriedenheit („ZUF 8 Score“) ($p < 0,001$ bzw. $p = 0,005$). Wenn „Veränderung Körpergewicht (vor–nach)“ um eine Einheit steigt (d. h. die Patienten nehmen ein Kilo ab), steigt die Zufriedenheit um 1,506 Punkte. Je übergewichtiger die Patienten und Patientinnen vor der Reha waren, desto größer ist die Zufriedenheit mit der Reha-Maßnahme (pro BMI-Einheit Zugewinn von 0,777 Punkten beim Zufriedenheitsscore).

Koeffizienten^a

Modell	Nicht standardisierte Koeffizienten		Standardisierte Koeffizienten	T	Sig.	95,0% Konfidenzintervalle für B	
	Regressionskoeffizient B	Std.-Fehler	Beta			Untergrenze	Obergrenze
3 (Konstante)	-9,171	7,676		-1,195	,235	-24,427	6,085
Veränderung Körpergewicht (vor-nach)	1,506	,113	,827	13,382	,000	1,282	1,729
Body Mass Index vor Reha	,777	,271	,177	2,868	,005	,238	1,315

a. Abhängige Variable: ZUF 8 Score

Abb. 3.11: Geschätzte Regressionskoeffizienten des endgültigen linearen Modells nach Anwendung der Rückwärtsselektion

3.2 Logistische Regression

Ist die **Zielgröße kategorial** skaliert, sollte die **logistische Regressionsanalyse** verwendet werden. Für den Spezialfall einer **binären Zielgröße** ist die **binäre logistische Regression** das richtige Verfahren, die ebenfalls häufig als logistische Regression bezeichnet wird.

Im Gegensatz zur linearen Regression wird hier kein linearer Zusammenhang zwischen Ziel- und Einflussvariable zugrunde gelegt. Stattdessen wird eine **Regressionsgleichung** aufgestellt, die den Effekt der Einflussgrößen x_1 bis x_p auf die Logit-Funktion der Zielvariablen Y : $L = \ln\left(\frac{p}{1-p}\right)$ untersucht. Dabei bezeichnet p die Wahrscheinlichkeit dafür, dass das Ereignis $Y=1$ vorliegt bzw. eintritt (d.h. $P(Y=1)=p$). Das führt für die Beobachtungseinheit i zu dem **Regressionsmodell**

$$L_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \varepsilon_i, \quad i = 1, \dots, n$$

Anders als bei der linearen Regression werden die Koeffizientenschätzer β_0 bis β_p nicht mit der Methode der kleinsten Quadrate, sondern iterativ mit der Maximum-Likelihood-Methode geschätzt. Die resultierenden Schätzer $\hat{\beta}_j$ mit $j = 1, \dots, p$ sind so erst mal nicht zu interpretieren. Berechnet man allerdings $e^{\hat{\beta}_j}$, erhält man damit das Odds Ratio (OR) für die Einflussgröße x_j . Wenn das OR gleich 1 ist, bedeutet das, dass x_j keinen Einfluss auf die Zielgröße hat. Ist das OR größer als 1, steigt das Risiko bzw. die Chance für $Y=1$; ist das OR kleiner als 1 sinkt das Risiko bzw. die Chance entsprechend.

Statt des Bestimmtheitsmaßes R^2 (wie bei der linearen Regression) dient hier die sog. Devianz als Maß für die Anpassungsgüte des Modells. Die Devianz liegt zwischen 0 und $+\infty$ (positiv unendlich) und ist bei einer perfekten Anpassung gleich 0. Dem logistischen Regressionsmodell liegen keine Verteilungsannahmen zugrunde. Das bedeutet, es müssen keine Annahmen überprüft werden, bevor eine logistische Regression durchgeführt werden kann (für weitere Informationen zur Anwendung der logistischen Regression in SPSS vgl. Fromm, 2012).

Maximum-Likelihood Methode

Odds Ratio

Devianz

Beispiel

Einbezug von Angehörigen in die Behandlung

An dieser Stelle soll die SPSS-Analyse des Beispiels aus Abschnitt 1.3 demonstriert werden. Dabei soll der Einfluss von Geschlecht, Alter, dem Raucherstatus vor der Reha und dem BMI vor der Reha auf die ärztliche Entscheidung, Angehörige in die Behandlung mit einzubeziehen, untersucht werden.

Zur binären logistischen Regression gelangt man über *Analysieren* → *Regression* → *binär logistisch*. In dem sich öffnenden Hauptfenster (siehe **Abb. 3.12**) zieht man die Zielgröße („Einbezug von Angehörigen“ – *v17_angehörige*) in das Feld *Abhängige Variable* und die Einflussgrößen in das Feld *Kovariaten*. Da hier eine Rückwärtsselektion durchgeführt werden soll, wählen wir bei *Methode* das Verfahren *Rückwärts: Bedingt*.

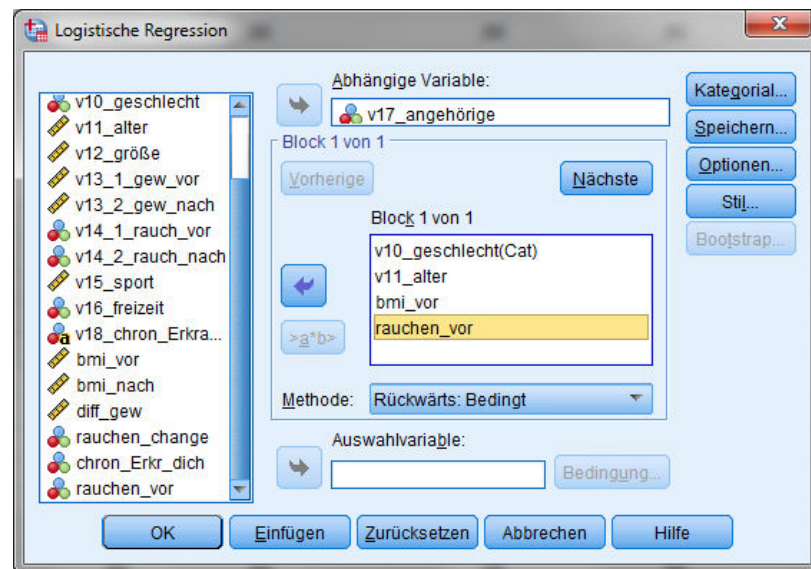


Abb. 3.12: Hauptfenster zur Auswahl der Variablen für die logistische Regression

Da SPSS nicht automatisch die Skalierung der Einflussgrößen berücksichtigt, wählen wir als nächstes das Feld *Kategorial* (siehe **Abb. 3.13**). Hier ziehen wir alle kategorialen Einflussgrößen (d.h. die Variablen *v10_geschlecht* und *rauchen_vor*) in das Feld *Kategoriale Kovariaten*. Hierbei wird automatisch die letzte Kategorie zur Referenzkategorie.

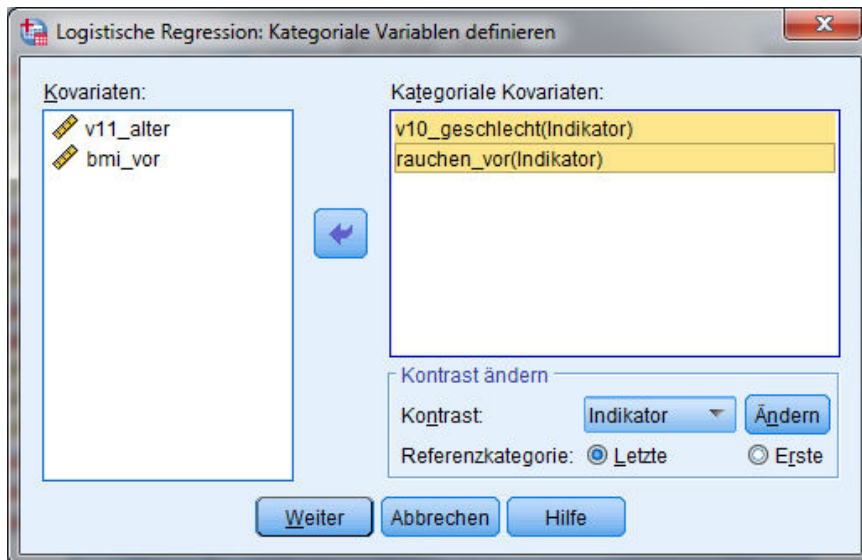


Abb. 3.13: Fenster zur Definition der kategorialen Einflussgrößen für die logistische Regression

Man kann diese Einstellung anpassen, muss aber darauf achten, nach der Anpassung das Feld *Ändern* anzuklicken. Erst dann wird die Anpassung übernommen. Nachdem man auf *Weiter* geklickt hat, kann man überprüfen, ob die zugrunde gelegte Skalierung jetzt stimmt. Bei den als kategorial modellierten Einflussgrößen steht dann „Cat“ in Klammern hinter dem Variablennamen (siehe Abb. 3.14).

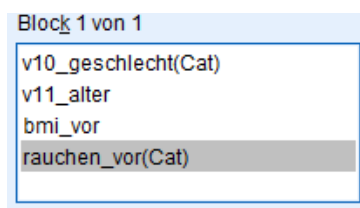


Abb. 3.14: Feld Kovariaten im Hauptfenster nach der Festlegung der kategorialen Einflussgrößen

Über das Feld *Optionen* (siehe Abb. 3.15) kann man u. a. einstellen, dass die Konfidenzintervalle für die geschätzten Odds Ratios mit ausgegeben werden (*Konfidenzintervall für $Exp(B)$*). Über *Weiter* und *OK* im Hauptfenster wird die logistische Regressionsanalyse durchgeführt.

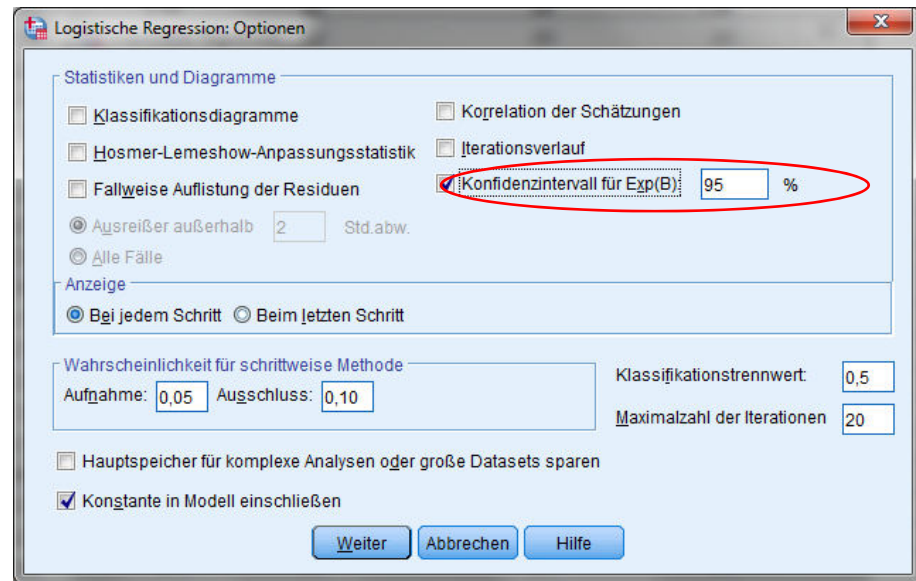


Abb. 3.15: Einstellen der Option „Ausgabe der Konfidenzintervalle“ bei der logistischen Regressionsanalyse

Ergebnisdarstellung in SPSS

Relativ zu Beginn der SPSS-Ausgabe steht die Tabelle „Codierungen kategorialer Variablen“ (siehe Abb. 3.15). Diese Informationen sind wichtig, um später die Regressionskoeffizienten richtig interpretieren zu können. Hier weist SPSS den Nichtrauchern und den Frauen jeweils die Kategorie (1) zu.

Codierungen kategorialer Variablen			
		Häufigkeit	Parameter-codierung
			(1)
rauchen_vor	0	64	1,000
	1	26	,000
Geschlecht	weiblich	64	1,000
	männlich	26	,000

Abb. 3.15: Kodierung der kategorialen Variablen rauchen_vor und Geschlecht in SPSS

Interpretation der Ergebnisse

Anschließend kommen die Ergebnisse der logistischen Regression. Da hier eine Rückwärtsselektion gewählt wurde, wurden mehrere Modelle gerechnet. Unter der Überschrift *Block 0: Anfangsblock* werden die Ergebnisse des kompletten Modells dargestellt. Weiter unten, unter der Überschrift *Block 1: Methode = Rückwärts Schrittweise (Konditional)*, sind dann die Ergebnisse der folgenden Modelle angegeben. Dabei bezeichnet *Schritt 1* usw. das jeweilige Modell. In unserem Beispiel wurden insgesamt vier Modelle gerechnet. Wir konzentrieren uns auf das endgültige, d. h. das vierte Modell. In der Tabelle „Modellzusammenfassung“ kann die Devianz unter der Bezeichnung *-2 Log-Likelihood* abgelesen werden (siehe Abb. 3.17). Man sieht, dass sich die Devianz (*-2 Log-Likelihood*) zwischen den verschiedenen Modellen kaum unterscheidet.

Modellzusammenfassung

Schritt	-2 Log-Likelihood	Cox & Snell R-Quadrat	Nagelkerkes R-Quadrat
1	114,242 ^a	,088	,119
2	114,314 ^a	,088	,118
3	115,698 ^a	,074	,099
4	117,013 ^b	,060	,081

a. Schätzung beendet bei Iteration Nummer 4, weil die Parameterschätzer sich um weniger als ,001 änderten.

b. Schätzung beendet bei Iteration Nummer 3, weil die Parameterschätzer sich um weniger als ,001 änderten.

Abb. 3.17: Tabelle mit Angaben zur Devianz (-2 Log-Likelihood) in den einzelnen Modellschritten

Weiterhin ist die Tabelle „Variablen in der Gleichung“ von Interesse (siehe Abb. 3.18). In der ersten Spalte sieht man, welche Variablen wann aus dem Modell entfernt wurden (oder bei der Vorwärtsselektion bzw. schrittweisen Selektion entsprechend auch hinzugefügt wurden). Im endgültigen Modell ist also nur noch der Raucherstatus vor der Reha (*rauchen_vor*) mit einem p-Wert (*Sig.*) von 0,020 enthalten. Das zugehörige Odds Ratio (*Exp(B)*) liegt bei 0,327. Dabei bezieht sich das Odds Ratio auf den Vergleich der Kategorie (*rauchen_vor(1)*) mit den anderen Kategorien. Im Beispiel *rauchen_vor(1)* gibt es nur eine Vergleichskategorie (Referenzkategorie): die Raucher (*rauchen_vor = 0*). Wie oben bereits erwähnt, bezeichnet die Kategorie (1) bei der Variable *rauchen_vor* die Nichtraucher. Das bedeutet, dass die Wahrscheinlichkeit, dass die Ärzte Angehörige in die Behandlung mit einbeziehen, bei Nichtrauchern (*rauchen_vor(1)*) deutlich geringer ist als bei Rauchern (genauer gesagt: Die Wahrscheinlichkeit für den Einbezug von Angehörigen ist 0,327-mal so groß). Das zugehörige Konfidenzintervall (95 % Konfidenzintervall für *EXP(B)*) geht von 0,127 bis 0,841.

Variablen in der Gleichung

		Regressions- koeffizient B	Standard- fehler	Wald	df	Sig.	Exp(B)	95% Konfidenzintervall für EXP(B)	
								Unterer Wert	Oberer Wert
Schritt 1 ^a	Geschlecht(1)	,661	,575	1,321	1	,250	1,937	,627	5,978
	Alter	,095	,074	1,659	1	,198	1,100	,951	1,272
	Body Mass Index vor Reha	,065	,242	,072	1	,788	1,067	,664	1,716
	rauchen_vor(1)	-1,085	,499	4,723	1	,030	,338	,127	,899
	Konstante	-7,348	9,015	,664	1	,415	,001		
Schritt 2 ^a	Geschlecht(1)	,592	,511	1,342	1	,247	1,807	,664	4,916
	Alter	,091	,072	1,591	1	,207	1,096	,951	1,263
	rauchen_vor(1)	-1,112	,490	5,156	1	,023	,329	,126	,859
	Konstante	-5,218	4,266	1,497	1	,221	,005		
Schritt 3 ^a	Alter	,081	,071	1,291	1	,256	1,084	,943	1,247
	rauchen_vor(1)	-1,084	,484	5,004	1	,025	,338	,131	,874
	Konstante	-4,216	4,140	1,037	1	,308	,015		
Schritt 4 ^a	rauchen_vor(1)	-1,117	,481	5,380	1	,020	,327	,127	,841
	Konstante	,470	,403	1,359	1	,244	1,600		

a. In Schritt 1 eingegebene Variablen: Geschlecht, Alter, Body Mass Index vor Reha, rauchen_vor.

Abb. 3.18: Tabelle mit den geschätzten Regressionskoeffizienten aus der logistischen Regression

Übungsaufgaben

- 3.1) Es soll der Einfluss von sportlicher Betätigung (in Minuten pro Tag), Alter und BMI vor der Reha auf die Gewichtsveränderung (in kg) untersucht werden. Überprüfen Sie, ob die Voraussetzungen für eine lineare Regression erfüllt sind. Wenn dem so ist, führen Sie die lineare Regression durch (u.U. mit Variablenselektion) und interpretieren Sie die Ergebnisse.
- 3.2) Es soll für verschiedene Variablen (Geschlecht, Alter, sportliche Betätigung, wahrgenommene Freizeitangebote und Gewichtsveränderung) untersucht werden, ob sie einen Einfluss auf die Zufriedenheit mit der Reha-Maßnahme haben.

Dichotomisieren Sie die Variable `v2_v9_Score` mit dem Median von 19 als Grenze. Führen Sie die logistische Regression durch (u.U. mit Variablenselektion) und interpretieren Sie die Ergebnisse.

Zusammenfassung

Wenn der Zusammenhang mehrerer Variablen gemeinsam untersucht werden soll, benötigt man multivariate Verfahren. Die gängigsten multivariaten Verfahren sind die Varianzanalyse und die Regressionsanalyse (hierbei insbesondere die lineare und die binäre logistische Regression. In diesem Studienbrief sollte gezeigt werden, wie die Verfahren in SPSS angewendet werden können.

Die Varianzanalyse kann für metrische Zielgrößen angewendet werden, wenn mindestens eine kategoriale Einflussgröße (Einflussfaktor) vorliegt. Mit Hilfe der Varianzanalyse kann insbesondere überprüft werden, ob Wechselwirkungen zwischen Einflussfaktoren bestehen. Ferner wird zwischen Varianzanalysen ohne Kovariablen (ANOVA) und mit Kovariablen (ANCOVA) unterschieden.

Die lineare Regression ist in erster Linie gedacht für die Untersuchung des Zusammenhangs zwischen metrischen Variablen. Es können aber auch Einflussfaktoren in das Modell aufgenommen werden.

Wenn die Zielgröße kategorial, insbesondere binär skaliert ist, ist die logistische Regression das Analyseinstrument der Wahl. Dabei können die Einflussgrößen sowohl kategorial als auch metrisch sein. Die geschätzten Regressionskoeffizienten entsprechen hier den Odds Ratios.

Das Ziel dieses Studienbriefs war, die Leserinnen und Leser zu befähigen, die vier grundlegenden multivariaten Verfahren (ANOVA, ANCOVA, lineare Regression und logistische Regression) selbständig in SPSS anzuwenden. Dazu gehört, zu wissen, wann welche Voraussetzungen überprüft werden müssen und wie das mit SPSS geht. Weiterhin gehört dazu, die Analysen durchführen zu können, und auch zu wissen, welche wichtigen Zusatzfunktionen (wie z. B. Ausgabe der Koeffizientenschätzer) zur Verfügung stehen. Und als Letztes, aber keinesfalls weniger Wichtiges, gehört dazu, die Ergebnisausgabe „lesen“ und interpretieren zu können.

Multivariate Verfahren

Varianzanalyse

Lineare Regression

Logistische Regression

Ziele des Studienbriefs

Glossar

Bestimmtheitsmaß (R^2): quadrierter Korrelationskoeffizient, der angibt, welcher Anteil der Streuung der Zielgröße durch die Einflussgrößen erklärt wird (Gütekriterium bei der linearen Regression).

Binär: Skalenniveau einer Variablen, die nur zwei Ausprägungen hat (z. B. Geschlecht mit den Ausprägungen „männlich“ und „weiblich“). Synonym für dichotom.

Devianz: Gütekriterium bei der logistischen Regression.

Dichotom: Skalenniveau einer Variablen, die nur zwei Ausprägungen hat (z. B. Geschlecht mit den Ausprägungen „männlich“ und „weiblich“).

Einfachregression: Regressionsanalyse mit einer Einflussgröße.

Einfaktorielle ANOVA: ANOVA mit einer Einflussgröße.

Einflussgröße: Variable, deren Einfluss auf die Zielgröße untersucht werden soll (auch unabhängige Variable genannt).

Einflussfaktor: kategoriale Einflussgröße.

Interaktion: Wechselwirkung zwischen Einflussfaktoren.

Konfirmatorisch: eine bestehende Auffassung/Hypothese empirisch bestätigend.

Kovariabel: nicht primär interessierende Einflussgröße.

Kovarianzanalyse (ANCOVA): ANOVA unter zusätzlicher Berücksichtigung von Kovariablen.

Levene-Test: statistischer Test, mit dem geprüft werden kann, ob die Varianzen von zwei oder mehreren Gruppen gleich sind.

Lineare Regression: Regressionsanalyse für Variablen mit linearem Zusammenhang.

Logistische Regression: Regressionsanalyse für binäre (dichotome) Zielgrößen.

MANOVA: multivariate Varianzanalyse.

Maximum-Likelihood-Methode: rekursives Verfahren zur Schätzung der Regressionskoeffizienten bei der logistischen Regression.

Mehrfaktorielle ANOVA: ANOVA mit mehreren Einflussgrößen.

Messwiederholung: mehrfaches Messen einer Zielvariablen.

Methode der Kleinsten Quadrate: Methode zur Schätzung der Regressionskoeffizienten bei der linearen Regression.

Multikollinearität: starke Abhängigkeiten zwischen den Einflussgrößen. Wechselseitige Abhängigkeit der unabhängigen Variablen (Einflussgrößen) in einem statistischen Modell. Perfekte Multikollinearität liegt vor, wenn die Werte einer (oder mehrerer) unabhängiger Variablen aus den anderen unabhängigen Variablen exakt vorhergesagt werden können.

Multiple Regression: Regressionsanalyse mit mehreren Einflussgrößen.

Multivariat: Verfahren, in denen mindestens drei Merkmale (Variablen) statistisch analysiert werden.

Odds Ratio: Chancenverhältnis. Maß für die Stärke des Unterschieds zwischen zwei Gruppen. Das Odds Ratio setzt die Odds (Chancen) der beiden Gruppen zueinander ins Verhältnis.

Power: $(1-\beta)$. Wahrscheinlichkeit eine Nullhypothese korrekterweise abzulehnen.

p-Wert: Wahrscheinlichkeit, dass der beobachtete Effekt bei Gültigkeit der Nullhypothese auftritt.

Referenzkategorie: in der Regression diejenige Kategorie einer erklärenden Variablen, die als Vergleichskategorie dient.

Regressionsanalyse: Zusammenhangsanalyse für zwei oder mehr Variablen.

Regressionsgleichung: Gleichungssystem, das den Zusammenhang zwischen Variablen quantifiziert.

Regressionskoeffizienten: konstante Werte in der Regressionsgleichung.

Repeated measures ANOVA: ANOVA unter der Berücksichtigung von Messwiederholungen.

Residuen: Fehlerterme bei der Regressionsanalyse.

Scheinkorrelation: scheinbarer Zusammenhang zwischen zwei Variablen.

Signifikanz: Überzufälligkeit eines Effekts.

Streuungszerlegung: Zerlegung der Gesamtstreuung in die Streuung innerhalb der Gruppen, die Streuung zwischen den Gruppen und die zufällige Streuung.

Syntax: Programmiersprache.

Variablenselektion: Reduktion der Einflussgrößen in der mehrfaktoriellen Analyse.

Varianzanalyse (ANOVA): Verallgemeinerung des t-Tests zum Vergleich mehrerer Gruppen.

Varianzhomogenität: Gleichheit der Varianzen verschiedener Gruppen.

Zielgröße: interessierende Variable deren Wert sich u. U. durch Einflussgrößen erklären lässt (auch abhängige Variable, Zielvariable genannt). Der Begriff ist relativ zu dem interessierenden Modell zu sehen: Eine Variable, die in einem Modell als Zielgröße fungiert, kann in einem anderen Modell lediglich eine Einflussgröße sein (und umgekehrt).

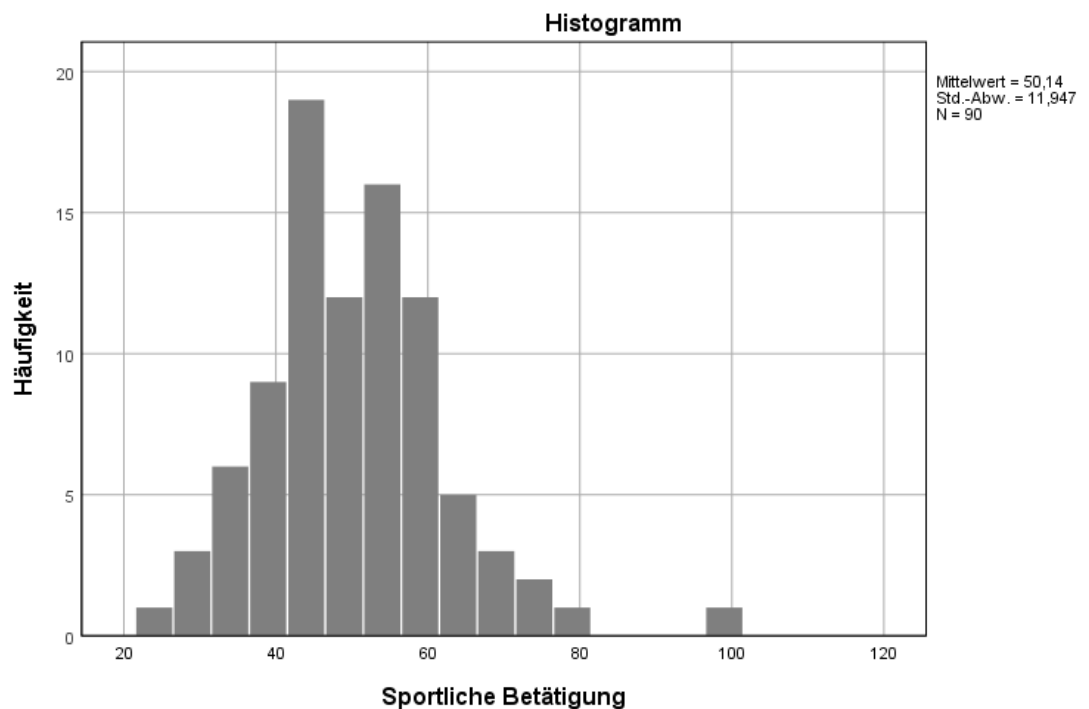
Lösungen zu den Übungsaufgaben

Lösungen zu Kapitel 1

- 1.1) Das Ergebnis ist statistisch signifikant, da der p-Wert=0,020 kleiner als 0,05 ist. Die Signifikanz kann auch daran abgelesen werden, dass die obere Grenze des Konfidenzintervalls für das Odds Ratio kleiner 1 ist. Der Odds Ratio-Wert drückt aus, dass Nichtraucher gegenüber Rauchern eine geringere Chance haben, dass ihre Angehörigen in die Therapie einbezogen werden. Die Chance von Nichtrauchern beträgt gegenüber der von Rauchern 0,327, d. h. nur 32,7%. Die Nullhypothese (Raucher und Nichtraucher unterscheiden sich nicht hinsichtlich der Einbeziehung von Angehörigen in die Therapie) ist daher abzulehnen.
- 1.2) Da die Zielgröße und auch die Einflussgrößen metrisch sind, ist die lineare Regression das geeignete Analyseverfahren. Das statistische Modell beinhaltet die folgenden Komponenten: die Zufriedenheit (Zielgröße) und als metrische Einflussgrößen: Anzahl der wahrgenommenen Freizeitangebote; durchschnittliche Zeit, die pro Tag Sport getrieben wird und die Veränderung des BMI.

Lösungen zu Kapitel 2

- 2.1) Die Fragestellung führt zu einem statistischen Modell mit einer metrischen Zielgröße (*v15_sport*) und zwei kategorialen Einflussgrößen (*v10_geschlecht*, *v17_angehörige*). Voraussetzungen 1 und 3 für die Durchführung einer Varianzanalyse sind damit erfüllt. Daher ist die mehrfaktorielle ANOVA mit Interaktionsterm die korrekte Analyseverfahren. Über *Analysieren* → *Allgemeines lineares Modell* → *Univariat* werden *v15_sport* als *Abhängige Variable* und *v10_geschlecht* sowie *v17_angehörige* als *Feste Faktoren* ausgewählt. Es wird ein *gesättigtes Modell* gewählt und der *Levene-Homogenitätstest* wird berechnet. Zur Überprüfung der Normalverteilung in der Gesamtpopulation wird ein Histogramm gezeichnet (es könnten ebenso deskriptive Kenngrößen oder ein Boxplot zur Überprüfung verwendet werden). Das Ergebnis des Histogramms zeigt, dass eine Normalverteilung angenommen werden kann.



Die Ergebnisse der ANOVA beinhalten den Levene-Test auf Varianzhomogenität. Da der p-Wert hier gleich 0,603 ist, kann auch Voraussetzung 4 als erfüllt angesehen werden.

Die Ergebnisse der ANOVA sind:

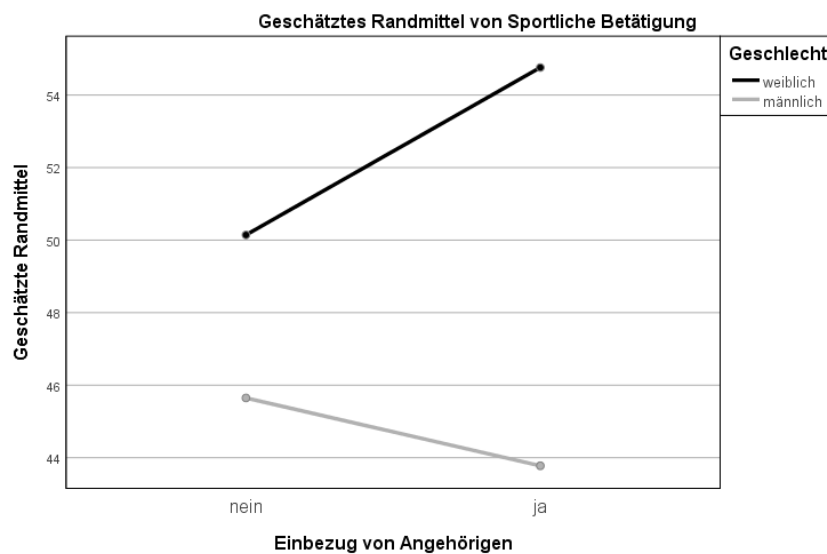
Tests der Zwischensubjekteffekte

Abhängige Variable: Sportliche Betätigung

Quelle	Quadratsumme vom Typ III	df	Mittel der Quadrate	F	Sig.
Korrigiertes Modell	1326,088 ^a	3	442,029	3,341	,023
Konstanter Term	162079,516	1	162079,516	1225,173	,000
v10_geschlecht	1028,060	1	1028,060	7,771	,007
v17_angehörige	32,376	1	32,376	,245	,622
v10_geschlecht * v17_angehörige	180,506	1	180,506	1,364	,246
Fehler	11377,034	86	132,291		
Gesamt	239005,000	90			
Korrigierte Gesamtvariation	12703,122	89			

a. R-Quadrat = ,104 (korrigiertes R-Quadrat = ,073)

Der p-Wert von 0,246 bei der Interaktion v10_geschlecht*v17_angehörige bedeutet, dass eine Wechselwirkung (Interaktion) nicht ausgeschlossen werden kann. Das *Profildiagramm* zeigt, dass Frauen, bei denen Angehörige in die Behandlung einbezogen werden, im Durchschnitt mehr Sport treiben als Frauen, deren Angehörige nicht einbezogen werden. Bei Männern ist der Effekt umgekehrt, d.h. Männer, bei denen die Angehörigen einbezogen werden, haben sich weniger sportlich betätigt. Daher empfiehlt es sich, den Effekt von Geschlecht und der Einbeziehung Angehöriger getrennt voneinander zu untersuchen.



Bei den resultierenden einfaktoriellen Varianzanalysen (d.h., es wird jeweils nur ein *Fester Faktor* bei der ANOVA berücksichtigt: v10_geschlecht bzw. v17_angehörige) zeigt sich, dass für Geschlecht als Einflussgröße die Voraussetzung der Varianzhomogenität erfüllt ist (p-Wert des Levene-Tests=0,803). Die ANOVA bestätigt den starken Einfluss des Geschlechts auf die durchschnittliche Zeit sportlicher Betätigung (p=0,008).

Bei der Einbeziehung Angehöriger ist die Voraussetzung der Varianzhomogenität nicht erfüllt (p-Wert des Levene-Tests=0,104). Dadurch kann man lediglich schlussfolgern, dass ein Zusammenhang beobachtbar ist (anhand des Profildiadgramms und der deskriptiven Kenngrößen). Der Unterschied zwischen den beiden Gruppen (Einbeziehung Angehöriger nein vs. ja) ist allerdings relativ klein (ca. 3,5 Minuten Mittelwertdifferenz; 52,16–48,67=3,49).

ONEWAY deskriptive Statistiken

Sportliche Betätigung

	N	Mittelwert	Std.-Abweichung	Std.-Fehler	95%-Konfidenzintervall für den Mittelwert		Minimum	Maximum
					Untergrenze	Obergrenze		
nein	52	48,67	9,513	1,319	46,02	51,32	28	71
ja	38	52,16	14,539	2,359	47,38	56,94	24	98
Gesamt	90	50,14	11,947	1,259	47,64	52,65	24	98

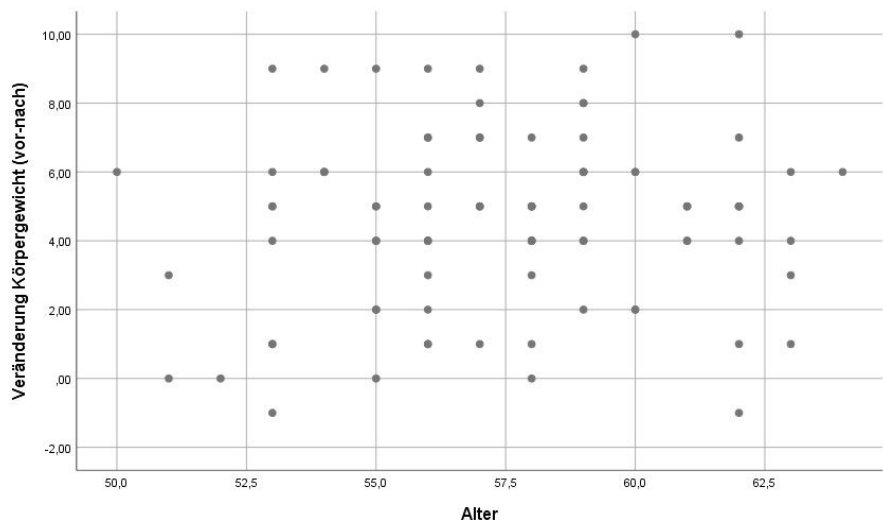
- 2.2) Da wir in Aufgabe 2.1) eine Wechselwirkung zwischen der Einbeziehung Angehöriger und dem Geschlecht aufgedeckt haben, wird die Kovariable „BMI vor der Reha“ in die jeweiligen einfaktoriellen Varianzanalysen aufgenommen. Damit erhält man zwei einfaktorielle ANCOVA's (d.h., es wird jeweils nur ein *Fester Faktor* bei der ANCOVA berücksichtigt: *v10_geschlecht* bzw. *v17_angehörige*)

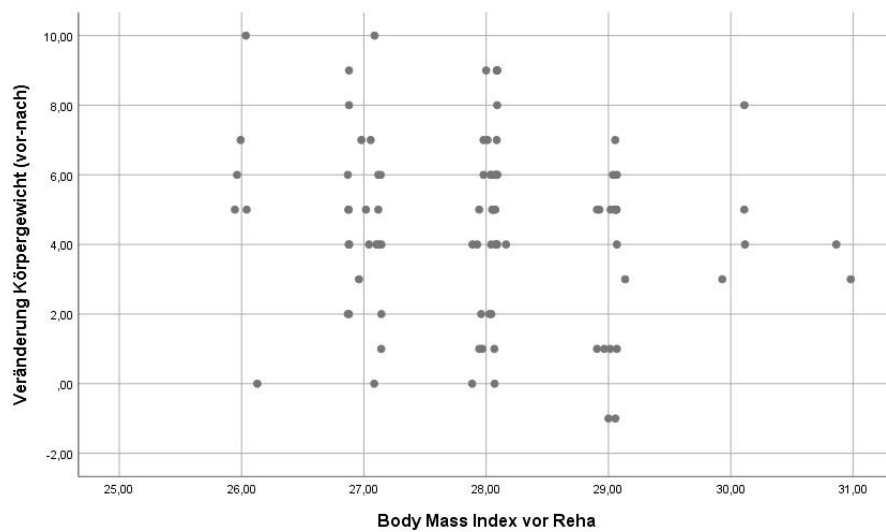
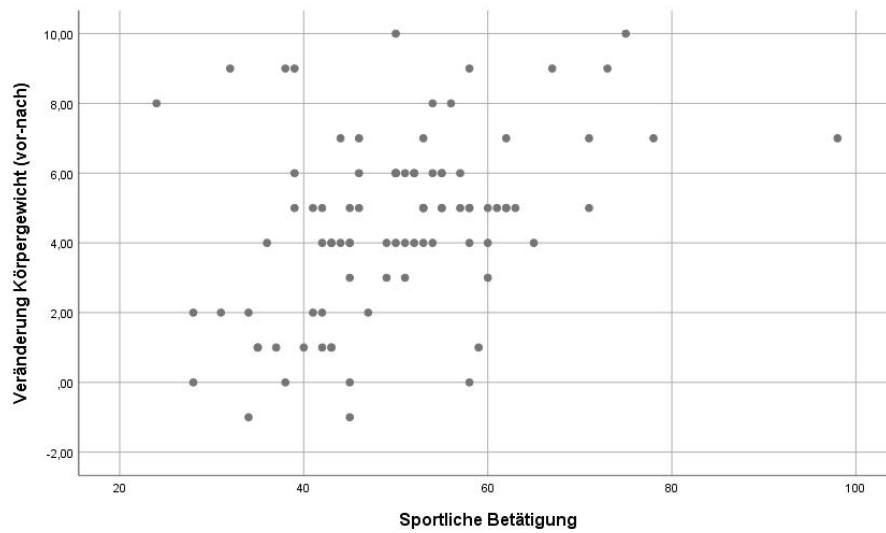
Bei der ANCOVA mit Geschlecht als Einflussgröße ist die Varianzhomogenität erfüllt (p-Wert des Levene-Tests=0,876). Der Einfluss des Geschlechts auf die durchschnittliche Zeit sportlicher Betätigung bleibt bestehen (p=0,036). Für den BMI vor der Reha zeigt sich hingegen kein Effekt (p=0,430).

Bei der ANCOVA mit Einbeziehung Angehöriger als Einflussgröße kann Varianzheterogenität weitgehend ausgeschlossen werden (p-Wert des Levene-Tests=0,280). Hier zeigt sich ein größerer Effekt des BMI vor der Reha als die Einbeziehung Angehöriger (p=0,075 vs. 0,167).

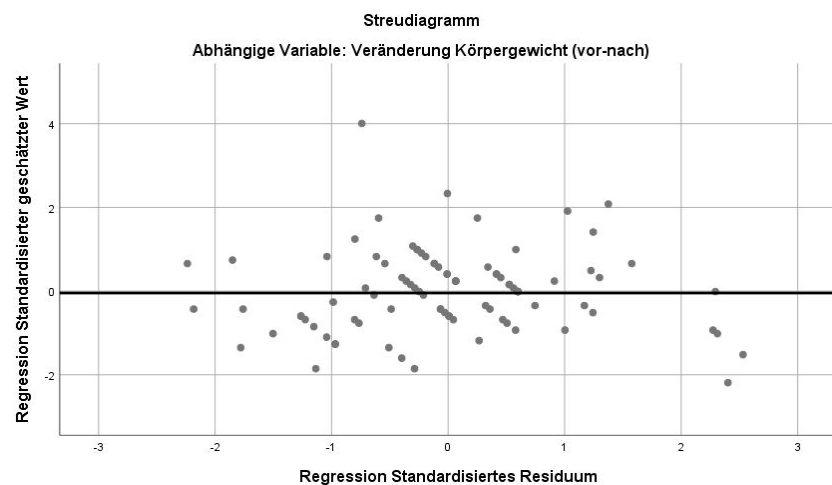
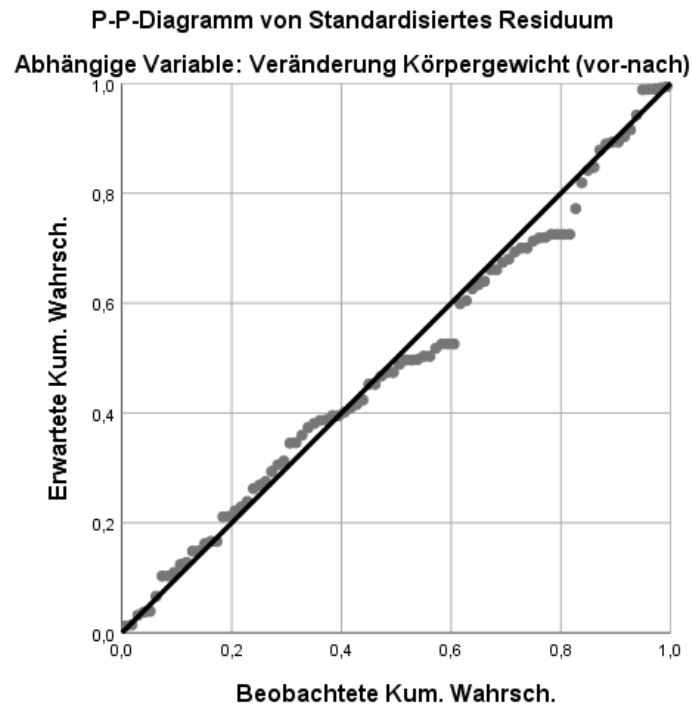
Lösungen zu Kapitel 3

- 3.1) Zuerst wird überprüft, ob ein linearer Zusammenhang zwischen den Einflussgrößen (*v11_alter*, *v15_sport*, *bmi_vor*) und der Zielgröße (*diff_gew*) besteht. Die Streudiagramme zeigen, dass diese Voraussetzung angenommen werden kann.





Das P-P-Diagramm zeigt lediglich leichte Abweichungen von der Diagonale. Das Streudiagramm (standardisierte Residuen vs. standardisierte Vorhersagewerte) zeigt eine Punktwolke, zentriert um die Null und ohne einen Trend hinsichtlich der Varianzen. Die Annahmen der Normalverteilung und Varianzhomogenität für die Varianzen sind also ebenfalls erfüllt. Somit liegen die Voraussetzungen für eine lineare Regression vor.



Anschließend wird die lineare Regression mit Rückwärtsselektion durchgeführt. Das endgültige Modell beinhaltet lediglich die tägliche Zeit sportlicher Betätigung (*v15_sport*). Das Bestimmtheitsmaß von 0,404 spricht für keine besonders gute Modellanpassung. Die Regressionskoeffizienten sind in der Tabelle „Koeffizienten“ zu finden:

Koeffizienten^a

Modell	Nicht standardisierte Koeffizienten		Standardisierte Koeffizienten	T	Sig.	95,0% Konfidenzintervalle für B	
	Regressionskoeffizient B	Std.-Fehler	Beta			Untergrenze	Obergrenze
1 (Konstante)	1,265	8,782		,144	,886	-16,192	18,722
Alter	,058	,079	,072	,726	,470	-,100	,216
Sportliche Betätigung	,083	,021	,388	3,882	,000	,041	,126
Body Mass Index vor Reha	-,149	,242	-,062	-,614	,541	-,630	,332
2 (Konstante)	-3,326	4,590		-,725	,471	-12,449	5,796
Alter	,063	,079	,079	,802	,424	-,093	,220
Sportliche Betätigung	,086	,021	,399	4,076	,000	,044	,127
3 (Konstante)	,253	1,079		,235	,815	-1,891	2,398
Sportliche Betätigung	,087	,021	,404	4,140	,000	,045	,128

a. Abhängige Variable: Veränderung Körpergewicht (vor-nach)

Das bedeutet, pro Minute, die im Durchschnitt pro Tag mehr Sport getrieben wird, wird im Durchschnitt 0,087 kg Gewicht während der Reha verloren.

- 3.2) In einem ersten Schritt muss die Zufriedenheit dichotomisiert werden. Entweder über *Transformieren* → *Umkodieren in andere Variable*, oder direkt über die Syntax (siehe *Syntaxdatei*).

Da für die logistische Regression keine Voraussetzungen zu erfüllen sind, kann die Analyse direkt durchgeführt werden. Bei Verwendung der bedingten Rückwärtsselektion erhalten wir am Ende ein Modell mit Geschlecht und Gewichtsveränderung als Einflussgrößen. Die Devianz ist mit einem Wert von 85,380 relativ niedrig, was für eine gute Modellanpassung spricht. Die Hauptergebnisse sind im Ergebnisfenster von SPSS in der Tabelle „Variablen in der Gleichung“ dargestellt:

Variablen in der Gleichung

		Regressions- koeffizientB	Standard- fehler	Wald	df	Sig.	Exp(B)	95% Konfidenzintervall für EXP(B)	
								Unterer Wert	Oberer Wert
Schritt 1 ^a	Geschlecht(1)	-1,862	,743	6,277	1	,012	,155	,036	,667
	Alter	-,050	,087	,339	1	,561	,951	,802	1,127
	Sportliche Betätigung	-,025	,028	,785	1	,376	,975	,922	1,031
	Anzahl wahrgenomme- ner Freizeitangebote			2,054	2	,358			
	Anzahl wahrgenomme- ner Freizeitangebote(1)	-1,100	,891	1,521	1	,217	,333	,058	1,911
	Anzahl wahrgenomme- ner Freizeitangebote(2)	-1,261	,891	2,003	1	,157	,283	,049	1,625
	Veränderung Körperge- wicht (vor-nach)	,921	,216	18,188	1	,000	2,513	1,645	3,837
	Konstante	2,049	5,093	,162	1	,687	7,763		
Schritt 2 ^a	Geschlecht(1)	-1,813	,733	6,116	1	,013	,163	,039	,686
	Sportliche Betätigung	-,027	,028	,928	1	,335	,973	,921	1,028
	Anzahl wahrgenomme- ner Freizeitangebote			1,998	2	,368			
	Anzahl wahrgenomme- ner Freizeitangebote(1)	-1,117	,892	1,567	1	,211	,327	,057	1,881
	Anzahl wahrgenomme- ner Freizeitangebote(2)	-1,233	,890	1,919	1	,166	,291	,051	1,668
	Veränderung Körperge- wicht (vor-nach)	,918	,215	18,191	1	,000	2,504	1,642	3,818
	Konstante	-,789	1,492	,280	1	,597	,454		
Schritt 3 ^a	Geschlecht(1)	-1,646	,699	5,538	1	,019	,193	,049	,760
	Sportliche Betätigung	-,031	,028	1,211	1	,271	,970	,918	1,024
	Veränderung Körperge- wicht (vor-nach)	,911	,210	18,869	1	,000	2,488	1,649	3,753
	Konstante	-1,712	1,319	1,684	1	,194	,181		
Schritt 4 ^a	Geschlecht(1)	-1,694	,696	5,921	1	,015	,184	,047	,719
	Veränderung Körperge- wicht (vor-nach)	,834	,191	19,029	1	,000	2,302	1,583	3,349
	Konstante	-2,899	,825	12,361	1	,000	,055		

a. In Schritt 1 eingegebene Variablen: Geschlecht, Alter, Sportliche Betätigung, Anzahl wahrgenommener Freizeitangebote, Veränderung Körpergewicht (vor-nach).

Die Ergebnisse der logistischen Regression (siehe Schritt 4) verdeutlichen, dass die Chance für Zufriedenheit mit der Reha-Maßnahme (d.h., dass ein Score-Wert von > 19 erreicht wird) bei den Frauen geringer ist als bei den Männern (Exp(B)=0,184). Mit jedem Kilo, dass ein Patient während der Reha-Maßnahme abnimmt, erhöht sich die Chance für Zufriedenheit (Exp(B)=2,302), d.h. mit jedem abgenommenen Kilo erhöht sich die Zufriedenheit um 2,302.

Literaturverzeichnis

- Bühl, A. (2012). *SPSS 20. Einführung in die moderne Datenanalyse* (13., aktualisierte Auflage). München: Pearson.
- Duller, C. (2013). *Einführung in die Statistik mit EXCEL und SPSS* (3. Auflage). Berlin, Heidelberg: Springer Gabler.
- Fahrmeir, L., Künstler, R. Pigeot, I. & Tutz, G. (2012). *Statistik. Der Weg zur Datenanalyse* (7. Auflage). Berlin, Heidelberg: Springer.
- Fromm, S. (2012). *Datenanalyse mit SPSS für Fortgeschrittene 2: Multivariate Verfahren für Querschnittsdaten*. Wiesbaden: Springer VS.
- Gaus, W. & Mücke, R. (2014). *Medizinische Statistik. Angewandte Biometrie für Ärzte und Gesundheitsberufe*. Stuttgart: Schattauer.
- Hartung, J., Elpelt, B. & Klösener, K.-H. (2005). *Statistik. Lehr- und Handbuch der angewandten Statistik* (14. Auflage). München: Oldenbourg.
- Janssen, J. & Laatz, W. (2013). *Statistische Datenanalyse mit SPSS. Eine anwendungsorientierte Einführung in das Basissystem und das Modul Exakte Tests* (8. Auflage). Berlin, Heidelberg: Springer.
- Rudolf, M. & Müller, J. (2012). *Multivariate Verfahren. Eine praxisorientierte Einführung mit Anwendungsbeispielen in SPSS*. Göttingen: Hogrefe.
- Schumacher, M. & Schulgen, G. (2008). *Methodik klinischer Studien – Methodische Grundlagen der Planung, Durchführung und Auswertung*. Berlin, Heidelberg: Springer.