

Übereinstimmungsmaße in der PROC FREQ – Option AGREE bei nichtsymmetrischen Tabellen – eine Macro-Lösung

Bettina Danner, Rainer Muehe

Abteilung Biometrie und Medizinische Dokumentation, Universität Ulm

Telefon: 0731 / 50-26891

eMail: bettina.danner@medizin.uni-ulm.de

Abstract

In dem Beitrag werden Probleme bei der Berechnung von Übereinstimmungsmaßen (Cohens Kappa, McNemar-Test) beschrieben, die bei der Nutzung der Option AGREE der PROC FREQ (seit Version 6.11) auftreten. Angegeben werden die Lösungsmöglichkeiten und das zugehörige SAS-Macro, welches eine Berechnung auch in den beschriebenen Problemsituationen erlaubt.

Einleitung

Sowohl in Studien zum Vergleich diagnostischer Instrumente als auch in gematchten Fall-Kontroll-Studien werden Übereinstimmungsmaße und davon abgeleitete statistische Verfahren sehr häufig benutzt. Zu diesen statistischen Verfahren gehören Cohens Kappa und der McNemar-Test [1].

Seit der Version 6.11 können die beiden Verfahren mit der Option AGREE in der Prozedur PROC FREQ mit SAS berechnet werden [2]. Dabei müssen die abhängigen Beobachtungen (z.B. in der Diagnostik zwei Messungen an ein und demselben Patienten) in zwei Variablen gespeichert sein.

1. Problem: nicht-quadratische Häufigkeitstabelle

PROC FREQ unterdrückt nichtbesetzte Zeilen und Spalten in der Häufigkeitstabelle. Sie lassen sich nicht mit irgendeiner Option als Null-besetzt darstellen. In dieser dann nicht-symmetrischen Tabelle werden keine Übereinstimmungsmaße mehr berechnet.

2. Problem: inhaltlich nicht-quadratische Häufigkeitstabelle

Sollte die Tabelle zwar von der Form, aber durch nichtbesetzte Merkmalsausprägungen nicht inhaltlich symmetrisch sein, so berechnet die Prozedur die Kenngrößen, ohne zu prüfen, ob die Ausprägungen übereinstimmen. Man bekommt falsche, inhaltlich nicht interpretierbare Ergebnisse ohne irgendeine Warnung.

3. Problem: BY-Variablen

Bei der Verwendung einer BY-Variablen werden nur die quadratischen Tafeln berechnet, unabhängig davon, ob diese inhaltlich gleich sind oder nicht. Damit stellt dieses Problem eine Kombination der unter 1. und 2. genannten Situationen dar.

Lösungsstrategie

Nachfolgend wird die allgemeine Lösungsstrategie aufgezeigt und an einem Beispiel dargestellt. Da diese Strategie sehr zeitaufwendig sein kann, wurde die Macro-Lösung programmiert.

Folgende Schritte sind durchzuführen bzw. laufen im Macro ab:

1. Auslesen der Kreuztabelle in eine Out-Datei, falls in der Ausgangsdatei Einzelbeobachtungen und keine Häufigkeiten enthalten sind
2. Einlesen mindestens einer Zelle für die fehlende Zeile oder Spalte (für jede Ausprägung der BY-Variablen) mittels CARDS, wobei ein sehr kleines Gewicht, z. B $1/10^{20}$ (1E-20), für die Zelhäufigkeit vergeben wird
3. Zusammensetzen der Ausgangsdatei / der Out-Datei und der unter 2. erstellten Datei mit dem SET-Befehl
4. erneuter Aufruf des Tests für die neue Datei unter Verwendung des WEIGHT-Statements für die Zelhäufigkeiten. Die Tabelle ist nunmehr quadratisch und der Test wird gerechnet. Die kleinen Zelhäufigkeiten haben keinen Einfluß auf das Ergebnis des Tests, allerdings sind die Spalten- und Zeilenprozent unter Umständen nicht mehr korrekt.

Beispiel:

Quadratische, aber inhaltlich nicht symmetrische Tafel	Quadratische und inhaltlich symmetrische Tafel (nach Macro)																																																																																																																																																																										
<p>TABLE OF GRUPPE_A BY GRUPPE_B</p> <table border="1"> <thead> <tr> <th>GRUPPE_A</th> <th colspan="2">GRUPPE_B</th> <th>Total</th> </tr> <tr> <th>Frequency</th> <th></th> <th></th> <th></th> </tr> <tr> <th>Percent</th> <th></th> <th></th> <th></th> </tr> <tr> <th>Row Pct</th> <th></th> <th></th> <th></th> </tr> <tr> <th>Col Pct</th> <th>2</th> <th>3</th> <th></th> </tr> </thead> <tbody> <tr> <td>1</td> <td>21</td> <td>14</td> <td>35</td> </tr> <tr> <td></td> <td>31.82</td> <td>21.21</td> <td>53.03</td> </tr> <tr> <td></td> <td>60.00</td> <td>40.00</td> <td></td> </tr> <tr> <td></td> <td>52.50</td> <td>53.85</td> <td></td> </tr> <tr> <td>2</td> <td>19</td> <td>12</td> <td>31</td> </tr> <tr> <td></td> <td>28.79</td> <td>18.18</td> <td>46.97</td> </tr> <tr> <td></td> <td>61.29</td> <td>38.71</td> <td></td> </tr> <tr> <td></td> <td>47.50</td> <td>46.15</td> <td></td> </tr> <tr> <td>Total</td> <td>40</td> <td>26</td> <td>66</td> </tr> <tr> <td></td> <td>60.61</td> <td>39.39</td> <td>100.00</td> </tr> </tbody> </table> <p>Bei Gruppe_A fehlt Merkmalsausprägung 3 Bei Gruppe_B fehlt Merkmalsausprägung 1 Die Tabelle ist quadratisch, aber inhaltlich nicht symmetrisch, dennoch wird die Statistik erstellt:</p> <p>STATISTICS FOR TABLE OF GRUPPE_A BY GRUPPE_B McNemar's Test</p> <p>----- Statistic = 0.758 DF = 1 Prob = 0.384</p> <p>Simple Kappa Coefficient</p> <p>----- Kappa = -0.013 ASE = 0.122 95% Confidence Bounds -0.251 0.225</p> <p>Sample Size = 66</p>	GRUPPE_A	GRUPPE_B		Total	Frequency				Percent				Row Pct				Col Pct	2	3		1	21	14	35		31.82	21.21	53.03		60.00	40.00			52.50	53.85		2	19	12	31		28.79	18.18	46.97		61.29	38.71			47.50	46.15		Total	40	26	66		60.61	39.39	100.00	<p>TABLE OF GRUPPE_A BY GRUPPE_B</p> <table border="1"> <thead> <tr> <th>GRUPPE_A</th> <th colspan="3">GRUPPE_B</th> <th>Total</th> </tr> <tr> <th>Frequency</th> <th></th> <th></th> <th></th> <th></th> </tr> <tr> <th>Percent</th> <th></th> <th></th> <th></th> <th></th> </tr> <tr> <th>Row Pct</th> <th></th> <th></th> <th></th> <th></th> </tr> <tr> <th>Col Pct</th> <th>1</th> <th>2</th> <th>3</th> <th></th> </tr> </thead> <tbody> <tr> <td>1</td> <td>1E-20</td> <td>21</td> <td>14</td> <td>35</td> </tr> <tr> <td></td> <td>0.00</td> <td>31.82</td> <td>21.21</td> <td>53.03</td> </tr> <tr> <td></td> <td>0.00</td> <td>60.00</td> <td>40.00</td> <td></td> </tr> <tr> <td></td> <td>100.00</td> <td>52.50</td> <td>53.85</td> <td></td> </tr> <tr> <td>2</td> <td>0</td> <td>19</td> <td>12</td> <td>31</td> </tr> <tr> <td></td> <td>0.00</td> <td>28.79</td> <td>18.18</td> <td>46.97</td> </tr> <tr> <td></td> <td>0.00</td> <td>61.29</td> <td>38.71</td> <td></td> </tr> <tr> <td></td> <td>0.00</td> <td>47.50</td> <td>46.15</td> <td></td> </tr> <tr> <td>3</td> <td>0</td> <td>0</td> <td>1E-20</td> <td>1E-20</td> </tr> <tr> <td></td> <td>0.00</td> <td>0.00</td> <td>0.00</td> <td>0.00</td> </tr> <tr> <td></td> <td>0.00</td> <td>0.00</td> <td>100.00</td> <td></td> </tr> <tr> <td></td> <td>0.00</td> <td>0.00</td> <td>0.00</td> <td></td> </tr> <tr> <td>Total</td> <td>1E-20</td> <td>40</td> <td>26</td> <td>66</td> </tr> <tr> <td></td> <td>0.00</td> <td>60.61</td> <td>39.39</td> <td>100.00</td> </tr> </tbody> </table> <p>STATISTICS FOR TABLE OF GRUPPE_A BY GRUPPE_B Test of Symmetry</p> <p>----- Statistic = 47.000 DF = 3 Prob = 0.001</p> <p>Kappa Coefficients</p> <table border="1"> <thead> <tr> <th>Statistic</th> <th>Value</th> <th>ASE</th> <th colspan="2">95% Confidence Bounds</th> </tr> </thead> <tbody> <tr> <td>Simple Kappa</td> <td>0.004</td> <td>0.042</td> <td>-0.078</td> <td>0.087</td> </tr> <tr> <td>Weighted Kappa</td> <td>-0.000</td> <td>0.000</td> <td>-0.000</td> <td>-0.000</td> </tr> </tbody> </table> <p>Sample Size = 66</p>	GRUPPE_A	GRUPPE_B			Total	Frequency					Percent					Row Pct					Col Pct	1	2	3		1	1E-20	21	14	35		0.00	31.82	21.21	53.03		0.00	60.00	40.00			100.00	52.50	53.85		2	0	19	12	31		0.00	28.79	18.18	46.97		0.00	61.29	38.71			0.00	47.50	46.15		3	0	0	1E-20	1E-20		0.00	0.00	0.00	0.00		0.00	0.00	100.00			0.00	0.00	0.00		Total	1E-20	40	26	66		0.00	60.61	39.39	100.00	Statistic	Value	ASE	95% Confidence Bounds		Simple Kappa	0.004	0.042	-0.078	0.087	Weighted Kappa	-0.000	0.000	-0.000	-0.000
GRUPPE_A	GRUPPE_B		Total																																																																																																																																																																								
Frequency																																																																																																																																																																											
Percent																																																																																																																																																																											
Row Pct																																																																																																																																																																											
Col Pct	2	3																																																																																																																																																																									
1	21	14	35																																																																																																																																																																								
	31.82	21.21	53.03																																																																																																																																																																								
	60.00	40.00																																																																																																																																																																									
	52.50	53.85																																																																																																																																																																									
2	19	12	31																																																																																																																																																																								
	28.79	18.18	46.97																																																																																																																																																																								
	61.29	38.71																																																																																																																																																																									
	47.50	46.15																																																																																																																																																																									
Total	40	26	66																																																																																																																																																																								
	60.61	39.39	100.00																																																																																																																																																																								
GRUPPE_A	GRUPPE_B			Total																																																																																																																																																																							
Frequency																																																																																																																																																																											
Percent																																																																																																																																																																											
Row Pct																																																																																																																																																																											
Col Pct	1	2	3																																																																																																																																																																								
1	1E-20	21	14	35																																																																																																																																																																							
	0.00	31.82	21.21	53.03																																																																																																																																																																							
	0.00	60.00	40.00																																																																																																																																																																								
	100.00	52.50	53.85																																																																																																																																																																								
2	0	19	12	31																																																																																																																																																																							
	0.00	28.79	18.18	46.97																																																																																																																																																																							
	0.00	61.29	38.71																																																																																																																																																																								
	0.00	47.50	46.15																																																																																																																																																																								
3	0	0	1E-20	1E-20																																																																																																																																																																							
	0.00	0.00	0.00	0.00																																																																																																																																																																							
	0.00	0.00	100.00																																																																																																																																																																								
	0.00	0.00	0.00																																																																																																																																																																								
Total	1E-20	40	26	66																																																																																																																																																																							
	0.00	60.61	39.39	100.00																																																																																																																																																																							
Statistic	Value	ASE	95% Confidence Bounds																																																																																																																																																																								
Simple Kappa	0.004	0.042	-0.078	0.087																																																																																																																																																																							
Weighted Kappa	-0.000	0.000	-0.000	-0.000																																																																																																																																																																							

Das SAS-Macro: AGREE_NS.MAC

<p style="text-align: center;">%MACRO agree_ns (datname,varlist1,varlist2,n,wt,wherevar,m,byvar,j);</p>	<p style="text-align: center;">PROC FREQ/AGREE (McNemar Test/Bowker Test auf Symmetrie, Kappa-Koeffizient) - wenn Tafel nicht quadratisch und symmetrisch ist</p>
<pre> OPTIONS NONUMBER NOCENTER NODATE PS=63 LS=88; %Macro agree_ns (datname,varlist1,varlist2,n,wt,wherevar,m,byvar,j); TITLE2; TITLE3; DATA zzzz; %DO k=1 %TO &n %BY 1; call symput ("a&k",scan(&varlist1,&k," ")); call symput ("b&k",scan(&varlist2,&k," ")); %END; RUN; %IF &byvar^=. %THEN %DO; PROC SORT DATA=&datname OUT=file; BY &byvar; RUN; DATA file; SET file; BY &byvar; IF FIRST.&byvar THEN FIRST=1; DATA file; SET file; IF FIRST=1 THEN x +1; RUN; PROC FREQ DATA=file; %IF &wt^=. %THEN %DO; TABLE &byvar*x/MISSING; WEIGHT &wt; %END; %ELSE %DO; TABLE &byvar*x/MISSING; %END; TITLE2 "Tabelle mit ursprünglicher (in &byvar) und neuer"; TITLE3 "(in x) Codierung der BY-Variablen für &datname"; RUN; %DO i=1 %TO &j %BY 1; DATA d_&i; SET file; IF x=&i THEN OUTPUT d_&i; RUN; %END; PROC DATASETS; DELETE file; RUN; QUIT; %END; </pre>	<p>datname Dateiname</p> <p>varlist1 Variablenliste 1 (z.B. Fälle) - gleiche Länge wie varlist2</p> <p>varlist2 Variablenliste 2 (z.B. Kontrollen) - gleiche Länge wie varlist1</p> <p>n Anzahl der zu testenden Variablen von varlist1 bzw varlist2</p> <p>wt WEIGHT-Variable - kann auf missing gesetzt werden</p> <p>wherevar WHERE-Variable - kann auf missing gesetzt werden</p> <p>m enthält Operationszeichen und Wert für WHERE-Variable, (z.B. m= ^=3)</p> <p>byvar BY-Variable - enthält Anzahl der Ausprägungen der BY-Variablen, bzw. der ersten j Ausprägungen der aufsteigend sortierten BY-Variablen, für die Berechnung durch geführt werden soll</p> <p>TITLE1 - kommt aus aufrufendem SAS-Programm</p> <p>Liest alle Variablen aus varlist1 und varlist2 ein gezielter Zugriff auf alle Variablen aus varlist1 über &a&i gezielter Zugriff auf alle Variablen aus varlist2 über &b&i</p> <p>Sortiert Datei &datname nach &byvar, wenn diese Variable gesetzt wurde</p> <p>Markiert jede neue Ausprägung von &byvar mit 1 in SAS-interner Variable First</p> <p>Zählt die Anzahl der Ausprägungen von &byvar - ggf. Umwandlung von char in num</p> <p>Gegenüberstellung der ursprünglichen (in byvar gespeicherten) und neuen Klassierung (in x gespeicherten) Ausprägung der BY-Variablen</p> <p>Teilt ursprüngliche Datei &datname in &j Einzeldateien, eine für jede Ausprägung von &byvar - es entstehen die Dateien d_1 bis d_&j</p> <p>Löschen der Hilfsdatei file aus PROC SORT</p>

<pre> %ELSE %DO; %LET j=%EVAL (1); DATA d_1; SET &datname; RUN; %END; %DO k=1 %TO &n %BY 1; %DO i=1 %TO &j %BY 1; PROC FREQ DATA=d_&i; %IF &wherevar^=. %THEN %DO; WHERE &wherevar&m; %IF &wt^=. %THEN %DO; TABLE &a&k*&b&k/ OUT=xxx&k&i NOROW NOCOL; WEIGHT &wt; %END; %ELSE %DO; TABLE &a&k*&b&k/ OUT=xxx&k&i NOROW NOCOL; %END; TITLE2 "PROC FREQ (AGREE) für Datei &datname, mit WHERE = &wherevar mit &m"; TITLE3 "mit Gruppierung nach BY = &byvar mit neu codierter Ausprägung &i"; %END; %ELSE %DO; %IF &wt^=. %THEN %DO; TABLE &a&k*&b&k/ OUT=xxx&k&i NOROW NOCOL; WEIGHT &wt; %END; %ELSE %DO; TABLE &a&k*&b&k/ OUT=xxx&k&i NOROW NOCOL; %END; TITLE2 "PROC FREQ (AGREE) für Datei &datname, ohne WHERE"; TITLE3 "mit Gruppierung nach BY = &byvar mit neu codierter Ausprägung &i"; %END; RUN; DATA yy&k&i; SET xxx&k&i; IF &a&k^=&b&k THEN count_n=1e-20; IF &a&k>&b&k THEN f_=&b&k; IF &a&k<&b&k THEN f_=&b&k; IF &a&k>&b&k THEN k_=&a&k; IF &a&k<&b&k THEN k_=&a&k; RUN; DATA yy1&k&i yy2&k&i; SET yy&k&i; IF &a&k^=f_ THEN OUTPUT yy1&k&i; IF &b&k^=k_ THEN OUTPUT yy2&k&i; RUN; DATA yy1&k&i; SET yy1&k&i; KEEP f_ &b&k count_n; DATA yy1&k&i; SET yy1&k&i; RENAME count_n=count_f_=&a&k; DATA yy2&k&i; SET yy2&k&i; KEEP k_ &a&k count_n; </pre>	<p>Wenn keine &byvar gesetzt wurde, wird ursprüngliche Datei &datname nach d_1 geschrieben</p> <p>Schleife über die Anzahl der Variablen in den Variablenlisten (&n) Schleife über die Anzahl der Ausprägungen von &byvar</p> <p>Die Kreuztabellen werden am Bildschirm angezeigt, die Teststatistik aber unterdrückt, da an dieser Stelle quadratische aber nicht-symmetrische Tafeln vorkommen können, für die eine falsche Berechnung erscheinen würde</p> <p>Neue Variable count_n wird mit sehr kleinem Wert gesetzt, wenn Inhalt von Variable a und b ungleich ist Neue Variable f_, wenn Inhalt von a und b nicht gleich ist; f_ bekommt Inhalt von b Neue Variable k_, wenn Inhalt von a und b nicht gleich ist; k_ bekommt Inhalt von a</p> <p>Trennen der Datei in zwei Teildateien, wenn</p> <p>Inhalt von a ungleich dem Inhalt von f_ Inhalt von b ungleich dem Inhalt von k_</p> <p>In der Datei yy1&k&i werden drei Variablen behalten</p> <p>Umbenennung von zwei Variablen</p> <p>In Datei yy2&k&i werden drei Variablen behalten</p>
---	--

<pre> DATA yy2&k&i; SET yy2&k&i; RENAME count_n=count k_=&b&k; RUN; DATA yyy&k&i; SET yy1&k&i yy2&k&i; RUN; DATA test&k&i; SET xxx&k&i yyy&k&i; WHERE count^=.; RUN; PROC SORT DATA=test&k&i OUT=test&k&i NODUP; BY &a&k &b&k; DATA test&k&i; SET test&k&i; BY &a&k &b&k; IF FIRST.&b&k THEN first=1; RUN; DATA test&k&i; SET test&k&i; IF first=1 THEN OUTPUT test&k&i; RUN; PROC FREQ DATA=test&k&i; TABLE &a&k*&b&k/AGREE NOPRINT NOCOL NOROW; WEIGHT count; TITLE2; TITLE3; RUN; PROC DATASETS; DELETE xxx&k&i yyy&k&i yy1&k&i yy2&k&i test&k&i aaa&k&i; RUN; QUIT; %END; %END; PROC DATASETS; %DO i=1 %TO &j %BY 1; DELETE d_&i; %END; PROC DATASETS; DELETE zzzz; RUN; QUIT; %MEND; </pre>	<p>Umbenennung von zwei Variablen</p> <p>Zusammenbinden der beiden Dateien yy1&k&i und yy2&k&i zu yyy&k&i</p> <p>Zusammenbinden der Datei yyy&k&i mit Out-Datei xxx&k&i zu test&k&i, bei denen count-Variable nicht missing</p> <p>Sortieren nach Kombination von &a&k und &b&k</p> <p>Markieren jeder neuen Ausprägung von &b&k in SAS-interner Variablen First mit 1</p> <p>Schreiben aller Datensätze mit first=1 nach test&k&i</p> <p>Anzeige der Teststatistik für die nunmehr symmetrische Kreuztabelle, die Darstellung der ergänzten Tabelle am Bildschirm wird unterdrückt, falls man diese sehen möchte NOPRINT-Anweisung löschen</p>
---	---

Literatur

- [1] I. Guggenmoos-Holzmann, K.-D. Wernecke: Medizinische Statistik. Blackwell Wissenschafts-Verlag, Berlin 1996.
- [2] SAS / STAT Software. Changes and Enhancements through Release 6.11 (55356). SAS Inst. Inc, Cary, NC 27513.