

# Logistische Regression: Die Analyse binärer Zielvariablen

brunner & hess software ag

Lorenz Gygax & Christian Schmidhauser

Version 1.0, März 2002

## 1 Einführung

Bei der logistischen Regression geht es um binäre Zielvariablen, d. h. solche Zielvariablen welche mit den Werten 0 oder 1 kodiert werden können. Potentiell gehören alle Daten mit zwei Gruppen in diese Sparte von Problemen. Es kann um die Ausprägung eines zweistufigen Merkmals gehen, um ausgefallene und noch funktionierende Maschinen, kranke und gesunde Menschen oder Tiere, um Lebende und Tote, um das Auftreten von Fehlern oder das Vorhandensein eines Merkmals.

**Übung 1.1** *Versuchen Sie sich einige konkrete Beispiele vorzustellen, wo ein Einsatz der logistischen Regression sinnvoll wäre.*

**Lösung 1.1** *Sie können sich sicher beliebig viele Beispiele zurechtlegen.*

Warum brauchen wir nun eine spezialisierte Methode? Wieso können wir nicht einfach unsere klassische Regression benutzen? Eine Regression mit einer Zielvariablen, die nur die Werte 0 oder 1 annehmen kann, kann als Schätzwerte "unmögliche" Werte kleiner 0 oder grösser 1 liefern. Keine der bisherigen Transformationen kann hier Abhilfe schaffen. Ein weiterer Grund ist, dass ein Modell, das die richtigen Annahmen über die Daten trifft, befriedigender ist.

**Übung 1.2** *Was sind die Voraussetzungen für eine Diskriminanzanalyse (eine gewöhnliche multiple Regression, bei der die Zielvariable nur die Werte 0 und 1 umfasst), welche Einschränkungen hat man dort? Was können wir mit der logistischen Regression gewinnen?*

**Lösung 1.2** *In der Diskriminanzanalyse müssen die Residuen normalverteilt sein und es können keine Interaktionen oder Funktionen (wie Quadrat, sin, etc.) von Variablen in die Gleichung aufgenommen werden.*

Es gibt zwei Hauptanwendungen für die logistische Regression:

**Zusammenhänge.** Wir können uns fragen, was für erklärende Variablen (z. B. kontinuierliche und kategoriale Baseline Charakteristika in einer medizinischen Untersuchung, wie Blutdruck oder Geschlecht) einen Einfluss auf die Wahrscheinlichkeit haben, auf eine Behandlung anzusprechen. Zudem können wir uns fragen wie dieser Einfluss aussieht und wie stark er ist.

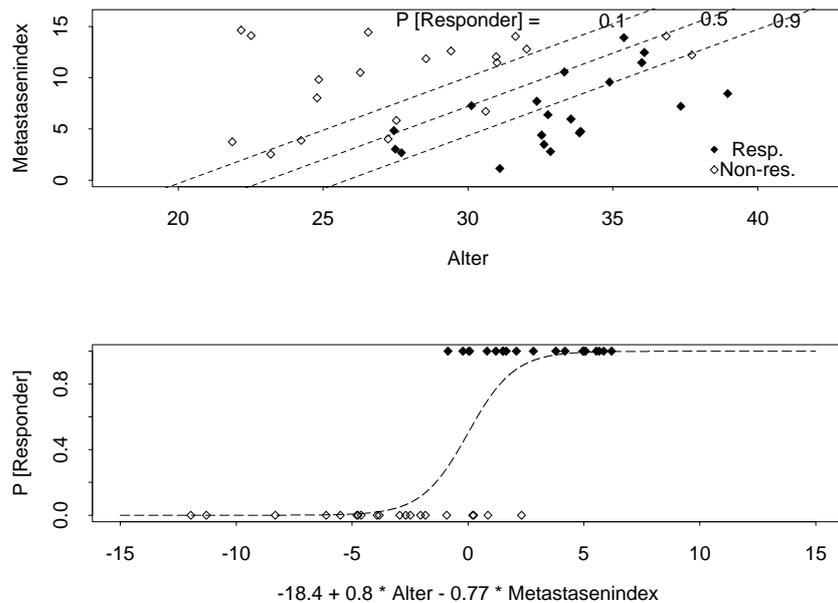


Abbildung 1: Responder in Abhängigkeit des Alters des Patienten und eines Metastasenindex (oben). Die Wahrscheinlichkeiten sind das Resultat der logistischen Regression. Unten: Die Daten und die angepasste Kurve entlang der Achse, welche durch die Parameter der logistischen Regression gegeben wird (dies entspricht einem Profil quer zu den Wahrscheinlichkeitslinien im oberen Teil).

**Klassifizierung.** Bei dieser Art von Fragestellung möchte man anhand einer Serie von kontinuierlichen und kategoriellen Variablen (erklärende Variablen) die Gruppenzugehörigkeit bestimmen. Z. B. können wir uns fragen, ob ein Patient mit bestimmten Eigenschaften auf eine medizinische Behandlung ansprechen wird.

## 2 Ein (simuliertes) Beispiel

Nehmen wir an, dass wir Responder in einer Krebstherapie untersuchen (wir nehmen ebenfalls an, dass die Art wie wir diese Responder definieren, vorgegeben ist). Dazu nehmen wir Daten zu zwei Variablen auf, von denen wir glauben, dass sie einen Einfluss darauf haben, ob jemand ein Responder wird: ein Metastasenindex (der den Zustand des Patienten an der Baseline definiert) und das Alter des Patienten. Responder haben wir mit 1, Non-responder mit 0 codiert (Abb. 1).

Bereits in dieser Graphik sehen wir, dass wir die Responder bei den älteren Patienten mit einem tiefen Metastasenindex finden.

## 3 Modellwahl

Wie wir bereits bei den linearen Modellen (Varianzanalyse und Regression) gesehen haben, müssen wir einen Weg finden, um zu entscheiden, welche erklärende Variablen wir in einem statistischen Modell berücksichtigen wollen.

In der logistischen Regression können wir grundsätzlich die gleichen Verfahren verwenden. Wir können von einem Modell ausgehen, welches alle uns interessierenden (vorhandenen) Vari-

ablen enthält, und Schritt für Schritt jene Variable ausschliessen, welche den höchsten p-Wert hat (stepwise backward). D. h. wir schliessen jene Variablen aus, welche in statistischem Sinne am wenigsten hergeben. Als Alternative können wir ein grosses Modell aus einem kleinen aufbauen (stepwise forward) oder (eine Teilmenge) aller möglichen Modelle untersuchen (allsubset Analyse). Auch hier werden wir feststellen, dass wir nicht ein einziges richtiges Modell finden, sondern eine Serie von Modellen, die alle eine ähnliche Qualität erreichen und somit (von der Statistik her) gleichwertig sind.

## 4 Schätzungen

In unserem Modell möchten wir als Zielvariable die Wahrscheinlichkeit dafür, dass wir eine 1 beobachten, in Abhängigkeit einer Funktion  $h$  der erklärenden Variablen setzen. Wir benutzen folgenden Ansatz:

$$\begin{aligned} P[Y = 1] &= h(x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}) \\ &= \tilde{h}(\alpha + \beta_1 x^{(1)} + \beta_2 x^{(2)} + \dots + \beta_m x^{(m)}), \text{ wobei} \end{aligned}$$

$$\tilde{h}(\eta) = \frac{e^\eta}{1 + e^\eta} = \frac{1}{1 + e^{-\eta}} \begin{cases} \xrightarrow{\eta \rightarrow \infty} 1 \\ \xrightarrow{\eta \rightarrow -\infty} 0 \end{cases}$$

Umgekehrt lässt sich schreiben:

$$g(P[Y = 1]) = \log\left(\frac{P[Y = 1]}{1 - P[Y = 1]}\right) = \alpha + \beta_1 x^{(1)} + \beta_2 x^{(2)} + \dots + \beta_m x^{(m)}$$

Die Funktion  $g$  wird “logit”-Funktion genannt und die Verwandtschaft mit der bisherigen Regression lässt sich an der rechten Seite der Gleichung leicht erkennen (wo wir wiederum eine Linearkombination unserer erklärenden Variablen finden). Die logit Funktion beschreibt auch gerade die sogenannten log-Odds, d. h. der Logarithmus des Quotienten aus der Wahrscheinlichkeit, dass wir einen Responder finden, geteilt durch die Wahrscheinlichkeit, dass wir einen Non-responder finden. Die Odds können dann durch  $e^{\alpha + \beta_1 x^{(1)} + \dots}$  geschätzt werden (dies brauchen wir weiter unten).

Um unsere Parameter zu schätzen benutzen wir die Maximum-Log-Likelihood Methode; dazu bauen wir zuerst eine Funktion, die die Wahrscheinlichkeit unserer Daten beschreiben soll:

$$LL = \log\left(\prod_{y_j=1} P[Y_j = 1] \prod_{y_j=0} (1 - P[Y_j = 1])\right) = \sum_{y_j=1} \log(\pi_j) + \sum_{y_j=0} \log(1 - \pi_j), \pi_j = \tilde{h}(\dots)$$

In Worten ausgedrückt, multiplizieren wir hier zuerst die Wahrscheinlichkeiten einer Serie von 0en und 1en, womit wir die Wahrscheinlichkeit der Sequenz unserer Zielvariablen berechnen. Von dieser Grösse ziehen wir den Logarithmus, was uns zu einfacheren Rechenoperationen führt (Summen statt Produkte).

Im Falle der logistischen Regression ist diese Gleichung noch etwas komplizierter, weil die Werte  $\pi_j$  für  $\tilde{h}(\eta)$  stehen. Diese Funktion wird nun für unsere vorgegebenen X-Werte maximiert und somit finden wir diejenigen Parameter unter denen unser Modell am besten zu den Daten passt. Diese Maximierung kann nicht analytisch gelöst werden und es kommen iterative numerische Verfahren zum Einsatz.

Wir erhalten aus diesem Verfahren Schätzgrössen anhand derer wir unser Modell interpretieren können. Wir betrachten sinnvollerweise die  $\beta_s$ , einen Standardfehler für die  $\beta_s$  und

$e^\beta$ . Diese Grössen sehen für unser Beispiel und das Modell mit beiden Variablen aber ohne Interaktion folgendermassen aus (vgl. Tabelle 1):

	$\beta$	s. e.	$e^\beta$
Alter	0.78	0.28	2.18
Metastasenindex	-0.77	0.29	0.46
Konstante ( $\alpha$ )	-18.40	6.75	

Die Konstante interessiert uns im Normalfall nicht und ist der Vollständigkeit halber aufgeführt (trotzdem sollte man sie zur Schätzung im Modell haben).

Anhand der  $\beta$ s können wir die Richtung des Zusammenhanges zwischen der geschätzten Wahrscheinlichkeit und unseren erklärenden Variablen sehen: zunehmendes Alter führt zu einer Erhöhung und zunehmender Metastasenindex zu einer Erniedrigung der Wahrscheinlichkeit ein Responder zu werden. Die naheliegende Interpretation ist, dass eine schwerwiegende Krankheit schwieriger zu heilen ist, und dass die Erkrankung je nach Alter eine andere ist. Junge Menschen haben eine schwerwiegendere Form der Krankheit oder eine die auf das geprüfte Medikament nicht gut anspricht.

Die Standardfehler geben einen ersten Hinweis darauf, wie wichtig die Variablen in statistischem Sinne sind. Variablen, bei welchen die Standardfehler relativ klein sind im Vergleich zur Schätzung der  $\beta$ s werden wichtig sein (d. h. die  $\beta$ s unterscheiden sich signifikant von 0). Ausserdem benötigen wir die Standardfehler, um Konfidenzintervalle zu berechnen.

**Einschub: Odds** Ein Odd ist der Quotient der Wahrscheinlichkeit eines Ereignisses geteilt durch die Wahrscheinlichkeit des Nicht-Eintreffens des Ereignisses. Somit ist das Odd einen Kopf bei einem Münzwurf zu erhalten gleich 1 (0.5/0.5) und eine sechs zu Würfeln gleich 0.2 (1/6 geteilt durch 5/6).

Die  $e^\beta$  lassen sich als Odds-ratios interpretieren und reflektieren somit die medizinische Relevanz. In unserem Falle heisst das, dass ein um ein Jahr älterer Patient 2.18 mal höhere Odds hat, ein Responder zu werden und eine Erhöhung um eine Einheit im Metastasenindex zieht eine Änderung der Odds, ein Responder zu werden, um den Faktor 0.46 nach sich.

Wie kommt man darauf? Schauen wir uns ein Beispiel an, das zwei Patienten vergleicht, welche ein Jahr Altersunterschied haben, aber in allen anderen Charakteristika übereinstimmen:

$$\begin{aligned} \text{Odds ratio} &= \frac{\text{Odds Patient alt}}{\text{Odds Patient jung}} = \frac{e^{\alpha + \beta_1 \cdot (\text{alter} + 1) + \beta_2 \cdot \text{metast}}}{e^{\alpha + \beta_1 \cdot \text{alter} + \beta_2 \cdot \text{metast}}} \\ &= \frac{e^\alpha \cdot e^{\beta_1 \cdot \text{alter}} \cdot e^{\beta_1 \cdot 1} \cdot e^{\beta_2 \cdot \text{metast}}}{e^\alpha \cdot e^{\beta_1 \cdot \text{alter}} \cdot e^{\beta_2 \cdot \text{metast}}} = e^{\beta_1} \end{aligned}$$

Schauen wir uns die Sache mit der Interpretation eines Odds-Ratios nochmals anhand zweier kategorischer Variablen an. Nehmen wir an, wir würden unser Modell der Responder mit den Variablen Geschlecht (0 = Männer, 1 = Frauen) und der Variable 'genereller Gesundheitszustand' (1 = schlecht, 2 = mittel, 3 = gut) erweitern. Zusätzlich bekämen wir also Resultate zu diesen Variablen:

	$\beta$	s. e.	$e^\beta$
Geschlecht (Frauen)	1.44	0.55	4.22
Gesundheitszustand			
schlecht	-3.01	0.87	0.05
mittel	-1.62	0.56	0.20

Dies bedeutet, dass die Odds der Frauen, Responder zu werden, 4.22 mal grösser sind als diejenigen der Männer. (Das heisst, wenn wir eine Frau und einen Mann mit den gleichen

Charakteristika vergleichen, sind die Chancen der Frau 4.22 mal grösser zu den Respondern zu gehören.)

Wenn wir mehr als zwei Kategorien haben, wird eine (wählbare) Kategorie als Vergleichskategorie gesetzt und die geschätzten Parameter beziehen sich auf diesen Vergleich. Im konkreten Fall heisst das, dass die Odds von Patienten mit mittlerem und schlechtem Gesundheitszustand um den Faktor 0.2, resp. um 0.05 verkleinert sind im Vergleich zu gutem Zustand.

Oft werden zu den  $e^\beta$  (also zu den Odds ratios) noch die 95% Konfidenzintervalle angegeben, was noch genauere Schlüsse zur Relevanz zulässt (vgl. Kapitel über Konfidenzintervalle).

Die Wahrscheinlichkeit, dass eine 30-jährige Frau mit Metastasenindex 7 und mittlerem Gesundheitszustand zu einem Responder wird, berechnet sich wie folgt:

$$P[Y = 1] = \frac{1}{1 + e^{-(-18.40 + 0.78 \cdot 30 - 0.77 \cdot 7 + 1.44 - 1.62)}} = 0.361$$

Für einen Mann mit den gleichen Merkmalen:

$$P[Y = 1] = \frac{1}{1 + e^{-(-18.40 + 0.78 \cdot 30 - 0.77 \cdot 7 + 0 - 1.62)}} = 0.118$$

Die Odds-Ration als:

$$\frac{0.361/0.639}{0.118/0.882} = 4.22$$

Bei der Interpretation heisst es wie immer aufzupassen. Was wir mit statistischen Methoden nachweisen können sind Korrelationen und Koinzidenz verschiedener Variablen. Auf kausale Zusammenhänge lässt sich nur durch ein sauberes Versuchsdesign schliessen.

## 5 Tests

Nun möchten wir natürlich noch berechnen, ob das Modell als Ganzes und die einzelnen Variablen signifikant zur Beschreibung der Daten beitragen. Für ersteres vergleichen wir das Gesamtmodell mit einem Minimalmodell (das besagt, dass wir unabhängig von den gemessenen Parametern eine konstante Wahrscheinlichkeit für einen Responder finden) und testen dann ob die einzelnen Parameter signifikant verschieden von Null sind.

Die Tests, die man hier anwenden kann sind sogenannte Likelihood-Ratio-Tests, d. h. man bildet Quotienten aus den Likelihoods oder Differenzen der Log-Likelihoods von einem kleinen (mit wenigen erklärenden Variablen) und einem grossen Modell (mit vielen erklärenden Variablen), welches das kleine Modell umfasst. Man kann zeigen, dass die doppelte Differenz  $\chi^2$  verteilt ist mit der Anzahl Freiheitsgraden, die der Differenz der Freiheitsgrade der Modelle entspricht. Da man die doppelten Log-Likelihoods braucht, werden von den Statistik-Programmen meist Devianzen ausgedrückt, die genau solche doppelten Log-Likelihoods sind.

Was tun wir nun konkret: In Tabelle 1 sehen wir die relevante Information einer logistischen Regression mit S-Plus. Wir sehen neben der Schätzung für die Grösse der Parameter (in der ersten Spalte) eine Schätzung deren Varianz und einen zugehörigen t-Wert. Des weiteren sehen wir die Nulldevianz, die dem kleinstmöglichen Modell entspricht, und die zum spezifischen Modell gehörende (Residual-) Devianz. Das kleinstmögliche Modell ist jenes, für das wir annehmen, dass alle unsere Variablen keinen Einfluss haben, und es somit eine konstante Wahrscheinlichkeit für die Beobachtung einer 1 gibt. Wir können nun zuerst testen, ob die Modelle überhaupt etwas erklären, indem wir jeweils die Residuen-Devianz von der Null-Devianz abziehen und den Wert mit der entsprechenden  $\chi^2_{fg}$  Verteilung mit  $fg$  Freiheitsgraden vergleichen:

Tabelle 1: Relevanter Computer-output von S-Plus zur logistischen Regression der Responderdaten (alter: Alter, metast: Metastasenindex)

**volles Modell:** mit Interaktion

	Value	Std. Error	t value
(Intercept)	-30.47725573	16.5371528	-1.8429567
alter	1.20699345	0.5915404	2.0404243
metast	0.70430774	1.5859311	0.4440973
alter:metast	-0.04742437	0.0521893	-0.9086989

Null Deviance: 55.45177 on 39 degrees of freedom  
 Residual Deviance: 20.87178 on 36 degrees of freedom

**volles Modell:** ohne Interaktion

	Value	Std. Error	t value
(Intercept)	-18.3956000	6.7485490	-2.725860
alter	0.7977799	0.2800214	2.848996
metast	-0.7688496	0.2906027	-2.645707

Null Deviance: 55.45177 on 39 degrees of freedom  
 Residual Deviance: 21.83712 on 37 degrees of freedom

**Modell nur mit einer Variablen:** Metastasenindex

	Value	Std. Error	t value
(Intercept)	1.8196151	0.81830766	2.223632
metast	-0.2194826	0.08961215	-2.449251

Null Deviance: 55.45177 on 39 degrees of freedom  
 Residual Deviance: 48.50085 on 38 degrees of freedom

**Modell nur mit einer Variablen:** Alter

	Value	Std. Error	t value
(Intercept)	-9.8019397	3.291214	-2.978214
alter	0.3200286	0.106097	3.016377

Null Deviance: 55.45177 on 39 degrees of freedom  
 Residual Deviance: 41.6674 on 38 degrees of freedom

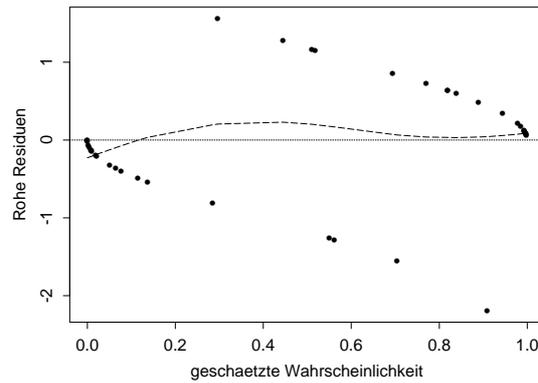


Abbildung 2: Rohe Residuen der logistischen Regression (gestrichelt: ein loess-Glätter).

volles Modell mit Interaktionen	$P[\chi_{fg=39-36}^2 \geq 55.45 - 20.87] < 0.0001$
volles Modell ohne Interaktionen	$P[\chi_{fg=39-37}^2 \geq 55.45 - 21.84] < 0.0001$
Modell mit nur einer Variablen: Metastasenin.	$P[\chi_{fg=39-38}^2 \geq 55.45 - 48.50] = 0.0084$
Modell mit nur einer Variablen: Alter	$P[\chi_{fg=39-38}^2 \geq 55.45 - 41.67] = 0.0002$

Wir sehen an den Signifikanzen, dass alle Modelle mehr bringen, als wenn wir eine konstante Wahrscheinlichkeit annehmen würden. Um zu sehen, welches Modell nun genügt, müssen wir zwischen den Modellen vergleichen:

ohne Interaktion versus mit Interaktion	$P[\chi_{fg=37-36}^2 \geq 21.84 - 20.87] = 0.326$
nur Metastasenin. versus mit beiden Variablen	$P[\chi_{fg=38-37}^2 \geq 48.50 - 21.84] < 0.0001$
nur Alter versus mit beiden Variablen	$P[\chi_{fg=38-37}^2 \geq 41.67 - 21.84] < 0.0001$

Die Nicht-Signifikanz des ersten Testes sagt uns, dass es nichts bringt, wenn wir noch die Interaktion ins Modell nehmen. Die beiden andern Signifikanzen zeigen jedoch, dass es beide der Variablen braucht. Diese Information kann man (asymptotisch) auch an den t-Werten in Tabelle 1 sehen (der Wald-Test ist ein anderer häufig benutzter asymptotischer Test).

## 6 Residuen

Auch in diesen Modellen sollten wir die Residuen anschauen, aber es ist nicht mehr so klar, was wir erwarten sollten. Auch gibt es keine durch Diskussion ausgereiften und implementierten Methoden. Ein wichtiger Plot ist noch immer der Tukey-Anscombe-Plot, bei dem wir die Residuen versus den geschätzten Wert auftragen. Um zu sehen, ob der Erwartungswert ungefähr null beträgt, sollte man unbedingt einen Glätter einzeichnen, da Artefakte in diesen Plots entstehen: die Residuen liegen auf zwei Geraden (Abb. 2).

## 7 Cut-off value

Um eine Entscheidung zu treffen, ob wir einen Responder oder einen Non-responder vor uns haben, müssen wir in den geschätzten Wahrscheinlichkeiten einen Cut-off value setzen. D. h. wir müssen uns für einen Wert entscheiden, wo wir die Überschreitung vom Non-responder zum Responder erwarten. In anderen Worten, wenn ein neuer Patient eine vom Modell geschätzte

Wahrscheinlichkeit über dem Cut-off value erhält, denken wir, er wird ein Responder. Dieser Wert wird oft bei 0.5 gesetzt, ist aber grundsätzlich beliebig wählbar. Wir können uns vorstellen, dass wir ein sehr potentes Heilmittel haben, bei dem durchschnittlich 80% der Patienten geheilt werden. Es gibt aber vielleicht einen Unterschied zwischen Männern und Frauen. Wenn wir den Cut-off value bei 0.5 belassen, dann prognostizieren wir für jeden Patienten einen Responder. Wenn wir den Cut-off value jedoch auf 0.8 setzten, dann fallen z. B. die Männer mehrheitlich unter den Cut-off und die Frauen mehrheitlich darüber. So sehen wir den Geschlechtsunterschied auch wenn die Responderrate an sich hoch ist.

## 8 Literaturhinweise

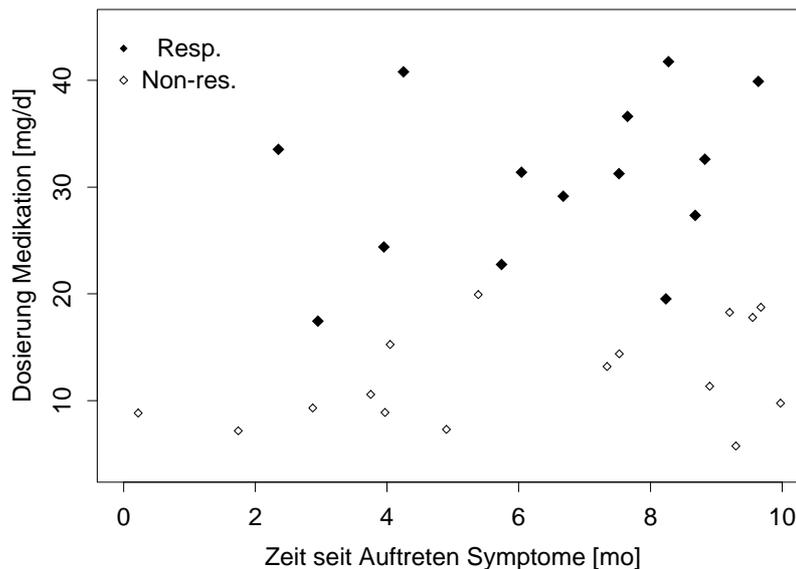
Weitere Informationen über Generalisierte Lineare Modelle finden sich in McCullagh and Nelder (1989) und speziell im Bezug auf S-plus in Venables and Ripley (1997):

McCullagh, P. and Nelder, J. A. (1989). *Generalized Linear Models*. Chapman and Hall, London.

Venables, W. N. and Ripley, B. D. (1997). *Modern Applied Statistics with S-PLUS*. Springer, New York, second edition.

## 9 Eine Übung

**Übung 9.1** Wir untersuchen erneut die Responder in einer klinischen Studie. Dabei betrachten wir die Variablen Dosierung der Medikation (*dos*) und die Zeit seit Auftreten der Symptome (*zeit*) als mögliche Prädiktoren für einen Behandlungserfolg. Die Daten sind in der folgenden Graphik dargestellt. Zeichnen Sie die Linie ein, wo Sie in etwa erwarten, dass die Wahrscheinlichkeiten einen Responder oder Non-responder zu sehen gleich sind. Was haben Sie nun für eine Erwartungen an die logistische Regression?



Hier folgt nun eine Auflistung des relevanten Computer-outputs. Welche Modelle und welche Variablen sind signifikant? Wie lässt sich das Interpretieren? Wurden Ihre Erwartungen von oben erfüllt?

mit Interaktion

	Value	Std. Error	t value
(Intercept)	-7.09179695	9.3122856	-0.7615528
zeit	-2.12331232	3.1251751	-0.6794219
dos	0.50733209	0.5764607	0.8800809
zeit:dos	0.08855432	0.1687055	0.5249049

Null Deviance: 41.4554 on 29 degrees of freedom  
Residual Deviance: 6.762643 on 26 degrees of freedom

ohne Interaktion

	Value	Std. Error	t value
(Intercept)	-12.3859016	8.6651672	-1.429390
zeit	-0.5137501	0.4054912	-1.266982
dos	0.8137422	0.4976112	1.635297

Null Deviance: 41.4554 on 29 degrees of freedom  
Residual Deviance: 7.019384 on 27 degrees of freedom

nur mit zeit

	Value	Std. Error	t value
(Intercept)	-0.42159164	0.9302550	-0.4532001
zeit	0.04560032	0.1350489	0.3376578

Null Deviance: 41.4554 on 29 degrees of freedom  
 Residual Deviance: 41.34084 on 28 degrees of freedom

nur mit dos

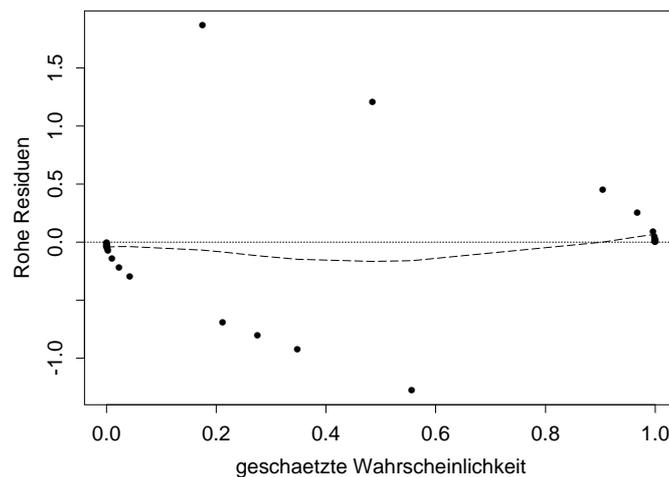
	Value	Std. Error	t value
(Intercept)	-14.0025294	8.0586952	-1.737568
dos	0.7135532	0.4233168	1.685625

Null Deviance: 41.4554 on 29 degrees of freedom  
 Residual Deviance: 8.979398 on 28 degrees of freedom

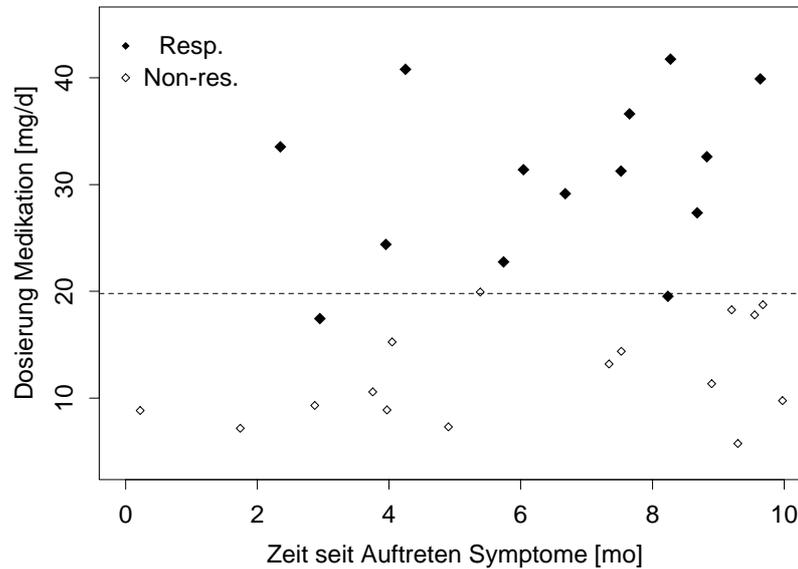
Einige kritische Werte der  $\chi^2_{fg}$ -Verteilung:

$fg$	$p = 0.05$	$p = 0.01$	$p = 0.001$
1	3.84	6.63	10.83
2	5.99	9.21	13.82
3	7.82	11.35	16.27

Was sagen Sie zu den Residuen?



**Lösung 9.1** Die Responder sind insbesondere im Bereich der hohen Dosen zu finden. Die Dauer der Symptomatik scheint keinen grossen Einfluss zu haben. Die Linie gleicher Wahrscheinlichkeit für einen Responder und einen Non-responder liegt bei etwa 20:



Es ergeben sich folgende Signifikanzen für die einzelnen logistischen Regressionen:

volles Modell mit Interaktionen	$P[\chi_{fg=29-26}^2 \geq 41.46 - 6.76] < 0.001$
volles Modell ohne Interaktionen	$P[\chi_{fg=29-27}^2 \geq 41.46 - 7.02] < 0.001$
Modell nur mit Zeit	$P[\chi_{fg=29-28}^2 \geq 41.46 - 41.34] > 0.05$
Modell nur mit Dosierung	$P[\chi_{fg=29-28}^2 \geq 41.46 - 8.98] < 0.001$

Alle Modelle sind signifikant ausser demjenigen, das nur die Variable Zeit enthält. Ein erster Hinweis darauf, dass diese Variable wohl keinen grossen Einfluss hat.

ohne Interaktion versus mit Interaktion	$P[\chi_{fg=27-26}^2 \geq 7.02 - 6.76] > 0.05$
nur Zeit versus mit beiden Variablen	$P[\chi_{fg=28-27}^2 \geq 41.34 - 7.02] < 0.001$
nur Dosierung versus mit beiden Variablen	$P[\chi_{fg=28-27}^2 \geq 8.98 - 7.02] > 0.05$

Wir sehen, dass es keine Verbesserung bringt, wenn wir zur Variablen Dosierung noch die Variable Zeit dazu nehmen. Offensichtlich unterscheiden sich die Responder nur dadurch von den Non-respondern, dass sie mit höheren Dosierungen behandelt wurden. Die Dosierung erreicht eine Odds-Ratio von  $e^{0.0456} = 1.047$  per mg, resp. eine Odds-ration von  $e^{10 \cdot 0.0456} = 1.578$  per 10 mg Dosissteigerung.

Die Residuen sehen recht gut aus (wie es sich für simulierte Daten gehört), denn der Glätter ist in seiner ganzen Länge sehr nahe an Null.