

1. Konferenz der SAS-Benutzer in Forschung und Entwicklung (KSFE)

an der

Humboldt-Universität zu Berlin

20./21. Februar 1997

Erklärte Variabilität bei Logistischer Regression unter Verwendung von SAS

Martina MITTLBÖCK und Michael SCHEMPER

Abteilung für Klinische Biometrie

Institut für Medizinische Computerwissenschaften der Universität Wien

Adresse: Spitalgasse 23, A-1090 Wien, Österreich

Fax: + 43 / 1 / 40400 / 6687

E-mail: Martina.Mittlboeck@akh-wien.ac.at

Zusammenfassung

Mehrere Maße für den Anteil an erklärter Variabilität werden in der Literatur vorgeschlagen. Im SAS werden für die Logistische Regression (PROC LOGISTIC) standardmäßig mehrere nichtparametrische Übereinstimmungsmaße zwischen den beobachteten Werten und den geschätzten Wahrscheinlichkeiten ausgegeben. Auf Anforderung werden auch 2 Bestimmtheitsmaße (R^2) ausgegeben. Wir haben die von SAS verwendeten und weitere in der Literatur empfohlenen Maße untersucht und ihre Eigenschaften näher betrachtet.

1. Einführung

Der quadrierte multiple Korrelationskoeffizient R^2 , auch Bestimmtheitsmaß genannt, ist ein aus der klassischen Regressionsanalyse bekanntes Maß. R^2 wird als Anteil der Variabilität der abhängigen Variable, die durch unabhängige Variablen im Modell erklärt werden können, definiert.

Für das Allgemeine Lineare Modell existiert nur ein geeignetes Maß, das multiple R^2 , um erklärte Variabilität zu quantifizieren; es gibt jedoch mehrere äquivalente Definitionen¹ dafür. Wendet man diese Definitionen auf die Logistische Regression an, so unterscheiden sich die Ergebnisse manchmal gravierend. Weiters wurden in der Literatur auch R^2 -Maße empfohlen, die nur auf die Logistische Regression anwendbar sind.

Im folgenden werden mehrere Maße vorgestellt und ihre Eigenschaften diskutiert.

2. Systematische Präsentation der Maße

Die vorgestellten Maße werden in 3 Gruppen zusammengefaßt: 1) quadrierte Korrelation zwischen beobachteten Werten und geschätzten Wahrscheinlichkeiten; 2) proportionale Reduktion der Streuung und 3) Maße basierend auf der Likelihoodfunktion.

Seien (y_i, x_i) , $i=1, \dots, n$, Beobachtungspaare, wobei $y_i=0$ oder 1 die abhängige Variable und x_i der entsprechende Kovariablenvektor darstellen. Die Schätzung der Logistischen Regression ist $\text{Prob}(y_i=1|x_i) = \hat{p}_i = \exp(\hat{\beta}x_i) / [1 + \exp(\hat{\beta}x_i)]$, wobei $\hat{\beta}$ der Vektor mit Parameterschätzungen ist, und $\text{Prob}(y_i = 1) = \bar{p} = \sum_i y_i / n$.

2.1. R^2 -Maße, basierend auf der quadrierten Korrelation von y und \hat{p}

i) *Quadrierte Pearson Korrelation* (r^2)

$$r = \frac{\sum_i (y_i - \bar{p})(\hat{p}_i - \bar{p})}{\sqrt{\sum_i (y_i - \bar{p})^2 \sum_i (\hat{p}_i - \bar{p})^2}}$$

ii) *Quadrierte Spearman Korrelation* (r_s^2)

$$r_s = \frac{\sum_i (R(y_i) - \bar{R})(R(\hat{p}_i) - \bar{R})}{\sqrt{\sum_i (R(y_i) - \bar{R})^2 \sum_i (R(\hat{p}_i) - \bar{R})^2}}$$

wobei $R(z)$ der Rang von z und $\bar{R} = (n+1)/2$ ist.

iii) *Quadriertes Kendall's τ_a* (τ_a^2)

$$\tau_a = \frac{\sum_{i < j} \text{sign}(y_j - y_i) \text{sign}(\hat{p}_j - \hat{p}_i)}{[n(n-1)/2]} \quad \text{mit } \text{sign}(z) = \begin{cases} 1 & \text{wenn } z > 0 \\ 0 & \text{wenn } z = 0 \\ -1 & \text{wenn } z < 0 \end{cases}$$

iv) *Quadriertes Kendall's* τ_b (τ_b^2)

$$\tau_b = \frac{\sum_{i < j} \text{sign}(y_j - y_i) \text{sign}(\hat{p}_j - \hat{p}_i)}{\sqrt{\sum_{i < j} \text{sign}^2(y_j - y_i) \sum_{i < j} \text{sign}^2(\hat{p}_j - \hat{p}_i)}}$$

v) *Quadriertes Somers' D* $D_{\hat{p}y}$ ($D_{\hat{p}y}^2$)

$$D_{\hat{p}y} = \frac{\sum_{i < j} \text{sign}(y_j - y_i) \text{sign}(\hat{p}_j - \hat{p}_i)}{\sum_{i < j} \text{sign}^2(y_j - y_i)}$$

vi) *Quadriertes Goodman und Kruskal's* γ (γ^2)

$$\gamma = \frac{\sum_{i < j} \text{sign}(y_j - y_i) \text{sign}(\hat{p}_j - \hat{p}_i)}{\left[\sum_{i < j} \text{sign}^2(y_j - y_i) \text{sign}^2(\hat{p}_j - \hat{p}_i) \right]}$$

r , r_s , τ_b kann man im SAS² unquadriert berechnen, indem man die beobachteten Werte (y) und die entsprechenden geschätzten Wahrscheinlichkeiten (\hat{p}) mit OUTPUT in eine Datei schreibt und mit PROC CORR berechnet. Die Maße τ_a , $D_{\hat{p}y}$ und γ werden unquadriert automatisch von der PROC LOGISTIC als Maß für den Zusammenhang von geschätzten Wahrscheinlichkeiten und beobachteten Responses ausgegeben.

Manche diese Maße haben offensichtliche, unerwünschte Eigenschaften, z.B. Kendall's τ_a kann nicht einmal bei perfekter Vorhersage den geforderten Wert von eins erreichen. Hingegen erreicht Goodman und Kruskal's γ den Wert eins, auch wenn die beobachteten Werte nicht vollständig bestimmbar sind.

Jedoch sind alle ordinalen Assoziationsmaße, wie r_s , τ_b , τ_a , $D_{\hat{p}y}$ und γ , für die Logistische Regression ungeeignet, da die Logistische Regression ein parametrisches Modell ist. Eine monotone, nichtlineare Transformation einer unabhängigen Variable kann den Fit des Modells und die Parameterschätzungen ändern, und sollte daher auch einen Einfluß auf ein gutes Maß für erklärte Variabilität haben; ordinale Maße bleiben jedoch von solchen Transformationen unbeeinflusst.

2.2. R²-Maße, basierend auf der proportionalen Reduktion der Streuung von y

Die allgemeine Form dieser Maße ist $PEV = \left[\sum_i D(y_i) - \sum_i D(y_i|x_i) \right] / \sum_i D(y_i)$, wobei $D(y_i)$ und $D(y_i|x_i)$ ein Maß für die Distanz zwischen y_i und einem unbedingten bzw. bedingten (auf ein Modell und Kovariablenvektor x_i) zentralen Lokationsparameter ist. Die vier Maße dieses Kapitels unterscheiden sich in ihrer Wahl für $D(y_i)$ und $D(y_i|x_i)$.

vii) *Sums-of-squares* $R^2(R_{SS}^2)$

R_{SS}^2 verwendet die quadrierten Abweichungen zwischen beobachteten Werten und geschätzten Wahrscheinlichkeiten: $D(y_i) = (y_i - \bar{p})^2$ und $D(y_i|x_i) = (y_i - \hat{p}_i)^2$.

viii) *Gini's Konzentrationsmaß* (R_G^2)

Für die Logistische Regression verwendet dieses Maß³ $D(y_i) = \bar{p}(1 - \bar{p})$ und $D(y_i|x_i) = \hat{p}_i(1 - \hat{p}_i)$, was der erwarteten Varianz unter dem Logistischen Modell entspricht. Durch die explizite Verwendung der binomialen Varianz für R_G^2 muß jedoch angenommen werden, daß das Modell korrekt geschätzt wird.

R_{SS}^2 verwendet die beobachteten Differenzen, während R_G^2 die unter dem Modell erwarteten Distanzen verwendet. R_{SS}^2 , R_G^2 und r^2 sind asymptotisch äquivalent.

ix) *Entropie* $R^2(R_E^2)$

Verwendet man die Entropie⁴ der binomialen Verteilung bzw. die Devianz Residuen, so gilt $D(y_i) = -[y_i \log \bar{p} + (1 - y_i) \log(1 - \bar{p})]$ und $D(y_i|x_i) = -[y_i \log \hat{p}_i + (1 - y_i) \log(1 - \hat{p}_i)]$, was in $\sum_i D(y_i) = -\log L(0)$ und $\sum_i D(y_i|x_i) = -\log L(\hat{\beta})$ resultiert. $L(\hat{\beta})$ und $L(0)$ entsprechen den Likelihoods der mit Kovariablen und ohne Kovariablen gefitteten Modelle. Daher entspricht R_E^2 der proportionalen Reduktion des maximierten log-Likelihood.

Die Berechnung von R_{SS}^2 , R_G^2 und R_E^2 in SAS ist nicht einfach.

2.3. R²-Maße, basierend auf der Likelihoodfunktion

x) *Likelihood-Ratio* $R^2(R_{LR}^2)$

Das Maß⁵ $R_{LR}^2 = 1 - \left[L(0) / L(\hat{\beta}) \right]^{2/n}$ ist im Allgemeinen Linearen Modell identisch mit dem multiplen R^2 . Es wurde vorgeschlagen, dieses Maß auch auf verallgemeinerte Lineare Modelle wie die Logistische Regression anzuwenden. R_{LR}^2 kann jedoch bei Logistischer Regression niemals den bei perfekter Vorhersage erwünschten Maximalwert von eins annehmen; z.B. mit $\bar{p}=0.5$ erreicht R_{LR}^2 sein Maximum bei 0.75.

xi) *Likelihood-Ratio* R^2 modifiziert (R_{CU}^2)

Da R_{LR}^2 einen Wert von eins nicht erreichen kann, wurde vorgeschlagen⁶, $R_{CU}^2 = R_{LR}^2 / \max R_{LR}^2$ zu verwenden, wobei $\max R_{LR}^2 = 1 - [L(0)]^{2/n}$ der Maximalwert ist, der durch R_{LR}^2 erreicht werden kann. Obwohl R_{CU}^2 den erwünschten Maximalwert von eins erreichen kann, ist damit noch nicht gesichert, daß die Skalierung dazwischen adäquat ist.

Die Maße R_{LR}^2 , R_{CU}^2 werden bei SAS in der PROC LOGISTIC (ab V6.10) auf Anfrage ausgegeben, wenn man die Option RSquare im MODEL-Statement angibt.

3. Empirischer Vergleich der Maße

Das Verhalten und die Eigenschaften der beschriebenen Maße wurden untersucht und in einer Stichprobe von $n=50.000$ verglichen, die entsprechend dem Logistischen Modell mit einer kontinuierlichen Kovariablen (x) erzeugt wurde. Es wurden aus der Gleichverteilung $(0,1)$ systematisch die Werte für x genommen, und die Werte der abhängigen Variable y (1 und 0) wurden mit den Wahrscheinlichkeiten $\text{Prob}(y_i=1|x_i)=\exp(b_0+b_1x_i)/[1+\exp(b_0+b_1x_i)]$ und

$1-\text{Prob}(y_i=1|x_i)$ generiert. Die Werte der Parameter b_1 und $b_0=-b_1/2$ wurden so gewählt, daß der gesamte Wertebereich für Maße der erklärten Variabilität geeignet abgedeckt wurde, und daß $\bar{p}=0.5$ in der zugrundeliegenden Population war.

Da Vergleiche der untersuchten Maße mit dem etablierten multiplen R^2 des Allgemeinen Linearen Modells von Interesse sind, wurden die Stichproben für $R^2=0$ und 12 so generiert, daß sowohl die Logistische Regression als auch das Allgemeine Lineare Modell anwendbar waren. Dafür wurde die untere und obere Grenze der erklärenden Variable so eingeschränkt, daß $0.2 < E(\hat{p}) < 0.8$. Es kann gezeigt werden⁷, daß die Logistische Funktion innerhalb dieser Grenzen fast linear ist, und daß daher beide Modelle anwendbar sind. Wenn der Range von \hat{p} und entsprechend von x eingeschränkt wird, können nur kleinere Werte für R^2 erzielt werden. Daher werden nur Ergebnisse für $R^2=0$ und 0.12 berichtet, wobei 0.12 dem maximalen R^2 bei einer kontinuierlichen Kovariable entspricht, bis zu dem auch das Allgemeine Lineare Modell anwendbar ist.

R_{GLM}^2	r^2	R_{SS}^2	R_{G}^2	R_{E}^2	r_s^2	τ_a^2	τ_b^2	$D_{\hat{p}y}^2$	γ^2	R_{LR}^2	R_{CU}^2
0	0	0	0	0	0	0	0	0	0	0	0
12	12	12	12	9	12	4	8	16	16	12	16
-	25	25	25	21	24	7	16	36	36	23	32
-	50	50	50	46	44	12	29	70	70	43	61
-	75	75	75	72	57	16	38	92	92	58	83
-	100	100	100	100	77	25	54	100	100	75	100

Tabelle I: Ergebnisse der R^2 -Maße für Logistische Regression mit kontinuierlicher Kovariable und $n=50.000$

Die erwünschten Eigenschaften einer Übereinstimmung mit R_{GLM}^2 wird nur bei r^2 , r_s^2 , R_G^2 und R_{SS}^2 beobachtet. Für hohe Werte an erklärter Variabilität sind sich r^2 , R_{SS}^2 und auch R_G^2 ähnlich, während mit wachsendem R^2 die Maße r_s^2 , τ_a^2 und τ_b^2 niedrigere Werte als r^2 zeigen. Für vollständig erklärte Variabilität, wenn y und \hat{p} fast identisch sind, nehmen r_s^2 , τ_b^2 und τ_a^2 zu niedrige Werte an. Grund dafür ist, daß nichtparametrische Assoziationsmaße verschiedene Ränge für fast identische Werte von \hat{p} (nahe 0 und 1) annehmen. Somers' $D_{\hat{p}y}$ -quadriert und Goodman-Kruskal's γ -quadriert nehmen viel höhere Werte als r^2 und auch R_{GLM}^2 an. R_E^2 nimmt im Vergleich zu R_{GLM}^2 zu niedrige Werte an. R_{LR}^2 nimmt ebenfalls zu niedrige Werte an und kann, wie bereits diskutiert, nicht einmal einen Wert von eins bei perfekter Vorhersage erreichen; die Werte von R_{CU}^2 sind im Vergleich zu R_{GLM}^2 zu hoch.

4. Anpassung der Maße für erklärte Variabilität für kleine Stichproben

Wenn die Zahl der Kovariablen k im allgemeinen linearen Modell im Verhältnis zur Stichprobengröße n groß ist, so hat $R_{adj}^2 = 1 - \frac{SSE/(n-k-1)}{SST/(n-1)}$ die wünschenswerte Eigenschaft, daß $E(R_{adj}^2) = 0$ für $R^2=0$ in der zugrundeliegenden Population. Die Kritik einer Inflation von R^2 bei kleinen Stichproben kann durch die Verwendung von R_{adj}^2 vermieden werden. Obwohl eine gründliche Untersuchung von R_{adj}^2 in Zusammenhang mit der Logistischen Regression ($R_{SS,adj}^2$) fehlt, wäre eine Verwendung von $R_{SS,adj}^2$ in Logistischen Modellen analog zum Linearen Modell denkbar.

Eine Simulationsstudie untersucht, wie adäquat $R_{SS,adj}^2$ für die Stichprobengrößen $n=50$, 100, 200 und 1000 mit $k=1$, 5 und 10 unabhängigen dichotomen Kovariablen ist. Die Stichproben wurden, wie im vorigen Kapitel beschrieben, generiert, aber mit mehreren Kovariablen mit identischen Koeffizienten ($b_1=\dots=b_k$). Die Werte wurden so gewählt, daß $R_{SS}^2=0.50$ in der zugrundeliegenden Population.

n	Anzahl der Kovariablen (k)											
	1				5				10			
	R_{SS}^2	$R_{SS,adj}^2$	R_E^2	$R_{E,adj}^2$	R_{SS}^2	$R_{SS,adj}^2$	R_E^2	$R_{E,adj}^2$	R_{SS}^2	$R_{SS,adj}^2$	R_E^2	$R_{E,adj}^2$
50	0.52	0.51	0.42	0.40	0.57	0.53	0.52	0.44	0.67	0.59	0.62	0.47
100	0.51	0.51	0.41	0.40	0.53	0.51	0.47	0.43	0.57	0.52	0.51	0.43
200	0.50	0.50	0.40	0.40	0.52	0.51	0.45	0.44	0.54	0.52	0.48	0.44
1000	0.50	0.50	0.40	0.40	0.51	0.50	0.44	0.44	0.51	0.51	0.45	0.44

Tabelle II: Schätzungen für R_{SS}^2 , $R_{SS,adj}^2$, R_E^2 und $R_{E,adj}^2$ für eine zugrundeliegende erklärte Variabilität von $R_{SS}^2 = 0.50$ bzw. $R_E^2 = 0.400$, 0.434 und 0.438 für $k=1, 5$ und 10 aus 500 simulierten Studien (Mediane).

Die Ergebnisse von Tabelle II unterstreichen den Vorzug von $R_{SS,adj}^2$ über R_{SS}^2 , wenn der Quotient k/n groß ist, obwohl es scheint, daß die Korrektur durch die Freiheitsgrade für SSE in $R_{SS,adj}^2$ möglicherweise noch nicht groß genug ist. Das erfordert aber Bestätigung durch größere systematische Simulationsstudien.

Auch für R_E^2 wurde eine Korrektur vorgeschlagen⁸: $R_{E,adj}^2 = 1 - \frac{\log L(\hat{\beta}) - (k+1)/2}{\log L(\beta_0) - 1/2}$, die sich zufriedenstellend in Tabelle II verhält. Die Korrektur $(k+1)/2$ wird durch die χ_{k+1}^2 -Verteilung von $2[\log L(\hat{\beta}) - \log L(\beta_0)]$ unter $H_0: \beta_1, \dots, \beta_k = 0$ motiviert, wobei $\hat{\beta}$ ein $(k+1)$ -dimensionaler Vektor von Maximum-Likelihood Schätzungen und β_0 der wahre Intercept-Parameter sind. Daher ist $E[\log L(\hat{\beta}) - \log L(\beta_0)] = (k+1)/2$ und der erwartete Optimismus von $\log L(\hat{\beta})$ ist $(k+1)/2$.

Keine Korrekturen sind bekannt für die anderen Maße für erklärte Variabilität.

Wenn adjustierte- R^2 Maße für ein Modell berechnet werden, das durch schrittweise Regression erhalten wurde, dann sollte man k als die Zahl der Kovariablen wählen und nicht die Faktoren, die letztendlich im Modell übrigbleiben⁹.

5. Abschließende Beurteilungen und sonstige Anmerkungen

Um R^2 -Maße geeignet beurteilen zu können, benötigen wir Kriterien, die gute Maße erfüllen sollten. Unserer Ansicht nach sollten Maße für erklärte Variabilität für Logistische Regression folgende Eigenschaften haben:

- 1) Konsistenz mit dem Charakter der Logistischen Regression, d.h. nichtlineare monotone Transformationen von erklärenden Variablen sollten die Maße beeinflussen und lineare Transformationen sollten sie nicht beeinflussen.
- 2) Der Wertebereich für die Maße sollte $[0,1]$ sein, null bei vollkommenem Fehlen von Vorhersagbarkeit und eins bei vollständiger Vorhersagbarkeit.
- 3) Maße für erklärte Variabilität sollten gleich oder ähnlich sein, wenn die Daten alternativ mit dem Allgemeinen Linearen Modell analysiert werden können.

Die Übereinstimmung der Maße mit den gegebenen Kriterien wird in Tabelle III dargestellt. Wir glauben, daß alle ordinalen Assoziationsmaße nicht in Übereinstimmung mit dem Charakter der Logistischen Regression sind und daß, von den restlichen Maßen, die das Maximum von eins erreichen, nur r^2 , R_{SS}^2 und R_G^2 mit R_{GLM}^2 übereinstimmen und daher als geeignet angesehen werden können.

Es wurde argumentiert¹⁰, daß im Gegensatz zur Verwendung im Allgemeinen Linearen Modell, R_{SS}^2 in der Logistischen Regression nicht durch die Modellanpassung optimiert wird. Manche Statistiker bevorzugen daher Maße, die auf der Likelihoodfunktion basieren, aufgrund ihrer Übereinstimmung mit dem ML-Fitting Prozeß. Die Logistische Regression kann aber auch mittels gewichteter Kleinst-Quadrate-Schätzungen berechnet werden¹¹; dieser Ansatz ist zwar optimal bezüglich der Effizienz, die ungleiche Gewichtung der Residuen beeinträchtigt hingegen die intuitive Interpretation eines darauf basierenden Maßes der

erklärten Variabilität¹². r^2 , R_{SS}^2 und R_G^2 gewichten alle Residuen gleich stark, sie sind daher nicht optimiert durch die ML-Schätzungen, aber die daraus resultierenden Maße beeinträchtigen nicht ihre intuitive Interpretation.

Maße	entspricht dem Charakter der Logistischen Regression	Range [0 - 1]	Konsistenz mit R_{GLM}^2
r^2	ja	ja	ja
R_{SS}^2	ja	ja	ja
R_G^2	ja	ja	ja
R_E^2	ja	ja	nein
r_s^2	nein	nein	ja
τ_a^2	nein	nein	nein
τ_b^2	nein	nein	nein
D_{py}^2	nein	ja	nein
γ^2	nein	ja	nein
R_{LR}^2	ja	nein	nein
R_{CU}^2	ja	ja	nein

Tabelle III: Zusammenfassung der Eigenschaften der Maße für erklärte Variabilität

Die von SAS für die Logistische Regression angebotenen R^2 -Maße (R_{LR}^2 und R_{CU}^2) und ordinalen Assoziationsmaße (τ_a^2 , D_{py}^2 und γ^2) sind unserer Ansicht nach keine gute Wahl; vor allem vor der Interpretation der automatisch berechneten Werte von τ_a^2 , D_{py}^2 und γ^2 möchten wir abraten. Wir haben deshalb ein SAS-Macro geschrieben, das R_{SS}^2 , $R_{SS,adj}^2$, R_E^2 und $R_{E,adj}^2$ berechnet und leicht verwendet werden kann. Dieses Macro kann über FTP vom Host VM.AKH-WIEN.AC.AT mit dem User BIOMETRY (kein Passwort) bezogen werden.

Literatur

1. Kvalseth, T. O. 'Cautionary Note about R^2 ', *American Statistician*, **39**, 279-285 (1985).
2. 'The LOGISTIC Procedure' *SAS/STAT User's Guide*, Version 6, 4th Edition, 1071-1126 (1990).
3. Haberman, S. J. 'Analysis of dispersion of multinomial responses', *Journal of the American Statistical Association*, **77**, 568-580 (1982).
4. Theil, H. 'On the estimation of relationships involving qualitative variables', *American Journal of Sociology*, **76**, 103-154 (1970).
5. Cox, D. R. and Snell, E. J. *Analysis of binary data*, Chapman and Hall, London, 1989.
6. Cragg, J. G. and Uhler, R. 'The demand for automobiles', *Canadian Journal of Economics*, **3**, 386-406 (1970).
7. Cox, D. R. and Wermuth, N. 'A comment on the coefficient of determination for binary responses', *The American Statistician*, **46**, 1-4 (1992).
8. Mittlböck, M. and Schemper, M. 'Explained Variation for Logistic Regression', *Statistics in Medicine*, **15**, 1987-97(1996).
9. Rencher, A. C. and Pun, F. C. 'Inflation of R^2 in best subset regression', *Technometrics*, **22**, 49-53 (1980).
10. Agresti, A. *Categorical Data Analysis*, (p. 112), Wiley, New York, (1990).
11. Hosmer, D. W. Jr. and Lemeshow, S. *Applied logistic regression*, Wiley & Sons, New York, 1989.
12. Willet, J. B. and Singer, J. D. 'Another cautionary note about R^2 : Its use in weighted least-squares regression analysis', *The American Statistician*, **42**, 236 - 238 (1988).