

**SISS:**

**Schriftenreihe des Instituts für Sozialwissenschaften  
der Universität Stuttgart**

**No. 3 / 2010**

**Binär-logistische Regressionsanalyse.  
Grundlagen und Anwendung  
für Sozialwissenschaftler**

**Jochen Mayerl  
Dieter Urban**

**Universität Stuttgart  
Institut für Sozialwissenschaften  
Abteilung für Soziologie und  
empirische Sozialforschung (SOWI IV)  
70174 Stuttgart**



**SOWI**

ISSN 0945-9197

**SISS:  
Schriftenreihe  
des Instituts für Sozialwissenschaften  
der Universität Stuttgart: No. 3 / 2010**

---

Binär-logistische Regressionsanalyse.

Grundlagen und Anwendung für Sozialwissenschaftler

---

Jochen Mayerl  
Dieter Urban

**Universität Stuttgart  
Institut für Sozialwissenschaften  
Abteilung für Soziologie und  
empirische Sozialforschung (SOWI IV)  
70174 Stuttgart**

## Binär-logistische Regressionsanalyse. Grundlagen und Anwendung für Sozialwissenschaftler.

**Z U S A M M E N F A S S U N G:** Das Skript beschreibt die Durchführung von binär-logistischen Regressionsanalysen in den Sozialwissenschaften unter Verwendung des Statistik-Programmpakets SPSS. Dabei wird erläutert, warum für bestimmte Untersuchungen und für bestimmte Datenanalysen in der sozialwissenschaftlichen Forschung nicht die klassische OLS-Regressionanalyse sondern die logistische Regressionanalyse (mit Maximum-Likelihood Schätzverfahren) eingesetzt werden sollte. Es wird gezeigt, nach welcher internen Logik logistische Regressionsschätzungen verfahren, in welcher Weise diese Regressionsschätzungen mit SPSS durchgeführt werden können und wie die Ergebnisse von logistischen Regressionsanalysen zu interpretieren sind. Auch werden die häufigsten Probleme, die bei der Durchführung von logistischen Regressionsanalysen auftreten können, vorgestellt, und es werden Möglichkeiten zur Identifikation und Beseitigung dieser Probleme aufgezeigt.

## Binary logistic regression analysis. Basic principles and application for social scientists.

**A B S T R A C T:** This report describes how to conduct binary logistic regression analysis in social science research utilizing the statistical software package SPSS. It explains why certain social science problems and certain social science data should be analyzed by logistic regression (with maximum likelihood estimation techniques) and should not be analyzed by classical OLS-regression procedures. The report informs about the internal logic of logistic regression estimation, shows how to handle logistic regression modeling in SPSS, and demonstrates how to interpret the results of logistic regression estimations. In addition, the paper specifies some of the most frequent problems when estimating logistic regression models and gives some advice how to detect and solve these problems.

## 1 Einleitung<sup>1</sup>

Klassische Regressionsschätzungen werden nach der OLS-Methode durchgeführt. Denn die Ordinary-Least-Squares- bzw. die Kleinst-Quadrate-Schätzmethode ist dasjenige Verfahren, mit dem optimale Schätzwerte für die Koeffizienten der Regressionsgleichung ermittelt werden können. Die OLS-Schätzung kann optimale Schätzwerte mit BLUE-Eigenschaften errechnen (sog. „Best Linear Unbiased Estimation“, vgl. dazu Urban/Mayerl 2008), wenn die dafür geltenden Modellvoraussetzungen gegeben sind (z.B. die Abwesenheit von Heteroskedastizität bzw. von Streuungsungleichheit). In der OLS-Regressionsanalyse sind daher mögliche Verstöße gegen diese Voraussetzungen zu identifizieren und ggfs. Maßnahmen einzuleiten, um einige Modellverstöße zu beseitigen bzw. in ihren negativen Konsequenzen abzumildern (vgl. Urban/Mayerl 2008).

Was kann jedoch gemacht werden, wenn in einem Forschungsprojekt aufgrund der empirischen Datenlage von vornherein zu erkennen ist, dass die Voraussetzungen für eine BLUE-Schätzung nicht gegeben sind? Sollte dann dennoch eine OLS- Regressionsanalyse durchgeführt werden (vielleicht mit eingeschränkten Ansprüchen an die Qualität der zu ermittelnden Schätzwerte)? Nein, das wäre der falsche Weg. So etwas wurde zwar früher in der Forschungspraxis häufiger gemacht (etwa in Form von binären Regressionsanalysen mit dichotomen Y-Variablen), ist aber heutzutage nicht mehr notwendig, weil in fast allen Statistik-Softwarepaketen (so auch in SPSS) zusätzlich zur OLS-Schätzmethode ein weiteres Schätzverfahren enthalten ist: das sogenannte Maximum-Likelihood-Schätzverfahren (ML-Schätzung). Mit diesem Verfahren ist es möglich, auch dann qualitative ausreichende Regressionsschätzungen durchzuführen, wenn die zu analysierenden Daten den BLUE-Kriterien prinzipiell nicht entsprechen können und deshalb eine OLS-Regressionsschätzung nicht sinnvoll ist. Machen wir uns das an einem Beispiel deutlich:

In einer empirischen Studie sollen sozio-ökonomische Determinanten der Wahlentscheidung für oder gegen die politische Partei „ABC“ untersucht werden. So soll z.B. auch die Abhängigkeit der Wahlentscheidung pro oder contra ABC von der Höhe des sozialen Status der Wähler untersucht werden. Und für diese Untersuchung sei eine Regressionsanalyse durchzuführen. Es ist leicht zu erkennen, dass in diesem Falle die abhängige Variable im Regressionsmodell nur zwei Ausprägungen aufweisen kann: entweder wird die ABC-Partei gewählt ( $Y=1$ ), oder sie wird nicht gewählt ( $Y=0$ ). Die Y-Variable wäre also binär oder binomial skaliert. Und mit einer binären abhängigen Variablen können die wichtigen BLUE-Kriterien nicht erfüllt werden, so dass hier eine OLS-Regressions-

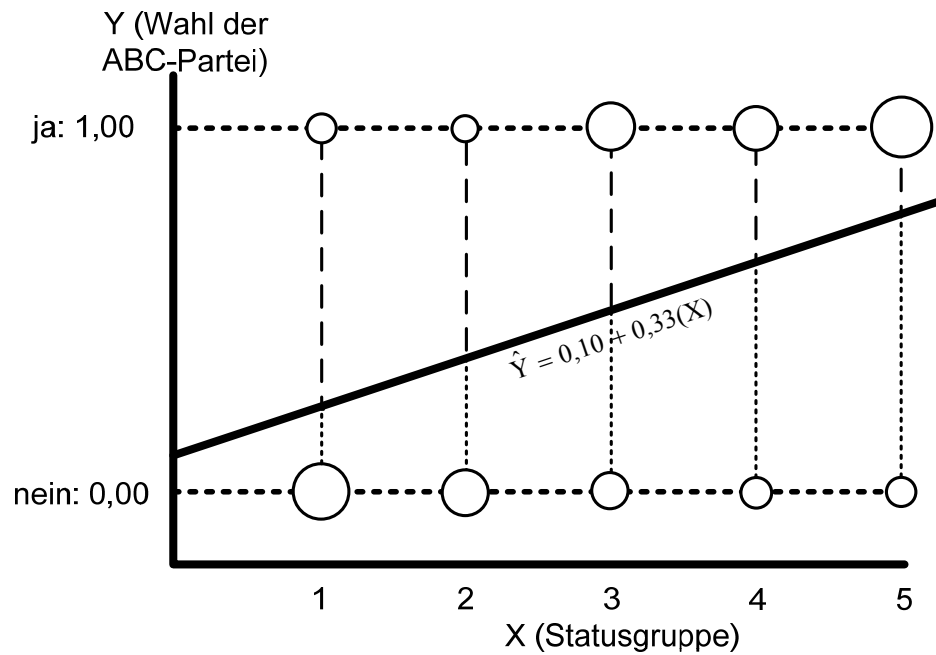
---

<sup>1</sup> Aus Gründen der sprachlichen Vereinfachung wird nachfolgend stets die maskuline Begriffsform verwendet. Die feminine Form gilt dabei als mit eingeschlossen.

schätzung nicht eingesetzt werden sollte. Denn in diesem Regressionsmodell können die Residuen niemals normalverteilt sein und auch die Streuungsgleichheit (Homoskedastizität) der Residuen ist in diesem Modell niemals zu erreichen.

In der folgenden Abbildung 1 wird verdeutlicht, warum BLUE-Kriterien (z.B. die Residuen-Normalverteilung und die Residuen-Streuungsgleichheit) bei binomialen abhängigen Variablen in der Regressionsschätzung nicht zu erfüllen sind. In der Abbildung symbolisieren große Kreise viele Befragte und kleine Kreise wenig Befragte. Die Befragten geben entweder an, die ABC-Partei gewählt zu haben, dann befinden sie sich in den oberen Kreisen (bzw. in den Kreisen auf der gestrichelten Linie  $Y=1$ ), oder sie geben an, die ABC-Partei nicht gewählt zu haben, dann befinden sie sich in den unteren Kreisen (bzw. in den Kreisen auf der gestrichelten Linie  $Y=0$ ). In welchem der oberen oder unteren Kreise sich ein jeder Befragte befindet, wird durch dessen Zugehörigkeit zu einer von fünf Statusgruppen bestimmt (die hier durch Gruppierung der Statuswerte gebildet wurden, um auf diese Weise die Abbildung zu vereinfachen). Mit den beschriebenen Datenpunkten wird nun mittels OLS-Schätzung die in der Abbildung eingezeichnete Regressionsgerade ermittelt. Kein Befragter liegt mit seinem  $Y$ -Wert auf dieser Regressionsgeraden, deshalb gibt es für jeden Befragten einen Residualwert als Differenz zwischen seiner Kästchenposition und dem dazugehörigen Schätzwert ( $\hat{Y}$ ) auf der Regressionsgeraden. Dieser ist in der Abbildung als senkrechte Residuallinie eingezeichnet. Das Besondere ist nun, dass es für alle Mitglieder einer jeden Statusgruppe immer nur einen von zwei Residualwerten gibt, denn entweder liegen sie auf der oberen horizontalen Linie von  $Y=1$  oder auf der unteren horizontalen Linie von  $Y=0$ . Und da es deshalb bei jedem  $X$ -Wert (bzw. bei jeder Statusgruppe) nur zwei verschiedene Residualwerte gibt, können diese auch nicht normalverteilt sein. Zudem müssen sich die Verteilungsmuster der Residuen zwischen den Statusgruppen zwangsläufig voneinander unterscheiden und können niemals identisch sein. Auch variieren die Residuen in Abhängigkeit vom Ausmaß der  $X$ -Variablen: oberhalb der Regressionslinie verkleinern sie sich mit Anstieg der  $X$ -Werte, und unterhalb der Regressionslinie werden sie mit Anstieg der  $X$ -Werte konstant größer. Es wird also schon durch einen kurzen Blick auf Abbildung 1 deutlich, dass bei einer Schätzung mit einer binomialen abhängigen  $Y$ -Variablen die BLUE-Kriterien nicht eingehalten werden können. Somit macht es auch keinen Sinn, bei einer solchen  $Y$ -Variablen eine Schätzung nach dem OLS-Verfahren durchzuführen. Stattdessen sollte in diesem Falle ein Maximum-Likelihood-Schätzverfahren eingesetzt werden. Die Grundprinzipien dieses Verfahrens werden im nächsten Abschnitt 2 verdeutlicht. Und im Anschluss daran wird gezeigt, wie das ML-Schätzverfahren zur Analyse binär-logistischer Regressionsmodelle eingesetzt werden kann.

Abbildung 1: Stilisierte lineare Regressionsschätzung mit binärer Y-Variable



## 2 Das Maximum-Likelihood Schätzverfahren

Das Maximum-Likelihood-Schätzverfahren ist intuitiv recht einfach zu verstehen, es erfordert jedoch rechentechnisch einigen Aufwand. Beginnen wir deshalb mit einer leicht fassbaren Beschreibung der allgemeinen Logik einer ML-Schätzung:

Wenn das ML-Schätzverfahren in der Regressionsanalyse eingesetzt wird, so werden mit dieser Schätzmethode die optimalen Schätzwerte für die  $\alpha$ - und  $\beta$ -Parameter einer Regressionsgleichung gesucht. Dabei gelten als optimale Schätzwerte für die  $\alpha$ - und  $\beta$ -Parameter diejenigen Werte, mit denen eine Regressionsgleichung zu konstruieren ist, die für bestimmte X-Werte die tatsächlich beobachteten Y-Werte einer bestimmten Stichprobe mit höchstmöglicher Wahrscheinlichkeit ermitteln kann.

Solche Parameter-Schätzwerte lassen sich leider nicht immer aus einem Gleichungssystem eindeutig ableiten. Die Schätzung muss deshalb oftmals iterativ erfolgen, d.h. die Schätzwerte werden mittels „trial and error“ in mehreren Schritten ermittelt und in jedem Schritt wird versucht, das Ergebnis des vorangegangenen Schrittes zu verbessern.<sup>2</sup> Erst wenn sich im Laufe eines solchen Prozesses die Wahrscheinlichkeit für eine möglichst exakte Regressionsschätzung der Stichprobenwerte nicht

<sup>2</sup> Nur in einfachen Modellen kann der ML-Schätzer direkt nach einer Formel berechnet werden (vgl. die nachfolgende Fußnote). In komplexeren Modellen, wie z.B. bei der in Abschnitt 3 vorgestellten logistischen Regressionsanalyse, muss das Maximum der Likelihoodfunktion iterativ ermittelt werden (dazu mehr auf den folgenden Seiten).

mehr steigern lässt, wird das iterative Vorgehen abgebrochen und es werden diejenigen Parameter-Schätzwerte, die im letzten Schritt erfolgreich ausprobiert wurden, als optimale Schätzwerte akzeptiert.

Im Laufe einer ML-Schätzung werden also bestimmte Parameter-Schätzwerte ausprobiert und wird danach geschaut, ob mit diesen Schätzwerten eine Regressionsschätzung durchzuführen ist, die mit hoher Wahrscheinlichkeit für bestimmte Werte der X-Prädiktoren die tatsächlichen Y-Werte einer Stichprobe ermitteln kann. Natürlich werden die im ersten Schritt geschätzten Werte noch nicht die optimale Regressionsschätzung ermöglichen, welche die empirischen Werte mit der maximal erreichbaren Wahrscheinlichkeit erbringen kann. Deshalb werden in einem anschließenden Schritt die zuvor ausprobierten Schätzwerte leicht verändert und dann überprüft, ob mit den neuen Schätzwerten die entsprechende Wahrscheinlichkeit zu steigern ist. Und dieses Vorgehen wird so lange wiederholt, bis die Wahrscheinlichkeit für eine gelungene Regressionsschätzung mit guten Schätzergebnissen nicht mehr bedeutsam gesteigert werden kann und deshalb die sogenannte „Konvergenz“ der Schätzung erreicht wurde.

Die ML-Schätzung sucht also nicht nach Parameter-Schätzwerten, die in einer Regressionsanalyse die kleinste quadrierte Residuensumme hervorbringen (wie es die OLS-Schätzung tut), sondern sie wählt im Zuge einer schrittweisen Annäherung diejenigen Koeffizienten als optimale Schätzwerte aus, die, unter der Annahme sie wären identisch mit den wahren Parametern in der Grundgesamtheit, die beobachteten Stichprobenwerte mit der größten Wahrscheinlichkeit hervorbringen können.

Wir wollen die zugrundeliegende Logik eines solchen Maximum-Likelihood-Schätzverfahrens (ML-Schätzung) an einem einfachen (didaktischen) Beispiel verdeutlichen:

Wenn in einer Stichprobe von insgesamt 10 befragten Personen 4 dieser Befragten angegeben hätten, bei einer zukünftigen Bundestagswahl die CDU wählen zu wollen, könnte die Forschungsfrage gestellt werden, welcher Prozentwert „ $\pi$ “ als der wahre Anteilswert aller CDU-Wähler in der Grundgesamtheit, also als Parameter der Population, vermutet werden sollte. Denn da es stets zufällige Stichprobenschwankungen und Messfehler gibt, muss der erfragte Anteilswert nicht auch automatisch mit dem wahren Wert identisch sein. Der wahre Wert müsste also ein Wert sein, der für die gesamte Population, für welche die Stichprobe repräsentativ gezogen wurde, gültig ist.

In diesem Beispiel könnte in einem ersten Schritt vermutet werden, dass ein Anteilswert von 10% der gesuchte, wahre Wert ist. Man könnte dann unter Verwendung dieses Wertes die Wahrscheinlichkeit für das erfragte Stichprobenergebnis aus der theoretischen Binomial-Verteilung (vgl. Kriz 1983: 90-92) ableiten:

$$\binom{n}{s} \pi^s (1-\pi)^{n-s} = \binom{10}{4} \times 0,1^4 \times 0,9^6 = 0,0112 \quad (1)$$

Nach Gleichung (1) ergäbe sich bei einem als wahr vermuteten Wert von 10% CDU-Wählern ( $\pi = 0,1$ ) das erfragte Ergebnis mit einer Wahrscheinlichkeit von 1,12%. Oder anders: nach Gleichung (1) wäre bei 100 Stichproben aus derselben Grundgesamtheit in nur einer Stichprobe das tatsächlich erfragte Ergebnis zu erwarten (wenn der wahre Anteilswert von CDU-Wählern bei 10% läge).

Natürlich würde ein Wahrscheinlichkeitswert von 1,12% nicht sehr überzeugend für einen geschätzten wahren Anteilwert von 10% sprechen. Aber wie oben angesprochen, ließe sich der als wahr vermutete Anteilswert auch schrittweise erhöhen. Auf diese Weise könnte ermittelt werden, welcher vermutete Anteilswert die beobachteten Werte mit der höchsten Wahrscheinlichkeit hervorbringt. Die folgende Tabelle 1 zeigt die nach Gleichung (1) berechneten Wahrscheinlichkeiten für verschiedene Anteilswerte.

Tabelle 1: Nach Gleichung (1) errechnete Wahrscheinlichkeiten für schrittweise erhöhte Schätzwerte des wahren CDU-Anteils

Schätzwert	Wahrscheinlichkeitswert (Likelihood-Wert: $L(\pi)$ )
10 %	1,12 %
20 %	8,81 %
30 %	20,01 %
40 %	25,08 %
50 %	20,51 %
60 %	11,15 %



Wie Tabelle 1 zeigt, hätte sich die Wahrscheinlichkeit, mit der unter 10 befragten Personen 4 CDU-Wähler anzutreffen sind, von 1,12% auf 8,81% erhöht, wenn als wahrer CDU-Anteil nicht 10% sondern 20% geschätzt worden wären. Und natürlich zeigt Tabelle 1 auch, dass der beste Schätzwert für den wahren CDU-Anteil ein Wert von 40% wäre, da sich dann der beobachtete Wähleranteil mit einer Wahrscheinlichkeit von 25,08% ergibt und dieser Wahrscheinlichkeitswert nicht mehr zu überbieten ist.<sup>3</sup>

In der Sprache der ML-Schätzmethode werden die in Tabelle 1 berichteten Wahrscheinlichkeitswerte als Likelihood-Werte (L) bezeichnet.<sup>4</sup> Der Wert „ $L(\pi)$ “ ist der Likelihood-Wert für die Beobachtung bestimmter empirischer Werte unter der Voraussetzung, dass in der Population der Wert eines bestimmten Parameters (hier: der geschätzte Prozentanteil „ $\pi$ “) gilt. Eine Gleichung wie die Gleichung (1) wird dementsprechend als „Likelihood-Funktion“ bezeichnet. Durch sie wird bestimmt, in welcher Weise sich der Wert „L“ in Abhängigkeit von einem bestimmten Parameterschätzwert (hier: „ $\pi$ “) verändert. Der mit der ML-Schätzmethode gesuchte optimale Schätzwert für einen bestimmten Populationsparameter ist der maximale Likelihood-Schätzwert (maximum likelihood estimator = mle). Das ist derjenige Schätzwert, der die Likelihood-Funktion maximiert. Denn mit diesem Schätzwert ist die Wahrscheinlichkeit am größten, dass aus einer Population, in der dieser Wert als wahrer Wert gilt, auch der beobachtete Stichprobenwert gezogen würde. Im oben benutzten Beispiel ist das  $\pi=0,4$  mit einer zu erwartenden Trefferquote von 25%.

Leider erfordert die rechentechnische Umsetzung einer ML-Schätzung in der Regressionsanalyse einigen formal-statistischen Aufwand. Im Folgenden sollen die wesentlichen Argumentationsschritte der ML-Schätzung am Beispiel einer binären logistischen Regressionsanalyse verdeutlicht werden. Weitergehende Erläuterungen zur Logik, Anwendung und Interpretation von binären logistischen Regressionsanalysen werden im daran anschließenden Abschnitt 3 vorgetragen.

In der ML-Schätzung eines binären logistischen Regressionsmodells gehen wir davon aus, dass  $\pi_i$  die wahre Wahrscheinlichkeit bezeichnet, mit der eine bestimmte Person das Ereignis „ $Y_i=1$ “ reali-

---

<sup>3</sup> Dieses Ergebnis ist nicht zufällig identisch mit der beobachteten Prozentzahl in der Stichprobe. Der ML-Schätzwert einer Prozentzahl ist immer identisch mit dem beobachteten Stichprobenwert. Zudem hätte sich dieser Wert auch ohne Iteration direkt aus der  $L(\pi)$ -Funktion ableiten lassen. Dazu wäre allein das Maximum der Funktion zu ermitteln gewesen, denn wenn die erste Ableitung der Funktion gleich null gesetzt und nach  $\pi$  aufgelöst wird, ergibt sich daraus der Schätzwert „ $x/n$ “.

<sup>4</sup> Für Wahrscheinlichkeiten gilt, dass die Summe aller Wahrscheinlichkeiten von Ereignissen, die sich gegenseitig ausschließen, 100% betragen muss. Berechnet man jedoch die Summe der Wahrscheinlichkeitswerte aller Schätzwerte, die in einem iterativen Schätzverfahren möglich sind, so bekommt man schnell Zahlen, die gegen unendlich gehen können. Schon die Summe der nur sechs Wahrscheinlichkeitswerte (bzw. Likelihood-Werte) in Tabelle 1 beträgt 96,61% und diese Liste ließe sich sehr schnell um weitere Schätzwerte mit weiteren Wahrscheinlichkeitswerten erweitern. Deshalb ist es sinnvoll, beim ML-Schätzverfahren von Likelihood-Werten anstatt von Wahrscheinlichkeitswerten zu sprechen.

siert (z.B. die politische Partei „CDU“ wählt). Dann bezeichnet „ $1-\pi_i$ “ die wahre Wahrscheinlichkeit, mit der eine bestimmte Person das Ereignis „ $Y_i=0$ “ realisiert (hier: die CDU nicht wählt). Insgesamt muss es für eine ML-Schätzung stets  $n_1$  Personen geben, die das Ereignis „ $Y_i=1$ “ realisieren und  $n_2$  Personen, die das Ereignis „ $Y_i=0$ “ realisieren. Da angenommen wird, dass jede Person ihren Y-Wert unabhängig von anderen Personen wählt, ergibt sich der Likelihood-Wert für die beobachtete Y-Verteilung im Sample „ $N=n_1+n_2$ “ aus der Multiplikation der Einzel-Wahrscheinlichkeiten in der Likelihood-Funktion:

$$L(\pi) = (\pi_1)(Y_1) \times (\pi_2)(Y_2) \times \dots \times (\pi_{n_1})(Y_{n_1}) \times (1 - \pi_{n_1+1})(Y_{n_1+1}) \times (1 - \pi_{n_1+2})(Y_{n_1+2}) \times \dots \times (\pi_{n_1+n_2})(Y_{n_1+n_2}) \quad (2)$$

Benutzt man in Gleichung (2) das mathematische Symbol „ $\Pi$ “, um das Produkt einer beliebigen Anzahl von Faktoren zu beschreiben, lässt sie sich folgendermaßen verkürzen:

$$L(\pi) = \left( \prod_{i=1}^{n_1} (\pi_i)(Y_i) \right) \times \left( \prod_{i=n_1+1}^{n_1+n_2} (1 - \pi_i)(1 - Y_i) \right) \quad (3)$$

Da es mathematisch einfacher ist, mit Summen als mit Produkten zu arbeiten, wird Gleichung (3) im nächsten Schritt logarithmiert, so dass die sogenannte „Log- Likelihood Funktion“ oder „ $LL(\pi)$ “ entsteht:

$$LL(\pi) = \left( \sum_{i=1}^{n_1} \ln(\pi_i)(Y_i) \right) + \left( \sum_{i=n_1+1}^{n_1+n_2} \ln(1 - \pi_i)(1 - Y_i) \right) \quad (4)$$

Natürlich ist in Gleichung (4) die Wahrscheinlichkeit „ $\pi$ “ nach wie vor unbekannt. Sie wird entsprechend der logistischen Grundgleichung (6.1), die in Abschnitt 3 noch ausführlich erläutert wird, berechnet:

$$\pi_i = P_i = \frac{e^{(\alpha + \sum \beta_k X_{ki})}}{1 + e^{(\alpha + \sum \beta_k X_{ki})}} \quad (5)$$

Nunmehr können in Gleichung (4) (in welche die Gleichung 5 eingesetzt wird), diejenigen Koeffizienten für die Parameter „ $\alpha$ “ und „ $\beta_k$ “ iterativ gesucht werden, die den Log-Likelihood-Wert

„ $LL(\alpha, \beta)$ “ maximieren.<sup>5</sup> Sind sie gefunden, erhält man damit auch gleichzeitig die ML-Schätzwerte „ $a$ “ und „ $b_k$ “ für das logistische Regressionsmodell (vgl. dazu Amemiya 1981; Dhrymes 1978: 335).

In der Praxis benutzen die meisten Iterationsverfahren den negativen Log-Likelihoodwert „ $-LL$ “ als Annäherungskriterium (SPSS benutzt sogar den zweifachen negativen LL-Wert: „ $-2 \times LL$ “), so dass das Maximum der Schätzung dort erreicht wird, wo der absolute Wert von  $-LL$  bzw. von  $-2 \times LL$  am geringsten ist. Die folgende Tabelle 2 zeigt das Protokoll einer Iteration in fünf Schritten. Es ist in der Tabelle leicht zu erkennen, dass zwischen dem vierten und fünften Iterationsschritt keine Veränderung des Likelihood-Wertes mehr eintritt und deshalb der Koeffizientenschätzwert von 0,492 als ML-Schätzer akzeptiert werden kann.

Tabelle 2: Beispiel eines Iterationsprotokolls

Iteration	-2 Log-Likelihood	Regressionskoeffizient b
Schritt 1	2648,125	0,343
Schritt 2	2607,061	0,468
Schritt 3	2606,116	0,491
Schritt 4	2606,115	0,492
Schritt 5	2606,115	0,492

Die ML-Schätzmethode erbringt Schätzwerte, die asymptotische<sup>6</sup> Eigenschaften aufweisen (vgl. auch Urban/Mayerl 2008: Kap. 3.1). Insbesondere sind die ML-Schätzer:<sup>7</sup>

- asymptotisch konsistent, d.h. umfangreichere Stichproben können Verzerrungen und unzulässig große Streuungen der Schätzwerte verringern (je größer die Stichprobe wird, umso kleiner wird die Wahrscheinlichkeit, dass geschätzte und wahre Parameter voneinander abweichen);
- asymptotisch effizient, d.h. dass die Varianz von vielfach wiederholten Schätzungen mit anderen Schätzverfahren nicht zu unterbieten ist;
- asymptotisch normalverteilt, d.h. die ML-Schätzwerte können in Signifikanz-Tests überprüft werden.

<sup>5</sup> Dazu muss die erste Ableitung der LL-Funktion für jeden Parameter berechnet und diese gleich null gesetzt werden. Da die so erhaltene Gleichung aber in  $\beta$  nicht linear ist, kann es für die Auflösung nach  $\beta$  auch keine analytische Lösung geben. Deshalb ist eine iterative Lösung nötig, in der schrittweise ein entsprechendes Maximum gesucht wird.

<sup>6</sup> Im Unterschied zu exakten Standards beziehen sich asymptotische Standards auf Verteilungsmerkmale, die erst bei einem gegen Unendlich konvergierenden Stichprobenumfang ihre Gültigkeit erlangen. Das bedeutet für die statistische Praxis, dass das Vorhandensein asymptotischer Merkmale nicht überprüft werden kann, sondern von der Hoffnung legitimiert werden muss, dass der gegebene Stichprobenumfang groß genug ist.

<sup>7</sup> Vgl. dazu Dhrymes 1978: 336ff.

### 3 Binär-logistische Regressionsanalyse

Eine Regressionsanalyse, die das vorgestellte ML-Schätzverfahren nutzt, ist die logistische Regression. Im Folgenden wird eine spezielle Form der logistischen Regression vorgestellt, die besonders häufig in der Forschungspraxis verwendet wird: die binär-logistische Regressionsanalyse.

Die binär-logistische Regression sollte immer dann eingesetzt werden, wenn die abhängige Variable eines Regressionsmodells nur zwei Ausprägungen hat, wenn also die Y-Variable binär bzw. binomial kodiert ist. Dies wäre z.B. dann der Fall, wenn in einem Regressionsmodell untersucht werden soll, ob die Absicht, eine bestimmte politische Partei zu wählen (z.B. die „CDU“), von der Links-Rechts-Selbsteinstufung oder der Mitgliedschaft in einer Gewerkschaft abhängt. Denn die Wahlabsicht zugunsten einer politischen Partei (Y) kann vorhanden sein (Y=1) oder kann nicht vorhanden sein (Y=0). Die Y-Variable hat also nur zwei mögliche Ausprägungen und ist damit eine binär kodierte abhängige Variable.

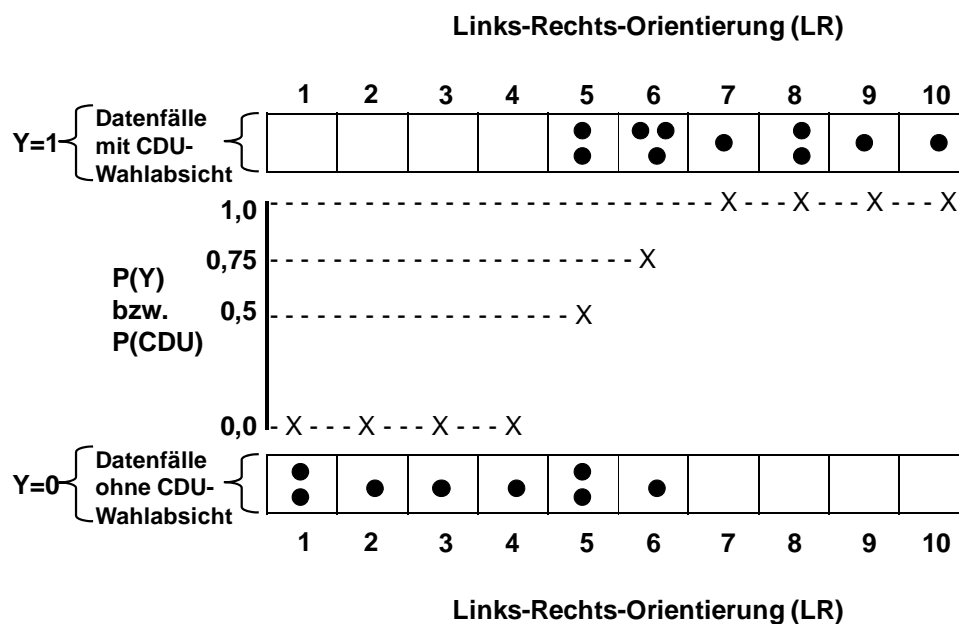
Binär kodierte abhängige Variablen erzeugen für die klassische OLS-Regression eine Fülle von Problemen. Einige davon wurden ganz zu Beginn angesprochen. Es wurde dort z.B. gezeigt, dass in einer Regression mit binärer Y-Variablen die Residuen nicht normalverteilt sein können und sie prinzipiell auch keine Homoskedastizität aufweisen (vgl. Urban/Mayerl 2008: Kapitel 3). Um diesen Problemen von vornherein aus dem Wege zu gehen, sollten Regressionsanalysen mit binomial skalierten abhängigen Variablen stets als binär-logistische Regressionsanalysen durchgeführt werden.

Was ist nun der Unterschied zwischen einer klassischen OLS-Regression und einer ML-basierten logistischen Regression? Im klassischen OLS-Regressionsmodell wird der zu erwartende Wert von Y durch eine Linearkombination von Prädiktoren geschätzt, so wie es die folgende Gleichung für zwei X-Variablen darstellt:  $\hat{Y}_i = a + b_1X_{1i} + b_2X_{2i}$

In der binär-logistischen Regression wird nicht der Vorhersagewert bzw. der Erwartungswert für Y geschätzt. Stattdessen wird der Vorhersagewert für die bedingte Wahrscheinlichkeit von Y berechnet. Die folgende Abbildung 2 verdeutlicht, wie die binär kodierte Variable Y(1/0) als bedingte Wahrscheinlichkeit  $P(Y|X)$  zu verstehen ist.

Abbildung 2 zeigt 18 Befragte als schwarze Punkte. Zehn Befragte geben als Wahlabsicht eine CDU-Wahl an ( $Y=1$ ) und befinden sich deshalb in der oberen Kästchenreihe, während die Befragten, die eine andere Partei zu wählen beabsichtigen ( $Y=0$ ), in der unteren Kästchenreihe zu finden sind. Jedes Kästchen markiert einen Wert auf der Skala der politischen Links-Rechts-Grundorientierung (LR) der befragten Wahlberechtigten. Die LR-Skalenwerte liegen zwischen 1 („extrem links“) und 10 („extrem rechts“). Jedem Befragten kann entsprechend seiner Lage in den LR-Kästchen ein  $P(Y)$ -Wert zugeordnet werden. So gibt es z.B. nur zwei Befragte, die einen LR-Wert von 1 aufweisen. Beide wollen nicht die CDU wählen ( $Y=0$ ). Daher ergibt sich für Personen mit  $LR=1$  ein Wahrscheinlichkeitswert  $P(Y=1)$  von 0,00 Prozent. Im Unterschied dazu zeigt die Abbildung bei einem LR-Wert von 5 insgesamt vier Befragte, von denen die Hälfte (zwei Personen) eine CDU-Wahl beabsichtigen ( $Y=1$ ). Somit beträgt  $P(Y=1|5)$  fünfzig Prozent oder 0,50.

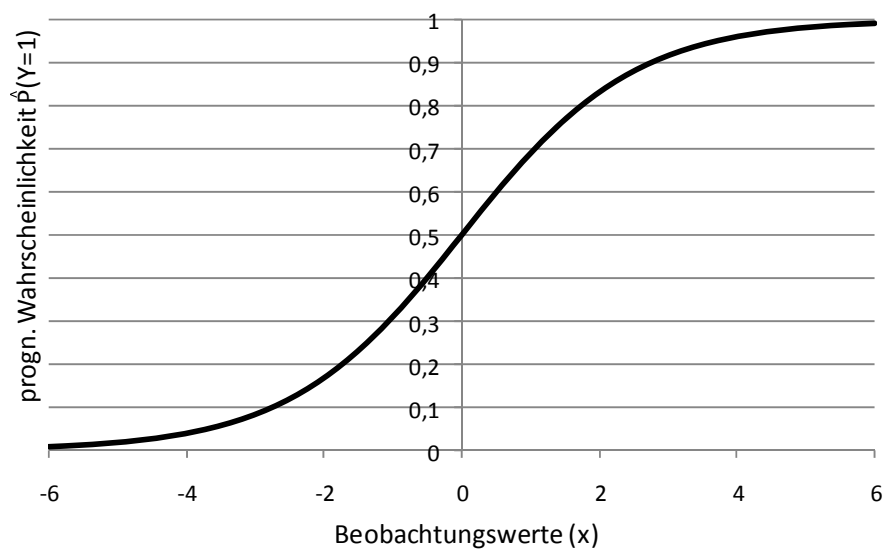
Abbildung 2: Transformation von  $Y(1/0)$  in  $P(Y)$  in Abhängigkeit von LR (fiktives Beispiel)



Nachdem, wie oben gezeigt, aus  $Y(1/0)$  die neue Variable  $P(Y|X)$  wurde, kann die Beziehung der abhängigen Variablen  $P(Y)$  und der unabhängigen Variablen  $X$  (im Beispiel: LR) regressionsanalytisch untersucht werden. Was dazu allerdings noch fehlt, ist die Festlegung einer mathematischen Funktion, welche die beiden Variablen miteinander verknüpfen soll. In der OLS-Regression ist das die lineare Funktion „ $a + bX$ “. Eine solche Funktionsbestimmung ist hier nicht möglich. Ein Grund (u.a.) dafür ist, dass die vorhergesagten  $P(Y)$ -Werte entsprechend der Prozentskala nicht größer als 1,0 und nicht kleiner als 0,0 werden dürfen (mehr dazu auch in Abschnitt 1 sowie in Urban 1993: 15-23). Dies kann jedoch durch Verwendung einer logistisch verlaufenden Funktionskurve sichergestellt werden. Folglich erhält die „logistische Regressions-

analyse“ ihren Namen dadurch, dass in dieser Regressionsanalyse die abhängige Variable mit der/den unabhängigen Variablen auf logistische Weise verknüpft wird (weshalb man in diesem Falle auch von einer logistischen „Link-Funktion“ spricht). Die folgende Abbildung 3 zeigt die allgemeine Form einer solchen logistischen Funktionskurve.

Abbildung 3: Allgemeine logistische Funktionskurve



Wie in Abbildung 3 gut zu erkennen ist, weist der kurvilineare, logistische Einflussverlauf im Unterschied zur linearen Einflussbeziehung keine konstanten Veränderungsrate bei der abhängigen Variablen auf. Steigt  $X$  um einen konstanten Betrag an, so sind die Steigungen in den  $P$ -Werten unterschiedlich groß, denn das Ausmaß der Steigung hängt von dem Startwert ab, von dem aus der  $X$ -Wert anwächst. Generell gilt für die logistische Einflussbeziehung: Veränderungen in den extremen Wahrscheinlichkeitswerten (nahe 0% und nahe 100%) sind sehr viel schwerer zu erreichen als Veränderungen im mittleren Wahrscheinlichkeitsbereich. Im mittleren Bereich impliziert nur eine kleine Veränderung in der/den unabhängigen Variablen weitreichende Veränderungen in den prognostizierten Wahrscheinlichkeitswerten. Im Unterscheid dazu bleiben gleichgroße Verschiebungen immer dort relativ konsequenzlos, wo sie von extremen  $X$ -Startwerten aus erfolgen.

Eine solche logistische Bestimmung des Verlaufs von Einflussbeziehungen ist sicherlich in vielen Forschungsbereichen wesentlich realistischer als eine lineare Bestimmung, die von der absoluten Konstanz der Veränderungsrate ausgeht und für die es unerheblich ist, von welchen Startwerten aus die Veränderungen erfolgen.

Generell wird eine logistische Link-Funktion durch die folgende Gleichung (6) beschrieben. Darin ist  $V$  diejenige Variable, welche die genaue Lage der Funktionskurve bestimmt. Sie legt die jeweilige Steigung und den jeweiligen Wendepunkt der Kurve fest.

$$\hat{P}_i = \frac{e^{(V_i)}}{1 + e^{(V_i)}} \quad (6)$$

In der logistischen Regressionsanalyse wird  $V$  (vgl. Gleichung 6) als Linearkombination der unabhängigen  $X$ -Variablen des zu untersuchenden Regressionsmodells bestimmt:

$$V_i = a + \sum b_k X_{ki} \quad (7)$$

so dass Gleichung (6) auch folgendermaßen zu schreiben ist:

$$\hat{P}_i = \frac{e^{(a + \sum b_k X_{ki})}}{1 + e^{(a + \sum b_k X_{ki})}} \quad (6.1)$$

In der logistischen Regressionsschätzung geht es nun darum, für die  $X$ -Variablen eines bestimmten Regressionsmodells mittels des oben erläuterten Maximum-Likelihood Schätzverfahrens diejenigen  $a$ - und  $b$ -Koeffizienten zu ermitteln, die für alle  $X$ -Werte und für alle Kombinationen von  $X$ -Werten solche  $P(Y)$ -Prognosewerte ergeben (genannt:  $(\hat{P}(Y))$ ), die möglichst gut mit den empirischen  $P(Y)$ -Werten übereinstimmen. Diese  $a$ - und  $b$ -Koeffizienten sollten die beste Schätzung der beobachteten  $P(Y)$  ermöglichen, d.h. sie sollten die Schätzung mit dem maximalen Likelihood-Wert liefern.

Generell betrachtet wird der Verlauf einer logistischen Funktion, wie sie von Gleichung (6.1) beschrieben wird, in folgender Weise von den Regressionskoeffizienten beeinflusst (vgl. dazu die Abbildungen in Urban 1993: 31):

- (-) Der konstante  $a$ -Koeffizient (in der Linearkombination von  $V$ ) verschiebt die logistische Kurve in der Horizontalen, ohne ihre Steigung zu verändern.
- (-) Höhere Werte der  $b$ -Koeffizienten (in der Linearkombination von  $V$ ) vergrößern die Veränderungsrate von  $\hat{P}(Y)$ , d.h. mit ihrem Anwachsen wird der Funktionsverlauf steiler.
- (-) Ein negatives Vorzeichen der  $b$ -Koeffizienten ändert den Ursprung des logistischen Funktionsverlaufs, der dann links oben bei der höchsten Wahrscheinlichkeit von  $\hat{P}(Y)$  beginnt und mit Anwachsen von  $V$  in Richtung  $\hat{P}(Y) = 0,00$  nach rechts unten verläuft.

Wir wollen dies anhand eines Beispiels veranschaulichen. Dabei soll, wie auch schon zuvor erläutert, in Form einer Regressionsanalyse untersucht werden, ob die Absicht, die politische Partei „CDU“ zu wählen, von der Links-Rechts-Selbsteinstufung (LR) und der Mitgliedschaft in einer Gewerkschaft (GEW) abhängt. Aufgrund einer ML-Schätzung ergeben sich für die Linearkombination der beiden unabhängigen Variablen LR und GEW folgende Werte:

$$V_i = -3,48 + 0,52LR_i - 0,75GEW_i \quad (8)$$

Die in Gleichung (8) ausgewiesenen Schätzwerte können in Gleichung (6.1) eingesetzt werden, so dass folgende logistische Regressionsschätzung für die erwarteten Wahrscheinlichkeiten  $\hat{P}(Y)$  entsteht:

$$\hat{P}_i = \frac{e^{(-3,48 + 0,52 \times LR_i - 0,75 \times GEW_i)}}{1 + e^{(-3,48 + 0,52 \times LR_i - 0,75 \times GEW_i)}} \quad (9)$$

Mit Gleichung (9) haben wir die Möglichkeit geschaffen, die Abhängigkeit der  $\hat{P}$ -Variablen von den Einflüssen der X-Prädiktoren zu analysieren. Und diese Abhängigkeit ist nicht mehr linear sondern nunmehr logistisch zu verstehen. Die Gleichung (9) beschreibt nämlich einen logistischen, d.h. S-förmigen Verlauf der Effekte aller Prädiktoren bzw. unabhängigen X-Variablen.

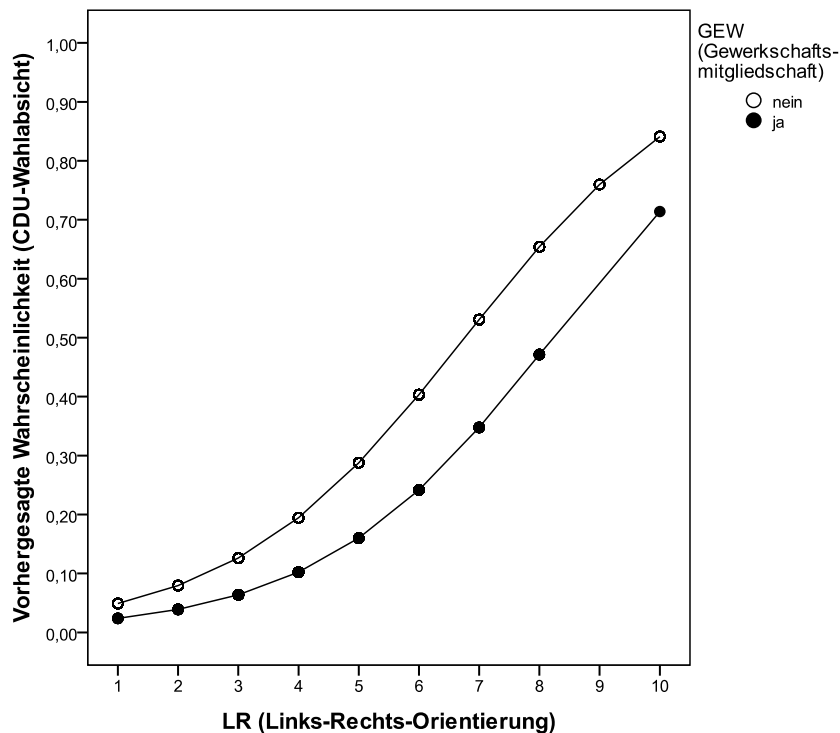
Die folgende Abbildung 4 zeigt den nach Gleichung (9) berechneten Verlauf zweier logistischer Funktionskurven für unser Beispiel zur Untersuchung der CDU-Wahlabsicht. Dabei wird der Verlauf der beiden Kurven nur für den Wertebereich LR=1 („extrem links“) bis LR=10 („extrem rechts“) abgedruckt.

Die Kurve in Abbildung 4 mit den hellen Datenpunkten betrifft den Funktionsverlauf ohne Gewerkschaftsmitgliedschaft (GEW=0), während die Linie mit den schwarzen Datenpunkten den Funktionsverlauf mit Gewerkschaftsmitgliedschaft (GEW=1) markiert. Deutlich ist die reduzierte Wahrscheinlichkeit einer CDU-Wahlabsicht bei vorhandener Gewerkschaftsmitgliedschaft zu erkennen. Insbesondere im Wertebereich von LR=5 bis LR=10 ist der Abstand der beiden Funktionskurven beträchtlich. Deutlich ist auch zu erkennen, dass die Steigung beider Funktionskurven erst ab einem Wert von LR=4 erheblich zunimmt und dann recht konstant bleibt. Wenn Personen eine Links-Rechts-Selbsteinstufung von LR=1 aufweisen (d.h. eine extrem linke politische Grundorientierung besitzen), haben geringfügige Einstellungsverschiebungen nach „rechts“ nur wenig Einfluss auf die



hier untersuchte Wahlabsicht. Die prognostizierte Wahrscheinlichkeit einer CDU-Wahlabsicht wächst dann zwar an, aber zunächst doch nur in recht unbedeutendem Maße. Deutlicher ist das Wachstum der Wahrscheinlichkeit bei Verschiebungen im mittleren Bereich beider Funktionskurven ausgeprägt. Dort ist jede Einstellungsverschiebung von einer LR-Einheit in Richtung „rechts“ mit einem ausgeprägten Anstieg der Wahrscheinlichkeit für eine CDU-Wahlabsicht verbunden.

Abbildung 4: Zwei empirisch geschätzte logistische Funktionskurven für GEW=0 und GEW=1



Die hier gezeigte Beispielgraphik (Abbildung 4) entstammt dem SPSS-Beispiel einer logistischen Regression, das im letzten Abschnitt noch ausführlich vorgestellt wird. Dort wird auch gezeigt, wie die vorhergesagte Wahrscheinlichkeit ( $\hat{P}_i$ ) als neue Variable erstellt werden kann. In SPSS kann dann die Graphik aus Abbildung 4 leicht mit dem folgenden Befehl erzeugt werden:

```
GRAPH /SCATTERPLOT(BIVAR)=LR WITH PRE_1 BY GEW.
```

Zusätzlich muss zur Einblendung der Verbindungslinie zwischen den Datenpunkten im Diagramm-Editor unter dem Menü-Punkt „Elemente“ die Option „Interpolationslinie“ ausgewählt werden.

Generell ist für jede logistische Regressionsanalyse zu empfehlen, die geschätzten Wahrscheinlichkeiten zumindest für bestimmte, für die Interpretation der Ergebnisse besonders aussagekräftige Kombinationen von X-Werten zu ermitteln und zu berichten. Für Dummy-Variablen lassen sich dabei schnell alle Kombinationen durchspielen. Bei metrischen X-Variablen sollte dies zumindest für Extremwerte und die Mittelkategorie gemacht werden.

In Gleichung (9) ist die Bedeutung der geschätzten Regressionskoeffizienten für das Ausmaß von  $\hat{P}$  nur schwer zu erkennen, weil die X-Variablen in einer logistischen Verbindung zur abhängigen

$\hat{P}$ -Variablen stehen. Jedoch kann die Interpretation der logistischen Regressionsschätzung dadurch erleichtert werden, dass Gleichung (9) dermaßen transformiert wird, dass sie auf der rechten Seite nur noch die Linearschätzung der Regressionskoeffizienten enthält. Auf diese Weise entsteht die folgende Gleichung (10):

$$\ln\left(\frac{\hat{P}_i}{1-\hat{P}_i}\right) = a + b_1 \times (\text{LR}_i) + b_2 \times (\text{GEW}_i) \quad (10)$$

In verallgemeinerter Form lässt sich Gleichung (10) auch als Gleichung (11) schreiben:

$$\ln\left(\frac{\hat{P}_i}{1-\hat{P}_i}\right) = a + \sum b_k \times X_{ki} \quad (11)$$

Auf der linken Seite der Gleichungen (10) und (11) befinden sich nun nicht mehr erwartete Wahrscheinlichkeiten von Y sondern sogenannte „Logits“ von Y. Diese haben dem Verfahren der logistischen Regressionsanalyse auch den Namen „Logit-Analyse“ gegeben. Die Logits sind das Ergebnis einer zweifachen Transformation von  $\hat{P}$ . Zum einen wird die Wahrscheinlichkeit von Y=1 in Verhältnis zu ihrer Komplementär-Wahrscheinlichkeit gesetzt, und zum anderen wird die daraus entstehende Verhältniszahl logarithmiert.<sup>8</sup>

Empirisch interpretieren kann man die Logits kaum noch. Man kann die Gleichungen (10) und (11) allerdings benutzen, um eine lineare Interpretation der geschätzten logistischen Regressionskoeffizienten vorzunehmen:

In Gleichung (11) benennt der zu schätzende Regressionskoeffizient „b“ die Einflussstärke und die Einflussrichtung für jede X-Variable auf die Logit-Variable, wobei dieser Einfluss im multivariaten Modell als kontrollierter bzw. partieller Effekt zu verstehen ist. Im formalen Sinne wird der partielle Regressionskoeffizient der logistischen Regression in gleicher Weise wie der partielle Regressionskoeffizient der klassischen OLS-Regression interpretiert (vgl. Urban/Mayerl 2008: Kapitel 2.3.1). In unserem Beispiel ergeben sich z.B. für die Regressionskoeffizienten in Gleichung (10) die folgenden Schätzwerte:

---

<sup>8</sup> Zur Erinnerung an die Schulmathematik: Der natürliche Logarithmus einer beliebigen Zahl „x“ ist gleich dem Exponenten „n“, mit dem die konstante Basiszahl „e“ (=2,718) zu potenzieren ist, um die gewählte Zahl „x“ wieder zurückzubekommen (also:  $\ln(x)=n$  und  $e^n=x$ ). Verständlicher wird das im Beispiel: Man nehme eine beliebige Zahl z.B. die Zahl „100“. Ihr natürlicher Logarithmus ist 4,605 oder:  $\ln 100 = 4,605$ , da Folgendes gilt:  $2,718^{4,605} = 100$  oder:  $e^n = x$ .

$$\ln\left(\frac{\hat{P}_i}{1-\hat{P}_i}\right) = -3,48 + 0,52 \times (\text{LR}_i) - 0,75 \times (\text{GEW}_i) \quad (12)$$

Die geschätzten Koeffizienten haben dann die folgende Bedeutung: 1.) Es gibt einen positiven partiellen Effekt von LR auf  $\text{Logit}(\hat{P})$ . Für jede Verschiebung auf der LR-Skala (mit Werten von 1 bis 10) um eine empirische Einheit nach „rechts“ (bzw. in Richtung des Wertes „10“) steigt der Logit ( $\hat{P}$ )-Wert um 0,52 Einheiten. 2.) Es gibt einen negativen partiellen Effekt von GEW auf  $\text{Logit}(\hat{P})$ . Eine Gewerkschaftsmitgliedschaft (dichotom gemessen:  $X_2 = 0 / 1$ ) lässt  $\text{Logit}(\hat{P})$  um einen Logit ( $\hat{P}$ )-Wert von 0,75 Einheiten absinken.

Ob der LR-Effekt oder der GEW-Effekt einen stärkeren Einfluss auf  $\text{Logit}(\hat{P})$  hat, lässt sich anhand der Werte von +0,52 und -0,75 nicht entscheiden. Denn bei den partiellen Regressionskoeffizienten der logistischen Regression handelt es sich um unstandardisierte Koeffizienten, die nur dann zum Vergleich der Effekte innerhalb eines Modells herangezogen werden dürfen, wenn die entsprechenden Prädiktoren in gleicher Weise skaliert sind.

Zum Vergleich der Einflussstärken unterschiedlich skalierte Prädiktoren innerhalb eines Modells wird in klassischen OLS-Regressionen häufig die Standardisierung der Regressionskoeffizienten eingesetzt (man beachte jedoch auch die damit verbundenen Probleme, vgl. Urban/Mayerl 2008: Kapitel 2.3.3). Die Standardisierung der Regressionskoeffizienten in der logistischen Regression ist mit den gleichen Problemen verbunden, weist aber zusätzlich auch noch die Problematik auf, dass die Standardabweichung der dem Modell zugrunde liegenden Logits nicht zu berechnen ist, da die Logitwerte empirisch nicht gemessen wurden (anders als die Y-Werte in der OLS-Regression). Deshalb wird in der Forschungspraxis der logistischen Regressionsanalyse häufig eine Teilstandardisierung der Regressionskoeffizienten durchgeführt. Dabei wird eine Multiplikation der Regressionskoeffizienten nur mit der Standardabweichung der jeweiligen X-Variablen durchgeführt:

$$b_k^* = b_k S_{X_k} \quad (13)$$

Im vorliegenden Beispiel betragen die Standardabweichungen  $S(\text{LR})=1,778$  und  $S(\text{GEW})=0,325$ . Durch Multiplikation der Koeffizienten aus Gleichung (12) mit diesen Standardabweichungen ergeben sich dann die teilstandardisierten partiellen Regressionskoeffizienten:  $b^*(\text{LR})=0,93$  und  $b^*(\text{GEW})=-0,24$ . Demnach hätten Veränderungen der LR-Variablen als Prädiktor für den Logitwert

von  $\hat{P}(\text{CDU})$  in etwa die vierfache (absolute) Stärke von Veränderungen der GEW-Variablen, wenn als Veränderungsmaß die jeweiligen Standardabweichungen benutzt werden. Oder anders gesagt: Eine Zunahme der LR-Variablen um eine Standardabweichung würde Logit(CDU) um 0,93 Einheiten anwachsen lassen, während die Zunahme von GEW um eine Standardabweichung den Logit(CDU)-Wert um 0,24 Einheiten absinken ließe. Unklar bleibt bei dieser Teilstandardisierung jedoch (wie auch bei der Standardisierung in der OLS-Regression), wie man sich Veränderungen um eine Standardabweichung bei einer dichotomen Variablen (wie z.B. bei GEW) mit den empirischen Werten 0 und 1 in inhaltlicher Weise vorstellen soll. Die Teilstandardisierung scheint sich deshalb eher zum Vergleich metrisch skaliertes Effekte in der logistischen Regression anzubieten.

### 3.1 Gewinnchancen (odds) und Effektkoeffizienten (odds ratios)

Da die Logitwerte einer zu erklärenden Y-Variablen in ihrer empirischen Bedeutung äußerst schwierig zu interpretieren sind, bleibt auch die empirische Bedeutung der geschätzten Regressionskoeffizienten eher im Unklaren. So geben diese zwar die Richtung eines entsprechenden Variablen-Einflusses exakt wieder, und auch die Signifikanz der geschätzten Koeffizienten kann wie üblich interpretiert werden (vgl. dazu auch noch Abschnitt 3.4). Jedoch stehen die Logitkoeffizienten in einer nur schwer zu begreifenden, logistischen Beziehung zum Ausmaß der durch X ausgelösten Veränderungen in  $\hat{P}(Y = 1)$ .

Leichter ist die Interpretation logistischer Regressionskoeffizienten, wenn diese auf Veränderungen der abhängigen Variablen  $\hat{P}(Y)$  in Form der so genannten „Gewinnchance“ bzw. der „odds“ bezogen werden. Als Gewinnchance (engl.: odds) wird der Quotient aus erwarteter Wahrscheinlichkeit für das Ereignis  $Y=1$  (z.B. für das Ereignis „CDU-Wahlabsicht“) und der erwarteten Wahrscheinlichkeit für das entsprechende Komplementär-Ereignis  $Y=0$  (hier: „keine CDU-Wahlabsicht“) verstanden (dabei gilt:  $\hat{P}(Y = 0) = 1 - \hat{P}(Y = 1)$ ):

$$\frac{\hat{P}_i(Y = 1)}{1 - \hat{P}_i(Y = 1)}$$

Um in der logistischen Regressionsanalyse vereinfachte Interpretationen mit der Gewinnchance als abhängiger Variablen durchführen zu können, muss Gleichung (11) in folgende Gleichung (14) umgeformt werden:

$$\frac{\hat{P}_i}{1 - \hat{P}_i} = e^{(a + \sum b_k X_{ki})} = e^a \times e^{(\sum b_k X_{ki})}$$

$$= e^a \times e^{(b_1 X_{1i})} \times e^{(b_2 X_{2i})} \times e^{(b_3 X_{3i})} \times \dots \times e^{(b_k X_{ki})} \quad (14)$$

In Gleichung (14) dient die Gewinnchance für  $\hat{P}(\text{CDU})$  als abhängige, zu erklärende Variable. Denn auf der linken Gleichungsseite steht das Wahrscheinlichkeitsverhältnis zwischen  $\hat{P}(\text{CDU})$  und  $\hat{P}(\text{nicht CDU})$ . Diese Verhältniszahl hat eine wesentlich eingängigere empirische Bedeutung als die Logitwerte einer CDU-Wahlabsicht: Der Wertebereich der Gewinnchance liegt zwischen 0,00 und  $+\infty$ . Die Gewinnchance hat einen Wert von 1,00, wenn die Wahrscheinlichkeiten für beide Handlungsalternativen gleich groß sind. Liegt der Wert der Gewinnchance über 1,00, ist die Wahrscheinlichkeit für eine CDU-Wahlabsicht größer als die Absicht, irgendeine andere Partei zu wählen. Liegt er unter 1,00, so ist die Wahrscheinlichkeit für eine CDU-Wahlabsicht schlechter als für eine andere Partei.

Die Beziehung zwischen Ereignis-Wahrscheinlichkeit, Ereignis-Gewinnchance und Ereignis-Logit wird in der folgenden Tabelle 3 anhand von drei verschiedenen Beobachtungsfällen verdeutlicht.

Tabelle 3: Beziehung zwischen Ereignis-Wahrscheinlichkeit, Ereignis-Gewinnchance und Ereignis-Logit

Fall-Nr.	$\hat{P}(\text{CDU})$	$\hat{P}(\text{nicht CDU})$	Chance(CDU)	Logit(CDU)
1	0,50	0,50	0,50:0,50 = 1,00	$\ln(1,00) = 0,00$
2	0,80	0,20	0,80:0,20 = 4,00	$\ln(4,00) = 1,39$
3	0,20	0,80	0,20:0,80 = 0,25	$\ln(0,25) = -1,39$

Wie Tabelle 3 für Fall-Nr. 1 ausweist, ist bei gleicher Wahrscheinlichkeit eines Ereignisses (CDU-Wahlabsicht) und des komplementären Nicht-Ereignisses (keine CDU-Wahlabsicht) die Chance einer CDU-Wahlabsicht gleich 1,00. Ein Chancenwert von 1,00 bedeutet mithin, dass beide Ereignisse gleich wahrscheinlich sind und erzeugt einen Logitwert von 0,00. Im zweiten Fall ist die Wahrscheinlichkeit für eine CDU-Wahlabsicht viermal so groß wie die Absicht, eine andere Partei zu wählen (0,80:0,20) und deshalb hat ihre Gewinnchance einen Wert von 4,00.

Ist die Wahrscheinlichkeit für die CDU-Wahlabsicht jedoch viermal kleiner (0,20:0,80), wie bei Fall-Nr. 3, so hat sie nur eine Gewinnchance von 0,25. Das ist verwirrend, denn obwohl die Wahrscheinlichkeit einmal viermal höher ist (0,80:0,20) und ein anderes Mal viermal geringer ist (0,20:0,80), entstehen unterschiedliche Chancenwerte, die intuitiv nicht vergleichbar sind (4,00 und

0,25). Diese Zahlen entstehen dadurch, dass die Gewinnchance zwar nach oben hin unendlich groß werden kann, jedoch nach unten hin einen festen Grenzwert von 0,00 aufweist (wir werden dies im Folgenden noch ausführlich erläutern und auch Interpretationshilfen mittels Kehrwertbildung vorstellen). Erst die Logitwerte, für die auch keine Untergrenze mehr gilt, können in beiden Fällen gleiche absolute Werte aufweisen (|1,39|).

Welche Bedeutung hat nun die Veränderung einer jeden X-Variablen für die Gewinnchance von  $Y=1$ ? Hierzu betrachtet man das *Verhältnis der Gewinnchancen* für einen beliebigen X-Wert im Vergleich zur Gewinnchance bei einem Anstieg von X um eine Einheit. Das Verhältnis der Gewinnchancen wird auch *odds ratio* genannt. Für ein bivariates Modell bedeutet dies:

$$\frac{\frac{\hat{P}_1}{1 - \hat{P}_1} \quad \text{(für den um eine Einheit erhöhten } X_{1i} \text{ - Wert "m + 1")}}{\frac{\hat{P}_1}{1 - \hat{P}_1} \quad \text{(für einen } X_{1i} \text{ - Wert "m")}} = \frac{e^a \times e^{(b_1[X_{1i}+1])}}{e^a \times e^{(b_1X_{1i})}} = \frac{e^{(b_1[X_{1i}+1])}}{e^{(b_1X_{1i})}} = \frac{e^{(b_1X_{1i})} \times e^{b_1}}{e^{(b_1X_{1i})}} = e^{b_1} \quad (15)$$

Nach Gleichung (15) kann man durch Entlogarithmierung des logistischen Regressionskoeffizienten (also durch Berechnung von  $e^b$  bzw. von  $2,718^b$ ) ein neues Maß zur Beschreibung der Einflussstärke der verschiedenen Prädiktoren in einem logistischen Regressionsmodell erhalten. Dieses Maß gibt in Form eines Multiplikationsfaktors die Veränderungen im Wahrscheinlichkeitsverhältnis der beiden Handlungsalternativen („CDU“ versus „nicht CDU“) an, wenn sich ein entsprechender X-Prädiktor um eine empirische Einheit vergrößert. Die Größe  $e^b$  (bzw. äquivalent auch „exp(b)“ geschrieben) ist der Multiplikationsfaktor für die Berechnung des neuen Wahrscheinlichkeitsverhältnisses, das durch Veränderung der dazugehörigen X-Variablen um eine empirische Einheit entsteht.

Was dies im Konkreten bedeutet, lässt sich an unserem CDU-Wahlbeispiel leicht verdeutlichen. Dort entsteht z.B. durch Entlogarithmierung von  $b_1=0,52$  ein Wert von  $\exp(b_1)=1,68$ . Folglich bedeutet ein Wert von „exp( $b_1$ ) = 1,68“, dass sich bei einer Verschiebung auf der Links-Rechts-Skala („LR“) um +1,00 Einheiten das Wahrscheinlichkeitsverhältnis zwischen einer CDU-Wahlabsicht und der Absicht, eine andere Partei zu wählen, auf das 1,68-fache der Chance, die vor der Verschiebung galt, erhöht.

Dieses Beispiel lässt sich auch anhand von Gleichung (15) verdeutlichen. Es werden dann die geschätzten logistischen Regressionskoeffizienten in die Gleichung (15) eingesetzt (z.B.  $b_{LR} = 0,52$ ), und es wird dann ein Anstieg der X-Variablen „LR“ um eine Einheit angenommen (hier: beispielhaft ein Anstieg von LR=4 auf LR=5; dasselbe gilt für beliebige LR-Skalenwerte und deren Verschiebung um eine Einheit). Auch hierbei zeigt sich, dass das Chancenverhältnis um das 1,68-fache ansteigt, wenn LR von 4 auf 5 anwächst und die Gewerkschaftsmitgliedschaft („GEW“) konstant bleibt:

$$\frac{\frac{\hat{P}_1}{1-\hat{P}_1} \text{ (für LR = 5)}}{\frac{\hat{P}_1}{1-\hat{P}_1} \text{ (für LR = 4)}} = \frac{e^{-3,48} \times e^{(0,52 \times [LR=5])} \times e^{(-0,75 \times GEW)}}{e^{-3,48} \times e^{(0,52 \times [LR=4])} \times e^{(-0,75 \times GEW)}} = \frac{e^{(0,52 \times 5)}}{e^{(0,52 \times 4)}} = e^{0,52} = 1,68$$

In gleicher Weise lässt sich auch der GEW-Regressionskoeffizient umrechnen. Durch Entlogarithmierung von  $b_2 = -0,75$  ergibt sich ein Wert von  $\exp(b_2) = 0,47$ . Der Wert von 0,47 ist der Multiplikationsfaktor mit dem sich das Wahrscheinlichkeitsverhältnis zwischen einer CDU-Wahlabsicht und der Absicht, eine andere Partei zu wählen, verändert, wenn Wähler einer Gewerkschaft angehören (bzw. wenn sich GEW=0 in GEW=1 verändert). In diesem Falle würde also die Chance für eine CDU-Wahlabsicht sinken, denn der Multiplikationsfaktor ist kleiner als 1,00. Nur Multiplikationsfaktoren von exakt 1,00 indizieren, dass eine X-Variable keinen Einfluss auf das Wahrscheinlichkeitsverhältnis zweier Handlungsalternativen ausübt.

Der hier vorgestellte Multiplikationsfaktor „ $\exp(b_k)$ “ wird in der logistischen Regressionsanalyse auch als *Effektkoeffizient* bezeichnet. Er beschreibt die Veränderung der Chance für das Ereignis  $Y=1$ , wenn sich im Regressionsmodell ein Prädiktor um eine empirische Einheit erhöht.<sup>9</sup>

Natürlich können Effektkoeffizienten auch durch Entlogarithmierung von teilstandardisierten Regressionskoeffizienten berechnet werden. Es entstehen dann teilstandardisierte Effektkoeffizienten. In unserem CDU-Beispiel entstehen dadurch folgende teilstandardisierte Effektkoeffizienten für LR und GEW:

$$\begin{aligned} b_{LR}^* &= 0,93 \quad \rightarrow \quad \exp(b_{LR}^*) &&= 2,53 \\ b_{GEW}^* &= -0,24 \quad \rightarrow \quad \exp(b_{GEW}^*) &&= 0,79 \quad \rightarrow \quad 1/\exp(b_{GEW}^*) = 1,27 \quad (-) \end{aligned}$$

<sup>9</sup> Da die Effektkoeffizienten multiplikativ wirken (vgl. Gleichung 14), müssen bei Veränderungen über mehrere X-Stufen hinweg auch dementsprechend viele Multiplikationen (und nicht etwa Additionen wie bei Regressionskoeffizienten) ausgeführt werden, um das Wahrscheinlichkeitsverhältnis auf der anvisierten X-Zielstufe zu berechnen.

Wie in dieser Aufstellung zu erkennen (unten rechts), empfiehlt es sich, für alle nicht- oder teilstandardisierten Effektkoeffizienten kleiner als 1,00 (also auch für den obigen Koeffizienten von 0,79) deren Kehrwert bei einem Vergleich zwischen mehreren Effektkoeffizienten zu benutzen. Der Kehrwert von 0,79 ist 1,27. Folglich ist der teilstandardisierte Effektkoeffizient von LR nicht ca. dreimal so stark wie der von GEW (2,53:0,79), sondern hat nur ca. zweimal die Stärke von GEW (2,53:1,27). Zur Kenntlichmachung der Kehrwertbildung kann, wie oben gezeigt, ein Minus-Zeichen in Klammern an den Kehrwert angehängt werden. Warum ist das sinnvoll?

Effektkoeffizienten haben zwei ungleich skalierte Wertebereiche, die oberhalb und unterhalb ihres neutralen Punkts von 1,00 liegen. Während der untere Bereich (welcher Verschiebungen der Wahrscheinlichkeitsverhältnisse zugunsten von  $Y=0$  signalisiert) zwischen 0,00 und 1,00 liegt, reicht der obere Bereich (welcher Verschiebungen der Wahrscheinlichkeitsverhältnisse zugunsten von  $Y=1$  ausdrückt) von 1,00 bis  $+\infty$ . Erst durch die Kehrwertbildung bei Koeffizienten kleiner als 1,00 wird auch die Begrenzung des unteren Bereichs aufgehoben, wodurch Koeffizienten aus beiden Bereichen hinsichtlich ihrer Größe miteinander vergleichbar werden. Entsprechend ist also auch der oben ausgewiesene unstandardisierte Effektkoeffizient von „ $\exp(b_2)=\exp(-0,75)=0,47$ “ in den Kehrwert von „ $1/0,47=2,13$ “ zu überführen. Dieser Kehrwert ist dann wie folgt inhaltlich zu interpretieren: Wähler, die keine Gewerkschaftsmitglieder sind, haben eine 2,13-fach höhere Chance, die CDU zu wählen, als Wähler, die Gewerkschaftsmitglieder sind.<sup>10</sup>

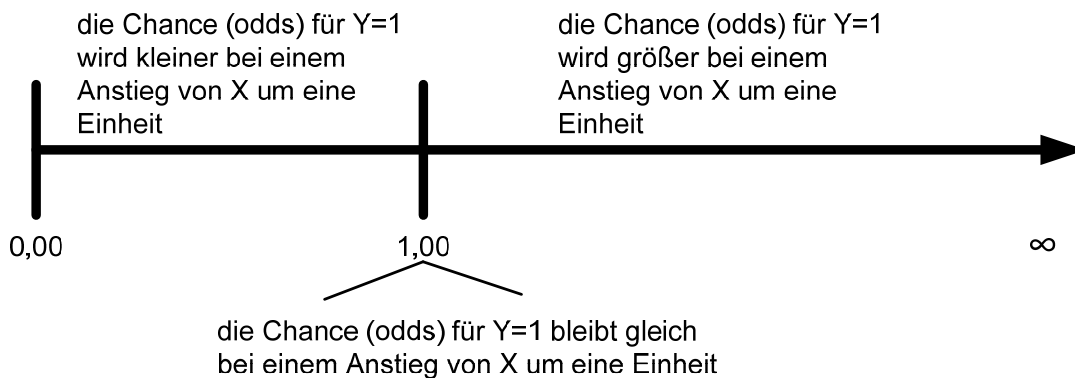
Die folgende Abbildung 5 verdeutlicht die asymmetrische Skalierung des nichtstandardisierten Effektkoeffizienten.

---

<sup>10</sup> Bei der inhaltlichen Interpretation des Effektkoeffizienten ist stets darauf zu achten, dass dieser nicht als Multiplikationsfaktor von einfachen Wahrscheinlichkeiten, sondern von Chancen bzw. Wahrscheinlichkeitsverhältnissen (d.h. von der Wahrscheinlichkeit für  $Y=1$  im Verhältnis zur Gegenwahrscheinlichkeit) zu verstehen ist. Und ebenso ist darauf zu achten, dass der Multiplikationsfaktor stets für die Erhöhung in der X-Variablen um eine Einheit gilt (vgl. hierzu auch Gleichung 15).



Abbildung 5: Die Skalierung des Effektkoeffizienten



Aufgrund seiner Interpretationsmöglichkeit als Multiplikationsfaktor sowie der Eigenschaft, dass der Effektkoeffizient („exp(b)“) unabhängig von den Startwerten der X-Variablen ist (d.h. der Multiplikationsfaktor ist stets gleich hoch, egal ob die entsprechende X-Variable z.B. von „1“ auf „2“ oder von „6“ auf „7“ um eine Einheit ansteigt), empfehlen wir, neben den b-Koeffizienten sowie den prognostizierten Wahrscheinlichkeiten für bestimmte X-Wertekombinationen stets auch die Effektkoeffizienten zu berichten und zu interpretieren.

### 3.2 Modell-Evaluation

Ein empirisch geschätztes, binär-logistisches Regressionsmodell kann bedeutsame und signifikante Regressionskoeffizienten aufweisen und dennoch für die Forschungspraxis unbrauchbar sein, weil die Modellschätzung insgesamt betrachtet nicht nahe genug an die empirisch beobachteten Daten herankommt. Im Folgenden werden daher einige Maßzahlen zur Überprüfung der Modellgüte vorgestellt. Im Einzelnen werden folgende Verfahren zur Evaluation der logistischen Modellschätzung erläutert:

- der Likelihood-Ratio-Test,
- die Analyse von Pseudo-R<sup>2</sup>-Koeffizienten,
- die Klassifizierung prognostizierter Wahrscheinlichkeiten.

Im *Likelihood-Ratio-Test* wird überprüft, ob die ML-Modellschätzung unter Verwendung von X-Prädiktoren eine bedeutsam bessere Anpassung an die beobachteten Daten erreicht, als eine Modellschätzung, bei der nur die a-Konstante aber ansonsten keinerlei Prädiktoren zur Vorhersage von P(Y) benutzt wird. Mithin basiert dieser Test auf einem Vergleich der ML-Schätzungen von zwei logistischen Regressionsmodellen, nämlich dem Modell ohne Prädiktoreffekte (dem Null-Modell) und dem Modell mit spezifizierten Prädiktoreffekten (dem Prädiktoren-Modell).

Im Likelihood-Ratio-Test wird nicht die Signifikanz einzelner Modell-Effekte sondern die Signifikanz des logistischen Gesamt-Modells überprüft. Die dementsprechende Null-Hypothese formuliert die Bedeutungslosigkeit aller im Modell spezifizierten Effekte:

$$H_0: (\beta_1 = \beta_2 = \dots = \beta_k) = (b_1 = b_2 = \dots = b_k) = 0 \quad (16)$$

Wenn der Log-Likelihood-Wert des Null-Modells als  $LL_0$  bezeichnet wird, und  $LL_p$  der Log-Likelihood-Wert des kompletten Prädiktoren-Modells ist, wird für den Likelihood-Ratio-Test die G-Statistik berechnet nach:<sup>11</sup>

$$G = -2 \ln \left( \frac{L_0}{L_p} \right) = -2(LL_0 - LL_p) = (-2LL_0) - (-2LL_p) \quad (17)$$

Diese G-Statistik kann mit Hilfe eines Chi-Quadrat-Tests geprüft werden. Denn der G-Wert ist asymptotisch chi-quadrat-verteilt und hat so viele Freiheitsgrade, wie es im Prädiktoren-Modell X-Variablen gibt. Im Test ist zu überprüfen, ob der  $LL_p$ -Wert signifikant kleiner als der  $LL_0$ -Wert ist. Denn nach der Null-Hypothese (Gleichung 16) gibt es keinen Unterschied zwischen dem Null- und dem Prädiktoren-Modell, da die im Null-Modell fehlenden Effekte sowieso bedeutungslos sind.

In unserem Beispiel zur binär-logistischen Analyse der CDU-Wahlabsicht ist  $-2LL_0=3066,599$  und ist  $-2LL_p=2685,712$ . Die G-Statistik beträgt demnach 380,887 und ist bei zwei Freiheitsgraden und entsprechend dem Verlauf der theoretischen Chi-Quadrat-Verteilung auf einem Niveau von 0,000 statistisch signifikant (vgl. dazu Abschnitt 3.4). Mithin muss das Prädiktoren-Modell im Vergleich zum Null-Modell als statistisch bedeutsame Verbesserung der Modellanpassung bzw. des Modellfits an die beobachteten Stichprobendaten gewertet werden.

Die Analyse von *Pseudo-R<sup>2</sup>-Koeffizienten* bietet eine weitere Möglichkeit, die Güte einer logistischen Modellschätzung zu bewerten. Üblicherweise werden als Pseudo-R<sup>2</sup>-Maßzahlen die folgenden beiden Koeffizienten berechnet:

- Cox & Snell R<sup>2</sup>
- Nagelkerkes R<sup>2</sup>

---

<sup>11</sup> SPSS präsentiert nicht den reinen LL-Wert einer Modellschätzung, sondern multipliziert diesen Wert mit „-2“ so dass der „-2LL“-Wert entsteht. Da SPSS den ursprünglichen LL-Wert als negativen Wert schätzt, ist der -2LL-Wert stets positiv. Je größer er ist, umso schlechter ist die Anpassung des geschätzten Modells. Der Grund dafür, dass SPSS den -2LL-Wert benutzt, liegt darin, dass dieser Wert für ein Modell, das nur den a-Koeffizienten aber keine Prädiktoren enthält, identisch ist mit der Summe der Abweichungsquadrate für dieses Modell in der OLS-Regression (vgl. Urban/Mayerl 2008: Kapitel 3.5).

Beide Koeffizienten sind nicht analog zum Determinationskoeffizienten „ $R^2$ “ in der OLS-Regression zu interpretieren. Sie sagen nichts über den Anteil ausgeschöpfter Varianz in einer Regressionsschätzung aus. Insofern ist ihre Bezeichnung als  $R^2$ - oder Pseudo- $R^2$ -Koeffizienten höchst unglücklich. Beide Koeffizienten sind reine Fit-Indizes, die den Grad der relativen Anpassung einer Regressionsschätzung an die beobachteten Stichprobenwerte durch Vergleich der Log-Likelihood-Werte von Null-Modell ( $LL_0$ ) und Prädiktoren-Modell ( $LL_P$ ) ermitteln. Ihre Zahlenwerte sind so zu interpretieren, dass diese den Prozentanteil berichten, um den der Schätzerfolg des Null-Modells (gemessen im  $LL_0$ -Wert) durch den Schätzerfolg des Prädiktoren-Modells (gemessen im  $LL_P$ -Wert) verbessert werden kann. Die Pseudo- $R^2$ -Koeffizienten sind somit ein modell-relatives Gütemaß. Sie vergleichen nur die Schätzergebnisse von zwei logistischen Regressionsmodellen.

Der Cox & Snell  $R^2$ -Koeffizient wird berechnet mit Hilfe der G-Statistik (vgl. Gleichung 17) in Form von:

$$R^2_{CS} = 1 - \exp\left[-\frac{G}{n}\right] \quad (18)$$

In unserer Modellschätzung zur CDU-Wahlabsicht (mit  $n = 2474$  Befragten) beträgt  $R_{CS}^2 = 0,14$ . Demnach kann die Schätzung von  $\hat{P}(CDU)$  um 14% verbessert werden, wenn zur Schätzung nicht das Null-Modell sondern das Prädiktoren-Modell benutzt wird. Leider kann der Koeffizient das Maximum von 1,00 nicht erreichen und ist deshalb nur zurückhaltend zu interpretieren.

Der Pseudo- $R^2$ -Koeffizient von Nagelkerke standardisiert  $R_{CS}^2$  und kann deshalb auch einen maximalen Wert von 1,00 erreichen:

$$R^2_N = \frac{R^2_{CS}}{R^2_{CS_{max}}} = \frac{R^2_{CS}}{1 - \exp(-(-2LL_0/n))} \quad (19)$$

In unserer Modellschätzung zur CDU-Wahlabsicht beträgt  $R_N^2 = 0,20$  und liegt damit deutlich über  $R_{CS}^2 = 0,14$ . Aufgrund seiner Standardisierung ist  $R_N^2$  immer größer als  $R_{CS}^2$ . Für unsere logistische Regression bedeutet ein  $R_N^2$  von 0,20, dass entsprechend der Logik des Pseudo- $R^2$ -Koeffizienten von Nagelkerke die Schätzung von  $\hat{P}(CDU)$  um 20% verbessert werden kann, wenn zur Schätzung nicht das Null-Modell sondern das Prädiktoren-Modell eingesetzt wird.

Die Anpassungsgüte einer logistischen Regressionsschätzung kann auch durch eine *Klassifizierung prognostizierter Wahrscheinlichkeiten* überprüft werden. Dieses Verfahren ist zwar intuitiv leicht nachvollziehbar, ist aber dennoch nur mit Vorsicht einzusetzen. Im Verfahren selbst wird ermittelt, wie viele der Personen, die in der Befragung eine Wahlabsicht zugunsten der CDU berichtet haben (CDU-Wahlabsicht = 1), und wie viele der Personen, die eine Wahlabsicht zugunsten einer anderen Partei berichtet haben (CDU-Wahlabsicht = 0), mit dem geschätzten Regressionsmodell aufgrund ihres jeweiligen LR- und GEW-Wertes als potentieller CDU-Wähler oder Nicht-CDU-Wähler „richtig“ erkannt werden konnten. Dazu wurden alle Personen mit einem im Modell geschätzten  $\hat{P}$ -Wert zwischen 0,50 und 1,00 als mögliche CDU-Wähler klassifiziert, und alle Personen mit einem geschätzten  $\hat{P}$ -Wert unterhalb von 0,50 als mögliche Nicht-CDU-Wähler eingestuft. In unserem Analysebeispiel können insgesamt 72,9% aller Befragten richtig klassifiziert werden (vgl. dazu die in Abschnitt 3.4 abgedruckte SPSS-Ausgabetable). Je höher der Anteil richtig klassifizierter Befragter ist, umso größer ist der Modellfit und damit auch die Modellgüte der entsprechenden Regressionsschätzung.

Zur Bewertung der Trefferquote einer Regressionsschätzung (hier: 72,9%) kann der Prognoseerfolg des Prädiktoren-Modells mit dem Prognoseerfolg des Null-Modells verglichen werden. Wie erinnerlich, werden im Null-Modell keine Prädiktoreneffekte geschätzt, sondern nur nach einer Konstanten gesucht, die für das Regressionsmodell den höchsten Likelihood-Wert erbringt. In unserem Beispiel beträgt der Prognoseerfolg des Null-Modells 68,9%. Die Trefferquote des Prädiktoren-Modells liegt also nur 4% oberhalb des ohne zusätzliche Dateninformation geschätzten Null-Modells. Das spricht auf den ersten Blick nicht für die Qualität des Prädiktoren-Modells. Jedoch muss bei der Interpretation dieser Zahlen berücksichtigt werden, dass die Klassifizierung prognostizierter Wahrscheinlichkeiten ein nicht sehr sensitives bzw. ein eher grobes Verfahren zur Bewertung des Modellfits ist. Denn bei der Zuordnung von Ereignissen ( $Y=0/1$ ) zu Wahrscheinlichkeiten ( $\hat{P}$ ) bleiben alle Informationen im Schätzergebnis außer der  $\hat{P}$ -Ausprägung „ $<0,5/\geq 0,5$ “ unberücksichtigt. Auf diese Weise führt z.B. sowohl ein Schätzergebnis von 0,51 als auch von 0,98 zur gleichen Prognose von  $Y=1$ . Bei der Wahrscheinlichkeitsklassifikation zur Ermittlung der Modellgüte wird also nicht die gesamte zur Verfügung stehende Schätzinformation ausgenutzt und insofern eine metrische Skala auf eine Nominal-Skala reduziert. Dementsprechend ist eine Bewertung des Modellfits mittels einer Klassifizierung prognostizierter Wahrscheinlichkeiten nur sehr zurückhaltend und stets in Kombination mit anderen Fit-Indizes vorzunehmen.

### 3.3 Problemdiagnostik im logistischen Regressionsmodell

Die ML-basierte, binär-logistische Regressionsanalyse ist störanfällig. So kann u.U. die Maximum-Likelihood-Schätzung nicht konvergieren und zu keinem endgültigen Schätzergebnis führen. Auch kann die ML-Schätzung zwar u.U. konvergieren, dabei jedoch eher zweifelhafte Parameterschätzwerte erbringen. Aber auch bei normal ablaufender ML-Schätzung können u.U. die Ergebnisse einer logistischen Regressionsanalyse fehlerbelastet sein, weil mit den ausgewerteten Daten bestimmte Modellannahmen der Analyse (z.B. die Linearität der Beziehung zwischen dem Logit von  $\hat{P}$  und den X-Prädiktoren) nicht einzuhalten sind. Im Folgenden soll deshalb auf einige der häufigsten Probleme bei Durchführung von logistischen Regressionsanalysen aufmerksam gemacht werden, und es sollen Tipps zum Umgang mit evtl. vorhandenen Problemen gegeben werden.

#### *Fallzahl*

Die zu analysierende Stichprobe sollte möglichst umfangreich sein. Da die ML-Schätzer asymptotisch konsistent und effizient sind, steigt auch die Qualität der Schätzergebnisse mit zunehmender Fallzahl (vgl. dazu Abschnitt 2). Wie groß die Fallzahl sein sollte, um robuste Schätzwerte zu erhalten, ist nur schwer zu ermitteln und hängt u.a. vom benutzten Schätzalgorithmus sowie von diversen Modelleigenschaften (z.B. von der Prädiktoren-Anzahl) und von gegebenen Datenstrukturen ab (z.B. von der Skalierung der Prädiktoren). Geringe Fallzahlen erhöhen die Wahrscheinlichkeit von Null-Zellen (s.u.) sowie von vollständiger Separation (s.u.) und reduzieren die Teststärke des Modells (vgl. Urban/Mayerl 2008: Kapitel 3.4.3). Eine Faustregel zur Überprüfung der Fallzahl bei kategorialen Modellvariablen besteht darin, in SPSS zwischen den Variablen bivariate Kreuztabellen erstellen zu lassen („CROSSTABS /TABLES=x1 BY x2 /CELLS=COUNT EXPECTED.“). In diesen Tabellen sollten die erwarteten Zelhäufigkeiten größer als 1,00 sein und nicht mehr als 20% sollten kleiner 5,00 sein.

#### *Y-Werte Verteilung*

Jeder der beiden Werte der abhängigen, binomial skalierten Y-Variablen sollte zumindest von einer kleinen Personengruppe der Stichprobe gewählt worden sein (Daumenregel: mindestens von ca. 10% aller Fälle). Wird ein Y-Wert zu selten beobachtet d.h. ist die Varianz von Y sehr gering, kann die ML-Schätzung u.U. nicht konvergieren oder erzeugt unplausible Schätzwerte (vgl. die Erläuterung zu „unvollständige Information“). Zur Diagnose kann auch hier der unter Punkt „Fallzahl“ beschriebene Kreuztabellen-Test (diesmal mit Y- und X-Variablen) eingesetzt werden. Bei extrem ungleicher Y-Werte Verteilung muss ein anderes Design für Stichprobenziehung und

Datenanalyse gewählt werden (z.B. ein gematchtes Fall-Kontroll-Design, vgl. Hosmer/Lemeshow 1989: 187ff).

### *unvollständige Information*

Für alle Werte-Kombinationen, die unter den im Regressionsmodell vertretenen Variablen möglich sind, sollte es eine ausreichende Anzahl von empirisch beobachteten Fällen geben. Wenn z.B. die beiden Variablen „Geschlecht“ (1/0) und „Parteimitgliedschaft“ (1/0) in einem Modell vorkommen, sollte es zu allen vier möglichen Wertekombinationen (also zu: 1-1, 1-0, 0-1, 0-0) zumindest einige Beobachtungsfälle geben. Das Risiko von Null-Wertekombinationen ist durch Inspektion von uni- und bivariaten Häufigkeitstabellen zu erkennen. Diese sollten keine Null-Zellen oder sehr schwach besetzte Zellen aufweisen. Zur Diagnose kann auch hier der unter Punkt „Fallzahl“ beschriebene Kreuztabellen-Test eingesetzt werden. Sind viele Zellen mit geringen Fallzahlen vorhanden, sollten diese durch Zusammenlegung von Variablenkategorien eliminiert werden. Schon eine Null-Zelle kann dazu führen, dass die ML-Schätzung nicht konvergiert oder fehlerhafte Schätzwerte ausgibt. Der beste Hinweis auf ML-Schätzprobleme ist ein übermäßig groß geschätzter Standardfehler, der sich sehr deutlich vom Schätzwert für den b-Koeffizienten unterscheidet (z.B.:  $b=12,3$ ;  $SE=140,2$ ). Sollte dies der Fall sein, kommt als eine mögliche Ursache die Existenz von Null-Zellen in Betracht.

### *vollständige Separation*

Es sollte keine „vollständige Separation“ im Datensatz für eine binär-logistischen Regressionsanalyse bestehen, d.h. die Werte  $Y=1$  oder  $Y=0$  sollten nicht ausschließlich in Kombination mit bestimmten X-Werten in der Stichprobe auftreten (z.B. sollte es keine X-Dummies geben, die für  $Y=1$  immer nur den Wert  $D=0$  aufweisen). Ansonsten wären die betreffenden X-Variablen perfekte Prädiktoren für jedes Y-Ereignis und die ML-Schätzung könnte keine maximalen Likelihood-Werte für die zu schätzenden Parameter finden. Begünstigt wird eine vollständige Separation durch kleine Fallzahlen (s.o.), durch eine stark ungleiche (unbalancierte) Y-Werteverteilung (s.o.) und durch eine große Anzahl von X-Variablen. Um das Risiko einer vollständigen Separation schon vor Beginn der Modellschätzung zu erkennen, sollten noch vor der Modellanalyse bivariate Kreuztabellen zwischen den einzelnen X-Variablen und der Y-Variablen erstellt werden. Ist darin keine vollständige Separation zu erkennen, so ist diese zwar nicht mit Sicherheit auszuschließen, kann dann aber noch in der Modellschätzung durch unplausibel hohe Schätzwerte (insbesondere, aber nicht nur, für einen oder mehrere Standardfehler) identifiziert werden.

### *Multikollinearität*

Es sollten keine übermäßig starken linearen Beziehungen zwischen den unabhängigen X-Variablen des logistischen Regressionsmodells bestehen (und erst recht keine vollständigen Kollinearitäten vorhanden sein). Wenn dem so ist, wird die ML-Schätzung immer falsche, sehr hohe Standardfehler und manchmal auch unzutreffend hohe Regressionskoeffizienten erbringen. Zum Test auf Multikollinearität sollte mit allen Variablen des logistischen Regressionsmodells eine OLS-Regression und sollten die von der OLS-Regressionsanalyse bekannten Tests durchgeführt werden (vgl. Urban/Mayerl 2008: Kapitel 4.5). Bei diagnostizierter Multikollinearität können dann die üblichen Strategien zur Beseitigung der Kollinearität der X-Variablen eingesetzt werden (u.a. Bildung von Indices, Mittelwertzentrierung, Ausschluss von X-Variablen, vgl. Urban/Mayerl 2008).

### *Linearität*

In der logistischen Regression sollte eine lineare Beziehung zwischen dem Logit-Wert von  $\hat{p}$  und den X-Prädiktoren bestehen (vgl. Gleichung 11), wobei dies insbesondere für kontinuierliche bzw. metrisch-skalierte X-Variablen gilt. In SPSS können die Logits durch die Anweisung „COMPUTE LOGIT = LN (PRE\_1 / (1 - PRE\_1))“ berechnet werden. Dabei ist „PRE\_1“ die von SPSS ausgegebene  $\hat{p}$ -Variable (vgl. Abschnitt 3.4). Sodann kann mit der Y-Variablen „LOGIT“ und den X-Variablen „LR, GEW“ eine OLS-Regression durchgeführt werden und die üblichen Verfahren zur Diagnose von Nicht-Linearität eingesetzt werden (vgl. Urban/Mayerl 2008: Kapitel 4.3). Eine weitere Strategie besteht darin, zusätzlich zu jeder metrischen X-Variablen eine Interaktionsvariable durch Multiplikation der Original-X-Variablen mit ihrem natürlichen Logarithmus zu bilden. Werden diese Interaktionsvariablen ergänzend in das Regressionsmodell hinein genommen, signalisieren signifikante Interaktionseffekte nicht-lineare Beziehungen zwischen der Logit-Variablen und den metrischen X-Prädiktoren.

### *Residuen-Unabhängigkeit*

Residuen-Unabhängigkeit bedeutet, dass die Werte einer bestimmten X-Variablen nicht über die Beobachtungsfälle hinweg miteinander korrelieren und somit keine Autokorrelation besteht. Wenn, wie unter Punkt „Linearität“ dargestellt, in der logistischen Regression die Variable „LOGIT“ berechnet wird, kann eine OLS-Regression mit den Variablen „LOGIT, LR, GEW“ durchgeführt werden und können die von der OLS-Regressionsanalyse bekannten Verfahren zur Diagnose von Autokorrelation eingesetzt werden (vgl. Urban/Mayerl 2008: Kapitel 4.7). Autokorrelation kann die Schätzwerte für die Standardfehler fehlerhaft klein werden lassen und so falsche Signifikanztestresultate erzeugen. Um dies zu korrigieren, kann bei vorhandener Multikollinearität

die Wald-Statistik (vgl. Abschnitt 3.4) durch Multiplikation mit einem Varianz-Inflationsfaktor modifiziert werden (vgl. Tabachnick/Fidell 2007: 444).

### *Ausreißer*

Ausreißer sind Fälle, bei denen sich beobachteter Y-Wert und prognostizierter  $\hat{P}$ -Wert überdurchschnittlich stark unterscheiden. Wenn der Anteil von Ausreißerfällen an der Gesamtzahl aller analysierten Fälle hoch ist, hat das geschätzte Modell einen schlechten Modellfit. Auch die Modellschätzwerte können durch einen hohen Anteil von Ausreißerfällen verzerrt werden. SPSS stellt in der logistischen Regressionsanalyse nach entsprechender Anforderung (vgl. Abschnitt 3.4) eine Vielzahl von Kennwerten zur Ausreißeranalyse zur Verfügung (u.a.: DFBETA, LEVER, RESID, SRESID, ZRESID, LRESID; zur Interpretation dieser Statistiken vgl. Urban/Mayerl 2008: Kapitel 4.1.1). Am häufigsten wird in der logistischen Regressionsanalyse das standardisierte „Pearson Residuum“ eingesetzt (in SPSS: „ZRESID“). Es wird berechnet als:  $R_{\text{Pearson}_i} = (Y_i - \hat{P}_i) / \sqrt{[\hat{P}_i (1 - \hat{P}_i)]}$  und sollte zwischen -2 und +2 liegen. Anders als bei der OLS-Regression sollte bei der logistischen Regressionsanalyse sehr sorgfältig überlegt werden, ob Fälle mit großen Residualwerten aus der Analyse ausgeschlossen werden müssen. Denn ihr Y-Wert beträgt ja immer nur „1“ oder „0“ und ein großes Residuum kommt nur dadurch zustande, dass die geschätzten X-Effekte für diese Fälle nicht zutreffen.

### *unbeobachtete Heterogenität und Koeffizientenvergleich*

Ein häufig übersehenes Problem bei der Interpretation der Ergebnisse von logistischen Regressionen wird durch Einflüsse unbeobachteter Heterogenität verursacht. Grundsätzlich meint unbeobachtete Heterogenität, dass wichtige X-Variablen zur „Erklärung“ der Y-Variablen im Regressionsmodell nicht berücksichtigt wurden. Anders als bei der OLS-Regression können bei der logistischen Regression jedoch auch diejenigen nicht-berücksichtigten X-Variablen bedeutsame Auswirkungen auf die Modellschätzungen von b und exp(b) haben, die zwar einen Effekt auf Y ausüben, aber *nicht* mit den X-Variablen des Regressionsmodells korrelieren ( $r = 0,00$ ). Dabei gilt, dass die Richtung dieser Verzerrung stets bekannt ist: b und exp(b) werden immer *unterschätzt* bei unkorrelierter unbeobachteter Heterogenität. Entsprechend führt die Berücksichtigung von weiteren X-Variablen im Regressionsmodell, die zwar einen Einfluss auf Y ausüben, mit den anderen X-Variablen aber unkorreliert sind, zu einer Erhöhung der b- und exp(b)-Koeffizienten der übrigen X-Variablen im Modell (vgl. zum Nachweis z.B. Mood 2010). *Innerhalb* eines Regressionsmodells können b und exp(b) zwar weiterhin als *relative* Maße zum Vergleich der Einflussstärke von gleich skalierten X-Variablen verwendet werden, da die Höhe der unbeobachteten Heterogenität innerhalb eines Modells für alle X-Variablen gleich hoch ist. Die absolute Höhe von b und exp(b) ist jedoch



stets durch unkorrelierte unbeobachtete Heterogenität nach unten verzerrt. Vergleiche von  $b$  und  $\exp(b)$  *zwischen* Gruppen (z.B. ein Modellvergleich zwischen Frauen und Männern), zwischen verschiedenen Modellen mit unterschiedlichen X-Variablen oder zwischen verschiedenen Zeitpunkten sind hingegen mit großer Vorsicht zu betrachten, da diese Vergleiche die (nicht testbare!) empirische Gültigkeit der Annahme voraussetzen, dass die unbeobachtete Heterogenität in allen Modellen bzw. Gruppen bzw. zu allen Zeitpunkten gleich hoch ist. Hat man für diese Annahme keine triftigen theoretischen Argumente, so raten wir von einem einfachen Vergleich von  $b$  und  $\exp(b)$  zwischen Gruppen, Modellen mit unterschiedlichen X-Variablen oder verschiedenen Zeitpunkten ab. In der Literatur werden verschiedene Verfahren vorgeschlagen, mit deren Hilfe man versuchen kann, einen unverzerrten Koeffizientenvergleich zwischen unterschiedlichen Modellen, Zeitpunkten oder Gruppen durchzuführen (z.B. Allison 1999; Mood 2010). Möchte man dennoch mit guten Gründen einen herkömmlichen Gruppenvergleich durchführen, so sollten zumindest so viele bedeutsame X-Variablen wie möglich als Kontrollvariablen in das Modell aufgenommen werden, um den potentiellen Einfluss ungleicher unbeobachteter Heterogenität zu reduzieren. Mit den gruppenspezifischen unstandardisierten  $b$ -Koeffizienten (und deren Standardfehlern) könnte dann genau dasselbe Testverfahren zum Gruppenvergleich durchgeführt werden, das auch bei der OLS-Regression verwendet wird (vgl. Urban/Mayerl 2008: Kapitel 5.2.2).

In diesem Abschnitt wurden einige der schwerwiegendsten und häufigsten Probleme von ML-basierten binär-logistischen Regressionsanalysen vorgestellt. Ausführlicher werden diese und noch weitere Probleme sowie diesbezügliche Diagnosestrategien beschrieben u.a. in: Fox 1984: 320-332; Hosmer/Lemeshow 1989: 149-170; Menard 2002: 67-90; Vittinghoff et al. 2005: 188-195.

### **3.4 SPSS-Beispiel: Binär-logistische Regression**

Im Folgenden werden wir die Durchführung einer binär-logistischen Regressionsanalyse mit dem Statistikprogramm SPSS erläutern. Dabei konzentrieren wir uns auf die Interpretation der wichtigsten Tabellen, die SPSS bei einer logistischen Regressionsanalyse ausgibt.

Im SPSS-Beispiel werden die Daten und Variablen analysiert, die auch in den vorausgegangenen Abschnitten zur Erläuterung der logistischen Regression benutzt werden. Die binomial (binär) skalierte, abhängige Y-Variable „CDU-Wahlabsicht“ wird im Folgenden als „cdu“ bezeichnet (1=CDU-Wahlabsicht, 0=keine CDU-Wahlabsicht). Die unabhängigen X-Prädiktoren sind: „Links-Rechts-Selbsteinstufung“ (Variable: „lr“, 1=links ... 10=rechts) und „Gewerkschaftsmitgliedschaft“ (Variable „gew“, 1=ja, 0=nein).

Wird die Syntax-Steuerung von SPSS benutzt, so sind zumindest folgende Befehlszeilen für eine binär-logistische Regressionsanalyse einzugeben:

```
GET FILE = 'xxx.SAV'.
LOGISTIC REGRESSION VARIABLES cdu
/METHOD = ENTER gew lr
/PRINT= ITER(1)
/SAVE = PRED ZRESID DFBETA LEVER RESID SRESID ZRESID LRESID.
SAVE OUTFILE = 'xxx.SAV' / KEEP = ALL.
```

Mit dem Zusatz „/SAVE PRED ZRESID“ werden, analog zum Syntax-Befehl der OLS-Regression (vgl. Urban/Mayerl 2008: Kap. 4), zwei neue Modellvariablen gebildet: mit „PRED“ erhält man die durch die X-Variablen prognostizierte Wahrscheinlichkeit  $\hat{P}$  für die Wahlabsicht der Partei CDU (neue Variable mit dem Titel „PRE\_1“), und mit „ZRESID“ wird das standardisierte Pearson Residuum in einer neuen Variablen (mit dem Titel „ZRE\_1“) abgespeichert. Damit können die in Abschnitt 3.3 angesprochenen Berechnungen durchgeführt werden. Mit einer Erweiterung dieser Befehlszeile in „/SAVE = PRED ZRESID DFBETA LEVER RESID SRESID LRESID“ können weitere Kennwerte zur Ausreißeranalyse (vgl. Abschnitt 3.3) angefordert werden. Mit dem Befehl „GET FILE“ lässt sich ein Datensatz öffnen und mit „SAVE OUTFILE“ mit den neu generierten Variablen in einem neuen Datensatz abspeichern.

Im Folgenden werden Hinweise zur Interpretation nicht aller, aber der wichtigsten SPSS-Ausgabetafeln bei einer binär-logistischen Regressionsanalyse gegeben.

## BLOCK 0: Anfangsblock

### Iterationsprotokoll

Iteration	-2 Log-Likelihood	Koeffizienten
		Constant
Schritt 0 1	3067,432	-,757
2	3066,599	-,796
3	3066,599	-,796

Im „Iterationsprotokoll“ für den „Anfangsblock“ wird mitgeteilt, dass die Maximum-Likelihood-Schätzung des Null-Modells (Modell ohne Prädiktoren, nur mit a-Konstante) drei Schritte benötigte und dann beendet wurde, weil sich der Schätzwert für den a-Parameter (Constant) nur noch um einen Betrag kleiner als 0,001 veränderte (vgl. dazu Abschnitt 2). Der  $-2LL_0$ -Wert beträgt nach dem dritten und letzten Iterationsschritt 3066,599.

**BLOCK 1: Methode = Einschluß****Iterationsprotokoll**

Iteration	-2 Log-Likelihood	Koeffizienten		
		Constant	GEW	LR
Schritt 1 1	2719,034	-2,607	-,449	,374
2	2686,329	-3,363	-,694	,495
3	2685,712	-3,477	-,750	,514
4	2685,712	-3,480	-,752	,515
5	2685,712	-3,480	-,752	,515

Im „Iterationsprotokoll“ für den „Block 1“ wird mitgeteilt, dass die Maximum-Likelihood-Schätzung des Prädiktoren-Modells (Modell mit a-Konstante und zwei Prädiktoren) fünf Schritte benötigte und dann beendet wurde, weil sich danach die Schätzwerte für die drei Parameter nur noch um Beträge kleiner als 0,001 veränderten (vgl. dazu Abschnitt 2). Der  $-2LL_P$ -Wert beträgt somit 2685,712.

**Omnibus-Tests der Modellkoeffizienten**

	Chi-Quadrat	df	Sig.
Schritt 1 Schritt	380,887	2	,000
Block	380,887	2	,000
Modell	380,887	2	,000

In der Tabelle „Omnibus-Tests der Modellkoeffizienten“ wird das Ergebnis des Likelihood-Ratio-Tests ausgegeben (vgl. dazu Abschnitt 3.2). Der Wert „380,887“ in der Zeile „Modell“ ist der G-Wert, der durch die Subtraktion „ $(-2LL_0) - (-2LL_P)$ “ entsteht (vgl. Gleichung 17). Dieser Wert ist im Chi-Quadrat-Test hoch signifikant. Der Test hat zwei Freiheitsgrade (df), weil sich Null-Modell und Prädiktoren-Modell um zwei zu schätzende Parameter (für die Effekte von LR und GEW) unterscheiden.

**Modellzusammenfassung**

Schritt	-2 Log-Likelihood	Cox & Snell R-Quadrat	Nagelkerkes R-Quadrat
1	2685,712	,143	,201

Die Tabelle „Modellzusammenfassung“ berichtet noch einmal den -2LL-Wert des Prädiktoren-Modells sowie die Werte der beiden Pseudo-R<sup>2</sup>-Koeffizienten zur Beurteilung des Modellfits (vgl. dazu Abschnitt 3.2).

### Klassifizierungstabelle

Beobachtet	Vorhergesagt		
	cdu		Prozentsatz der Richtigen
	0	1	
Schritt 1 cdu 0	1535	170	90,0
1	501	268	34,9
Gesamtprozentsatz			72,9

Die „Klassifizierungstabelle“ weist den Anteil von Personen aus, die mit dem Prädiktoren-Modell hinsichtlich ihrer CDU-Wahlabsicht richtig klassifiziert werden konnten (vgl. dazu Abschnitt 3.2). Das sind insgesamt 72,9% aller 2474 Befragten. Dabei konnten die Befragten mit einem CDU-Wert von 0 (insgesamt 1535+170=1705 Personen) treffsicherer bestimmt werden (90,0 %) als die Befragten mit einem CDU-Wert von 1 (insgesamt 501+268=769 Personen), die „nur“ zu 34,9% richtig erkannt wurden. Der Wert von 72,9% liegt über dem Anteil von 68,9% richtig prognostizierten Personen aus der Klassifizierungstabelle des Null-Modells (hier nicht abgedruckt) (vgl. dazu Abschnitt 3.2).

### Variablen in der Gleichung

	Regressi- onskoeffi- zient B	Standard- fehler	Wald	df	Sig.	Exp(B)
Schritt 1 GEW	-,752	,172	19,044	1	,000	,471
LR	,515	,031	275,137	1	,000	1,673
Konstante	-3,480	,180	374,988	1	,000	,031

Die Tabelle „Variablen in der Gleichung“ berichtet die Schätzwerte der unstandardisierten partiellen Regressionskoeffizienten (B) und Effektkoeffizienten (Exp(B)) im Prädiktoren-Modell (vgl. dazu Abschnitt 3). Alle Schätzwerte sind hoch signifikant. Die Wald-Statistik ergibt sich aus:  $B^2/SE^2$ . Die Wald-Statistik ist nicht t-verteilt, sondern Chi<sup>2</sup>-verteilt. Einen t-Test kann man aber leicht „per Hand“ durchführen, indem man entweder  $|B/SE|$  ermittelt oder die Wurzel aus der Wald-Statistik zieht. Die Interpretation ist ansonsten mit dem t-Test identisch (vgl. Urban/Mayerl 2008: Kapitel 3.3.1 und 3.4.1).

## Literaturverzeichnis

- Allison, P.D., 1999: Comparing Logit and Probit Coefficients Across Groups. *Sociological Methods & Research* 28(2):186-208.
- Amemiya, T., 1981: Qualitative Response Models: A Survey. *Journal of Economic Literature* 19: 1483-1536.
- Dhrymes, P.J., 1978: *Introductory Econometrics*. New York et al.: Springer.
- Fox, J., 1984: *Linear Statistical Models and Related Methods: With Applications to Social Research*. New York: Wiley.
- Hosmer, D.W. / Lemeshow, S., 1989: *Applied Logistic Regression*. New York: Wiley.
- Kriz, J., 1983: *Statistik in den Sozialwissenschaften. Einführung und kritische Diskussion* (5. Auflage). Opladen: Westdeutscher Verlag.
- Menard, S., 2002: *Applied Logistic Regression Analysis* (2. Auflage). Thousand Oaks, CA: Sage.
- Mood, C., 2010: Logistic Regression: Why We Cannot Do What We Think We Can Do, and What We Can Do About It. *European Sociological Review* 26(1): 67-82.
- Tabachnick, B.G. / Fidell, L.S., 2007: *Using Multivariate Statistics* (5. Auflage). Boston: Pearson Education.
- Urban, D., 1993: *LOGIT-Analyse: Statistische Verfahren zur Analyse von Modellen mit qualitativen Response-Variablen*. Stuttgart: Fischer.
- Urban, D. / Mayerl, J., 2008: *Regressionsanalyse: Theorie, Technik und Anwendung* (3. Auflage). Wiesbaden: VS Verlag.
- Vittinghoff, E. / Shiboski, S.C. / Glidden, D.V. / McCulloch, C.E., 2005: *Regression Methods in Biostatistics: Linear, Logistic, Survival, and Repeated Measures Models*. New York: Springer.

## **SISS: Schriftenreihe des Instituts für Sozialwissenschaften der Universität Stuttgart bisher sind erschienen:**

- No. 1/1994 "Vertrauen" - soziologisch betrachtet. Ein Beitrag zur Analyse binärer Interaktionssysteme.  
Peter Antfang, Dieter Urban
- No. 2/1994 Report on the German Machine Tool Industry.  
Frank C. Englmann, Christian Heyd, Daniel Köstler, Peter Paustian  
with the assistance of Susanne Baur and Peter Bergmann
- No. 3/1994 Neue württembergische Rechtstatsachen zum Unternehmens- und Gesellschaftsrecht.  
Udo Kornblum
- No. 4/1994 Rechtstatsachen zum Unternehmens- und Gesellschaftsrecht aus den neuen Bundesländern.  
Udo Kornblum
- No. 1/1995 Die Bedeutung Neuronaler Netze in der Ökonomie.  
Hermann Schnabl
- No. 2/1995 Regionale Strukturprobleme.  
Sammelband der Beiträge zum Symposium vom 13. und 14. Oktober 1994.  
Frank C. Englmann (Hrsg.)
- No. 3/1995 Latent Attitude Structures Directing the Perception of New Technologies.  
An Application of SEM-Methodology to the Construction of Attitude Measurement Models Related to Technologies of Prenatal Genetic Engineering and Testing.  
Dieter Urban
- No. 4/1995 Handbuch zur empirischen Erhebung von Einstellungen/Kognitionen zur Bio- und Gentechnologie (inklusive Diskette)  
(zweite, überarbeitete und erweiterte Auflage)  
Uwe Pfenning, Dieter Urban, Volker Weiss
- No. 5/1995 Social Indicators in a Nonmetropolitan County: Testing the Representativeness of a Regional Nonrandom Survey in Eastern Germany.  
Dieter Urban, Joachim Singelmann
- No. 1/1996 Jugend und Politik im Transformationsprozeß. Eine Fallstudie zur Stabilität und Veränderung von politischen Einstellungen bei ostdeutschen Jugendlichen zwischen 1992 und 1995.  
Dieter Urban, Joachim Singelmann, Helmut Schröder
- No. 2/1996 Einstellungsmessung oder Einstellungsgenerierung? Die Bedeutung der informationellen Basis bei Befragten für die empirische Rekonstruktion von Einstellungen zu gentechnischen Anwendungen.  
Martin Slaby
- No. 1/1997 Gentechnik: „Fluch oder Segen“ versus „Fluch und Segen“.  
Bilanzierende und differenzierende Bewertungen der Gentechnik in der öffentlichen Meinung.  
Dieter Urban und Uwe Pfenning

(Fortsetzung ...)

## **SISS: Schriftenreihe des Instituts für Sozialwissenschaften der Universität Stuttgart bisher sind erschienen (Fortsetzung):**

- No.2/1997 Die soziale Vererbung von Ausländer“feindlichkeit“. Eine empirische Längsschnittanalyse der intra- und intergenerativen Transmission von sozialen Einstellungen.  
Dieter Urban und Joachim Singelmann
- No. 3/1997 Politische Sozialisation im Transformationsprozeß: Die Entwicklung demokratiebezogener Einstellungen von ostdeutschen Jugendlichen und deren Eltern zwischen 1992 und 1996.  
Barbara Schmidt, Dieter Urban, Joachim Singelmann
- No.1/1998 Bewertende Einstellungen zur Gentechnik: ihre Form, ihre Inhalte und ihre Dynamik. Kurzbericht zu Ergebnissen des Forschungsprojektes „Einstellungen zur Gentechnik“.  
Dieter Urban, Uwe Pfenning, Joachim Allhoff
- No.2/1998 Technikeinstellungen: gibt es die überhaupt? Ergebnisse einer Längsschnittanalyse von Bewertungen der Gentechnik.  
Dieter Urban
- No.3/1998 Zur Interaktion zwischen Befragten und Erhebungsinstrument. Eine Untersuchung zur Konstanz des Meinungsurteils von Befragten im Interviewverlauf.  
Martin Slaby
- No.1/1999 Role Models and Trust in Socio-Political Institutions: A Case Study in Eastern Germany, 1992-96.  
Joachim Singelmann, Toby A. Ten Ayck, Dieter Urban
- No.1/2000 Die Zufriedenheit von Stuttgarter Studierenden mit ihrer Lebens- und Wohnsituation. Erste deskriptive Ergebnisse einer sozialwissenschaftlichen Studie zu allgemeinen und bereichsspezifischen Zufriedenheiten der Studierenden des Campus Vaihingen und des Campus Hohenheim.  
Projektgruppe Campus: Slaby, M.; Grund, R.; Mayerl, J.; Noak, T.; Payk, B.; Sellke, P.; Urban, D.; Zudrell, I.
- No.2/2000 Längsschnittanalysen mit latenten Wachstumskurvenmodellen in der politischen Sozialisationsforschung.  
Dieter Urban
- No.1/2001 Unser „wir“ - ein systemtheoretisches Modell von Gruppenidentitäten.  
Jan A. Fuhse
- No.2/2001 Differentielle Technikakzeptanz, oder: Nicht immer führt die Ablehnung einer Technik auch zur Ablehnung ihrer Anwendungen.  
Eine nutzentheoretische und modell-statistische Analyse.  
Martin Slaby, Dieter Urban
- No.3/2001 Religiosität und Profession. Longitudinale Analysen zur Entwicklung des religiösen Selbstbildes bei Erzieherinnen.  
Heiko Lindhorst
- No.4/2001 Ist Glück ein affektiver Sozialindikator subjektiven Wohlbefindens?  
Dimensionen des subjektiven Wohlbefindens und die Differenz zwischen Glück und Zufriedenheit.  
Jochen Mayerl

(Fortsetzung ...)

**SISS: Schriftenreihe des Instituts für Sozialwissenschaften der Universität Stuttgart  
bisher sind erschienen (Fortsetzung):**

- No.1/2002 Risikoakzeptanz als individuelle Entscheidung.  
Zur Integration der Risikoanalyse in die nutzentheoretische  
Entscheidungs- und Einstellungsforschung.  
Martin Slaby, Dieter Urban
- No.2/2002 Vertrauen und Risikoakzeptanz. Zur Relevanz von Vertrauen  
bei der Bewertung neuer Technologien.  
Martin Slaby, Dieter Urban
- No.3/2002 Probleme bei der Messung individueller Veränderungsdaten.  
13 empirisch und methodisch induzierte Effekte, die es schwierig machen,  
Veränderungen von generalisierten Bewertungen zu ermitteln.  
Dieter Urban
- No.1/2003 Systeme, Netzwerke, Identitäten. Die Konstitution sozialer Grenzziehungen  
am Beispiel amerikanischer Straßengangs.  
Jan A. Fuhse
- No.2/2003 Können Nonattitudes durch die Messung von Antwortreaktionszeiten ermittelt werden?  
Eine empirische Analyse computergestützter Telefoninterviews.  
Jochen Mayerl
- No.1/2004 Erhöht ein Opfer-Täter-Zyklus das Risiko, Sexualstraftaten als pädosexuelle Straftaten zu  
begehen? Ergebnisse einer ereignisanalytischen Pilotstudie  
Dieter Urban, Heiko Lindhorst
- No.1/2005 Persönliche Netzwerke in der Systemtheorie  
Jan A. Fuhse
- No.2/2005 Analyzing cognitive processes in CATI-Surveys with response latencies:  
An empirical evaluation of the consequences of using different  
baseline speed measures.  
Jochen Mayerl, Piet Sellke, Dieter Urban
- No.1/2006 Ist Bildung gleich Bildung? Der Einfluss von Schulbildung auf ausländerablehnende  
Einstellungen in verschiedenen Alterskohorten.  
Klaus Hadwiger
- No.2/2006 Zur soziologischen Erklärung individuellen Geldspendens.  
Eine Anwendung und Erweiterung der Theory of Reasoned Action unter Verwendung von  
Antwortlatenzzeiten in einem Mediator-Moderator-Design.  
Jochen Mayerl
- No.1/2007 Antwortlatenzzeiten in TRA-Modellen. Zur statistischen Erklärung von (Geld)-  
Spendenverhalten.  
Dieter Urban, Jochen Mayerl
- No.1/2008 Berufseintritt und Berufssituation von Soziologieabsolventen der Universität Stuttgart.  
Deskriptive Ergebnisse einer Absolventenbefragung aus dem Jahr 2007.  
Jochen Mayerl, Dieter Urban
- No.1/2010 Der Bystander-Effekt in alltäglichen Hilfsituationen:  
Ein nicht-reaktives Feldexperiment.  
Katrin Alle, Jochen Mayerl

(Fortsetzung ...)



**SISS: Schriftenreihe des Instituts für Sozialwissenschaften der Universität Stuttgart  
bisher sind erschienen (Fortsetzung):**

- No.2/2010 Das Working-Poor-Problem in Deutschland.  
Empirische Analysen zu den Ursachen von Armut trotz Erwerbstätigkeit.  
Leonie Hellmuth, Dieter Urban
- No.3/2010 Binär-logistische Regressionsanalyse.  
Grundlagen und Anwendung für Sozialwissenschaftler.  
Jochen Mayerl, Dieter Urban