Analyzing Survey Data Using SUDAAN[®] Release 7.5

by

Gayle S. Bieler gbmac@rti.org Research Triangle Institute

and

Rick L. Williams willy@rti.org Research Triangle Institute

Prepared for: Joint Statistical Meetings Continuing Education Workshop August, 1997 Analyzing Survey Data Using SUDAAN Release 7.5 was written by Gayle S. Bieler and Rick L. Williams

Copyright 1997 by Research Triangle Institute P.O. Box 12194 Research Triangle Park, NC 27709-2194

All rights reserved. No part of this publication may be reproduced or transmitted by any means without permission from the publisher.

SUDAAN and RTI are trademarks of the Research Triangle Institute. SAS is a trademark of SAS Institute, Inc. SPSS is a trademark of SPSS, Inc.

How to contact us:

Voice:	919-541-6236
Fax:	919-541-7431
Email:	sudaan@rti.org

Visit SUDAAN on the Web!

www.rti.org/patents/sudaan/sudaan.html

SUDAAN Software for the Statistical Analysis of Correlated Data

Analyzing Survey Data Using SUDAAN® Release 7.5

by

Gayle S. Bieler gbmac@rti.org

and

Rick L. Williams willy@rti.org Research Triangle Institute

Prepared for: Joint Statistical Meetings Continuing Education Workshop August, 1997



Research Triangle Institute PO Box 12194 Research Triangle Park, NC 27709-2194

Analyzing Survey Data Using SUDAAN Release 7.5

Table of Contents

About Survey Data	-13 -30
New Features in SUDAAN 7.5 Summary of Enhancements Variance Estimation Methods: Taylor linearization, Jackknife and BRR MULTILOG procedure (Release 7.0) GEE for Efficient Parameter Estimation R-Square in Logistic Regression REFLEVEL statement 94-7 EFFECTS statement 101-7 LSMEANS statement 121-7 Design Effects	-79 -83 -92 93 100 120 132
Additional Example Using SUDAAN Release 7.5: Evaluation of Project DARE	163
References	168

Analyzing Survey Data Using SUDAAN® Release 7.5

ABSTRACT

In the social sciences and public health, researchers often analyze survey data which were collected via a complex sampling design. Such survey designs often include stratification and cluster sampling (*e.g.*, sampling by geographic clusters) in one or more stages, where the clusters may be sampled with unequal probabilities. Such designs complicate the statistical analysis since the observations are not independent and identically distributed (*iid*). Failure to account for the design in the statistical analysis typically result in underestimated standard errors and false positive test results.

Unlike standard statistical packages, SUDAAN is specifically designed to handle non-*iid* observations drawn from finite populations. SUDAAN offers a powerful set of analytic tools for linear regression, logistic regression, multinomial logistic regression, proportional hazards modelling, and descriptive data analysis.

This seminar will review the statistical methods used in SUDAAN and demonstrate its use via a series of examples from the social sciences and public health. The basic concept throughout all SUDAAN procedures is to use consistent variance estimators for statistics derived from complex samples (*e.g.*, means, proportions, odds ratios, regression coefficients), without imposing strict distributional assumptions, and treating the intracluster correlation as a nuisance parameter. SUDAAN is currently the only statistical package to offer three well-known methods for variance estimation in sample surveys: Taylor linearization, BRR, and Jackknife.

This workshop will highlight many of the new features in SUDAAN Release 7.5 that are of particular interest to survey researchers, including: 1) BRR and Jackknife variance estimation for descriptive statistics and regression modelling; 2) Four choices for computing *design effects*; 3) *User-friendly contrast statements* and a *reference level statement* for specifying the reference cells of categorical covariates in all regression procedures; 4) *Least squares means* estimation for linear regression; 5) a more useful R^2 statistic based on the log-likelihood for logistic regression; and 6) better compatibility with well-known software packages (*SAS-Callable versions* for SUN Solaris and Win 95, and *reading SPSS datasets*). Attendees should be familiar with the basics of survey sampling and analysis, as well as fitting linear and non-linear regression models.

What is a Sample Survey?

- A study involving a subset (or sample) of individuals selected from a larger population with known probabilities of selection
- Measurements are aggregated over all sample members to obtain summary statistics (*e.g.*, means, proportions, totals, or ratios) for the sample
- Extrapolations made to the entire population (estimates of population parameters)

What is a Census?

 All individuals in a population are selected for measurement (summary stats are not extrapolations)

Main Advantages of Sampling:

- Reduced Cost
- Greater Speed

Many Surveys are Purely Descriptive:

 Estimation of summary statistics often the primary objective; hypothesis testing a secondary objective

Target Population

Entire set of individuals to which findings are to be extrapolated Individual members of the population whose characteristics are to be measured are called *population elements*

Example:

Select a nationally representative sample of students from the US population

Sample Design:

Stratified, Multi-Stage Nested Design

- 1) Divide the country into 4 *regions* (*strata*)
- Obtain a comprehensive list of schools in each region; Select a sample of *schools* from each region, according to a known probability sampling scheme, such as:
 - simple random sampling,
 - probability proportional to size sampling (PPS),
 - certainty sampling (probability of selection = 1)
- 3) Obtain a comprehensive list of *students* within each sample school;

Select a sample of *students* from each school according to a known probability sampling scheme (as above), or select all students from the school (probability of selection = 1)

Stratification

 Selection of sampling units (population elements) from mutually exclusive and exhaustive subpopulations

STUDENT SAMPLE: Strata = REGION

Regions were *not* randomly selected; they were chosen in advance

Independent samples of schools chosen within each stratum

- Different sampling methods can be used in different strata
- Obtain an estimate for the population as a whole by aggregating the individual stratum estimates over all strata
- Can *reduce variance* of sample statistics (e.g., average GPA, average height) if strata are chosen efficiently (i.e., if strata are homogeneous wrt the variable of interest).

Clustering

Problems with direct element sampling if:

- There exists no sampling frame for the population elements (*e.g.*, no master list of students in US from which to select a sample)
- The population elements are scattered over a wide area in which case direct element sampling will result in a scattered sample
 - field costs prohibitive

Solution:

- Use cluster sampling
 - Population elements are aggregated into larger units (*clusters*) for which complete lists are available

e.g., schools

■ Use multistage designs

Clustering

 Subunits (students) are selected into the sample from clusters or primary sampling units (schools)

Student Sample:

Primary Sampling Unit (*PSU* or *Cluster*) = *School* Students clustered within schools;

- Usually positive correlation within clusters (*i.e.*, students within schools are more alike than across schools, so they tend to respond similarly)
- Variance of sample statistics (e.g., average GPA) is typically *increased* under cluster sampling

How Much ?

Clustering

Design Effect

Describes the change in variance of an estimated statistic due to clustering

Estimated as the ratio of variance under the cluster design vs. a simple random sample of the same size (*i.e.*, independence), or via an analytic expression:

Design Effect =
$$\frac{V(\hat{\theta})_{CLUSTER}}{V(\hat{\theta})_{SRS}} = 1 + \rho(m-1)$$
,

m = average cluster size (students)

 ρ = intra-cluster correlation coefficient (measure of association within the cluster)

If
$$\rho = 0$$
:

No correlation (DEFF = 1)

If $\rho = 1$:

Perfect positive correlation (*e.g.*, everyone responds the same)

DEFF = m (cluster size)

If $0 < \rho < 1$:

Some degree of correlation (units respond similarly) 1 < DEFF < m

Clustering

- For a sample of a given size, as the cluster size and the intracluster correlation increase, the variance is also *increased*.
- Another way to think about clustering:

Loss of precision Reduction in effective sample size

Effective sample size < number of observations (students) > number of clusters (schools)

Multi-Stage Sampling

- Cluster sampling in 2 or more stages:
- Selection Order:
 - Stage 1: CountiesStage 2: Schools from sampled countiesStage 3: Students from sampled schools

or

- Stage 1: Schools (stratified by region)
- Stage 2: Homeroom classes (stratified by grade)
- Stage 3: Students.

Unequal Weighting

- Relates to *Probability* or *Population-Based Sampling*
 - Every element in the target population has a known, non-zero probability of being included in the sample
- *Unequal weighting* results when sample members (*e.g.*, students) selected with unequal probabilities
 - Oversampling certain subpopulations, such as the elderly, the poor, Hispanics, or Native Americans.
- Each sample member has a sampling weight associated with their data

Sampling weight = inverse of selection probability

Refers to number of individuals in target population that the sample member represents

Weights needed for unbiased estimation of population parameters (findings are then generalizable to a finite population of interest).

- Downside: Variability in sampling weights can lead to inefficiency, meaning loss of power and wider confidence intervals.
- Variance of sample statistics usually *increased* if weights are highly variable.

Nonlinear Statistics

Most survey statistics are *not* simple linear functions of the data, but rather ratios of random variables

Linear Statistic: the weighted total

$$x_w = \sum_{i=1}^n w_i x_i$$

Nonlinear Statistic: the weighted mean, proportion, etc.

$$\bar{x}_{w} = \frac{\sum_{i=1}^{n} w_{i} x_{i}}{\sum_{i=1}^{n} w_{i}} \qquad w_{i} = \frac{1}{p_{i}} = sampling weight for sample member i$$

$$x_{i} = outcome of interest for sample member i$$

- \bar{x}_{w} is a weighted mean if the response variable is continuous, or a weighted proportion if x_{i} is coded 0 (characteristic absent) and 1 (characteristic present)
- The denominator $\sum w_i$ is not a fixed quantity but rather an *estimate* of the population size, which varies from sample to sample when the weights are unequal.
- Non-standard techniques required for variance estimation: SUDAAN offers *Taylor series linearization* and *replication methods* (*BRR and Jackknife*)

Without Replacement Sampling

- Units selected into the sample do not have another chance of being selected
- In practice, almost all sampling is done without replacement, but can often be ignored in the analysis
- If you account for it in analysis, the variance of sample statistics is *decreased* when the sampling fractions (*e.g.*, proportion of schools selected from each region) become large
- Sampling fraction:

 $\frac{n}{N} = \frac{number \ units \ selected \ into \ sample}{number \ units \ in \ population}$

In other words, the more you know about a population, the smaller the variance of sample statistics

• Why is it ignored most of the time?

Accounting for without-replacement sampling makes variance estimation slightly more complicated, since you must know the sampling fractions within each of the firststage strata

There is little efficiency to be gained when sampling fractions are small.

Applications in Epidemiologic Studies

Longitudinal Studies Repeated Measures Studies

Multiple events, such as hospital visits or illness episodes, are observed over time on each subject.

Example 1:

Relationship between MDI (Mental Development Index) measurements and umbilical cord blood lead levels in children (Waternaux, et al, *JASA*, 1989)

MDI measurements recorded at 6, 12, and 18 months of age for each child

Example 2:

Logistic regression of the propensity of daily asthma attacks on the average daily level of total suspended particulates in the air (Korn and Whittemore, 1979, *Biometrics*)

Daily asthma measurements and other time-dependent covariates recorded on each person in a sample of adults and children (up to 34 weeks of daily measurements)

Applications in Behavioral Research

Examples:

- School-based evaluations of substance-abuse prevention programs in the student population (observations are on students nested within schools)
- Evaluation of Project DARE (Ennett et al, 1994, Addictive Behaviors; Norton, et al., 1996, Journal of Consulting and Clinical Psychology)

Multi-Stage Sample Surveys

- Data are obtained via a complex survey design (Cluster sampling in 1 or more stages; clusters may be sampled with differing probabilities)
- Practical advantages to multi-stage design (e.g., sampling by geographic clusters):
 - Not always feasible to enumerate the population of interest (sample frame)
 - Reduces cost of data collection (travel)
- Design-based methods of analysis:
 - Weighting of the data for unbiased estimates
 - *Linearization* and *replication methods* to estimate variances
- Examples: NHANES, NHIS, BRFSS, NHSDA

Sudaan vs. Other Software (SAS[®], SPSS[®], ...)

SAS, SPSS, etc

SUDAAN

Simple random sampling; Infinite populations	SAMPLE SELECTION ASSUMPTIONS	Complex probability sampling schemes; finite populations
Known probability distributions (normal, binomial)	DISTRIBUTIONAL ASSUMPTIONS	No strict distributional assumptions
Linear statistics <i>only</i>	RANDOM VARIABLE ASSUMPTIONS	Functions of linear statistics

RESULT: *SAS* yields unbiased point estimates if you include appropriate weights, but variance estimates *wrong* (usually *underestimated*) due to clustering.

Test statistics have inflated Type I error rates (reject null hypothesis more often than nominally specified, *i.e.*, false positives)

SUDAAN yields consistent variance estimates for sample statistics (*e.g.*, means, totals, proportions, ratios, regression coefficients) needed for unbiased inference.

An Example

WIC Mothers and Infants:	Two-stage clustered design				
	Strata = Region PSU = WIC local agencies				
Sample Size:	953				
Population Size:	Approximately 506,000 WIC participants				
Outcome of Interest:	Initiated breastfeeding				
ESTIMATE:	Percentage breastfed their infant				
COMPARISON DOMAINS:	Race groups (white vs. non-whites)				

(Results Follow)

An Example

SAS Results

Association Between Breast-feeding and Mother's Race

Sampling Weights Sum to Population Size

	TABLE BFEED (Br		ng Initia	ation)	
	MOMRACE (Mother R Frequency Col Pct White Other				
	Yes		141964 56.40		
	No		109748 43.60		
	Total	254567	251712	506279	
STA	ATISTICS F	OR TABLE	OF BFEED	BY MOMRACE	
Statist	cic	DF	' Valu	le Prob)
Chi-Squ	lare	1	1462.54	4 0.001	-
Sample	Size = 50	6279			

An Example

SAS Results

Association Between Breast-feeding and Mother's Race

Weights Normalized to Sum to Sample Size: NORMWGT = WEIGHT * (953 / 506,279)

Normalized Weights TABLE OF BFEED BY MOMRACE BFEED (Breastfeeding Initiation) MOMRACE (Mother Race) Frequency Col Pct White Other To Yes 244.57 267.23 51 51.04 56.40 No 234.61 206.59 44 48.96 43.60	
BFEED (Breastfeeding Initiation) MOMRACE (Mother Race) Frequency Col Pct White Other To 	
MOMRACE (Mother Race) Frequency Col Pct White Other To 	
Frequency Col Pct White Other To 	
Col Pct White Other To 	
51.04 56.40 No 234.61 206.59 44	cal
 No 234.61 206.59 44	L.8
	L.2
Total 479.187 473.813	953
STATISTICS FOR TABLE OF BFEED BY MOM	RACE
Statistic DF Value	Prob
Chi-Square 1 2.753	
Sample Size = 953	0.097

An Example

SUDAAN Results

Association Between Breast-feeding and Mother's Race

Weights Sum to Population Size

 Date:
 07-17-97
 Research Triangle Institute
 Page : 2

 The
 10:12:22
 The
 The
 Table : 1 Time: 16:13:03 The CROSSTAB Procedure Number of observations read : 953 Weighted count : 506279 Denominator degrees of freedom : 21 Variance Estimation Method: Taylor Series (WR) STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS | **Mother Race** | Total | **White | Non-White** | Breastfeeding Initiation _____

 Sample Size
 953
 480
 473

 Population
 506279
 254567
 251712

 Column Percent
 100.00
 100.00
 100.00

 Std Error
 0.00
 0.00
 0.00

 Design Effect
 .
 .
 .

 Total -----

 Sample Size
 522
 249
 273

 Population
 271893
 129929
 141964

 Column Percent
 53.70
 51.04
 56.40

 Yes 3.10 | Std Error 3.18 4.99 Design Effect 3.74 1.97 | 4.90 _____

 Sample Size
 431
 231
 200

 Population
 234386
 124638
 109748

 Column Percent
 46.30
 48.96
 43.60

 Std Error
 2.10
 2.18
 4.99

 NO

 Std Error
 3.10
 3.18
 4.99

 Design Effect
 3.58
 1.91
 4.64

 4.99

L

Why SUDAAN?

An Example

SUDAAN Results

Association Between Breast-feeding and Mother's Race

Weights Sum to Population Size

Date: 07-17-97 Time: 16:13:03		Research Triangle Institute The CROSSTAB Procedure					
Variance Estimation	n Method: Taylor Sen	ries (WR)					
Chi Square Test of Independence for Breastfeeding Initiation and Mother Race							
STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS							
			-				
		 I	-				
	Chi-Square	0.9403					
	DF	1					
	P-Value	0.3432					
			-				

An Example

SUMMARY OF RESULTS

Package	Method	%Breastfeed (SE): White	%Breastfeed (SE): Non-White	Chi- Square	P-value
SAS	Weighted	51.04 (0.10)	56.40 (0.10)	1462.5	0.001
	Weights Normalized (Sum to Sample Size)	51.04 (2.28)	56.40 (2.28)	2.75	0.097
SUDAAN	Weighted	51.04 (3.18)	56.40 (4.99)	0.94 *	0.343

NOTE:

SAS standard errors are calculated as $\sqrt{p(1-p)/n}$, where *n* is the sum of the weights, and *p* is the proportion breastfeeding.

Why Did We Bother Developing SUDAAN?

Intra-Cluster Correlation

- Potential for clustermates to respond similarly (genetic and environmental influences)
- Experimental units from the same cluster are not statistically independent
- Usually results in *overdispersion*, or extra-variation in the responses beyond what would be expected under independence
- Other standard statistical packages (e.g., SAS[®], SPSS[®]) do not uniformly address the correlated data problem in all analytical procedures

SUDAAN uses correlated data methods for:

- Regression modelling
- Estimating and analyzing: Means, medians, percentages, percentiles, odds ratios and relative risks, and ratios of random variables
- Chi-square tests in contingency tables
- Cochran-Mantel-Haenszel tests in contingency tables

Multivariate Responses (Clustered Data)

Notation

Data

$$(y_{ij}, \boldsymbol{x}_{ij}), \quad j = 1, \dots, m_i$$
$$i = 1, \dots, n$$
$$N = \sum_i m_i = total \ sample \ size$$

Responses

$$\mathbf{y}_{i} = (y_{i1}, y_{i2}, \dots, y_{im_{i}})$$

Covariates

$$\boldsymbol{x}_{ij} = (x_{ij1}, x_{ij2}, ..., x_{ijp})$$

This is the clustered data situation covered by SUDAAN

Assumptions: Independence Vs. Clustered Data

Independence

$$\boldsymbol{Y} = \begin{bmatrix} y_1 \\ \cdot \\ \cdot \\ \cdot \\ y_N \end{bmatrix} \qquad \qquad \boldsymbol{V}(\boldsymbol{Y}) = \sigma^2 \boldsymbol{I}_N = \begin{bmatrix} \sigma^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 \\ 0 & 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

Observations independent, constant variance

Clustered Data (SUDAAN):

$$\mathbf{Y} = \begin{bmatrix} y_{11} \\ \vdots \\ y_{1m_1} \\ \vdots \\ y_{n1} \\ \vdots \\ y_{nm_n} \end{bmatrix} \quad n \text{ clusters of } m_i \text{ observations } (N = \sum_{i=1}^n m_i)$$

$$\text{Unequal observations per cluster} = m_i$$

$$\text{Example: } n \text{ litters with } m_i \text{ pups per litter}$$

Assumptions: Independence Vs. Clustered Data

Clustered Data (SUDAAN):

$$\boldsymbol{V}(\boldsymbol{Y}) = \begin{vmatrix} V_1 & 0 & 0 & \cdots & 0 \\ 0 & V_2 & 0 & \cdots & 0 \\ 0 & 0 & V_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & V_n \end{vmatrix}$$
Cluster-Correlated Data
Block-Diagonal by Cluster
 V_i is an $m_i \ge m_i$ matrix

$$\boldsymbol{V_i} = \begin{bmatrix} \sigma_{(i)1}^2 & \sigma_{(i)12} & \sigma_{(i)13} & \cdots & \sigma_{(i)1m} \\ \sigma_{(i)21} & \sigma_{(i)2}^2 & \sigma_{(i)23} & \cdots & \sigma_{(i)2m} \\ \vdots & & \ddots & \vdots \\ \vdots & & \ddots & \vdots \\ \sigma_{(i)m1} & \sigma_{(i)m2} & \sigma_{(i)m3} & \cdots & \sigma_{(i)m}^2 \end{bmatrix}$$

- V_i is an $m_i \ge m_i$ variance covariance matrix of observations in the *i*-th cluster
- *No assumptions on structure* of V_i (could be unstructured, multi-level, AR(1), exchangeable, etc.)
- Observations independent between clusters, completely arbitrary correlation structure within clusters

Independence Vs. Clustered Data: Fitting Linear Regression Models

Standard Situation: Linear Regression

$$\boldsymbol{Y} = \begin{bmatrix} y_1 \\ \cdot \\ \cdot \\ \cdot \\ \cdot \\ y_N \end{bmatrix} \quad \begin{array}{l} E(\boldsymbol{Y}) = \boldsymbol{X}\boldsymbol{\beta} \\ \boldsymbol{V}(\boldsymbol{Y}) = \sigma^2 \boldsymbol{I}_N \\ \text{Independent obs, constant variance} \end{array}$$

Standard Solution to Normal Equations:

$$\boldsymbol{b} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{Y}$$

 $Var(\boldsymbol{b}) = \hat{\sigma}^2 (\boldsymbol{X}'\boldsymbol{X})^{-1} \qquad \hat{\sigma}^2 = \text{Mean Square Error}$

This variance formula only holds when: $V(\mathbf{Y}) = \sigma^2 \mathbf{I}_N$

Independence Vs. Clustered Data: Fitting Linear Regression Models

How is SUDAAN different?

$$\boldsymbol{V}(\boldsymbol{Y}) = \boldsymbol{V}_{\boldsymbol{Y}} = \begin{bmatrix} V_1 & 0 & 0 & \cdots & 0 \\ 0 & V_2 & 0 & \cdots & 0 \\ 0 & 0 & V_3 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & V_n \end{bmatrix}$$

Cluster-Correlated Data Block-Diagonal by Cluster V_i is an $m_i \ge m_i$ matrix

 $\boldsymbol{b} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{X}'\boldsymbol{Y}$

Use *robust variance formula* to estimate:

 $Var(b) = V_{b}$ Estimates each element separately

KEY POINT:

 $V_{h} \neq \hat{\sigma}^{2} (X'X)^{-1}$ due to cluster-correlated data

Independence Vs. Clustered Data Fitting Linear Regression Models

Null Hypothesis:

$$H_0: C\beta = 0$$

C is a contrast matrix of rank r

General Form for Test Statistic:

$$Q = (\boldsymbol{C}\boldsymbol{b})^{\prime} [\boldsymbol{C} Var(\boldsymbol{b}) \boldsymbol{C}^{\prime}]^{-1} (\boldsymbol{C}\boldsymbol{b})$$

Standard Situation

$$Q = (Cb)^{\prime} \left[\hat{\sigma}^2 C (X^{\prime} X)^{-1} C^{\prime} \right]^{-1} (Cb)$$
$$= \frac{r \cdot MS_{H_0}}{MS_{error}} \sim r F_{r, N-r}$$

Standard computing formula used by most software packages

SUDAAN Test Statistic:

$$Q = (Cb)^{\prime} \left[CV_b C^{\prime} \right]^{-1} (Cb)$$

Does not reduce to any simple computing formula

SUDAAN Software Package

Software for Statistical Analysis of Correlated Data

- Single program, written in the C language, consisting of a family of statistical procedures
- As easy to use as SAS!
 - Uses a SAS-like interface
 - Accepts SAS data sets as input
- Two Modes of Operation:
 - 1) SAS-Callable (Win 95, SUN/Solaris, VAX/VMS, IBM/MVS)
 - 2) Stand-Alone (many platforms, including Windows)
- SPSS Users: Release 7.5 reads SPSS files

SUDAAN Procedures

DESCRIPTIVE PROCEDURES

CROSSTAB

Computes frequencies, percentage distributions, odds ratios, relative risks, and their standard errors (or confidence intervals) for user-specified crosstabulations, as well as chi-square tests of independence and the Cochran-Mantel-Haenszel chi-square test for stratified two-way tables.

DESCRIPT

Computes estimates of means, totals, proportions, percentages, geometric means, quantiles, and their standard errors; also computes standardized estimates and tests of single degree-offreedom contrasts among levels of a categorical variable.

RATIO

Computes estimates and standard errors of generalized ratios of the form $\Sigma y / \Sigma x$, where x and y are observed variables; also computes standardized estimates and tests single-degree-of-freedom contrasts among levels of a categorical variable.

REGRESSION PROCEDURES

REGRESS

Fits linear regression models and performs hypothesis tests concerning the model parameters. Uses *GEE* to efficiently estimate regression parameters, with robust and model-based variance estimation.

LOGISTIC

Fits logistic regression models to binary data and computes hypothesis tests for model parameters; also estimates odds ratios and their 95% confidence intervals for each model parameter.

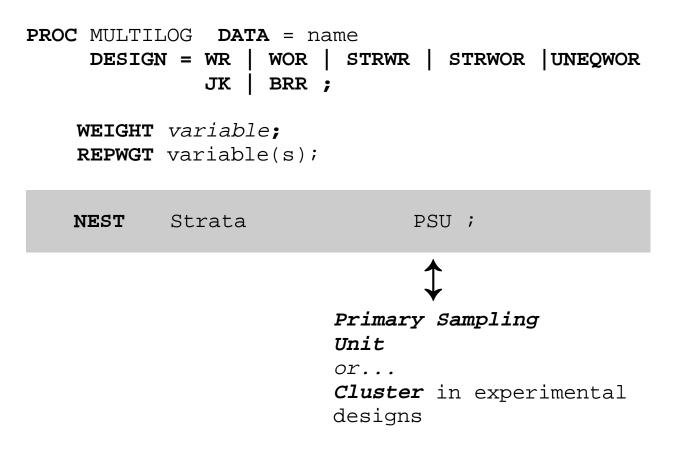
MULTILOG

Fits logistic and multinomial logistic regression models to ordinal and nominal categorical data and computes hypothesis tests for model parameters; estimates odds ratios and their 95% confidence intervals for each model parameter; uses *GEE* to efficiently estimate regression parameters, with robust and model-based variance estimation.

SURVIVAL

Fits discrete and continuous proportional hazards models to failure time data; also estimates hazard ratios and their 95% confidence intervals for each model parameter.

Elements of a SUDAAN Procedure



For Regression Modelling:

MODEL dependent = independent ;

DRUGUSE = AGE SEX RACE ;

For Descriptive Statistics:

VAR response_variables ;

TABLE categorical effects (e.g., RACE) ;

Enhancements to SUDAAN Release 7.5

Replication Methods for Robust Variance Estimation

- Jackknife
- Balanced Repeated Replication (BRR)

Enhancements of GEE Capabilities

- Exchangeable correlations in linear regression (as already in logistic and multinomial logistic since Release 7.0)
- Robust (default) and model-based variances in GEE applications

Regression Enhancements

- REFLEVEL statement to change the reference level for categorical covariates
- User-friendly contrast statement (EFFECTS) for testing simultaneous regression effects, simple effects in interaction models, and more
- R-square (Cox and Snell, 1989) in logistic regression
- Least Squares Means (LSMEANS) statement in linear regression
- MULTILOG Procedure for multinomial logistic regression (7.0)

SAS-Callable Platforms

- Windows 95
- SUN/Solaris

Now reads SPSS files (in addition to SAS and ASCII)

Three Variance Estimation Methods in SUDAAN

Basic Concept Behind All

1) Use *consistent estimators* of the parameters

e.g., Means, Proportions, Percentages, Odds Ratios, Regression Coefficients

Can even estimate the correlation structure and improve the efficiency of β

Intracluster correlation treated as a nuisance parameter

- 2) *Robust variance estimators* ensure consistent variance estimates and valid inferences:
 - Taylor linearization / GEE
 - Jackknife (new in Release 7.5)
 - BRR (new in Release 7.5)
 - Without imposing strict distributional assumptions about the response of interest

Taylor Linearization Approach

Two-Step Procedure for Variance Estimation:

1) Use Taylor series linearization to approximate functions of linear statistics (e.g., ratios of random variables)

Example: Prevalence of drug use

$$\hat{p} = \frac{\sum_{i=1}^{n} \sum_{j=1}^{m_{i}} w_{ij} y_{ij}}{\sum_{i=1}^{n} \sum_{j=1}^{m_{i}} w_{ij}} = \frac{Estimated number drug users}{Estimated population size}$$

Find *linear approximation* to this nonlinear statistic (Kendall and Stuart, 1973);

Design-specific variance formulas available for *linear* statistics.

Woodruff (1971):

- Equivalent computational procedure using Taylor series *linearized values*
- Each observational unit gets a linearized value for a particular statistic.
- 2) Compute design-specific variance of the linearized values

Taylor Linearization Approach

Design-Specific Variance: Choice of Sample Designs

■ DESIGN=WOR

Equal probability without-replacement sampling at *each stage* (finite population corrections)

■ DESIGN=UNEQWOR

Unequal probability without-replacement sampling at *first stage* (Yates-Grundy-Sen variance estimator)

■ DESIGN=WR

With-replacement sampling at first stage (this is referred to as the *between-cluster variance estimator*)

- *Most common choice*, as long as low sampling fractions at first stage
- Allows for any sample design within each PSU (*e.g.*, additional stages of sampling, with equal or unequal probabilities of selection)
- Stratification allowed with all designs, even if the sample is not clustered

Between-Cluster Variance Estimator (DESIGN=WR)

Goal is to estimate $Var(\hat{\theta})$:

$$\hat{\theta} = F(X,Y)$$
 where X and Y are linear statistics
 $Z_{ij} = Linearized$ value of $\hat{\theta}$ for unit-ij
 $= (\partial F_X)x_{ij} + (\partial F_Y)y_{ij}$

For a proportion, $\hat{p} = \frac{Y}{X}$,

$$Z_{ij} = w_i (y_{ij} - \hat{p}) / \sum_{i=1}^n \sum_{j=1}^{m_i} w_{ij}$$

$$Z_i = \sum_{j=1}^{m_i} Z_{ij}$$
 PSU Totals

$$\bar{Z} = \frac{1}{n} \sum_{i=1}^{n} Z_i$$
 Mean of PSU Totals

$$\hat{Var}(\hat{\theta}) = \frac{n}{n-1} \sum_{i=1}^{n} (Z_i - \bar{Z})^2$$

Logistic Regression Model:

$$p_{ijk}(\boldsymbol{\beta}) = \Pr(y_{ijk} = 1 | \boldsymbol{x}_{ijk}, \boldsymbol{\beta}) = \left[1 + \exp(-\boldsymbol{x}'_{ijk}\boldsymbol{\beta})\right]^{-1}$$

where:

i = stratum; j = PSU or cluster; k = observation within the cluster

$$\mathbf{x}_{ijk} = (1, x_{1,ijk}, ..., x_{q,ijk})' =$$
 vector of regression effects (stratum-,
PSU-, and observation-specific)

$$\boldsymbol{\beta} = (\boldsymbol{\beta}_0, \boldsymbol{\beta}_1, ..., \boldsymbol{\beta}_q)^{\prime}$$
 = vector of unknown regression
coefficients

 $y_{ijk} = \begin{cases} 1, \text{ outcome present} \\ 0, \text{ outcome absent} \end{cases}$

Overview of Implicit Taylor Linearization Method (using between-cluster variance estimator)

- 1) Find solutions to weighted pseudo-likelihood equations (identical to SAS PROC LOGISTIC)
- 2) Application of Taylor linearization for implicitly-defined parameter vectors in conjunction with a *between-cluster* variance estimation formula (Binder, 1983)

Yields consistent estimator for $Var(\hat{\beta})$

Maximize the Log-Likelihood

Weighted Score Equations:

$$\boldsymbol{U}(\boldsymbol{\beta}) = \frac{\partial \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{i} \sum_{j} \sum_{k} w_{ijk} \boldsymbol{x}'_{ijk} y_{ijk} - \sum_{i} \sum_{j} \sum_{k} w_{ijk} \boldsymbol{x}'_{ijk} p_{ijk} (\boldsymbol{\beta})$$

Solve via iteration: $U(\beta) = 0 \Rightarrow \hat{\beta}$

$$\hat{p}_{ijk} = \left[1 + \exp\left(-\boldsymbol{x}_{ijk}^{\prime} \hat{\boldsymbol{\beta}}\right)\right]^{-1}$$

Binomial-based estimates are asymptotically normally distributed and consistent, even under cluster sampling. Standard regression coefficient estimates are robust to violations of model assumptions.

Weighted Sample Information Matrix:

$$J = -\left[\frac{\partial^2 \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}^2}\right] = \sum_i \sum_j \sum_k \boldsymbol{x}'_{ijk} \boldsymbol{x}_{ijk} w_{ijk} \hat{d}_{ijk} ,$$

where $\hat{d}_{ijk} = \hat{p}_{ijk} (1 - \hat{p}_{ijk})$

Under Cluster Sampling (Intracluster correlation *≠* 0)

 $\hat{Var}(\hat{\boldsymbol{\beta}}) \neq J^{-1}$ Not a Consistent Estimator (Biased)

Taylor Linearization for Implicitly-Defined Parameter Vectors (*Robust or Sandwich Estimator*, Binder 1983)

$$\hat{Var}(\hat{\boldsymbol{\beta}}) = (\boldsymbol{J}^{-1}) \hat{Var}[\hat{\boldsymbol{U}}(\hat{\boldsymbol{\beta}})] (\boldsymbol{J}^{-1})^{\prime}$$

where

$$\hat{U}(\hat{\beta}) = \frac{\partial \log L(\beta)}{\partial \beta} \Big|_{\beta = \hat{\beta}}$$
 Estimating Equations (Score Function)
$$J = \frac{\partial \hat{U}(\hat{\beta})}{\partial \hat{\beta}}$$
 Sample Information Matrix

Outside Term

$$J^{-1}$$
 is the *model-based* (or *naive*) variance estimate

Inside Term

$$\hat{Var}[\hat{U}(\hat{\boldsymbol{\beta}})]$$
 is the *design-specific variance correction*

Estimate $Var[\hat{U}(\hat{\beta})]$ Using Between-Cluster Variance:

1) Score equations are simple linear functions of the observations

$$\hat{U}(\hat{\boldsymbol{\beta}}) = \sum_{i} \sum_{j} \sum_{k} \hat{U}(\boldsymbol{Z}_{ijk}; \boldsymbol{\beta})$$

Linearized variate vector for the *ijk*-th unit:

$$\boldsymbol{Z}_{ijk} = w_{ijk} \boldsymbol{x}'_{ijk} (y_{ijk} - \hat{\boldsymbol{p}}_{ijk})$$

2) Compute *between-PSU* within-stratum variance estimate for a vector of linear statistics:

Accumulations of linearized variate vectors at PSU level

$$\mathbf{Z}_{ij} = \sum_{k} \mathbf{Z}_{ijk}$$
, $k = 1, \dots, m_{ij}$

Form Between-PSU Within-Stratum Mean Square Matrix

$$S_z = \sum_i n_i S_{zi}$$
, $n_i = \#$ PSU's in stratum i

With sample mean squares and cross-products matrix:

$$S_{zi} = \sum_{j} (Z_{ij} - \overline{Z_{i}}) (Z_{ij} - \overline{Z_{i}})' / (n_{i} - 1)$$
$$\overline{Z_{i}} = \sum_{j} Z_{ij} / n_{i}.$$

Estimated Cluster Covariance Matrix for $\hat{\beta}$:

$$\hat{Var}(\hat{\boldsymbol{\beta}}) = (\boldsymbol{J}^{-1}) \boldsymbol{S}_{z} (\boldsymbol{J}^{-1})^{\prime}$$

Null Hypothesis:

 $H_0: C\beta = 0$ vs. $H_1: C\beta \neq 0$

C =Contrast Matrix

Wald Test Statistic:

 $\chi^{2} = \left[C\hat{\beta}\right]^{/} \left[C\hat{Var}(\hat{\beta})C^{/}\right]^{-1} \left[C\hat{\beta}\right]$ ~ χ^{2}_{c} , where c = rank of C

Small-Sample Modifications to Wald Chi-Square:

- Wald chi-square too liberal when DF associated with the hypothesis is large compared to the DF available for estimating variance of regression coefficients (#clusters-#strata) (Thomas and Rao, 1987)
- Satterthwaite-corrected Chi-Square (Rao and Scott, 1987)
- Adjusted Wald F-statistic (Folsom, 1974; Fellegi, 1980)

Implicit Linearization Method Also Used For:

- 1) Proportional Hazards Model (Cox Regression) Binder, 1992
- 2) Ordinary Linear Regression: parameter vector *explicitly* defined

$$\hat{\boldsymbol{\beta}} = (\boldsymbol{X}^{\prime} \boldsymbol{W} \boldsymbol{X})^{-1} \boldsymbol{X}^{\prime} \boldsymbol{W} \boldsymbol{Y}$$

Where W is a diagonal matrix with diagonal elements equal to the sample member weights.

Estimating Equations (Normal Equations):

$$\hat{\mathbf{U}}(\hat{\boldsymbol{\beta}}) = X'WX\hat{\boldsymbol{\beta}} - X'WY$$
$$\mathbf{J} = \frac{\partial \hat{\mathbf{U}}(\hat{\boldsymbol{\beta}})}{\partial \hat{\boldsymbol{\beta}}} = X'WX$$

Linearized Variate Vector for $\hat{U}(\hat{\beta})$

$$Z_{ijk} = \left[x_{ijk}' (x_{ijk} \hat{\beta} - y_{ijk}) \right] w_{ijk}$$

Replication Methods for Complex Survey Data

Quenouille (1956): Reducing bias in estimation Tukey (1958): Approximate confidence intervals

Start With Given Point Estimator:

Descriptive statistics (*e.g.*, means, proportions) Regression parameter vectors

- Use consistent estimators of location parameters
- Assumes with-replacement sampling of PSUs (same as the *between-cluster* estimator)

Prevalence of drug use in a complex sample survey:

$$\hat{p} = \frac{\sum_{i=1}^{n} \sum_{j=1}^{m_{i}} w_{ij} y_{ij}}{\sum_{i=1}^{n} \sum_{j=1}^{m_{i}} w_{ij}} = \frac{Estimated number drug users}{Estimated population size}$$

An estimate based on all PSU's *except the k-th* is as follows:

$$\hat{p}_{(k)} = \frac{\sum_{i \neq k}^{n} \sum_{j=1}^{m_{i}} w_{ij} y_{ij}}{\sum_{i \neq k}^{n} \sum_{j=1}^{m_{i}} w_{ij}}$$

Jackknife Variance Estimate for \hat{p} **:**

$$\hat{\sigma}_{JK}^2 = \frac{n-1}{n} \sum_{k=1}^n \left[\hat{p}_{(k)} - \hat{p}_{(.)} \right]^2$$

where $\hat{p}_{(.)}$ is the average of the Jackknife estimates:

$$\hat{p}_{(.)} = \frac{\sum_{k=1}^{n} \hat{p}_{(k)}}{n}$$

٠

Covariance of Regression Parameters

Start With Given Point Estimator $\hat{\beta}$:

Estimated parameter vector obtained by naively assuming the observations within a cluster are independent

Solution to any score estimating equation of the form

$$\mu(\hat{\beta}) = \sum_{i=1}^{n} \mu_{i}(\hat{\beta}) = 0$$

where $\mu_i(\hat{\beta})$ is the contribution to the "score" vector from the *i*-th cluster.

Example

Logistic score equations under binomial likelihood

$$\boldsymbol{U}(\boldsymbol{\beta}) = \frac{\partial \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = \sum_{i} \sum w_{ij} \boldsymbol{x}_{ij}' y_{ij} - \sum_{i} \sum_{j} w_{ij} \boldsymbol{x}_{ij}' p_{ij}(\boldsymbol{\beta})$$

Regression Parameters (continued)

As long as the model for the marginal mean is correctly specified, the MLE $\hat{\beta}$ is asymptotically consistent and normally distributed

Jackknife Variance Estimator For $\hat{\beta}$

$$Var_{JK}(\hat{\boldsymbol{\beta}}) = \left(\frac{n-p}{n}\right) \sum_{i=1}^{n} (\hat{\boldsymbol{\beta}}_{-i} - \hat{\boldsymbol{\beta}}_{.})(\hat{\boldsymbol{\beta}}_{-i} - \hat{\boldsymbol{\beta}}_{.})$$

where

p = number of parameters in the model,

 $\hat{\boldsymbol{\beta}}_{-i}$ = estimate of $\boldsymbol{\beta}$ obtained by deleting the m_i observations in PSU *i* and solving the estimating equations via the Newton-Raphson algorithm, and

$$\hat{\boldsymbol{\beta}}_{.}$$
 = is the average of the $\hat{\boldsymbol{\beta}}_{-i}$.

PSU's are removed sequentially and with-replacement

JK variance estimator is consistent for estimating the asymptotic variance of $\hat{\beta}$

Assumptions and Validity for Taylor Linearization and Jackknife

- PSUs are statistically independent
- No strict distributional assumptions for the response of interest
- Yields consistent estimates of the variance as the number of PSUs tends to infinity
- Method is valid for any underlying intra-PSU correlation structure, as long as PSUs are statistically independent

Balanced Repeated Replication (BRR)

BRR Variance Estimation for Complex Sample Surveys

McCarthy, PJ 1966, Vital and Health Statistics, 2(14), NCHS 1969, Review of the International Statistical Institute 37, 239-264.
Wolter, KM 1985, Introduction to Variance Estimation. Springer-Verlag.

Usually assumes PSUs selected with-replacement

Allows for any unbiased sampling method within PSUs

BRR Variance Estimation

How Does It Work?

Balanced Half Samples

- Assume two (2) PSUs are selected with-replacement from each of *L* strata (more than 2 selections can be made, but is more complicated to explain)
- Form *G* half-sample replicates, where each half-sample is formed by selecting one of the two PSUs from each stratum based on a Hadamard matrix (Plackett and Burman, 1946)
- Let $\hat{\theta}_g$ be the estimate of the parameter based on the *g*-th half-sample

$$Var(\hat{\theta}) = \frac{1}{G} \sum_{g=1}^{G} (\hat{\theta}_g - \hat{\theta})^2$$

where

 $\hat{\theta}$ = estimate based on the full sample

BRR Weights: REPWGT *variables*;

- Variables whose values are the BRR replicate weights for each sampled individual (nonnegative or missing)
- Many survey data bases will supply BRR or replicate weights.
 SUDAAN assumes the weights are supplied.

Let $w_{gi} = \begin{cases} replicate weight for sample unit-i in half-sample g \\ 0, if sample unit-i NOT in half-sample g \end{cases}$

• Use
$$w_{gi}$$
 to estimate $\hat{\theta}_g$

For example, the *Total* estimate from replicate-*g*:

$$\hat{Y}_g = \sum_{i=1}^n w_{gi} y_i$$

Possible to develop special weights to account for *without replacement* sampling. Need to consult with a statistician to develop such weights.

Example Comparing the Three Approaches

These data are taken from a one-year longitudinal study of infant feeding practices of participants in the *Special Supplemental Nutrition Program for Women*, *Infants, and Children (WIC)*. A national sample of 42 local agencies (sites) was selected at the first stage and implicitly stratified by region of the country and state within region. Local agencies were paired to form strata. A sample of about 22 pregnant women or new mothers participating in the WIC Program were then selected from each local agency. The participants were interviewed 9 times during each infant's first year of life to gain a complete picture of the feeding patterns of WIC infants. The data consist of one record per WIC respondent.

We use these data to demonstrate the three variance estimation methods in SUDAAN (Taylor linearization, Jackknife, and Balanced Repeated Replication, or BRR). We first estimate descriptive statistics on baby's birth weight and mother's breastfeeding status. Then, we fit a logistic regression model to the incidence of breastfeeding initiation. Point estimates of means, proportions, and regression coefficients are equivalent for all three approaches. Variance estimates are similar in most situations. This example does not point to favoring one method vs. another for variance estimation.

Example Comparing the Three Approaches

Descriptive Statistics

Proportion	~ -		S	tandard Erro	ors
Initiating Breastfeeding	Sample Size	Percentage	Taylor	Jackknife	BRR
TOTAL	953	54 %	3.1 %	3.1 %	3.1 %
White	480	51 %	3.2 %	3.2 %	3.2 %
African American	225	32 %	4.5 %	4.5 %	4.7 %
Latina	190	83 %	3.5 %	3.5 %	3.5 %
Other	58	63 %	12.2 %	12.8 %	15.0 %

Logistic Regression

Breastfeeding = baby's weight + sex + race + education + marital status

		Tayle	or	Jackki	nife	BRI	R
Effect	df	Chi-square	P-value	Chi-square	P-value	Chi-square	P-value
Race	3	39.9	0.000	38.7	0.000	36.5	0.000
Education	2	14.2	0.001	14.0	0.001	11.9	0.003
Marital Status	1	36.4	0.000	35.6	0.000	29.9	0.000
Sex	1	0.5	0.463	0.5	0.464	0.5	0.467
Baby's Weight	1	8.3	0.004	8.0	0.005	7.9	0.005

DESCRIPT Programming Statements for Taylor Linearization (DESIGN=WR)

SUDAAN Software for the Statistical Analysis of Correlated Data Copyright Research Triangle Institute May 1997 Beta Test Release 7.5 1 PROC DESCRIPT DATA="WIC" FILETYPE=SAS DESIGN=WR DEFT2 MERGEHI; NEST STRATUM SITE; 2 WEIGHT ANALWGT1; 3 VAR BRFDINIT BABYWGT; 4 5 SUBGROUP RACEMOM EDUC MRTLSTAT BABYSEX; б LEVELS 4 3 2 2; 7 SETENV LABWIDTH=28 COLSPCE=1 COLWIDTH=10 LINESIZE=78 DECWIDTH=4 PAGESIZE=60; PRINT NSUM="SAMPLE SIZE" WSUM="POPULATION SIZE" MEAN SEMEAN="S.E." 8 DEFFMEAN="DESIGN EFFECT" / STYLE=NCHS NSUMFMT=F6.0 WSUMFMT=F10.0 DEFFMEANFMT=F6.2 SEMEANFMT=F7.4; TITLE " " "STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS" " " 9 "TAYLOR LINEARIZATION VARIANCE ESTIMATION" " "; Opened SAS data file C:\TERA\EXAMPLES\WIC.SSD for reading. Number of observations read : 953 Weighted count : 506279 Denominator degrees of freedom : 21

Note the Strata and PSU variables *STRATUM* and *SITE* on the NEST statement, and the analysis weight variable ANALWGT1 on the WEIGHT statement. There are 953 WIC participants on the file, summing to an estimated 506,279 participants in the US population (this is a slight underestimate, since some sites were not included in this example).

ime: 14:41:15	Research The DE				Page : . Table : 1
ariance Estimation Met	hod: Taylor a	Series (WR)			
FUDY OF BREAST-FEEDING	PATTERNS AM	ONG WIC PART	ICIPANTS		
AYLOR LINEARIZATION VAL	RIANCE ESTIM	ATION			
Y: Mother Race.					
ariable					
Mother Race	SAMPLE	POPULATION			DESIGN
	SIZE	SIZE	Mean	S.E.	EFFECT
eastfeeding Initiation	n				
reastfeeding Initiation Total	n 953	506279	0.5370	0.0310	3.67
reastfeeding Initiation Total White	n 953 480	506279 254567	0.5370 0.5104	0.0310 0.0318	3.67 1.94
reastfeeding Initiation Total White African American	n 953 480 225	506279 254567 123217	0.5370 0.5104 0.3202	0.0310 0.0318 0.0445	3.67 1.94 2.05
Teastfeeding Initiation Total White African American Latina	n 953 480 225 190	506279 254567 123217 106306	0.5370 0.5104 0.3202 0.8324	0.0310 0.0318 0.0445 0.0348	3.67 1.94 2.05 1.65
eastfeeding Initiation Total White African American Latina Other	n 953 480 225 190	506279 254567 123217	0.5370 0.5104 0.3202 0.8324	0.0310 0.0318 0.0445 0.0348	3.67 1.94 2.05 1.65
eastfeeding Initiation Total White African American Latina Other by Weight (ozs.)	n 953 480 225 190 58	506279 254567 123217 106306 22189	0.5370 0.5104 0.3202 0.8324 0.6319	0.0310 0.0318 0.0445 0.0348 0.1220	3.67 1.94 2.05 1.65 3.71
Total White African American Latina Other Dy Weight (OZS.) Total	n 953 480 225 190 58 952	506279 254567 123217 106306 22189 505897	0.5370 0.5104 0.3202 0.8324 0.6319 116.6232	0.0310 0.0318 0.0445 0.0348 0.1220 0.9249	3.67 1.94 2.05 1.65 3.71 1.90
eastfeeding Initiation Total White African American Latina Other by Weight (ozs.) Total White	n 953 480 225 190 58 952 480	506279 254567 123217 106306 22189 505897 254567	0.5370 0.5104 0.3202 0.8324 0.6319 116.6232 118.6759	0.0310 0.0318 0.0445 0.0348 0.1220 0.9249 1.0306	3.67 1.94 2.05 1.65 3.71 1.90 1.17
reastfeeding Initiation Total White African American Latina Other by Weight (ozs.) Total White	n 953 480 225 190 58 952 480 225	506279 254567 123217 106306 22189 505897 254567	0.5370 0.5104 0.3202 0.8324 0.6319 116.6232 118.6759 108.6932	0.0310 0.0318 0.0445 0.0348 0.1220 0.9249 1.0306 1.6144	3.67 1.94 2.05 1.65 3.71 1.90 1.17 1.29

Here we see that breastfeeding initiation and baby's birth weight are both highest among Latina women and lowest among African American women. The standard errors are obtained through Taylor linearization.

Date: 07-07-97 Fime: 14:41:15		Triangle Ins CRIPT Proces			Page Table
Variance Estimation Me	thod: Taylor Se	eries (WR)			
STUDY OF BREAST-FEEDING	G PATTERNS AMON	NG WIC PARTI	ICIPANTS		
TAYLOR LINEARIZATION V	ARIANCE ESTIMA	FION			
by: Education.					
Variable					
Education	SAMPLE I	POPULATION			DESIGN
	SIZE S	SIZE	Mean	S.E.	EFFECT
Breastfeeding Initiation	 on				
Total		505920	0.5367	0.0310	3.67
< High School	368	212474	0.5187	0.0520	3.98
5		211345			
High School					
High School > High School	185	82101	0.6804	0.0460	
> High School	185	82101	0.6804	0.0460	
> High School		82101 505539			1.80
> High School Baby Weight (ozs.)	951		116.6045	0.9291	1.80 1.91
> High School Baby Weight (ozs.) Total	951 368	505539	116.6045 115.2054	0.9291	1.80 1.91 1.56

Breastfeeding initiation and baby's birth weight among WIC participants is highest among women with more than a high school education.

Date: 07-07-97 Time: 14:41:15		Triangle In: CRIPT Proced			Page Table	
Variance Estimation Method	: Taylor S	eries (WR)				
STUDY OF BREAST-FEEDING PA	TTERNS AMO	NG WIC PART	ICIPANTS			
TAYLOR LINEARIZATION VARIA	NCE ESTIMA	TION				
by: Marital Status.						
 Variable						
	SAMPLE	POPULATION			DESIGN	
		SIZE			EFFECT	
Breastfeeding Initiation						
Total	952	505897	0.5367	0.0310	3.68	
Currently Married	462	246325	0.6393	0.0317	2.02	
Not Currently Married	490	259572	0.4394	0.0330	2.16	
Baby Weight (ozs.)						
Total	952	505897	116.6232	0.9249	1.90	
	462	246325	119.0474	1.5696	2.55	
Currently Married	402					

Breastfeeding initiation and baby's birth weight are also higher among those currently married compared to those not currently married.

Date: 07-07-97 Time: 14:41:15		riangle Ins RIPT Procec			Page Table
Variance Estimation Me	thod: Taylor Ser	ries (WR)			
STUDY OF BREAST-FEEDIN	G PATTERNS AMON	G WIC PARTI	ICIPANTS		
TAYLOR LINEARIZATION V	ARIANCE ESTIMAT	ION			
by: Baby Sex.					
Variable					
Baby Sex	SAMPLE PO	OPULATION			DESIGN
	SIZE S		Mean		
Breastfeeding Initiati					
Total	953	506279	0.5370	0.0310	3.67
Воу	495	254670	0.5295	0.0366	2.66
Girl	458	251609	0.5446	0.0374	2.58
Baby Weight (ozs.)					
Baby Weight (ozs.) Total	952	505897	116.6232	0.9249	1.90
Baby Weight (ozs.) Total Boy	952 494		116.6232 118.5878		

Breastfeeding initiation is comparable for boy vs. girl babies.

10 PROC DESCRIPT DATA="WIC" FILETYPE=SAS DESIGN=JACKKNIFE MERGEHI;
11 NEST STRATUM SITE;
12 WEIGHT ANALWGT1;
13 VAR BRFDINIT BABYWGT;
14 SUBGROUP RACEMOM EDUC MRTLSTAT BABYSEX;
15 LEVELS 4 3 2 2;
16 SETENV LABWIDTH=28 COLSPCE=1 COLWIDTH=10 LINESIZE=78 DECWIDTH=4 PAGESIZE=60;
<pre>17 PRINT NSUM="SAMPLE SIZE" WSUM="POPULATION SIZE" MEAN SEMEAN="S.E." DEFFMEAN="DESIGN EFFECT" / STYLE=NCHS NSUMFMT=F6.0 WSUMFMT=F10.0 DEFFMEANFMT=F6.2 SEMEANFMT=F7.4;</pre>
18 TITLE " " "STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS" " " "JACKKNIFE VARIANCE ESTIMATION" " ";
Opened SAS data file C:\TERA\EXAMPLES\WIC.SSD for reading.
Number of observations read : 953 Weighted count : 506279 Denominator degrees of freedom : 21

For *DESIGN=JACKKNIFE*, we keep the NEST and WEIGHT statements as they were for Taylor linearization.

Date: 07-07-97 Fime: 14:41:15		Triangle In SCRIPT Proce			Page : Table :
Variance Estimation Met	hod: Jackkni:	fe			
STUDY OF BREAST-FEEDING	PATTERNS AM	ONG WIC PART	ICIPANTS		
JACKKNIFE VARIANCE ESTI	MATION				
by: Mother Race.					
Variable					
Mother Race	SAMPLE	POPULATION			DESIGN
	SIZE	SIZE	Mean	S.E.	EFFECT
Breastfeeding Initiation					
Total		506279	0.5370	0.0310	3.68
White	480	254567	0.5104	0.0318	1.94
African American	225	123217	0.3202	0.0449	2.07
African American Latina	225 190		0.3202 0.8324		
	190		0.8324	0.0352	1.68
Latina	190	106306	0.8324	0.0352	1.68
Latina Other	190	106306 22189	0.8324 0.6319	0.0352 0.1279	1.68 4.01
Latina Other Baby Weight (ozs.)	190 58	106306 22189 505897	0.8324 0.6319 116.6232	0.0352 0.1279 0.9253	1.68 4.01 1.90
Latina Other Baby Weight (ozs.) Total	190 58 952	106306 22189 505897 254567	0.8324 0.6319 116.6232	0.0352 0.1279 0.9253 1.0338	1.68 4.01 1.90 1.17
Latina Other Baby Weight (ozs.) Total White	190 58 952 480 225	106306 22189 505897 254567	0.8324 0.6319 116.6232 118.6759 108.6932	0.0352 0.1279 0.9253 1.0338 1.6389	1.68 4.01 1.90 1.17 1.33

Variance estimates for all Jackknife results are similar to Taylor linearization.

te: 07-07-97 me: 14:41:15	Research Tr The DESCE	-			Page Table	
IME• 14•41•15	THE DESCR	(IPI PIOCE)	Jure		Table	• 2
riance Estimation Me	thod: Jackknife					
UDY OF BREAST-FEEDIN	G PATTERNS AMONO	G WIC PART	ICIPANTS			
CKKNIFE VARIANCE EST	τμάτιον					
·: Education.						
ariable						
Education	SAMPLE PO	OPULATION			DESIGN	
			Mean	S.E.		
reastfeeding Initiati Total		505920	0.5367	0.0310	3.68	
< High School			0.5187			
High School	399					
> High School			0.6804			
aby Weight (ozs.)						
Total	951	505539	116.6045	0.9294	1.92	
	368		115.2054	1.3386	1.56	
< High School	500					
< High School High School			116.5428	1.3052	1.69	

Date: 07-07-97 Time: 14:41:15		Triangle In SCRIPT Procee			Page Table
Variance Estimation Method	: Jackknif	Ee			
STUDY OF BREAST-FEEDING PA	TTERNS AMO	ONG WIC PART	ICIPANTS		
JACKKNIFE VARIANCE ESTIMAT	ION				
by: Marital Status.					
Variable					
Marital Status	SAMPLE	POPULATION			DESIGN
		SIZE			EFFECT
Breastfeeding Initiation					
Total	952	505897	0.5367	0.0311	3.69
Currently Married	462	246325	0.6393	0.0318	2.02
Not Currently Married	490	259572	0.4394	0.0330	2.16
Baby Weight (ozs.)					
Total	952	505897	116.6232	0.9253	1.90
	462	246325	119.0474	1.5704	2.55
Currently Married	102				

Date: 07-07-97 Time: 14:41:15					Page Table	
Variance Estimation Me	ethod: Jackknife	1				
TUDY OF BREAST-FEEDI	NG PATTERNS AMON	G WIC PARTI	ICIPANTS			
JACKKNIFE VARIANCE ES	TIMATION					
by: Baby Sex.						
/ariable						
Baby Sex	SAMPLE P	OPULATION			DESIGN	
	SIZE S	IZE	Mean	S.E.	EFFECT	
reastfeeding Initiat:	ion					
Total	953	506279	0.5370	0.0310	3.68	
Воу	495	254670	0.5295	0.0366	2.66	
Girl	458	251609	0.5446	0.0374	2.58	
aby Weight (ozs.)						
aby Weight (ozs.) Total	952	505897	116.6232	0.9253	1.90	
Baby Weight (ozs.) Total Boy		505897 254288				

```
19 PROC DESCRIPT DATA="WIC" FILETYPE=SAS DESIGN=BRR MERGEHI;
20 WEIGHT ANALWGT1;
21 REPWGT RPL001--RPL024;
22 VAR BRFDINIT BABYWGT;
23 SUBGROUP RACEMOM EDUC MRTLSTAT BABYSEX;
24 LEVELS 4
                  3
                        2
                                 2;
25 SETENV LABWIDTH=28 COLSPCE=1 COLWIDTH=10 LINESIZE=78 DECWIDTH=4 PAGESIZE=60;
26 PRINT NSUM="SAMPLE SIZE" WSUM="POPULATION SIZE" MEAN SEMEAN="S.E."
         DEFFMEAN="DESIGN EFFECT" / STYLE=NCHS NSUMFMT=F6.0 WSUMFMT=F10.0
         DEFFMEANFMT=F6.2 SEMEANFMT=F7.4;
27 TITLE " " "STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS" " "
         "BRR VARIANCE ESTIMATION" " ";
Opened SAS data file C:\TERA\EXAMPLES\WIC.SSD for reading.
Number of observations read :
                                 953
                                         Weighted count : 506279
Denominator degrees of freedom :
                                   24
```

For *DESIGN=BRR*, we *remove the NEST statement* and include a statement for the known replicate weights (the REPWGT statement). There are 24 replicate weights in this study.

Date: 07-07-97 Cime: 14:41:15		-			-	
Variance Estimation Method	: BRR					
TUDY OF BREAST-FEEDING PA	TTERNS AMO	ONG WIC PART	ICIPANTS			
BRR VARIANCE ESTIMATION						
by: Mother Race.						
Jariable						
Mother Race	SAMPLE	POPULATION			DESIGN	
	SIZE	SIZE	Mean	S.E.	EFFECT	
Breastfeeding Initiation						
Total	953	506279	0.5370	0.0311	3.70	
White	480	254567	0.5104	0.0321	1.98	
African American	225	123217	0.3202	0.0469	2.27	
Latina	190	106306	0.8324	0.0351	1.67	
Other	58	22189	0.6319	0.1496	5.49	
Baby Weight (ozs.)						
Total	952	505897	116.6232	0.9215	1.88	
White	480	254567	118.6759	1.0957	1.32	
African American	225	123217	108.6932	1.6363	1.32	
	189	105924	120.3653	1.3577	1.36	
Latina				5.5170		

Variance estimates based on BRR are similar to Taylor linearization and Jackknife results. Of course, points estimates of the population mean are the same for all three methods.

Date: 07-07-97 Fime: 14:41:15		Triangle In CRIPT Proce			Page : Table :
Variance Estimation Me	thod: BRR				
TUDY OF BREAST-FEEDIN	G PATTERNS AMO	NG WIC PART	ICIPANTS		
BRR VARIANCE ESTIMATIO	N				
v: Education.					
-					
Education	SAMPLE	POPULATION			DESIGN
		SIZE			
Breastfeeding Initiati Total		505920	0 5367	0 0311	3 70
< High School		212474			
High School		211345			
> High School		82101			
aby Weight (ozs.)					
Total	951	505539	116.6045	0.9255	1.90
< High School	368	212474	115.2054	1.3756	1.65
	2.2.2	211245	116.5428	1.3217	1.73
High School	399	211343	110.0100		

Date: 07-07-97 Time: 14:41:15	Research Triangle Institute The DESCRIPT Procedure				
Variance Estimation Method	: BRR				
STUDY OF BREAST-FEEDING PA	TTERNS AM	ONG WIC PART	ICIPANTS		
BRR VARIANCE ESTIMATION					
by: Marital Status.					
Variable					
	SIZE	POPULATION SIZE	Mean	S.E.	DESIGN EFFECT
Breastfeeding Initiation					
Total	952	505897	0.5367	0.0311	3.71
Currently Married	462	246325	0.6393	0.0316	2.00
Not Currently Married	490	259572	0.4394	0.0338	2.27
Baby Weight (ozs.)					
Total	952	505897	116.6232	0.9215	1.88
Currently Married	462	246325	119.0474	1.6293	2.75
Not Currently Married	490	259572	114.3227	0.9802	1.17

Date: 07-07-97 Time: 14:41:15	Research Triangle Institute Page The DESCRIPT Procedure Table					
Variance Estimation M	lethod: BRR					
STUDY OF BREAST-FEEDI	NG PATTERNS AMO	NG WIC PART	ICIPANTS			
BRR VARIANCE ESTIMATI	ON					
by: Baby Sex.						
 Variable						
Baby Sex	SAMPLE	SAMPLE POPULATION			DESIGN	
	SIZE	SIZE			EFFECT	
Breastfeeding Initiat						
Total	953	506279	0.5370	0.0311	3.70	
Воу	495	254670	0.5295	0.0374	2.77	
Girl	458	251609	0.5446	0.0367	2.49	
Baby Weight (ozs.)						
Total	952	505897	116.6232	0.9215	1.88	
5	494	254288	118.5878	0.9277	1.01	
Воу						

LOGISTIC Modelling Based on Taylor Linearization

```
28 PROC LOGISTIC DATA="WIC" FILETYPE=SAS DESIGN=WR DEFT2 MERGEHI;
29 NEST STRATUM SITE;
30 WEIGHT ANALWGT1;
31 SUBGROUP RACEMOM EDUC MRTLSTAT BABYSEX;
32 LEVELS 4
                    3
                           2 ;
                         2
33 REFLEVEL RACEMOM=1 EDUC=1;
34 MODEL BRFDINIT = BABYWGT BABYSEX RACEMOM EDUC MRTLSTAT;
35 EFFECTS RACEMOM=(0 1 -1 0) / NAME="African Am Vs. Latina";
36 TEST WALDCHI;
37 SETENV COLSPCE=2 LABWIDTH=26 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
38 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" OR LOWOR UPOR
         DF="DF" WALDCHI="WALD CHI-SQ" WALDCHP="P-VALUE"
         / T_BETAFMT=F8.2 DEFTFMT=F6.2 SEBETAFMT=F8.6
          ORFMT=F5.2 LOWORFMT=F6.2 UPORFMT=F6.2
          DFFMT=F7.0 WALDCHIFMT=F8.2 ;
39 TITLE " " "STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS" " "
         "TAYLOR SERIES VARIANCE ESTIMATION";
NOTE: Terms in the MODEL statement have been rearranged
     to follow subgroup order.
Opened SAS data file C:\TERA\EXAMPLES\WIC.SSD for reading.
Number of observations read
                             :
                                    953 Weighted count: 506279
                                    951 Weighted count: 505539
Observations used in the analysis :
Observations with missing values :
                                     2 Weighted count:
                                                               740
Denominator degrees of freedom :
                                      21
Maximum number of estimable parameters for the model is 9
Number of zero responses : 431
Number of non-zero responses :
                               520
Parameters have converged in 4 iterations
R-Square for dependent variable BRFDINIT (Cox & Snell, 1989): 0.171348
```

In the logistic models, we want to see if baby's birth weight, sex, as well as the mother's race, education, and marital status significantly affect breastfeeding initiation. More than half the sample initiated breastfeeding (520 out of 951 non-missing responses). The SUBGROUP and

68 SUDAAN Release 7.5

LEVELS statements define the variables to be treated as categorical, and the REFLEVEL statement changes the default reference levels for two of the categorical covariates from the last level to the first level. The EFFECTS statement directly compares African American women to Latina women.

LOGISTIC Modelling Based on Taylor Linearization

Date: 07-07-97 Fime: 14:41:15		STIC Proce			Table :		
	1110 11091		aur c				
esponse variable BRFDINI	I: Breastfee	ding Initi	ation				
TUDY OF BREAST-FEEDING PA	ATTERNS AMON	IG WIC PARI	ICIPANTS				
AYLOR SERIES VARIANCE EST	TIMATION						
Independent Variables and							
Effects							
	BETA	S.E.	EFFECT	T:BETA=0	P-VALUE		
Intercept	-1.5928	0.427856	0.96	-3.72	0.0013		
Mother Race							
White	0.0000	0.000000					
African American	-0.5410						
Latina	1.7147	0.300657	1.80	5.70	0.0000		
Other	0.4155	0.547545	2.26	0.76	0.4563		
Education							
< High School	0.0000	0.000000					
High School	0.0565	0.210793		0.27	0.7911		
> High School	0.8533	0.298020	1.75	2.86	0.0093		
Marital Status							
Currently Married	0.6843	0.113357	0.57	6.04	0.0000		
Not Currently Married	0.0000	0.000000	•		•		
Baby Sex							
-	-0.1232	0.167831	1.33	-0.73	0.4709		
Girl	0.0000	0.000000					
Baby Weight (ozs.)	0 0000	0 002225	0 00	2.89	0 0000		

From the estimated regression coefficients we see immediately that significantly fewer African American women, but significantly more Latina women, initiated breastfeeding compared to white women. Also, having more than a high school education and being currently married both significantly improved the likelihood of breastfeeding. Finally, as baby's birth weight increased, the likelihood of breastfeeding was significantly increased.

Date: 07-07-97 Time: 14:41:15		-		Page : 2 Table : 1
11110- 14-41-13	THE LOG	FISTIC Proc	euure	
sponse variable BRFDINI	T: Breastfe	eding Init	iation	
UDY OF BREAST-FEEDING P	ATTERNS AMC	NG WIC PAR	TICIPANTS	
AYLOR SERIES VARIANCE ES	TIMATION			
· · · · · · · · · · · · · · · · · · ·				
ontrast		WALD		
	DF	CHI-SQ		
		144.66		
OVERALL MODEL	9	144.00	0.0000	
OVERALL MODEL MODEL MINUS INTERCEPT		135.98		
ODEL MINUS INTERCEPT	8		0.0000	
	8.	135.98	0.0000	
NODEL MINUS INTERCEPT INTERCEPT RACEMOM	8 3	135.98 39.92	0.0000	
NODEL MINUS INTERCEPT	8 3 2	135.98 39.92	0.0000 0.0000 0.0008	
NODEL MINUS INTERCEPT INTERCEPT RACEMOM EDUC	8 3 2 1	135.98 39.92 14.15	0.0000 0.0000 0.0008 0.0000	
NODEL MINUS INTERCEPT ENTERCEPT RACEMOM EDUC IRTLSTAT	8 3 2 1 1	135.98 39.92 14.15 36.44	0.0000 0.0000 0.0008 0.0000 0.4628	

Under Taylor linearization, mother's race, education, marital status, and baby's birth weight were all statistically significant. Also, the user-specified contrast comparing African American women to Latina women was statistically significant.

LOGISTIC Modelling Based on Taylor Linearization

STUDY OF BREAST-FEEDING PA					
	AIIERNS A	MONG WIC	2 PARTICIP	ANTS	
TAYLOR SERIES VARIANCE EST	IMATION				
Independent Variables and					
Effects		Lower			
		95%			
		Limit			
Intercept		0.08			
Mother Race					
White	1.00	1.00	1.00		
African American	0.58	0.36	0.95		
Latina		2.97			
Other	1.52	0.49	4.73		
Education					
< High School	1.00	1.00	1.00		
High School	1.06	0.68	1.64		
> High School	2.35	1.26	4.36		
Marital Status					
Currently Married	1.98	1.57	2.51		
Not Currently Married	1.00	1.00	1.00		
Baby Sex					
Воу		0.62			
Girl	1.00	1.00	1.00		
Baby Weight (ozs.)	1.01	1.00	1.02		

The estimated odds ratios and 95% confidence limits indicate that, for example:

- the odds of initiated breastfeeding are increased by more than five-fold for Latina women vs. white women
- the odds are reduced by half in African American women vs. white women
- the odds are approximately doubled for women who are currently married as well as for women with more than a high school education.

```
40 PROC LOGISTIC DATA="WIC" FILETYPE=SAS DESIGN=JACKKNIFE MERGEHI;
41 NEST STRATUM SITE;
42 WEIGHT ANALWGT1;
43 SUBGROUP RACEMOM EDUC MRTLSTAT BABYSEX;
44 LEVELS 4 3 2 2 ;
45 REFLEVEL RACEMOM=1 EDUC=1;
46 MODEL BRFDINIT = BABYWGT BABYSEX RACEMOM EDUC MRTLSTAT;
47 EFFECTS RACEMOM=(0 1 -1 0) / NAME="African Am Vs. Latina";
48 TEST WALDCHI;
49 SETENV COLSPCE=2 LABWIDTH=26 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
50 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" OR LOWOR UPOR
         DF="DF" WALDCHI="WALD CHI-SQ" WALDCHP="P-VALUE"
          / T_BETAFMT=F8.2 DEFTFMT=F6.2 SEBETAFMT=F8.6
          ORFMT=F5.2 LOWORFMT=F6.2 UPORFMT=F6.2
          DFFMT=F7.0 WALDCHIFMT=F8.2 ;
51 TITLE " " "STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS" " "
          "JACKKNIFE VARIANCE ESTIMATION";
NOTE: Terms in the MODEL statement have been rearranged
      to follow subgroup order.
Opened SAS data file C:\TERA\EXAMPLES\WIC.SSD for reading.
Number of observations read
                              :
                                     953 Weighted count: 506279
Observations used in the analysis :951Weighted count:505539Observations with missing values :2Weighted count:740
Denominator degrees of freedom :
                                       21
Maximum number of estimable parameters for the model is 9
Number of zero responses : 431
Number of non-zero responses :
                                520
Parameters have converged in 4 iterations
R-Square for dependent variable BRFDINIT (Cox & Snell, 1989): 0.171348
```

The results of logistic modelling using the Jackknife variance estimation method are very similar to Taylor linearization.

Date: 07-07-97 Time: 14:41:15		STIC Proce			Table :		
Response variable BRFDINI	: Breastfee	ding Initi	ation				
STUDY OF BREAST-FEEDING PA	ATTERNS AMON	IG WIC PART	ICIPANTS	5			
JACKKNIFE VARIANCE ESTIMAT							
JACKENIFE VARIANCE ESTIMA	IION						
Independent Variables and							
Effects			DESIGN				
				T:BETA=0			
Intercept				-3.68			
Mother Race							
White	0.0000	0.000000					
African American	-0.5410	0.235709	1.69	-2.30	0.0321		
Latina	1.7147	0.306057	1.96	5.60	0.0000		
Other	0.4155	0.601481	3.00	0.69	0.4972		
Education							
< High School	0.0000	0.000000					
High School	0.0565	0.215382	1.83	0.26	0.7955		
> High School	0.8533	0.302420	1.92	2.82	0.0102		
Marital Status							
Currently Married	0.6843	0.114777	0.59	5.96	0.0000		
Not Currently Married	0.0000	0.000000					
Baby Sex							
Воу	-0.1232	0.168455	1.35	-0.73	0.4725		
Girl	0.0000	0.000000			•		
				2.83	0 01 01		

Date: 07-07-97				
Time: 14:41:15	The LOG	ISTIC Proce	edure	Table : 1
Response variable BRFDINIT	Breastfe	eding Init:	iation	
STUDY OF BREAST-FEEDING PAT	TTERNS AMO	NG WIC PAR	FICIPANTS	
JACKKNIFE VARIANCE ESTIMAT	ION			
Contrast		WALD		
	DF	CHI-SQ		
OVERALL MODEL	 9	142.32	0.0000	
	8	134.01	0.000	
MODEL MINUS INTERCEPT	0	101.01		
MODEL MINUS INTERCEPT INTERCEPT				
INTERCEPT RACEMOM	3		0.0000	
INTERCEPT	3 2	38.67	0.0000 0.0009	
INTERCEPT RACEMOM EDUC MRTLSTAT	3 2 1	38.67 14.01	0.0000 0.0009 0.0000	
INTERCEPT RACEMOM EDUC	3 2 1 1	38.67 14.01 35.55	0.0000 0.0009 0.0000 0.4644	

	ATTERNS A	MONG WIC	PARTICIPAN	TS	
ACKKNIFE VARIANCE ESTIMA	TION				
Independent Variables and					
Effects		Lower	Upper		
		95%			
		Limit			
 Intercept		0.08			
Mother Race					
White	1.00	1.00	1.00		
African American		0.36			
Latina	5.55	2.94	10.50		
Other	1.52	0.43	5.29		
Education					
< High School	1.00	1.00 0.68	1.00		
High School	1.06	0.68	1.66		
> High School	2.35	1.25	4.40		
Marital Status					
Currently Married	1.98	1.56	2.52		
Not Currently Married	1.00	1.00	1.00		
Baby Sex					
Воу	0.88	0.62 1.00	1.25		
Girl	1.00	1.00	1.00		
Baby Weight (ozs.)	1.01	1.00	1.02		

```
52 PROC LOGISTIC DATA="WIC" FILETYPE=SAS DESIGN=BRR MERGEHI;
53 WEIGHT ANALWGT1;
54 REPWGT RPL001--RPL024;
55 SUBGROUP RACEMOM EDUC MRTLSTAT BABYSEX ;
56 LEVELS 4 3 2 2 ;
57 REFLEVEL RACEMOM=1 EDUC=1;
58 MODEL BRFDINIT = BABYWGT BABYSEX RACEMOM EDUC MRTLSTAT;
59 EFFECTS RACEMOM=(0 1 -1 0) / NAME="African Am Vs. Latina";
60 TEST WALDCHI;
61 SETENV COLSPCE=2 LABWIDTH=26 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
62 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" OR LOWOR UPOR
         DF="DF" WALDCHI="WALD CHI-SQ" WALDCHP="P-VALUE"
         / T_BETAFMT=F8.2 DEFTFMT=F6.2 SEBETAFMT=F8.6
          ORFMT=F5.2 LOWORFMT=F6.2 UPORFMT=F6.2
          DFFMT=F7.0 WALDCHIFMT=F8.2 ;
63 TITLE " " "STUDY OF BREAST-FEEDING PATTERNS AMONG WIC PARTICIPANTS" " "
         "BRR VARIANCE ESTIMATION";
NOTE: Terms in the MODEL statement have been rearranged
     to follow subgroup order.
Opened SAS data file C:\TERA\EXAMPLES\WIC.SSD for reading.
                            :
Number of observations read
                                   953 Weighted count: 506279
Observations used in the analysis : 951 Weighted count: 505539
Observations with missing values :
                                     2 Weighted count:
                                                             740
Denominator degrees of freedom :
                                     24
Maximum number of estimable parameters for the model is 9
Number of zero responses : 431
Number of non-zero responses :
                               520
Parameters have converged in 4 iterations
R-Square for dependent variable BRFDINIT (Cox & Snell, 1989): 0.171348
```

Logistic modelling results using BRR variance estimation methods are similar to those based on Taylor linearization and the Jackknife.

Time: 14:41:15	The LOGI	STIC Proce	dure		Table :		
Response variable BRFDINII	: Breastfee	ding Tniti	ation				
	. Dreabtree	aing initi	acton				
STUDY OF BREAST-FEEDING PA	TTERNS AMON	IG WIC PART	ICIPANTS				
BRR VARIANCE ESTIMATION							
Independent Variables and							
Effects							
	BETA	S.E.	EFFECT	T:BETA=0	P-VALUE		
 Intercept				-3.60			
Mother Race	1.5520	0.112512	0.90	5.00	0.0014		
White	0.0000	0.00000					
African American				-2.32			
Latina	1.7147	0.307665	1.98	5.57	0.0000		
Other				0.33			
Education							
< High School	0.0000	0.00000					
High School	0.0565	0.219845	1.90	0.26	0.7992		
> High School				2.74			
Marital Status							
Currently Married	0.6843	0.125232	0.71	5.46	0.0000		
Not Currently Married	0.0000	0.000000					
Baby Sex							
Воу	-0.1232	0.169507	1.37	-0.73	0.4742		
Girl	0.0000	0.00000		2.81			
Baby Weight (ozs.)		0.003416			0.0097		

Date: 07-07-97 Time: 14:41:15		Triangle In SISTIC Proc		Page : Table :
Response variable BRFDINI	: Breastfe	eding Init	iation	
STUDY OF BREAST-FEEDING PA	ATTERNS AMC	NG WIC PAR	FICIPANTS	
BRR VARIANCE ESTIMATION				
Contrast	DF	WALD CHI-SQ	P-VALUE	
OVERALL MODEL	 9	111.03	0.0000	
MODEL MINUS INTERCEPT	8	108.99	0.0000	
INTERCEPT				
RACEMOM	3	36.48		
EDUC	2	11.92	0.0026	
	1	29.86	0.0000	
MRTLSTAT		0 5 2	0 4672	
MRTLSTAT BABYSEX	1	0.53	0.40/2	
		0.53 7.89		

'ERNS A	MONG WIC	PARTICIPANTS	
	Lowor	Uppor	
0.20	0.08	0.51	
1.00	1.00	1.00	
0.58	0.36	0.94	
5.55	2.94	10.48	
1.52	0.11	21.15	
1.00	1.00	1.00	
1.06	0.67	1.67	
2.35	1.23	4.47	
1.98	1.53	2.57	
1.00	1.00	1.00	
0.88	0.62	1.25	
		1 00	
1.00	1.00	1.00	
	Odds Ratio 0.20 1.00 0.58 5.55 1.52 1.00 1.06 2.35 1.98 1.00	Lower 95% Ratio Limit 0.20 0.08 1.00 1.00 0.58 0.36 5.55 2.94 1.52 0.11 1.00 1.00 1.06 0.67 2.35 1.23 1.98 1.53 1.00 1.00	Lower Upper Odds 95% 95% Ratio Limit Limit 0.20 0.08 0.51 1.00 1.00 1.00 0.58 0.36 0.94 5.55 2.94 10.48 1.52 0.11 21.15 1.00 1.00 1.00 1.06 0.67 1.67 2.35 1.23 4.47 1.98 1.53 2.57 1.00 1.00 1.00 0.88 0.62 1.25 1.00 1.00 1.00

The MULTILOG Procedure

Multinomial Logistic Regression (*Release 7.0*)

- Generalized Logit Models
 - Nominal Outcomes

e.g., Type of health plan (A, B, C, D)

- Cumulative Logit Models
 - Ordinal Outcomes

e.g., Pain Relief: none, mild, moderate, complete relief

- "Proportional Odds Models"
- Binary Logistic is a special case of each
- Model-fitting Approach
 - Fits *marginal* or *population-averaged* models
 - Uses GEE to model the intracluster correlations and efficiently estimate regression coefficients

Applications in Pharmaceutical Research

Toxicology / Pre-Clinical Studies

Developmental Toxicity
 Severity of malformations recorded on fetuses clustered within litters (cluster = litter)

Clinical Trials

Repeated Measures Studies Multiple illness or adverse events per patient (cluster = patient)

Example

Repeated ordinal responses of pain relief over an 8-hour period in a randomized clinical trial of acute pain relief comparing placebo with 2 analgesics (Gansky, Koch, et al., 1994, Journal of Biopharmaceutical Statistics)

Cross-Over Studies

Subjects receive each treatment in sequence (cluster = patient)

Example

3-period, 3 treatment cross-over study (Snapinn and Small, 1986, Biometrics):

Investigational drug, aspirin, and placebo administered in sequence to headache sufferers

Patients rated each drug on scale of 1-4 according to amount of pain relief.

Generalized Logit Model

Y is a categorical response variable with *K* categories 1,2,...,*K* (nominal scale)

 $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip})^{\prime} =$ vector of explanatory variables for subject *i*

Model
$$\pi_k(x_i) = prob(Y_i = k | x_i)$$
 $k = 1, ..., K-1$

Generalized Logits Model (Agresti, 1990):

$$\log\left[\frac{\pi_k(\boldsymbol{x}_i)}{\pi_K(\boldsymbol{x}_i)}\right] = \boldsymbol{\beta}'_k \boldsymbol{x}_i \qquad k = 1, \dots, K-1$$

Separate parameter vector (intercepts and slopes) for *each* of the *K*-1 logit equations

 $\blacksquare \qquad \beta_K = 0.$

• $\exp(\beta_k) = \text{odds of being in category } k \text{ vs. } K \text{ (the last)}$ for each 1-unit increase in x

Cumulative Logit Model

Y is a categorical response variable with *K* categories 1,2,...,K ordinal scale: *e.g.*, none, mild, moderate, severe

 $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip})' =$ vector of explanatory variables for subject *i*

Model $F_k(x_i) = prob(Y \le k | x_i) =$ cum. prob. up to and including category k

McCullagh's (1980) Proportional Odds Model:

Cumulative Logits

$$\log\left[\frac{F_k(\boldsymbol{x}_i)}{1-F_k(\boldsymbol{x}_i)}\right] = \alpha_k + \boldsymbol{\beta}' \boldsymbol{x}_i \qquad k = 1, \dots, K-1$$

- Separate intercepts α_k , but a *common set of slopes* β, for k = 1,...,K-1
- β measures the effect of the covariates on the severity of response

Efficient Parameter Estimation

Efficiently Weight the Data to Estimate Regression Coefficients (β)

GEE Approach

(Longitudinal Data Analysis, Zeger and Liang, 1986):

- 1) Assume a Covariance Structure V_i to describe the relationship among observations within clusters, i=1,...,n
 - Mean / Variance Relationship: $V(y_{ij}) = g(\mu_{ij})$
 - Pairwise Correlation Model: $Corr(y_{ij}, y_{ik})$
- 2) Estimate Covariance Parameters
- 3) Weight Data Inversely Proportional to V_i to Estimate β

 V_i inserted into the usual estimating equations in order to weight the data efficiently

Efficiently Weight the Data to Estimate Regression Coefficients (β)

GEE Approach

(Longitudinal Data Analysis, Zeger and Liang, 1986):

$i = 1, \ldots, n$	Clusters
$j=1$,, m_i	Observational Units
$y_i = (y_{i1},, y_{im_i})$	Vector of responses
$\boldsymbol{\mu}_{i} = E(\boldsymbol{y}_{i}) = \boldsymbol{\mu}_{i}(\boldsymbol{\beta})$ $= (\boldsymbol{\mu}_{i1}, \dots, \boldsymbol{\mu}_{im_{i}})$	Vector of marginal means
$V_i(\boldsymbol{\alpha}) = Cov(y_i; \boldsymbol{\mu}_i, \boldsymbol{\alpha})$	Working Covariance matrix

"Generalized" Estimating Equations:

$$U(\boldsymbol{\beta}) = \sum_{i=1}^{n} \frac{\partial \boldsymbol{\mu}_{i}^{\prime}}{\partial \boldsymbol{\beta}} \boldsymbol{V}_{i}(\boldsymbol{\alpha})^{-1} (\boldsymbol{y}_{i} - \boldsymbol{\mu}_{i}) = \boldsymbol{0}$$

Working Covariance Structure

$$\boldsymbol{V}_i(\boldsymbol{\alpha}) = \boldsymbol{A}_i^{1/2} \boldsymbol{R}_i(\boldsymbol{\alpha}) \boldsymbol{A}_i^{1/2} \cdot \boldsymbol{\phi} \qquad \boldsymbol{V}$$
 is Block diagonal

$$A_i$$
 = diagonal matrix with diagonal elements equal to
the marginal variances of observational units
within clusters: $g(\mu_{i1}), \dots, g(\mu_{im_i})$

$$= \begin{bmatrix} g(\mu_{i1}) & 0 & 0 & 0 \\ 0 & g(\mu_{i2}) & 0 & 0 \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & g(\mu_{im_i}) \end{bmatrix}$$

Relationship Between Variance of y_{ij} and its Mean

 $Var(y_{ij}) = g(\mu_{ij}) \cdot \phi$

g is a known variance function, ϕ is an unknown scale parameter

Binary Responses

Marginal distribution of y_{ij} is Bernoulli

Therefore $Var(y_{ij}) = \mu_{ij}(1 - \mu_{ij})$ and $\phi = 1$.

Choices for Working Correlation Matrices

$R_i(\alpha)$ is the "Working" Correlation Matrix for y_i

 $\alpha_{jk} = corr(y_{ij}, y_{ik})$

1) *Independent Working Correlation Matrix* (Identity matrix implies 0 pairwise correlation)

$$\boldsymbol{R}_{i}(\boldsymbol{\alpha}) = \boldsymbol{I} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

- Estimating equations reduce to familiar forms:
 - Normal equations for linear regression
 - Score equations for logistic regression
- Leads to standard regression coefficient estimates
- Consistent and asymptotically normal, regardless of whether or not the correlation structure is correctly specified
- This approach is offered in SUDAAN, and it is perfectly valid for estimating the *regression parameters*.

Choices for Working Correlation Matrices

2) Exchangeable

(equal pairwise correlations)

$$\boldsymbol{R}_{i}(\boldsymbol{\alpha}) = \begin{bmatrix} 1 & \rho & \rho & \rho \\ \rho & 1 & \rho & \rho \\ \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & 1 \end{bmatrix}$$

- SUDAAN offers this form as well
- Can improve *efficiency* of parameter estimates over the independence working assumption when working correlations are close to truth.

Robust Variance Estimate for GEE

$$\operatorname{Var}(\hat{\boldsymbol{\beta}}) = M_0^{-1} M_1 M_0^{-1}$$

where

$$M_{0} = \sum_{i=1}^{n} \frac{\partial \mu_{i}^{\prime}}{\partial \beta} V_{i}^{-1} \frac{\partial \mu_{i}}{\partial \beta}$$
$$M_{1} = \sum_{i=1}^{n} \frac{\partial \mu_{i}^{\prime}}{\partial \beta} V_{i}^{-1} Var(y_{i}) V_{i}^{-1} \frac{\partial \mu_{i}}{\partial \beta}$$

• M_0^{-1} (outside term) is called the *naive* or *model-based* variance (inverse of information matrix, appropriate when working assumption about covariance structure is correct)

Sensitive to violations of model assumptions!

- M_1 (middle term) serves as a *variance correction* when the covariance model is misspecified
- **Robust variance** is consistent even when $var(y_{ij}) \neq g(\mu_{ij}) \cdot \phi$ or $R_i(\alpha)$ is not the true correlation matrix of Y_i
- $Var(y_i)$ empirically estimated by $(y_i \hat{\mu}_i)(y_i \hat{\mu}_i)^{\prime}$
- SUDAAN offers the *robust* (default) and in Release 7.5 the *model-based* variance estimates (via the *SEMETHOD=MODEL* option)

Robust Variance Estimate for GEE

- Also referred to as *Sandwich Estimator* or *Variance Correction*
- Properly accounts for intracluster correlation
- Yields *consistent variance estimates*, even if correlation structure is misspecified (*e.g.*, by specifying "working" independence when the correlations are in fact exchangeable)

Huber (1967) Royall (1986) Binder (1983, 1992)

SYNTAX for GEE options in REGRESS and MULTILOG

PROC REGRESS MULTILOG ... R = Independent | Exchangeable RSTEPS = count SEMETHOD = ZEGER | BINDER | MODEL

R = Independent | Exchangeable

Specifies the "working" assumption for estimating the within-cluster correlation structure. The default assumption is independent working correlations. When R=exchangeable, the estimated exchangeable correlation matrix is available for printing.

RSTEPS = count

Specifies the maximum number of steps (iterating between estimated regression coefficients and correlations) used to fit the model. The default value is 0 and the default correlation structure is independent (R=independent). If you specify exchangeable correlations, the default value for the RSTEPS parameter is 1.

SEMETHOD = ZEGER / BINDER / MODEL

Specifies the method for computing standard errors of regression coefficients. *SEMETHOD=ZEGER* and *BINDER* both specify the full *robust* or *sandwich* variance estimator. For the REGRESS procedure, *ZEGER* and *BINDER* produce identical results. For the MULTILOG procedure, *ZEGER* and *BINDER* produce different results for responses with more than 2 levels. *SEMETHOD=MODEL* requests the *model-based* or *naive* standard error estimator, which is simply the outside of the sandwich estimator and is appropriate when the pairwise correlations within a cluster have been correctly specified.

What Does SUDAAN Model?

Marginal Models (Population-Averaged)

 Marginal mean of the multivariate outcomes as a function of the covariates:

 $F\left[E(y_{ij} \mid \boldsymbol{x}_{ij})\right] = \boldsymbol{x}_{ij}^{\prime} \boldsymbol{\beta}$

- Focus on how X causes Y, while acknowledging the dependence within clusters (as opposed to how one Y causes another)
- Describes relationship between covariates and response across clusters
- Intracluster correlation treated as nuisance parameter

References:

Zeger and Liang (1986) Liang and Zeger (1986) Zeger, Liang, and Albert (1988) Binder (1983, 1992)

R-Square for Logistic Regression

Proportion of Log-Likelihood Explained by the Model (Cox and Snell, 1989)

$$R^{2} = 1 - \left(\frac{L(\mathbf{0})}{L(\hat{\mathbf{\beta}})}\right)^{\frac{2}{n}}$$

where:

 $L(\mathbf{0})$ is the likelihood of the intercept-only model $L(\hat{\boldsymbol{\beta}})$ is the likelihood of the specified model, and *n* is the sample size.

R-Square for *Linear* **Regression:**

Simple correlation between observed and predicted response (based on the model).

REFLEVEL Statement

- Available in all modelling procedures
- Allows the user to change the definition of the *reference cell* for all categorical covariates.
- By *default*, the reference cell is the *last level* of each categorical covariate.

Syntax:

```
REFLEVEL variable_1 = reference_level_1
variable_2 = reference_level_2
{... variable_k = reference_level_k};
```

- Each variable_i must be defined on the SUBGROUP and LEVELS statements
- For SUBGROUP variables *not* on the REFLEVEL statement, the default reference level is still the *last* level.

The following example comes from the NHANES I Survey and its Longitudinal Follow-up Study conducted 10 years later. NHANES I (*National Health and Nutrition Examination Survey I*) was a multi-stage sample survey of over 14,000 adults in the US aged 25-74 years, with data collection taking place in 1971-1975. The epidemiologic follow-up took place in 1981-1984.

In this analysis, we wish to determine whether follow-up cancer status (*CANCER12*, 1=yes vs. 0=no) is associated with a measure of body iron stores at the initial exam (*B_TIBC*, total ironbinding capacity), while adjusting for age group at initial exam (*AGEGROUP*, 1=20-49, 2=50+) and smoking status (*SMOKE*, 1=current, 2=former, 3=never, 4=unknown).

First, we supply the results with the *default reference cells*, the last level of each categorical covariate, *i.e.*, SMOKE=4 (*unknown*) and AGEGROUP=2 (50+):

PROC MULTILOG DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2; 1 2 NEST Q_STRATA PSU1; WEIGHT B_WTIRON; 3 4 SUBGROUP CANCER12 AGEGROUP SMOKE; 5 LEVELS 2 2 4; 6 MODEL CANCER12 = B TIBC AGEGROUP SMOKE / CUMLOGIT; 7 SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60; 8 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0" P_BETA="P-VALUE" DF WALDCHI WALDCHP / T_BETAFMT=F8.2 DEFTFMT=F6.2 WALDCHIFMT=F8.2 DFFMT=F8.0; 9 TITLE "Default Reference Cell Model"; Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading. Number of observations read : 3290 Weighted count: 40570323 Observations used in the analysis : 3290 Weighted count: 40570323 0 Weighted count: Observations with missing values : 0 Denominator degrees of freedom : 35 Maximum number of estimable parameters for the model is 6 File C:\ADVANCED\IRONSUD.SSD contains 67 Clusters Maximum cluster size is 111 records Minimum cluster size is 15 records Independence parameters have converged in 5 iterations Sample and Population Counts for Response Variable CANCER12 Cancer : Sample Count 232 Population Count 1745695 No Cancer: Sample Count 3058 Population Count 38824628

DEFAULT Reference Cell Parameterization

Date:05-29-97Research Triangle InstitutePage: 1Time:14:16:21The MULTILOG ProcedureTable: 1 Variance Estimation Method: Robust (Binder, 1983) Working Correlations: Independent Link Function: Cumulative Logit Response variable CANCER12: Cancer Status (1/2) Default Reference Cell Model ______ Independent Variables DESIGN and Effects BETA S.E. EFFECT T:BETA=0 P-VALUE _____
 Intercept
 -0.8618
 0.6605
 0.94
 -1.30
 0.2004

 Total Iron Binding Capacity
 -0.0024
 0.0018
 1.10
 -1.29
 0.2052
 Age Cohort 20-49 yrs.-2.25250.33431.89-6.740.0000**50+ yrs.0.00000.0000**... Smoking Status -0.5858 0.2771 0.77 -2.11 0.0417 Current -0.9418 0.2922 0.84 -3.22 0.0027 Former -0.4998 0.2743 0.85 -1.82 0.0770 0.0000 0.0000 . . Never Unknown _____

Here, each smoking group is automatically compared to the *unknown* smoking status (SMOKE=4), which may not be very meaningful.

DEFAULT Reference Cell Parameterization

```
Date: 05-29-97 Research Triangle Institute
                                                                       Page : 2
Time: 14:16:21
                            The MULTILOG Procedure
                                                                        Table : 1
Variance Estimation Method: Robust (Binder, 1983)
Working Correlations: Independent
Link Function: Cumulative Logit
Response variable CANCER12: Cancer Status (1/2)
Default Reference Cell Model
_____
                          Degrees P-value
Contrast
                          of Wald Wald
                          Freedom ChiSq ChiSq
_____

      OVERALL MODEL
      6
      708.28
      0.0000

      MODEL MINUS INTERCEPT
      5
      64.47
      0.0000

      B_TIBC
      1
      1.67
      0.1967

      AGEGROUP
      1
      45.39
      0.0000

      SMOKE
      3
      10.60
      0.0141

                                  3 10.60 0.0141
SMOKE
_____
MULTILOG used
 CPU time : 12.74 seconds
Elapsed time : 13 seconds
  Virtual memory : 2.84 MB
```

Here we see that *Age group* and *Smoking status* are significantly associated with follow-up cancer status, but *Total iron-binding capacity* is not (*p*=0.1967).

Using the REFLEVEL Statement

Next, using the REFLEVEL statement, we re-define the reference cells to be the *first level* of each categorical variable. Note the only differences in the results are in the estimates of the regression coefficients, where the expected value of the response for each level of the categorical covariate(s) is now compared to the user-specified *first* level instead of the last. The main effects tests remain unchanged.

```
PROC MULTILOG DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2;
10
11 NEST Q_STRATA PSU1;
12 WEIGHT B_WTIRON;
13 REFLEVEL AGEGROUP=1 SMOKE=1;
14 SUBGROUP CANCER12 AGEGROUP SMOKE;
15 LEVELS 2
                     2
                              4;
  MODEL CANCER12 = B_TIBC AGEGROUP SMOKE / CUMLOGIT;
16
17
  SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
18 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" DF WALDCHI WALDCHP / T_BETAFMT=F8.2 DEFTFMT=F6.2
         WALDCHIFMT=F8.2 DFFMT=F8.0;
19 TITLE "Using the REFLEVEL Statement";
Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading.
                             : 3290
Number of observations read
                                            Weighted count: 40570323
                                          Weighted count: 40570323
Observations used in the analysis :
                                    3290
Observations with missing values : 0
                                          Weighted count:
Denominator degrees of freedom :
                                      35
Maximum number of estimable parameters for the model is 6
File C:\ADVANCED\IRONSUD.SSD contains 67 Clusters
Maximum cluster size is 111 records
Minimum cluster size is 15 records
Independence parameters have converged in 5 iterations
Sample and Population Counts for Response Variable CANCER12
 Cancer : Sample Count 232 Population Count
                                                       1745695
 No Cancer: Sample Count 3058
                                     Population Count 38824628
```

Using the REFLEVEL Statement

Date: 05-29-97 R Time: 14:16:21	esearch Tr The MULT	-		te	Page : Table :
	1110 110111		cedure		
Variance Estimation Method:	Robust (Bi	inder, 19	983)		
Working Correlations: Indepe	endent				
Link Function: Cumulative Lo	git				
Response variable CANCER12:	Cancer Sta	atus (1/2	2)		
Using the REFLEVEL Statement	:				
Independent Variables			DESIGN		
and Effects	BETA				
Intercept					
Total Iron-Binding Capacity	-0.0024	0.0018	1.10	-1.29	0.2052
Age Cohort					
20-49 yrs.	0.0000	0.0000	•	•	•
50+ yrs.				6.74	
Smoking Status					
-	0.0000	0.0000	•	•	•
Current	0.0000 -0.3560				
Current	-0.3560	0.2716	1.16		0.1985

Now each smoking group is compared to the *current* smokers (SMOKE=1), and we see immediately that *current smokers* are not significantly different from *former smokers* (p=0.1985) nor from those who have *never smoked* (p=0.7330).

Using the REFLEVEL Statement

```
Date: 05-29-97 Research Triangle Institute
                                                                   Page : 2
                            The MULTILOG Procedure
Time: 14:16:21
                                                                   Table : 1
Variance Estimation Method: Robust (Binder, 1983)
Working Correlations: Independent
Link Function: Cumulative Logit
Response variable CANCER12: Cancer Status (1/2)
Using the REFLEVEL Statement
_____
                        Degrees P-value
Contrast
                        of Wald Wald
                        Freedom ChiSq ChiSq
-----

        OVERALL MODEL
        6
        708.28
        0.0000

        MODEL MINUS INTERCEPT
        5
        64.47
        0.0000

        B_TIBC
        1
        1.67
        0.1967

        AGEGROUP
        -
        -
        -

AGEGROUP
                               1 45.39 0.0000
                          3 10.60 0.0141
SMOKE
-----
MULTILOG used
 CPU time : 13.2 seconds
  Elapsed time : 14 seconds
  Virtual memory : 2.88 MB
```

The tests of main effects are the same, no matter which groups are designated as the reference cells.

EFFECTS Statement

• Available in all modeling procedures

Simplifies the following hypothesis testing situations:

- Testing multiple main effects and/or interactions simultaneously (*e.g.*, testing chunk interaction effects);
- Testing general linear contrasts (*e.g.*, pairwise comparisons, trends) for a specific variable(s) in the model by only specifying contrast coefficients for the variable(s) of interest;
- Testing main effects in the presence of interactions. If the model contains factors A, B, and their interaction A*B, the user can obtain the:
 - 1) *Simple effect* of A, which is the effect of variable A tested within a given level of variable B, and
 - 2) *Main effects* of A, which are averaged over the levels of B.

Syntax:

```
EFFECTS term(s) / [ NAME = "label" ] [ DISPLAY ]
[ REFLEVEL | AVERAGE |
VARIABLE_NAME = value ] ;
```

where *term(s)* are name of effect(s) (single variables or/and interactions) on the MODEL statement, which may include contrast matrices.

EFFECTS Statement Options

NAME = "*label*"

Assigns a label to the contrast. Default is "*Effect_nn*", where *nn* is the *nn*-th EFFECT statement in the procedure

DISPLAY

Prints the contrast coefficients

REFLEVEL, AVERAGE, VARIABLE_NAME = value

Tells SUDAAN how to test the effects of covariates in the model when they are interacted with other effects in the model.

Example:

MODEL Y = A B A*B;

To test the effect of A (which may be either continuous or categorical), the user has three options:

REFLEVEL (default)

Tests the effect of A when B (and all other variables A is interacted with) are set to their reference levels.

AVERAGE

Tests the effect of A *averaged over the interaction effect*, with proportional weighting over each level of B (Graubard and Korn, 1997). The contrast coefficient vector contains the weighted proportion of subjects in the *j*-th category of the *i*-th SUBGROUP variable.

VARIABLE_NAME = value

Similar to the REFLEVEL option, except here *the user chooses the level of B within which to test the effect of A*. This option is used to carry out what are commonly known as "simple effects," in which an effect A is to be tested within a specific level of B, other than the reference cell.

Using the NHANES I Study and its longitudinal follow-up (see the REFLEVEL statement examples for details), we evaluate the effects of body iron stores at initial exam (B_TIBC , continuous), age group at initial exam (AGEGROUP, 1=20-49, 2=50+), and smoking status (SMOKE, 1=current, 2=former, 3=never, 4=unknown) on follow-up cancer status (CANCER12, 1=yes, 2=no).

The **EFFECTS statement** can be used to:

1) Test the combined effect of *Agegroup* and *Smoke*:

EFFECTS AGEGROUP SMOKE / NAME = "Combined Age, Smoke";

2) Compare *Smoke* Level 1 to Level 2 (the default reference level for *Smoke* is Level 4):

EFFECTS SMOKE = (-1 1 0 0) / NAME="Smoke 1 vs 2";

```
PROC MULTILOG DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2;
1
2
  NEST Q_STRATA PSU1;
3
  WEIGHT B_WTIRON;
  SUBGROUP CANCER12 AGEGROUP SMOKE;
4
  LEVELS 2 2
5
                            4;
  MODEL CANCER12 = B_TIBC AGEGROUP SMOKE / CUMLOGIT;
6
7
  EFFECTS AGEGROUP SMOKE / NAME = "Combined Age, Smoke";
8
  EFFECTS SMOKE=(-1 1 0 0) / NAME = "Smoke 1 vs 2";
9 SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
10 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" DF WALDCHI WALDCHP /
         T_BETAFMT=F8.2 DEFTFMT=F6.2 DFFMT=F8.0 WALDCHIFMT=F8.2;
11 TITLE "EFFECTS Statement Example";
NOTE: Terms in the MODEL statement have been rearranged
     to follow subgroup order.
Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading.
Number of observations read
                              : 3290 Weighted count: 40570323
Observations used in the analysis : 3290 Weighted count: 40570323
Observations with missing values : 0 Weighted count:
Denominator degrees of freedom :
                                     35
Maximum number of estimable parameters for the model is 6
File C:\ADVANCED\IRONSUD.SSD contains 67 Clusters
Maximum cluster size is 111 records
Minimum cluster size is 15 records
Independence parameters have converged in 5 iterations
Sample and Population Counts for Response Variable CANCER12
 Cancer:Sample Count232Population Count1745695No Cancer:Sample Count3058Population Count38824628
```

Date: 05-29-97	Research Ti The MULT	5		ce	Page : Table :
Variance Estimation Method: Working Correlations: Indep Link Function: Cumulative L	endent ogit				
Response variable CANCER12: EFFECTS Statement Example	Cancer Sta	atus (1/.	2)		
	BETA	S.E.			
Intercept				-1.30	
Age Cohort					
20-49 yrs.	-2.2525	0.3343	1.89	-6.74	0.0000
50+ yrs.	0.0000	0.0000	•		
Smoking Status					
	-0.5858	0.2771	0.77	-2.11	0.0417
Current			0 0 1	-3.22	0.0027
Current Former	-0.9418	0.2922	0.04		
				-1.82	
Former Never		0.2743	0.85	-1.82	0.0770

The MU Robust Ident git Cancer	n Trianglo JLTILOG P: (Binder, Status (:	1983)		Page : Table :	
Robust ndent git Cancer	(Binder,	1983)		Table :	Ţ
ndent git Cancer		·			
6	708.28	0.0000			
5	64.47	0.0000			
1	45.39	0.0000			
3	10.60	0.0141			
1	1.67	0.1967			
4	53.16	0.0000			
1	1.72	0.1899			
	eedom 6 5 1 3 1 4 1	Wald ChiSq 6 708.28 5 64.47 1 45.39 3 10.60 1 1.67 4 53.16 1 1.72	grees P-value Wald Wald eedom ChiSq 6 708.28 0.0000 5 64.47 0.0000 1 45.39 0.0000 3 10.60 0.0141 1 1.67 0.1967 4 53.16 0.0000 1 1.72 0.1899	Wald Wald cedom ChiSq 6 708.28 0.0000 5 64.47 0.0000 1 45.39 0.0000 3 10.60 0.0141 1 1.67 0.1967 4 53.16 0.0000 1 1.72 0.1899	Wald Wald eedom ChiSq 6 708.28 0.0000 5 64.47 0.0000 1 45.39 0.0000 3 10.60 0.0141 1 1.67 0.1967 4 53.16 0.0000 1 1.72 0.1899

The combined effect of *Age* and *Smoking Status* is statistically significant (p=0.0000). However, *current smokers* (SMOKE=1) are not significantly different (p=0.1899) from *former smokers* (SMOKE=2).

In this example, we evaluate the effects of body iron stores at initial exam (*TRFSAT*, 1 = high vs. 0=normal indicator), smoking status (*SMOKE*, 1=current, 2=former, 3=never, 4=unknown), age group at initial exam (*AGEGROUP*, 1=20-49 yrs, 2=50+ yrs), and various two-way interactions on a binary response, cancer status at follow-up (*CANCER1*, 1=yes vs. 0=no).

The **EFFECTS Statement** can be used to easily test simultaneous interaction effects (smoking by age group, smoking by indicator of body iron stores):

EFFECTS SMOKE*AGEGROUP SMOKE*TRFSAT / NAME="Chunk Interactions";

```
66 PROC LOGISTIC DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2;
67 NEST Q_STRATA PSU1;
68 WEIGHT B_WTIRON;
69 SUBGROUP SMOKE AGEGROUP;
70 LEVELS 4
                  2;
71 MODEL CANCER1 = TRFSAT SMOKE AGEGROUP SMOKE*AGEGROUP SMOKE*TRFSAT;
72 EFFECTS SMOKE*AGEGROUP SMOKE*TRFSAT / NAME = "Chunk Interactions";
73 SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
74 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" DF WALDCHI WALDCHP
         / SEBETAFMT=F8.5 DFFMT=F8.0 T_BETAFMT=F8.2 DEFTFMT=F6.2 WALDCHIFMT=F8.2;
75 TITLE "Using EFFECTS to Test Chunk Interactions";
NOTE: Terms in the MODEL statement have been rearranged
     to follow subgroup order.
Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading.
Number of zero responses : 3058
Number of non-zero responses :
                              232
Parameters have converged in 5 iterations
Number of observations read
                             : 3290 Weighted count: 40570323
Observations used in the analysis : 3290 Weighted count: 40570323
                                    0 Weighted count:
Observations with missing values :
                                                                   Ω
Denominator degrees of freedom :
                                      35
Maximum number of estimable parameters for the model is 12
R-Square for dependent variable CANCER1 (Cox & Snell, 1989): 0.046486
```

Using EFFECTS to Test Chunk Interactions

Response variable CANCER1	: Cancer S	tatus (0	/1)			
Using Effects to Test Chu	ink Interac	tions				
Independent Variables and	 l					-
Effects			DESIGN			
				T:BETA=0		2
Intercept				-5.92)
Smoking Status						
Current	-0.6159	0.37457	0.97	-1.64	0.1090)
Former	-1.6133	0.33255	0.65	-4.85	0.000)
Never	-0.5606	0.35346	0.93	-1.59	0.1217	7
Unknown	0.0000	0.00000	•	•	•	
Age Cohort						
20-49 yrs.	-3.8676	0.84072	0.31	-4.60	0.0001	L
50+ yrs.	0.0000	0.00000				
High Transferrin						
Saturation $(0/1)$	0.1745	0.52386	0.72	0.33	0.7411	L
Smoking Status, Age Cohor	rt					
Current, 20-49 yrs.	1.4407	1.03113	0.41	1.40	0.1711	L
Current, 50+ yrs.	0.0000	0.00000		•		
Former, 20-49 yrs.	2.2305	1.05117	0.44		0.0410)
Former, 50+ yrs.	0.0000	0.00000	•		•	
Never, 20-49 yrs.	1.5366	1.03999	0.44	1.48	0.1485	5
Never, 50+ yrs.	0.0000	0.00000	•		•	
Unknown, 20-49 yrs.	0.0000	0.00000	•		•	
Unknown, 50+ yrs.	0.0000					
Smoking Status, High						
Transferrin Saturation						
Current	-0.1905	0.56612	0.58	-0.34	0.7385	5
Former	1.1955	0.69445	0.94	1.72	0.0940)
Never	-0.1575	0.50445	0.52	-0.31	0.7568	3

Using EFFECTS to Test Chunk Interactions

Date: 04-04-97 Time: 15:55:41				Page : 2 Table : 1
Response variable CANCER	R1: Cancer	Status (0	/1)	
Using EFFECTS to Test Cl	unk Intera	ctions		
Contrast	-			
	Freedom	Wald ChiSq	ChiSq	
OVERALL MODEL	12			
MODEL MINUS INTERCEPT	11	101.61	0.0000	
INTERCEPT			•	
SMOKE				
AGEGROUP		•	•	
TRFSAT				
SMOKE * AGEGROUP	3	4.96	0.1749	
TRFSAT * SMOKE	3	6.02	0.1105	
Chunk Interactions	6	21.21	0.0017	

The combined interaction effect is statistically significant (p=0.0017). To test the same hypothesis using the CONTRAST statement, we would specify the following 12-row contrast matrix. The number of rows equals the number of regression coefficients to be tested in the contrast, with 1's in the columns corresponding to those regression coefficients. All other columns for intercept and main effects are 0's.

CONTRAST 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 **1** 0 0 0 0 0 0 Ω 0 0 0 0 0 0 0 0 0 **1** 0 0 0 0 0 0 0 0 0 **1** 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 **1** 0 0 / NAME="CHUNK INTERACTIONS";

Comparison to the CONTRAST Statement

```
62 PROC LOGISTIC DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2;
63 NEST Q STRATA PSU1;
64 WEIGHT B_WTIRON;
65 SUBGROUP SMOKE AGEGROUP;
66 LEVELS 4
                2;
67 MODEL CANCER1=TRFSAT SMOKE AGEGROUP SMOKE*AGEGROUP SMOKE*TRFSAT;
0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0
           0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0
           0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0
           0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0
           0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0
           0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0
           0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0 0
           0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0 0
           0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 0
           / NAME="CHUNK INTERACTIONS";
69 SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
70 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
        P_BETA="P-VALUE" DF WALDCHI WALDCHP / SEBETAFMT=F8.5 T_BETAFMT=F8.2
        DEFTFMT=F6.2 WALDCHIFMT=F8.2 DFFMT=F8.0;
71 TITLE "Using CONTRAST to Test Chunk Interactions";
Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading.
Number of zero responses
                        : 3058
Number of non-zero responses :
                             232
Parameters have converged in 5 iterations
Number of observations read
                           : 3290 Weighted count: 40570323
Observations used in the analysis : 3290 Weighted count: 40570323
Observations with missing values :
                                  0
                                       Weighted count:
Denominator degrees of freedom :
                                   35
Maximum number of estimable parameters for the model is 12
R-Square for dependent variable CANCER1 (Cox & Snell, 1989): 0.046486
```

Comparison to the CONTRAST Statement

Response variable CANCER1:	Cancer S	tatus (O	/1)		
			(_)		
Using CONTRAST to Test Chu	ınk Intera	ctions			
Independent Variables and Effects			DEGION		
FILECUS	BETA	S.E.	DESIGN EFFECT	T:BETA=0	P-VALUE
Intercept	-1.6135	0.27254	0.72	-5.92	0.0000
High Transferrin Saturation (0/1)	0 1745	0 5000	0 70	0 22	0 7/11
Saturation (0/1) Smoking Status	U.1/45	0.52386	0.72	0.33	U./411
Current	-0 6150	0.37457	0 97	-1 64	0.1090
Former		0.37457			0.0000
Never		0.35346			0.1217
Unknown		0.00000		-1.59	
Age Cohort	0.0000	0.00000	•	·	•
20-49 yrs.	-3.8676	0.84072	0.31	-4.60	0.0001
50+ yrs.		0.00000			
Smoking Status, Age Cohort				•	-
Current, 20-49 yrs.		1.03113	0.41	1.40	0.1711
Current, 50+ yrs.					
Former, 20-49 yrs.				2.12	
Former, 50+ yrs.					
Never, 20-49 yrs.	1.5366	1.03999		1.48	
Never, 50+ yrs.	0.0000	0.00000			
Unknown, 20-49 yrs.					
Unknown, 50+ yrs.	0.0000	0.00000			
Smoking Status, High					
Transferrin Saturation					
(0/1)					
Current				-0.34	
Former	1.1955	0.69445	0.94	1.72	0.0940
Never	-0.1575	0.50445	0.52	-0.31	0.7568
Unknown	0.0000	0.00000			•

1 1

Comparison to the CONTRAST Statement

Date: 03-27-97 Time: 14:25:00		n Triangl DGISTIC P:		Page : 2 Table : 1
Response variable CANCE	R1: Cancer S	Status (O	/1)	
Using CONTRAST to Test	Chunk Intera	actions		
Contrast	Degrees	Wald		
	Freedom	ChiSq	ChiSq	
OVERALL MODEL	12			
MODEL MINUS INTERCEPT	11	101.61	0.0000	
INTERCEPT	•	•	•	
TRFSAT	•	•	•	
SMOKE	•		•	
AGEGROUP		•		
SMOKE * AGEGROUP	-	4.96		
TRFSAT * SMOKE		6.02		
CHUNK INTERACTIONS	6	21.21	0.0017	
,				
LOGISTIC used CPU time : 29.2	7 gogonda			
Elapsed time : 30 s				
Virtual memory : 2.23				

The results are the same as for the EFFECTS statement, with the simultaneous interactions being statistically significant.

In this example, we evaluate the effect of smoking status (*SMOKE*, 1=*current*, 2=*former*, 3=*never*, 4=*unknown*) on a binary response, cancer status at follow-up (*CANCER1*, 1=*yes* vs. 0=*no*) under the following conditions:

- 1) When Age Group=1 (20-49 yrs),
- 2) When Age Group=2(50 + yrs),
- 3) When Age Group is at its reference level (level 2=50+ yrs),
- 4) Averaged over the interaction cells with Age Group.

The **EFFECTS statement** can be used to easily test these hypotheses:

```
EFFECTS SMOKE / AGEGROUP=1 NAME = "SMOKE in AGEGROUP=1";
EFFECTS SMOKE / AGEGROUP=2 NAME = "SMOKE in AGEGROUP=2";
EFFECTS SMOKE / REFLEVEL NAME = "SMOKE in Age Reference Level";
EFFECTS SMOKE / AVERAGE NAME = "SMOKE Averaged Over
Interaction";
```

```
76 PROC LOGISTIC DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2;
77 NEST Q_STRATA PSU1;
78 WEIGHT B_WTIRON;
79 SUBGROUP AGEGROUP SMOKE;
80 LEVELS 2
                   4;
81 MODEL CANCER1 = TRFSAT AGEGROUP SMOKE AGEGROUP*SMOKE;
82 EFFECTS SMOKE / AGEGROUP=1 NAME="Smoke Effect in Age=20-49";
83 EFFECTS SMOKE / AGEGROUP=2 NAME="Smoke Effect in Age=50+";
84 EFFECTS SMOKE / REFLEVEL NAME="Smoke Effect at Age Reference Level";
85 EFFECTS SMOKE / AVERAGE NAME="Smoke averaged over interaction";
86 SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
87 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P BETA="P-VALUE" DF WALDCHI WALDCHP
         /SEBETAFMT=F8.5 DFFMT=F8.0 T_BETAFMT=F8.2 DEFTFMT=F6.2 WALDCHIFMT=F8.2;
88 TITLE "Using EFFECTS to Test Simple Effects;
NOTE: Terms in the MODEL statement have been rearranged
     to follow subgroup order.
Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading.
Number of zero responses : 3058
Number of non-zero responses :
                              232
Parameters have converged in 5 iterations
Number of observations read
                             : 3290 Weighted count: 40570323
Observations used in the analysis : 3290 Weighted count: 40570323
                                    0 Weighted count:
Observations with missing values :
                                                                  0
Denominator degrees of freedom :
                                      35
Maximum number of estimable parameters for the model is 9
R-Square for dependent variable CANCER1 (Cox & Snell, 1989): 0.043642
```

Time: 15:55:41	The LO	GISTIC P	rocedure	2	Ta	able :
Response variable CANCER	1: Cancer S	tatus (0	/1)			
Using EFFECTS to Test Si	mple Effect	s				
Independent Variables an	d					
Effects				T:BETA=0	P-VALUE	
Intercept		0.25187		-6.65	0.0000	
Age Cohort						
20-49 yrs.	-3.8681	0.84493	0.31	-4.58	0.0001	
50+ yrs.	0.0000	0.00000				
Smoking Status						
Current	-0.6625	0.31953	0.91	-2.07	0.0455	
Former	-1.1591					
Never	-0.6030				0.0554	
Unknown	0.0000	0.00000	•		•	
High Transferrin						
Saturation (0/1)		0.20980	1.19	1.91	0.0650	
Age Cohort, Smoking Stat						
20-49 yrs., Current				1.38		
20-49 yrs., Former						
20-49 yrs., Never				1.47	0.1515	
20-49 yrs., Unknown		0.00000		•	•	
50+ yrs., Current	0.0000			•	•	
50+ yrs., Former		0.00000		•	•	
50+ yrs., Never		0.00000	•	•	•	
50+ yrs., Unknown	0.0000	0.00000				

Using EFFECTS to Test Simple Effec	ts			
		Wald	P-value	
Fre	edom	ChiSq	ChiSq	
OVERALL MODEL			0.0000	
MODEL MINUS INTERCEPT			0.0000	
INTERCEPT				
AGEGROUP				
SMOKE				
TRFSAT	1	3.63	0.0567	
AGEGROUP * SMOKE	3	5.25	0.1547	
Smoke Effect in Age=20-49	3	1.66	0.6466	
Smoke Effect in Age=50+	3	11.15	0.0110	
Smoke Effect at Age Reference Leve	13	11.15	0.0110	
Smoke averaged over interaction	3	0.36	0.9491	

Note that the test for *"Smoke Effect in Age=50+"* is equivalent to *"Smoke in Age Reference Level."* Here we see that:

- 1) There is a marginally significant interaction between age and smoking on follow-up cancer status (p=0.1547). SUDAAN computes this test automatically, without the need for the EFFECTS statement.
- 2) There is no significant effect of smoking on cancer status when age group=20-49 yrs. (p=0.6466), although the regression coefficients on the previous page (provided automatically by SUDAAN) and the EFFECTS statement here indicates a significant smoking effect when age is at its reference level (50+ yrs., p=0.0110).
- 3) There is no significant effect of smoking when smoking is averaged over its interaction with age (p=0.9302).

Now the same results via the CONTRAST statement:

118 SUDAAN Release 7.5

Comparison to the CONTRAST Statement

```
72 PROC LOGISTIC DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2;
73 NEST Q_STRATA PSU1;
74 WEIGHT B_WTIRON;
75 SUBGROUP AGEGROUP SMOKE;
76 LEVELS 2
                   4;
77 MODEL CANCER1 = TRFSAT AGEGROUP SMOKE AGEGROUP*SMOKE;
78 CONTRAST 0 0 0 0 -1 0 0 1 -1 0 0 1 0 0 0 0
            0 0 0 0 -1 0 1 0 -1 0 1 0 0 0 0
            0 0 0 0 -1 1 0 0 -1 1 0 0 0 0 0 0
            / NAME="SMOKE IN AGE=1";
79 CONTRAST 0 0 0 0 -1 0 0 1 0 0 0 0 -1 0 0 1
            0 0 0 0 -1 0 1 0 0 0 0 0 -1 0 1 0
            0 0 0 0 -1 1 0 0 0 0 0 0 0 -1 1 0 0
            / NAME="SMOKE IN AGE=2";
80 SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
81 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" DF WALDCHI WALDCHP / SEBETAFMT=F8.5 DFFMT=F8.0
         T_BETAFMT=F8.2 DEFTFMT=F6.2 WALDCHIFMT=F8.2;
82 TITLE "Testing Simple Effects via the CONTRAST Statement";
Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading.
Number of zero responses
                         : 3058
Number of non-zero responses :
                              232
Parameters have converged in 5 iterations
Number of observations read
                              : 3290
                                          Weighted count: 40570323
Observations used in the analysis : 3290 Weighted count: 40570323
Observations with missing values :
                                     0
                                          Weighted count:
                                                                   Ω
Denominator degrees of freedom :
                                      35
Maximum number of estimable parameters for the model is 9
R-Square for dependent variable CANCER1 (Cox & Snell, 1989): 0.043642
```

Comparison to the CONTRAST Statement

Research Triangle Institute Date: 03-27-97 Page : 1 Time: 14:25:00 The LOGISTIC Procedure Table : 1 Response variable CANCER1: Cancer Status (0/1) Testing Simple Effects Via the CONTRAST Statement _____ Independent Variables and Effects DESIGN BETA S.E. EFFECT T:BETA=0 P-VALUE _____ -1.6762 0.25187 0.79 -6.65 0.0000 Intercept High Transferrin 0.3997 0.20980 1.19 1.91 0.0650 Saturation (0/1)Age Cohort -3.8681 0.84493 0.31 -4.58 0.0001 20-49 yrs. 0.0000 0.00000 . 50+ yrs. • • Smoking Status -0.6625 0.31953 0.91 -2.07 0.0455 -1.1591 0.34790 1.05 -3.33 0.0020 Current Former -0.6030 0.30426 0.91 -1.98 0.0554 0.0000 0.00000 Never 0.0000 0.00000 . Unknown . . Age Cohort, Smoking Status ge Cohort, Smoking Status20-49 yrs., Current1.42901.034430.411.380.175920-49 yrs., Former2.23991.041730.432.150.038520-49 yrs., Never1.53451.046520.451.470.151520-49 yrs., Unknown0.00000.00000...50+ yrs., Current0.00000.00000...50+ yrs., Never0.00000.00000...50+ yrs., Unknown0.00000.00000...50+ yrs., Unknown0.00000.00000... ______

Comparison to the CONTRAST Statement

Date: 03-27-97 Time: 14:25:00		h Trianglo OGISTIC P:		Page : 2 Table : 1
Response variable CANCER1	: Cancer	Status (0	/1)	
Testing Simple Effects Vi	a the CON	IRAST Sta	tement	
Contrast				
		Wald		
		ChiSq	-	
OVERALL MODEL		859.59		
MODEL MINUS INTERCEPT	8	89.15	0.0000	
INTERCEPT	•			
TRFSAT	1	3.63	0.0567	
AGEGROUP		•		
SMOKE				
AGEGROUP * SMOKE			0.1547	
SMOKE IN AGE=1	3	1.66	0.6466	
SMOKE IN AGE=2	3	11.15	0.0110	

LSMEANS Statement

- Available in the *linear regression procedure* (REGRESS).
- Produces "least squares" or "adjusted means" for any number of categorical covariates in the model.
- List one or more categorical effects from the right-hand-side of the MODEL statement. *Continuous variables are not allowed* on the LSMEANS statement.
- The keyword *INTERCEPT* specifies an overall least-squares mean, when the model contains an intercept.

Syntax:

```
LSMEANS [INTERCEPT] effect(s) / [ALL] [DISPLAY] ;
```

ALL

Requests least-squares means for *all effects* on the right-hand side of the MODEL statement.

DISPLAY

Requests least squares means contrast coefficients.

Construction of the LSMEANS Contrast

- SUDAAN calculates *contrast coefficients* that are the weighted means of each covariate to be adjusted for in the model, using all observations for which there are no missing independent or dependent variable values.
- Contrast coefficients corresponding to the levels of the *categorical covariates* (appearing on the SUBGROUP statement) are the weighted numbers of individuals in each category of the covariate. Sample member weights are provided by the variable specified on the WEIGHT statement. If weights are all equal to one (*e.g.*, via the keyword _ONE_), unweighted means are used.
- The set of contrast coefficients are vector-multiplied by the estimated regression coefficients.

The following example illustrates the construction of the LSMEANS contrast.

Data:

NHANES I Survey and its Longitudinal Follow-up Study.

Question:

Is smoking status at initial exam (*SMOKE*, where 1=current vs. 2=former, 3=never, 4=unknown) associated with a measure of body iron stores at the initial exam (*B_TIBC*, or total iron-binding capacity), while adjusting for age at initial exam?

LSMEANS

We request the least squares means of the response B_TIBC , total ironbinding capacity, within levels of *SMOKE*, adjusted for age at initial exam (first as categorical, then as a continuous covariate). The data are weighted by the variable B_WTIRON .

SUDAAN Programming Statements Demonstrating the Construction of the LSMEANS Contrast for Categorical Covariates

PROC REGRESS DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2; 1 2 NEST Q_STRATA PSU1; WEIGHT B_WTIRON; 3 4 SUBGROUP AGEGROUP SMOKE; LEVELS 2 4; 5 б MODEL B TIBC = SMOKE AGEGROUP; 7 LSMEANS SMOKE / DISPLAY; SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60; 8 9 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0" P_BETA="P-VALUE" DF WALDCHI WALDCHP / LSMEANS=ALL T_BETAFMT=F8.2 DEFTFMT=F6.2 DFFMT=F8.0 WALDCHIFMT=F8.2; 10 TITLE "LSMEANS With Categorical Covariate"; NOTE: Terms in the MODEL statement have been rearranged to follow subgroup order. Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading. Number of observations read : 3290 Weighted count: 40570323 Observations used in the analysis : 3290 Weighted count: 40570323 Observations with missing values : 0 Weighted count: 0 Denominator degrees of freedom : 35 Maximum number of estimable parameters for the model is 5 File C:\ADVANCED\IRONSUD.SSD contains 67 clusters Maximum cluster size is 111 records Minimum cluster size is 15 records Weighted mean response is 354.580621

Estimated Regression Coefficients for the Model

Date: 05-29-97 Research Triangle Institute Page : 4 Time: 15:28:17 The REGRESS Procedure Table : 1 Variance Estimation Method: Robust (Binder, 1983) Working Correlations: Independent Link Function: Identity Response variable B_TIBC: TOTAL IRON BINDING CAPACITY LSMEANS With Categorical Covariate _____ Independent Variables and Effects DESIGN BETA S.E. EFFECT T:BETA=0 P-VALUE _____ **352.8876** 3.8547 1.09 91.55 0.0000 Intercept Age Cohort 7.22101.89681.123.810.00050.00000.0000.... 20-49 yrs. 50+ yrs. Smoking Status **-7.5062** 3.7690 0.95 -1.99 0.0543 Current -1.67544.26361.25-0.390.6967-0.92613.82841.06-0.240.81030.00000.0000... Former Never Unknown _____

Least Squares Means Contrast Coefficients:

Smoking Status and Age Group

Since we want to estimate the least squares means of the response within each level of smoking status (a 4-level variable), SUDAAN will produce four rows of contrast coefficients. The first row of the matrix will produce the adjusted means for SMOKE=*current*, the second row is for SMOKE=*former*, and so on. The contrast coefficients for *smoking status* are 1's and 0's, indicating the level of interest. Since we are adjusting for *age group* as a categorical covariate, the age group coefficients are the weighted (weight = $b_w tiron$) proportion of people in each of the two categories.

Date: 05-29-97 Research Triangle Institute Page : 1 Table : 1 Time: 15:28:17 The REGRESS Procedure Variance Estimation Method: Robust (Binder, 1983) Working Correlations: Independent Link Function: Identity Response variable B_TIBC: TOTAL IRON BINDING CAPACITY LS Means Contrast _____ Age Cohort Age Cohort 20-49 yrs. Intercept 50+ yrs. _____ Smoking Status Current 1.000 0.603 0.397 0.603 1.000 0.397 Former 1.000 0.603 Never 0.397 Unknown 1.000 0.603 0.397

Age Group Contrast Coefficients

Least Squares Means Contrast Coefficients:

Smoking Status Coefficients

The contrast coefficients for *smoking status* are 1's and 0's, indicating the level of interest in each row.

```
Date: 05-29-97
           Research Triangle Institute
                                             Page : 2
Time: 15:28:17
                  The REGRESS Procedure
                                             Table : 1
Variance Estimation Method: Robust (Binder, 1983)
Working Correlations: Independent
Link Function: Identity
Response variable B_TIBC: TOTAL IRON BINDING CAPACITY
LS Means Contrast
_____
           Smoking Status Smoking Status Smoking Status Smoking Status
            Current Former Never Unknown
_____
Smoking Status
                1.0000.0000.0000.0001.0000.0000.0000.0001.0000.0000.0000.000
 Current
                                                 0.000
 Former
                                                0.000
                                                0.000
 Never
                                      0.000
 Unknown
                 0.000
                           0.000
                                                 1.000
 _____
```

Least Squares Means Results

Age Group as Categorical Covariate

This table shows the *estimated least-squares means*, with standard errors that are adjusted for clustering and stratification (via the NEST statement and DESIGN=WR option on the PROC statement).

Date: 05-29-97 Time: 15:28:17		Triangle EGRESS Pr			Page Table	
Variance Estimation Me Working Correlations: Link Function: Identit Response variable B_TI	Independent Y					
LSMEANS With Categoric	al Covariate					
Least-Square Means				P-value		
			T-Test LSM=0	T-Test		
				T-Test		
Least-Square Means	LS Mean	Mean		T-Test LSM=0		
Least-Square Means Smoking Status	LS Mean 349.7372	Mean 2.1938	LSM=0	T-Test LSM=0 		
Least-Square Means Smoking Status Current	LS Mean 349.7372 355.5680	Mean 2.1938 2.2920	LSM=0 	T-Test LSM=0 0.0000 0.0000		

Least Squares Means Contrast Coefficients:

Age at Exam as Continuous Covariate

Now we show how the contrast is formed when age is modelled as a *continuous* covariate.

```
11 PROC REGRESS DATA="C:\\ADVANCED\\IRONSUD" FILETYPE=SAS DESIGN=WR DEFT2;
12 NEST Q_STRATA PSU1;
13 WEIGHT B_WTIRON;
14 SUBGROUP SMOKE;
15 LEVELS 4;
16 MODEL B_TIBC = SMOKE AGEXAM;
17 LSMEANS SMOKE / DISPLAY;
18 SETENV COLSPCE=1 LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60;
19 PRINT BETA="BETA" SEBETA="S.E." DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
         P_BETA="P-VALUE" DF WALDCHI WALDCHP /
         LSMEANS=ALL T_BETAFMT=F8.2 DEFTFMT=F6.2 DFFMT=F8.0 WALDCHIFMT=F8.2;
20 TITLE "LSMEANS With Continuous Covariate";
Opened SAS data file C:\ADVANCED\IRONSUD.SSD for reading.
                             : 3290 Weighted count: 40570323
Number of observations read
Observations used in the analysis : 3290 Weighted count: 40570323
Observations with missing values : 0 Weighted count:
                                                                  0
Denominator degrees of freedom :
                                      35
Maximum number of estimable parameters for the model is 5
File C:\ADVANCED\IRONSUD.SSD contains 67 clusters
Maximum cluster size is 111 records
Minimum cluster size is 15 records
Weighted mean response is 354.580621
```

Estimated Regression Coefficients for the Model

Age at Exam as Continuous Covariate

Date: 05-29-97		5				age : 3
Time: 15:28:17	The RI	EGRESS Pi	cocedure	ē	Ta	able : 1
Variance Estimation Working Correlations Link Function: Ident Response variable B_	s: Independent tity		·	ITY		
LSMEANS With Continu	uous Covariate					
Independent Variable	es and					
Independent Variable Effects		S F	DESIGN	 Т:вета=0	D-VALUE	
-		S.E.		T:BETA=0	P-VALUE	
-			EFFECT			
Effects	BETA		EFFECT			
Effects Intercept	BETA	4.9483	EFFECT 1.19	74.86	0.0000	
Effects Intercept Smoking Status	BETA 	4.9483 3.7812	EFFECT 1.19 0.95	74.86 -2.14	0.0000	
Effects Intercept Smoking Status Current	BETA 370.4372 -8.0845 -2.0617	4.9483 3.7812 4.2763	EFFECT 1.19 0.95 1.26	74.86 -2.14	0.0000 0.0396 0.6327	
Effects Intercept Smoking Status Current Former	BETA 370.4372 -8.0845 -2.0617	4.9483 3.7812 4.2763 3.8930	EFFECT 1.19 0.95 1.26 1.09	74.86 -2.14 -0.48 -0.39	0.0000 0.0396 0.6327 0.6989	

Least Squares Means Contrast Coefficients:

Age at Exam as Continuous Covariate

When age at initial exam is modelled as a continuous covariate, its single contrast coefficient is the weighted mean of *AGEXAM* (45.706 years). The contrast coefficients for Smoking status are the same as previously.

Date: 05-29-97	Research Triangle Institute	Page : 2
Time: 15:28:17	The REGRESS Procedure	Table : 1
Working Correlation Link Function: Iden	-	
Response variable i		
LS Means Contrast		
-	Age at Exam	
-		
LS Means Contrast		
LS Means Contrast	Age at Exam	
LS Means Contrast	Age at Exam 45.706	

Least Squares Means Results with Age as Continuous Covariate

This table shows the *estimated least-squares means*, with standard errors that are adjusted for clustering and stratification (via the NEST statement and DESIGN=WR option on the PROC statement), when Age is modelled as a continuous covariate.

Date: 05-29-97 Time: 15:28:17	Research Triangle Institute The REGRESS Procedure					5 1
Variance Estimation M Working Correlations Link Function: Ident: Response variable B_1	: Independent ity		·			
LSMEANS With Continue	ous Covariate					
LSMEANS With Continue	ous Covariate			P-value		
			T-Test LSM=0	T-Test		
				T-Test		
Least-Square Means	LS Mean	Mean		T-Test LSM=0		
Least-Square Means Smoking Status	LS Mean 349.6539	Mean 	LSM=0	T-Test LSM=0 		
Least-Square Means Smoking Status Current	LS Mean 349.6539 355.6767	Mean 2.2333 2.2900	LSM=0 156.5668	T-Test LSM=0 0.0000 0.0000		

Design Effects

Numerator

Variance calculated according to the user-specified sample design option, the working correlation structure specified (independent or exchangeable), and the standard error method (robust vs. model-based).

Denominator = SRS Variance

Calculated according to the type of design effect requested on the PROC statement (see below: DEFT1, DEFT2, DEFT3, DEFT4). DEFT4 is the default, and leads to SRSCOV calculated as the model-based variance under the user-specified correlation structure (independent vs. exchangeable).

Design Effect	Measures Variance Inflation Due to:	Default?
DEFT1	Stratification (or blocking), Clustering, Unequal Weighting, and <i>Oversampling</i> Assumes that total sample size is fixed	No; This is the original one; Request on PROC Statement
DEFT2	Stratification (or blocking), Clustering, and Unequal Weighting Assumes that subgroup sample sizes are fized	No; Request on PROC statement
DEFT3	Stratification (or blocking), Clustering Assumes that subgroup sample sizes are fixed	No; Request on PROC Statement
DEFT4	Stratification (or blocking), Clustering, and Unequal Weighting: <i>Model-based</i> SRS variance (this is the standard software variance when no weights involved) <i>Good for experimental designs</i>	Yes

Evaluation of a Drug Abuse Prevention Program: Project DARE

- Ennett, Rosenbaum, Flewelling, Bieler, *et al* (1994)
 Norton, Bieler, Ennett, and Zarkin (1996)
- Longitudinal evaluation of the DARE program (Drug Abuse Resistance Education) in northern and central Illinois
- Semester-long drug-use prevention program for upper elementary school students (5th and 6th graders)
- Convenience sample of 36 schools (clusters) representing urban, suburban, and rural areas, randomly assigned to DARE and control conditions
- Data represent responses from students immediately before and after program implementation (Waves 1 and 2)
- 1,525 students present at both waves of data collection
- Outcome: Initiation of cigarette smoking by Wave 2 (includes only those students reporting no lifetime use in wave 1)

$$y = \begin{cases} 0, & \text{if student did not initiate smoking by follow-up} \\ 1, & \text{if student initiated smoking by follow-up} \end{cases}$$

Question: Does the DARE program reduce the incidence of adolescent cigarette smoking (at least during the intervention)?

Fitting GEE Logistic Regression Models in MULTILOG

Evaluation of a Drug Abuse Prevention Program (Project DARE)

Experimental studies of the effect of prevention programs on substance use are often based on nested cohort designs, in which intact social groups or clusters of individuals are randomized to treatment conditions, and individuals within the clusters are followed over time as a cohort to evaluate the effects of treatment. The units of assignment may be schools, communities, or worksites, but the units of observation are the students, community residents, or workers. Because they are exposed to a common set of circumstances, students within the same school tend to be positively correlated with one another. This positive intracluster correlation implies that the observational units are no longer statistically independent. Unless the intracluster correlation that results from the sampling design is accounted for in the statistical analysis, estimated standard errors of the treatment effects will generally be underestimated, leading to inflated Type I error rates and false-positive tests of treatment effects (Murray and Hannan, 1990; Moskowitz, Malvin, Schaeffer, and Schaps, 1984; Donner, 1982; Donner, Birkett, and Buck, 1981).

Illustrative data for this example were collected as part of a longitudinal evaluation of Project DARE (Drug Abuse Resistance Education) on substance abuse outcomes in Illinois (Ennett, Rosenbaum, Flewelling, Bieler, Ringwalt, and Bailey, 1994). The DARE curriculum is a semester-long drug use prevention program for late elementary school students. Respondents for the study were originally obtained in 1990 from the fifth and sixth grades of 36 schools representative of rural, urban, and suburban areas in the state of Illinois. Within each metropolitan status stratum, 6 pairs of schools (matched on various demographic characteristics) were randomly assigned to DARE and control conditions.

Researchers collected data immediately before and after program implementation (Waves 1 and 2) and have collected three additional waves at annual intervals since then. Analyses reported here draw only on data from Waves 1 and 2. The sample includes students for whom complete information is available on the variables of interest in both waves (N = 1525, 85% of Wave 1 sample). Students answered a self-administered questionnaire that took approximately 35 minutes to complete. The questions concerned substance use, attitudes toward drugs, self-esteem, and peer-resistance skills.

In this example we analyze a single dependent variable that is representative of outcome measures used to evaluate drug use prevention programs. At each Wave of data collection, students were asked whether they had ever smoked cigarettes. The binary dependent variable relates to the initiation of cigarette use between Waves 1 and 2 (coded 1 if the adolescent initiated cigarette use; 2 = otherwise). The desired effect is a negative correlation with DARE (coded 1 = adolescent exposed to DARE, 2 = not exposed). The sample for initiation analysis is limited to students who reported no lifetime use at Wave 1.

We report results for the covariate of primary interest, exposure to the DARE program, as well as the following background characteristics (with 8 degrees of freedom): grade in school, sex, race/ethnicity, family composition, and metropolitan status. Respondents included 34% fifth and 66% sixth-grade students; approximately half were male. The sample was 51% white, 24%

African American, 9% Hispanic, and 16% "other". The majority (65%) lived with both parents in the same household. Fewer respondents lived in rural areas (26%) compared with suburban (38%) and urban (36%) areas.

We used SUDAANs MULTILOG procedure to fit a logistic regression model to the binary response variable of interest via the GEE model-fitting method, under both independent and exchangeable working correlations. The independence working assumption here amounts to ordinary logistic regression. The use of the variance correction (standard in SUDAAN) yields valid results in the presence of intracluster correlation. In fact, the robust variance estimate ensures that the results are robust to any misspecification of the correlation structure. We also provide results using the model-based variance estimates. In Table 1, we compare the GEE/SUDAAN results to SAS PROC LOGISTIC, which currently fits ordinary logistic regression but naively makes no correction for intracluster correlation and instead considers the observations statistically independent.

Using SUDAAN, the DARE program is shown to have a significant negative effect on the initiation of cigarette use, regardless of the working assumptions about the correlation structure (p=0.0369 under working independence; p=0.0216 under exchangeability). The estimated intracluster correlation under exchangeability is 0.0206. Use of a robust variance estimate ensures that the results of statistical analyses are valid no matter what the true correlation structure is. In this example, the exchangeability assumption appears to be correct, since results using the robust and model-based variance estimates were essentially the same. The advantage of modelling the correlation structure (*e.g.*, through exchangeability) is its potential to improve efficiency and hence increase the power of statistical analyses.

The incidence of cigarette use during the intervention was significantly lower among students who participated in DARE (9.5% observed for DARE vs. 15.4% for controls). As seen in Table 1, naively ignoring the intracluster correlation as in SAS PROC LOGISTIC leads to a much more significant treatment effect (p=0.0069). The observed design effect for DARE was 1.75, which indicates almost a doubling in the variance of the estimated treatment effect under cluster randomization.

Exposure Group 1 = Control 2 = DARE	School ID (Cluster)	Student ID (unit of observation)	Y = cigarette initiation 1= yes 2 = no
1	1	1	2
1	1	2	1
1	1	3	2
1	2	1	2
1	2	2	2
2	10	1	2
2	10	2	1
2	20	1	1
2	20	2	1
2	30	1	1

Structure of the DARE Data

N = 1,525 records on the file

(1,525 students clustered within 36 schools)

Evaluation of the DARE Effect on Cigarette Initiation
Via Logistic Regression Modelling

		Working Correlations			
		Independent (Ordinary Logistic Regression)		Exchangeable	
Variable	Statistic	No Variance Correction	Variance Correction	No Variance Correction	Variance Correction
Initiation of	β	-0.5225	-0.5225	-0.5825	-0.5825
Cigarette Use	SE	0.1821	0.2408	0.2433	0.2422
By Wave 2	Observe d DEFF		1.75		1.77
	Z- statistic	-2.87	-2.17	-2.39	-2.41
	P-value	0.0069	0.0369	0.0221	0.0216

Working Correlations:	Software:		
Independent (Ordinary Logistic Regression)			
No variance correction:	SAS Logistic		
Variance Correction (robust variance):	SUDAAN Multilog		
Exchangeable			
No variance correction [model-based (naive) variance]:	SUDAAN Multilog		
Variance Correction (robust variance):	SUDAAN Multilog		

MULTILOG Programming Statements and Options

The following sets of programming statements fit different versions of a logistic model in SUDAAN PROC MULTILOG. The **DATA** option on the **PROC** statement specifies a SAS data set as input. Since there is no **DESIGN** option specified on the PROC statement, SUDAAN is using the default **DESIGN=WR** (with-replacement) option for variance estimation.

In the accompanying output, we fit the following types of GEE logistic regression models:

1) SEMETHOD=ZEGER and R=INDEPENDENT

Implements the GEE model-fitting technique under an *independent "working" assumption* and Zeger and Liang's (1986) *robust* variance estimator. This model is sometimes referred to as ordinary logistic regression <u>with</u> a variance correction. Note that for binary outcomes, SEMETHOD=ZEGER is equivalent to SEMETHOD=BINDER.

2) SEMETHOD=MODEL and R=INDEPENDENT

This amounts to ordinary logistic regression <u>without</u> a variance correction, which yields the same results as SAS PROC LOGISTIC. Literally, this combination implies an *independent "working" assumption* and a *model-based* or *naive* variance estimator. The variance estimator is naive in the sense that it computes variances as if the independence working assumption were correct.

3) SEMETHOD=ZEGER and R=EXCHANGEABLE

Implements the GEE model-fitting technique under *exchangeable "working" correlations* and Zeger and Liang's (1986) *robust* variance estimator.

4) SEMETHOD=MODEL and R=EXCHANGEABLE

We compare the results from the robust variance estimator (*SEMETHOD=ZEGER*) to the *model-based*, or *naive*, variance assumption (*SEMETHOD=MODEL*). When R=exchangeable is specified in conjunction with *SEMETHOD=MODEL*, variances are then computed as if the *exchangeable "working" correlation* assumption were correct.

The **NEST** statement indicates that SCHOOL is the cluster variable. The **WEIGHT** statement indicates equal sampling weights of 1.0 for each student on the file.

In MULTILOG, the **SUBGROUP** statement contains the dependent variable and all covariates that are to be modelled as categorical covariates (with level values of 1, 2, ..., K), where the maximum number of levels (*K*) appears on the **LEVELS** statement.

The **MODEL** statement specifies the categorical dependent variable INTCIG12 on the left of the "=" sign (with levels 1 and 2), and regressors on the right. For binary responses, the **CUMLOGIT** (cumulative logit) and **GENLOGIT** (generalized logit) links specify the same logistic regression model.

The **TEST** statement specifies that we want the Wald chi-square statistic to be the default for testing main effects, interactions, and user-defined contrasts.

Descriptive Statistics for Initiation of Cigarette Smoking in the DARE Study

DES Var	SCRIPTIVE STATISTICS FOR THE DARE DATA riable Sample Design DARE Program Size Percent STDERR Effect itiation of Cigarette Use: Yes
DES	
Var	
	riance Estimation Method: Taylor Series (WR)
	te: 04-28-97Research Triangle InstitutePage :me: 13:30:35The DESCRIPT ProcedureTable :
Num Den	ened SAS data file c:\tera\examples\DARE.SSD for reading. mber of observations read : 1525 Weighted count : 1525 nominator degrees of freedom : 35
	TITLE "DESCRIPTIVE STATISTICS FOR THE DARE DATA";
10	PRINT NSUM PERCENT SEPERCENT="STDERR" DEFFPCT="Design Effect" / NSUMFMT=F6.0 PERCENTFMT=F7.2 STYLE=NCHS;
9	SETENV LABWIDTH=30 COLWIDTH=6 DECWIDTH=2;
8	CATLEVEL 1;
7	VAR INTCIG12;
б	TABLES DARE;
5	LEVELS 2;
4	SUBGROUP DARE;
3	WEIGHT _ONE_;
	NEST _ONE_ SCHOOL;
2	

These results indicate that 15.4% of students not receiving DARE initiated cigarette smoking during the time of the intervention, compared to 9.5% of those exposed to DARE. The standard errors estimated by SUDAAN use a between-cluster variance formula and are therefore adjusted for clustering. The design effects indicate that the variances of the percentages are more than doubled under cluster randomization. Is the observed difference statistically significant, after adjustment for other covariates? The MULTILOG procedure will be used to find out.

GEE With Independent "Working" Correlations Robust Variance Estimator

```
12 PROC MULTILOG DATA="c:\\tera\\examples\\DARE" FILETYPE=SAS
                  SEMETHOD=ZEGER R=INDEPENDENT;
13 NEST _ONE_ SCHOOL;
14 WEIGHT _ONE_;
  SUBGROUP DARE FIFTH SEX RACE OTHFAM AREA INTCIG12;
15
                                        3
                                             2;
16 LEVELS
             2
                  2
                        2
                                 2
                            4
17 MODEL INTCIG12 = DARE FIFTH SEX RACE OTHFAM AREA / CUMLOGIT;
18 TEST WALDCHI;
  SETENV LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60 COLSPCE=2;
19
20 PRINT BETA="BETA" SEBETA="STDERR" DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
          P_BETA="P-VALUE" / RISK=ALL TESTS=DEFAULT
          BETAFMT=F8.6 SEBETAFMT=F8.6 T_BETAFMT=F8.2 WALDCHIFMT=F6.2
          WALDCHPFMT=F7.4 DEFTFMT=F6.2 DFFMT=F7.0
          ORFMT=F5.2 LOWORFMT=F6.2 UPORFMT=F6.2;
21 TITLE "MULTILOG Logistic Regression Model for the DARE Evaluation Study"
          "Ennett, et al, 1994";
Opened SAS data file c:\tera\examples\DARE.SSD for reading.
Number of observations read
                                 :
                                     1525
                                              Weighted count:
                                                                  1525
Observations used in the analysis :
                                     1188
                                             Weighted count:
                                                                  1188
Observations with missing values :
                                      337
                                              Weighted count:
                                                                  337
Denominator degrees of freedom
                                 :
                                       35
Maximum number of estimable parameters for the model is 10
File c:\tera\examples\DARE.SSD contains
                                         36 Clusters
Maximum cluster size is 153 records
Minimum cluster size is 11 records
Independence parameters have converged in 4 iterations
Sample and Population Counts for Response Variable INTCIG12
  Yes: Sample Count
                         145
                                 Population Count
                                                        145
 No : Sample Count
                        1043
                                 Population Count
                                                       1043
```

Here we see that there are 1,525 students (1 record/student) on the file, and that 1,188 were used in the analysis (337 students deleted due to missing values on one or more MODEL statement variables). There are 36 clusters (schools), with cluster sizes ranging from 11 to 153. Overall, 145 students reported having initiated cigarette use during the intervention, while 1043 did not.

GEE With Independent "Working" Correlations Robust Variance Estimator

Time: 13:30:35	The MUL	TILOG Proc	edure		Table :
	_				
Variance Estimation Meth		Zeger-Lian	g, 1986)		
orking Correlations: In	-				
Jink Function: Cumulativ Response variable INTCIG	-	on of circo	retta ITa	۹	
Cobourse variable INICIG	120 IIIILIALI	on or crda	LELLE US		
MULTILOG Logistic Regres	sion Model f	or the DAR	E Evalua	tion Study	
INTITION DOGIDETE REGIES	SION HOUEL L	OI CHE DAN	L DVarua	cion beday	
			DESIGN		
and Effects	RETA	STDEFF		T:BETA=0	P-VALIIF
NTCIG12 (cum-logit)					
	-1.84755	0.465860	1.99	-3.97	0.0003
DARE Program					
Yes	-0.52248	0.240770	1.75	-2.17	0.0369
No		0.000000			
Frade in School					
5th Grade	-0.50020	0.249380	1.25	-2.01	0.0527
6th Grade	0.00000	0.000000	•		•
SEX					
Male	0.084014	0.159940	0.78	0.53	0.6027
Female	0.00000	0.000000	•	•	•
ACE					
Black		0.378610			0.1977
Hispanic	0.095132	0.467026	1.46	0.20	0.8398
Other	0.493601		2.23	1.17	0.2494
White	0.00000	0.000000	•		•
amily Situation					
Non-Traditional				2.47	0.0187
Traditional	0.00000	0.000000	•	•	•
REA					
Rural				-0.20	
Suburban			1.77	-0.69	0.4918
	0.00000				

This first table contains the estimated regression coefficient vector, the estimated robust standard errors, design effects, t-statistics, and p-values for testing H₀: β =0. The CUMLOGIT option estimates only one model intercept in the case of a binary outcome, and is equivalent to the GENLOGIT option. The treatment effect (DARE) is observed to significantly reduce the incidence of cigarette initiation (p=0.0369) using the GEE-independent approach, after adjusting for other covariates in the model. Other than the treatment effect, only family situation is a statistically significant covariate (p=0.0187). The observed design effect for the treatment parameter is 1.75, indicating a 75% increase in variance due to cluster randomization.

GEE With Independent "Working" Correlations Robust Variance Estimator

Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure Page : 2 Time: 13:30:35 The MULTILOG Procedure Table : 1 Variance Estimation Method: Robust (Zeger-Liang, 1986) Working Correlations: Independent Link Function: Cumulative Logit Response variable INTCIG12: Initiation of Cigarette Use MULTILOG Logistic Regression Model for the DARE Evaluation Study Ennett, et al, 1994 _____ Degrees P-value Contrast of Wald Wald Freedom ChiSq ChiSq _____ 10 313.24 0.0000 OVERALL MODEL 9 28.27 0.0009 MODEL MINUS INTERCEPT 1 4.71 0.0300 1 4.02 0.0449 DARE FIFTH 1 0.28 0.5994 SEX 3 1.90 0.5930 RACE 1 6.08 0.0137 OTHFAM 2 0.60 0.7420 AREA _____

This table contains the statistical significance of all main effects, interactions, and user-defined contrasts. The Wald chi-square test (from the TEST statement) is used to evaluate these effects.

GEE With Independent "Working" Correlations Robust Variance Estimator

Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure Page : 3 Table : 1 Variance Estimation Method: Robust (Zeger-Liang, 1986) Working Correlations: Independent Link Function: Cumulative Logit Response variable INTCIG12: Initiation of Cigarette Use MULTILOG Logistic Regression Model for the DARE Evaluation Study Ennett, et al, 1994 _____ Independent Variables Lower Upper and Effects Odds 95% 95% Ratio Limit Limit _____ INTCIG12 (cum-logit) 0.16 0.06 0.41 Intercept 1: Yes DARE Program 0.59 0.36 0.97 Yes 1.00 1.00 1.00 No Grade in School 0.61 0.37 1.01 1.00 1.00 1.00 5th Grade 6th Grade SEX Male 1.09 0.79 1.50 Female 1.00 1.00 1.00 RACE 1.640.763.551.100.432.841.640.703.851.001.001.00 Black Hispanic Other White Family Situation
 Non-Traditional
 1.52
 1.08
 2.15

 Traditional
 1.00
 1.00
 1.00
 AREA Rural 0.92 0.41 2.07 0.78 0.37 Suburban 1.62 1.00 1.00 1.00 Urban _____ MULTILOG used CPU time : 12.3 seconds Elapsed time : 13 seconds Virtual memory : 1.52 MB

This final table contains the estimated odds ratios and their 95% confidence limits for each regression coefficient in the model. We see that the negative regression coefficient for DARE corresponds to an odds ratio for smoking initiation of 0.59, indicating a protective effect of the DARE program (the odds are reduced by around 40% in the DARE group). Again, each regression coefficient is adjusted for all others in the model.

GEE with Independent "Working" Correlations Model-Based (Naive) Variance Estimator

Below are the results obtained under working independence using the *model-based* or *naive variance-covariance matrix* of the estimated regression coefficients. The model-based variance is the M_0^{-1} matrix, or the outside portion of the robust variance estimate: $M_0^{-1} = [D'V^{-1}D]^{-1}$, where $D = \partial \pi_i / \partial \beta$ is the vector of first partial derivatives of the response probabilities π_i with respect to the regression coefficients β . In this case, the naive variance estimate is computed as *if the independent working correlation assumption were correct*. In other words, these are the results that would be obtained if clustering were ignored altogether. Although it is not recommended for analysis of clustered data, we are showing it to demonstrate the effects of clustering. We use the *SEMETHOD=MODEL* option on the PROC statement to obtain the model-based results.

```
PROC MULTILOG DATA="c:\\tera\\examples\\DARE" FILETYPE=SAS
2.2
                 SEMETHOD=MODEL R=INDEPENDENT;
23 NEST _ONE_ SCHOOL;
24 WEIGHT _ONE_;
25
  SUBGROUP DARE FIFTH SEX RACE OTHFAM AREA INTCIG12;
26
  LEVELS 2
                 2
                       2
                           4
                                2
                                       2
                                            2;
   MODEL INTCIG12 = DARE FIFTH SEX RACE OTHFAM AREA / CUMLOGIT;
27
  TEST WALDCHI;
28
29 SETENV LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 MAXIND=4 LINESIZE=78
          PAGESIZE=60 COLSPCE=2;
30 PRINT BETA="BETA" SEBETA="STDERR" DEFT="DESIGN EFFECT" T_BETA="T:BETA=0"
          P_BETA="P-VALUE" / RISK=ALL TESTS=DEFAULT
          BETAFMT=F8.6 SEBETAFMT=F8.6 T BETAFMT=F8.2 WALDCHIFMT=F6.2
          WALDCHPFMT=F7.4 DEFTFMT=F6.2 DFFMT=F7.0
          ORFMT=F5.2 LOWORFMT=F6.2 UPORFMT=F6.2;
31 TITLE "MULTILOG Logistic Regression Model for the DARE Evaluation Study"
          "Model-Based Variance Estimation";
Opened SAS data file c:\tera\examples\DARE.SSD for reading.
Number of observations read
                              : 1525
                                            Weighted count:
                                                                1525
                                          Weighted count:
Observations used in the analysis : 1188
                                                                1188
Observations with missing values : 337 Weighted count:
                                                                 337
Denominator degrees of freedom :
                                      35
Maximum number of estimable parameters for the model is 10
```

File c:\tera\examples\DARE.SSD contains 36 Clusters
Maximum cluster size is 153 records
Minimum cluster size is 11 records
Independence parameters have converged in 4 iterations
Sample and Population Counts for Response Variable INTCIG12
Yes: Sample Count 145 Population Count 145
No : Sample Count 1043 Population Count 1043

GEE with Independent "Working" Correlations Model-Based (Naive) Variance Estimator

Time: 13:30:35	ING MOL	TILOG Proc	.euure		Table :
Variance Estimation Meth	od: Model-Ba	sed (Naive)		
Working Correlations: In					
Link Function: Cumulativ	-				
Response variable INTCIG	12: Initiati	on of Ciga	rette Us	е	
MULTILOG Logistic Regres	sion Model f	or the DAR	E Evalua	tion Study	
Model-Based Variance Est	imation				
Independent Variables			DESIGN		
and Effects				T:BETA=0	P-VALUE
 INTCIG12 (cum-logit)					
Intercept 1: Yes	_1 Q1755	0 220101	1 00	_5 50	0 0000
DARE Program	1.04/33	0.000404	1.00	ور. ر –	0.0000
Yes	-0 500/0	0 182076	1 00	-2.87	0 0060
No		0.182078			0.0009
Grade in School	0.000000	0.000000	•	•	•
5th Grade	-0 50020	0.223481	1 00	-2.24	0 0317
6th Grade		0.223401			
SEX	2.200000		•	•	•
Male	0.084014	0.181161	1.00	0.46	0.6457
Female	0.000000				
RACE					
Black	0.497135	0.283813	1.00	1.75	0.0886
Hispanic	0.095132	0.386261	1.00	0.25	0.8069
Other		0.282336		1.75	0.0892
White		0.000000			
Family Situation					
Non-Traditional	0.420841	0.192760	1.00	2.18	0.0358
Traditional	0.000000	0.000000		•	
AREA					
Rural	-0.07878	0.319795	1.00	-0.25 -0.92	0.8069
Suburban			1.00	-0.92	0.3613
	0.00000	0.000000			

Here we see that the estimated regression coefficients are the same as previously, but the estimated standard errors using the model-based approach under independence are much smaller than with the robust variance estimator. The effects of DARE (p=0.0069), family situation (p=0.0358), and grade in school (p=0.0317) are all statistically significant. These standard error estimates are overly optimistic (naive), computed as if the data were truly independent. Therefore, these results are not valid for the data at hand. They merely demonstrate the

consequences of ignoring the experimental design. The design effects are all equal to 1.0, since both numerator and denominator values are the same.

GEE with Independent "Working" Correlations Model-Based (Naive) Variance Estimator

Time: 13:30:35			e Institute rocedure		Page : 2 Table : 1
Variance Estimation Meth Working Correlations: In		ased (Na	ive)		
Link Function: Cumulativ	re Logit				
Response variable INTCIG	12: Initiat	ion of C	igarette Use	2	
MULTILOG Logistic Regres	sion Model	for the	DARE Evaluat	tion Study	7
Model-Based Variance Est	imation				
	Degrees		P-value		
	Degrees of	Wald	P-value Wald		
Contrast	Degrees of Freedom	Wald ChiSq	P-value Wald ChiSq		
	Degrees of Freedom	Wald ChiSq	P-value Wald ChiSq		
Contrast	Degrees of Freedom 10	Wald ChiSq 470.28	P-value Wald ChiSq 0.0000		
Contrast OVERALL MODEL	Degrees of Freedom 10 9	Wald ChiSq 470.28 30.41	P-value Wald ChiSq 0.0000		
Contrast OVERALL MODEL MODEL MINUS INTERCEPT	Degrees of Freedom 10 9 1	Wald ChiSq 470.28 30.41 8.23	P-value Wald ChiSq 0.0000 0.0004		
Contrast OVERALL MODEL MODEL MINUS INTERCEPT DARE	Degrees of Freedom 10 9 1 1	Wald ChiSq 470.28 30.41 8.23 5.01	P-value Wald ChiSq 0.0000 0.0004 0.0041		
Contrast OVERALL MODEL MODEL MINUS INTERCEPT DARE FIFTH	Degrees of Freedom 10 9 1 1 1	Wald ChiSq 470.28 30.41 8.23 5.01 0.22	P-value Wald ChiSq 0.0000 0.0004 0.0041 0.0252		
Contrast OVERALL MODEL MODEL MINUS INTERCEPT DARE FIFTH SEX	Degrees of Freedom 10 9 1 1 1 3	Wald ChiSq 470.28 30.41 8.23 5.01 0.22 4.64	P-value Wald ChiSq 0.0000 0.0004 0.0041 0.0252 0.6428		

This table contains the main effects tests computed as if the naive assumption of independence were true. The Wald chi-square test is used to evaluate the null hypotheses.

GEE with Independent "Working" Correlations Model-Based (Naive) Variance Estimator

Date: 04-28-97Research Triangle InstitutePage : 3Time: 13:30:35The MULTILOG ProcedureTable : 1 Variance Estimation Method: Model-Based (Naive) Working Correlations: Independent Link Function: Cumulative Logit Response variable INTCIG12: Initiation of Cigarette Use MULTILOG Logistic Regression Model for the DARE Evaluation Study Model-Based Variance Estimation _____ Independent Variables Lower Upper and Effects Odds 95% 95% Ratio Limit Limit Ratio Limit Limit -----INTCIG12 (cum-logit) Intercept 1: Yes 0.16 0.08 0.31 DARE Program 0.59 0.41 0.86 Yes 1.00 1.00 1.00 No Grade in School 0.61 0.39 0.95 5th Grade 1.00 1.00 1.00 6th Grade SEX 1.09 0.75 1.57 Male Female 1.00 1.00 1.00 RACE 1.640.922.931.100.502.411.640.922.911.001.001.00 Black Hispanic Other White Family Situation
 Non-Traditional
 1.52
 1.03
 2.25

 Traditional
 1.00
 1.00
 1.00
 AREA 0.92 0.48 1.77 0.78 0.45 1.35 Rural Suburban Urban 1.00 1.00 1.00 _____ MULTILOG used CPU time : 6.81 seconds Elapsed time : 7 seconds Virtual memory : 1.52 MB

GEE with Exchangeable "Working" Correlations Robust Variance Estimator

This next set of SUDAAN programming statements fits the logistic regression model via the GEE model-fitting technique, under the assumption of exchangeable "working" correlations (R=exchangeable) and using a robust variance estimator. All other programming statements remain the same as previously.

```
32 PROC MULTILOG DATA="c:\\tera\\examples\\DARE" FILETYPE=SAS
                 SEMETHOD=ZEGER R=EXCHANGE;
33 NEST _ONE_ SCHOOL;
34 WEIGHT _ONE_;
35 SUBGROUP DARE FIFTH SEX RACE OTHFAM AREA INTCIG12;
36 LEVELS
            2
                 2
                       2
                           4
                                2
                                       3
                                            2;
37 MODEL INTCIG12 = DARE FIFTH SEX RACE OTHFAM AREA / CUMLOGIT;
38 TEST WALDCHI;
  SETENV LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60 COLSPCE=2;
39
40 PRINT BETA="BETA" SEBETA="STDERR" T_BETA="T:BETA=0" P_BETA="P-VALUE"
          / RISK=ALL TESTS=DEFAULT RHOS=ALL BETAFMT=F8.6 SEBETAFMT=F8.6
            T BETAFMT=F8.2 WALDCHIFMT=F6.2 WALDCHPFMT=F7.4 DFFMT=F7.0
            ORFMT=F5.2 LOWORFMT=F6.2 UPORFMT=F6.2;
41 TITLE "MULTILOG Logistic Regression Model for the DARE Evaluation Study"
         "Ennett, et al, 1994";
Opened SAS data file c:\tera\examples\DARE.SSD for reading.
Number of observations read
                              : 1525 Weighted count:
                                                                1525
Observations used in the analysis : 1188 Weighted count:
                                                                1188
Observations with missing values :
                                     337
                                             Weighted count:
                                                                 337
Denominator degrees of freedom
                                 :
                                       35
Maximum number of estimable parameters for the model is 10
File c:\tera\examples\DARE.SSD contains 36 Clusters
Maximum cluster size is 153 records
Minimum cluster size is 11 records
Independence parameters have converged in 4 iterations
Step 1 parameters have converged in 6 iterations.
Sample and Population Counts for Response Variable INTCIG12
  Yes: Sample Count
                        145
                                Population Count
                                                       145
 No : Sample Count
                        1043
                                Population Count
                                                      1043
```

By default, SUDAAN fits the 1-step GEE estimates (Lipsitz, et al., 1994). Here we see that the independence betas (the starting estimates for GEE exchangeable) have converged in 4 iterations, and the Step 1 GEE parameter estimates (under exchangeable working correlations) have converged in 6 iterations.

GEE with Exchangeable "Working" Correlations Robust Variance Estimator

```
Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure
                                                     Page : 1
Time: 13:30:35
                                                     Table : 1
Variance Estimation Method: Robust (Zeger-Liang, 1986)
Working Correlations: Exchangeable
Link Function: Cumulative Logit
Response variable INTCIG12: Initiation of Cigarette Use
MULTILOG Logistic Regression Model for the DARE Evaluation Study
Ennett, et al, 1994
Correlation Matrix
_____
Initiation of Cigarette Initiation of Cigarette Use
 Use
                         Yes
_____
                      0.0206
Yes
  _____
```

This table contains the estimated correlation matrix, which has only one parameter because the response is binary. We see that the estimated intracluster correlation is 0.0206. This value will be used in estimating the final regression parameters.

Note that although the intracluster correlation is small, the cluster sizes in this study are large enough to cause almost a doubling in the variance of estimated regression coefficients (deff=1.75 for the DARE effect in the working independence model with robust variance estimate).

GEE with Exchangeable "Working" Correlations Robust Variance Estimator

Date: 04-28-97 Research Triangle Institute Page : 2 Time: 13:30:35 The MULTILOG Procedure Table : 1 Variance Estimation Method: Robust (Zeger-Liang, 1986) Working Correlations: Exchangeable Link Function: Cumulative Logit Response variable INTCIG12: Initiation of Cigarette Use MULTILOG Logistic Regression Model for the DARE Evaluation Study Ennett, et al, 1994 _____ INTCIG12 (cum-logit), Independent Variables and Effects BETA STDERR T:BETA=0 P-VALUE _____ INTCIG12 (cum-logit) Intercept 1: Yes -1.88017 0.449771 -4.18 0.0002 DARE Program -0.58250 0.242184 -2.41 0.0216 Yes 0.000000 0.000000 No . . Grade in School 5th Grade -0.46289 0.221616 -2.09 0.0441 0.000000 0.000000 . . 0.000000 0.000000 6th Grade . . SEX 0.087569 0.157590 0.56 0.5820 0.000000 0.000000 . . Male Female RACE 0.5088010.3677071.380.17520.2778010.4124050.670.50500.5180410.4279641.210.2342 Black Hispanic Other White 0.000000 0.000000 . . Family Situation
 Non-Traditional
 0.436618
 0.173405
 2.52
 0.0165

 Traditional
 0.000000
 0.000000
 .
 .
 AREA -0.067640.378772-0.180.8593-0.261650.371397-0.700.48580.0000000.000000.. Rural Suburban Urban

In this example, the treatment effect (DARE) has become slightly more significant (p=0.0216) under exchangeability, as the parameter estimate (-0.5825) has increased compared to independence (-0.5225). The variance estimate has also increased, but only slightly. Nevertheless, the overall conclusions are qualitatively the same as for independent working correlations with a robust variance estimate.

GEE with Exchangeable "Working" Correlations Robust Variance Estimator

```
Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure
                                                                  Page : 3
                                                                   Table : 1
Variance Estimation Method: Robust (Zeger-Liang, 1986)
Working Correlations: Exchangeable
Link Function: Cumulative Logit
Response variable INTCIG12: Initiation of Cigarette Use
MULTILOG Logistic Regression Model for the DARE Evaluation Study
Ennett, et al, 1994
_____
                       Degrees P-value
of Wald Wald
Contrast
                         Freedom ChiSq ChiSq
_____

        OVERALL MODEL
        10
        312.35
        0.0000

        MODEL MINUS INTERCEPT
        9
        29.80
        0.0005

        DARE
        0
        0
        0
        0

                               1 5.78 0.0162
DARE
                                1 4.36 0.0367
FIFTH
SEX
                                1 0.31 0.5784
RACE
                               3 1.93 0.5867
OTHFAM
                               1 6.34 0.0118
                                2 0.58 0.7497
AREA
_____
```

Here we see the model main effects, under the exchangeable option and a robust variance estimate. All of the effects have become slightly more significant compared to working independence with a robust variance estimate. This should not be taken as a general result for exchangeability vs. working independence. Studies have shown that modelling the correlations tend to yield greater power for detecting within-cluster covariates (Neuhaus, 1993; Lipsitz, Fitzmaurice, Orav, and Laird, 1994), such as sex, race, and family status in the current example. Cluster-level covariates, such as the DARE effect, seem not to benefit as much from modelling the correlation structure.

GEE with Exchangeable "Working" Correlations Robust Variance Estimator

Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure Page : 4 Time: 13:30:35 The MULTILOG Procedure Table : 1 Variance Estimation Method: Robust (Zeger-Liang, 1986) Working Correlations: Exchangeable Link Function: Cumulative Logit Response variable INTCIG12: Initiation of Cigarette Use MULTILOG Logistic Regression Model for the DARE Evaluation Study _____ Independent Variables Lower Upper and Effects Odds 95% 95% Ratio Limit Limit _____ 0.56 0.34 0.91 1.00 1.00 1.00
 No

 Grade in School

 5th Grade
 0.63
 0.40
 0.99

 5th Grade
 1.00
 1.00
 1.00
 No 1.09 0.79 1.50 Male 1.00 1.00 1.00 Female RACE
 Black
 1.66
 0.79
 3.51

 Hispanic
 1.32
 0.57
 3.05

 Other
 1.68
 0.70
 4.00

 White
 1.00
 1.00
 1.00
 Family Situation
 Non-Traditional
 1.55
 1.09
 2.20

 Traditional
 1.00
 1.00
 1.00
 AREA 0.93 0.43 2.02 0.77 0.36 1.64 Rural Suburban Urban 1.00 1.00 1.00 _____ MULTILOG used CPU time : 30.15 seconds Elapsed time : 31 seconds Virtual memory : 1.58 MB

The estimated odds of initiating smoking by Wave 2 is now 0.56 under exchangeability, vs. 0.59 under working independence. It should also be noted that for binary outcomes, GEE results under the GENLOGIT and CUMLOGIT options are identical, since they specify the same model.

GEE with Exchangeable "Working" Correlations Model-Based (Naive) Variance Estimator

Below are results from the exchangeable correlation model using the *model-based* or *naive variance-covariance matrix* of the estimated regression coefficients. The model-based variance is the M_0^{-1} matrix, or the outside portion of the robust variance estimate: $M_0^{-1} = [D'V^{-1}D]^{-1}$, where $D = \partial \pi_i / \partial \beta$ is the vector of first partial derivatives of the response probabilities π_i with respect to the regression coefficients β . In this case, the naive variance estimate is computed *assuming that the exchangeable "working" correlation assumption were correct*. Since that is close to truth for students clustered within schools, we will see that results are essentially the same as with the robust variance estimator.

```
42 PROC MULTILOG DATA="c:\\tera\\examples\\DARE" FILETYPE=SAS
                 SEMETHOD=MODEL R=EXCHANGE;
43 NEST _ONE_ SCHOOL;
44 WEIGHT _ONE_;
   SUBGROUP DARE FIFTH SEX RACE OTHFAM AREA INTCIG12;
45
                       2 4
                                2
                                     3
46
  LEVELS 2
                 2
                                            2;
  MODEL INTCIG12 = DARE FIFTH SEX RACE OTHFAM AREA / CUMLOGIT;
47
48 TEST WALDCHI;
49
  SETENV LABWIDTH=25 COLWIDTH=8 DECWIDTH=4 LINESIZE=78 PAGESIZE=60 COLSPCE=2;
50
  PRINT BETA="BETA" SEBETA="STDERR" T_BETA="T:BETA=0"
          P_BETA="P-VALUE" / RISK=ALL TESTS=DEFAULT RHOS=ALL
          BETAFMT=F8.6 SEBETAFMT=F8.6 T_BETAFMT=F8.2 WALDCHIFMT=F6.2
          WALDCHPFMT=F7.4 DFFMT=F7.0
          ORFMT=F5.2 LOWORFMT=F6.2 UPORFMT=F6.2;
51 TITLE "MULTILOG Logistic Regression Model for the DARE Evaluation Study"
          "Model-Based Variance Estimation";
Opened SAS data file c:\tera\examples\DARE.SSD for reading.
Number of observations read
                              :
                                             Weighted count:
                                                                 1525
                                     1525
Observations used in the analysis :
                                    1188
                                             Weighted count:
                                                                 1188
                                    337
Observations with missing values :
                                             Weighted count:
                                                                 337
Denominator degrees of freedom :
                                     35
Maximum number of estimable parameters for the model is 10
File c:\tera\examples\DARE.SSD contains
                                         36 Clusters
Maximum cluster size is 153 records
```

Minimum cluster size is 11 records Independence parameters have converged in 4 iterations Step 1 parameters have converged in 6 iterations. Sample and Population Counts for Response Variable INTCIG12 Yes: Sample Count 145 Population Count 145 No : Sample Count 1043 Population Count 1043

GEE with Exchangeable "Working" Correlations Model-Based (Naive) Variance Estimator

Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure Page : 1 Time: 13:30:35 Table : 1 Variance Estimation Method: Model-Based (Naive) Working Correlations: Exchangeable Link Function: Cumulative Logit Response variable INTCIG12: Initiation of Cigarette Use MULTILOG Logistic Regression Model for the DARE Evaluation Study Model-Based Variance Estimation _____ Independent Variables and Effects BETA STDERR T:BETA=0 P-VALUE _____ INTCIG12 (cum-logit) Intercept 1: Yes -1.88017 0.369879 -5.08 0.0000 DARE Program -0.58250 0.243284 -2.39 0.0221 Yes 0.000000 0.000000 No . . Grade in School 5th Grade 6th Grade -0.46289 0.270334 -1.71 0.0957 0.000000 0.000000 . . SEX 0.087569 0.181267 0.48 0.6320 0.000000 0.000000 . . Male Female RACE 0.5088010.3041711.670.10330.2778010.3794790.730.46900.5180410.2875781.800.08030.0000000.0000000.000000 Black Hispanic Other White 0.000000 0.000000 . . Family Situation
 Non-Traditional
 0.436618
 0.193741
 2.25
 0.0306

 Traditional
 0.000000
 0.000000
 .
 .
 AREA -0.06764 0.375932 -0.18 0.8582 Rural -0.26165 0.344689 Suburban -0.76 0.4529 0.000000 0.000000 . Urban . _____

Here we have the *estimated regression coefficients* computed under exchangeability and the standard errors as if the exchangeable working assumption were correct. The standard errors are roughly the same as with the robust variance estimator for these data, indicating that the exchangeable correlation assumption is close to truth.

GEE with Exchangeable "Working" Correlations Model-Based (Naive) Variance Estimator

```
Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure
                                                         Page : 2
Time: 13:30:35
                      The MULTILOG Procedure
                                                         Table : 1
Variance Estimation Method: Model-Based (Naive)
Working Correlations: Exchangeable
Link Function: Cumulative Logit
Response variable INTCIG12: Initiation of Cigarette Use
MULTILOG Logistic Regression Model for the DARE Evaluation Study
Model-Based Variance Estimation
_____
                    Degrees P-value
Contrast
                      of Wald Wald
                 Freedom ChiSq ChiSq
_____

        OVERALL MODEL
        10
        262.19
        0.0000

        MODEL MINUS INTERCEPT
        9
        22.36
        0.0078

                           1 5.73 0.0167
DARE
FIFTH
                           1 2.93 0.0868
                           1 0.23 0.6290
SEX
                           3 4.06 0.2546
RACE
                           1 5.08 0.0242
OTHFAM
                           2 0.64 0.7249
AREA
_____
```

Here we have the *main effects tests* computed under exchangeability, using the model-based variance approach. Results are essentially the same as with the robust variance estimator.

GEE with Exchangeable "Working" Correlations Model-Based (Naive) Variance Estimator

Date:04-28-97Research Triangle InstituteTime:13:30:35The MULTILOG Procedure Page : 3 Table : 1 Variance Estimation Method: Model-Based (Naive) Working Correlations: Exchangeable Link Function: Cumulative Logit Response variable INTCIG12: Initiation of Cigarette Use MULTILOG Logistic Regression Model for the DARE Evaluation Study Model-Based Variance Estimation -----Independent Variables Lower Upper and Effects Odds 95% 95% Ratio Limit Limit Ratio Limit Limit -----INTCIG12 (cum-logit) Intercept 1: Yes 0.15 0.07 0.32 DARE Program 0.56 0.34 0.92 Yes 1.00 1.00 1.00 No Grade in School 0.63 0.36 1.09 5th Grade 6th Grade 1.00 1.00 1.00 SEX 1.09 0.76 1.58 Male 1.00 1.00 1.00 Female RACE 1.660.903.081.320.612.851.680.943.011.001.001.00 Black Hispanic Other White Non-Traditional 1.55 1.04 2.29 Traditional 1.00 1 00 Family Situation AREA 0.93 0.44 2.01 0.77 0.38 1.55 Rural Suburban Urban 1.00 1.00 1.00

References

Applications of SUDAAN and Related Techniques

Bieler, G. and Williams, R. (1995). Cluster sampling techniques in quantal response teratology and developmental toxicity studies. *Biometrics* **51**, 764-776.

Davies, G.M. (1994). Applications of sample survey methodology to repeated measures data structures in dentistry. *Institute of Statistics Mimeo Series* No. 2128T. University of North Carolina: Chapel Hill.

Donner, A. (1982). An empirical study of cluster randomization. *American Journal of Epidemiology* **11**, 283-286.

Ennett, S.T., Rosenbaum, D.P., Flewelling, R.L., Bieler, G.S., Ringwalt, C.L., and Bailey, S.L. (1994). Long-term evaluation of drug abuse resistance education. *Addictive Behaviors*, **19**, 113-125.

Fung, K.Y., Krewski, D., and Scott, A.J. (1994). Tests for trend in developmental toxicity experiments with correlated binary data. *Risk Analysis* **14**, 639-648.

Gansky, SA, Koch, GG, and Wilson, J. (1994). Statistical evaluation of relationships between analgesic dose and ordered ratings of pain relief over an eight-hour period. *Journal of Biopharmaceutical Statistics* **4**, 233-265.

Graubard, B.I. and Korn, E.L. (1994). Regression analysis with clustered data. *Statistics in Medicine* **13**, 509-522.

LaVange, L.M. and Koch, G.G. (1994). Analysis of repeated measures studies with multiple regression methods for sample survey data. Presented and the 1994 Drug Information Association Meetings, and submitted for publication.

LaVange, L.M., Keyes, L.L., Koch, G.G., and Margolis, P.A. (1994). Application of sample survey methods for modelling ratios to incidence densities. *Statistics in Medicine* **13**, 343-355.

Norton, E., Bieler, G., Zarkin, G., and Ennett, S. (1996). Analysis of prevention program effectiveness with clustered data using generalized estimating equations. *Journal of Consulting and Clinical Psychology* **64**, 919-926.

Rao, J. and Colin, D. (1991). Fitting dose-response models and hypothesis testing in teratological studies. In: <u>Statistics in Toxicology</u>, Krewski and Franklin, eds. NY: Gordon and Breach.

Rao, J. and Scott, A.J. (1992). A simple method for the analysis of clustered binary data. *Biometrics* **48**, 577-585.

Schmid, JE, Koch, GG, and LaVange, LE (1991). An overview of statistical issues and methods

of meta-analysis. Journal of Biopharmaceutical Statistics, 1, 103-120.

Williams, R.L. (1995). Product-limit survival functions with correlated survival times. *Lifetime Data Analysis* 1, 171-186.

Williams, R.L. and Bieler, G.S. (1993). Estimation of proportional hazards models for survival times with nested errors. <u>ASA Proceedings of the Biopharmaceutical Section</u>, 115-120.

Survey Sampling

Binder, D. (1992). Fitting Cox's proportional hazards models from survey data. *Biometrika* **79**, 139-147.

Binder, D. (1983). On the variances of asymptotically normal estimators from complex surveys. *International Statistical Review* **51**, 279-292.

Cochran, W.G. (1977). Sampling Techniques. Wiley, New York.

Folsom, R.E. (1974). National Assessment Approach to Sampling Error Estimation, <u>Sampling</u> <u>Error Monograph</u>. Prepared for the National Assessment of Educational Progress, Denver, CO.

Hansen, M.H., Hurwitz, W.N., and Madow, W.G. (1953). <u>Sample Survey Methods and Theory</u>, <u>Volume I: Methods and Applications</u>. NY: Wiley.

Kendall, M.G. and Stuart, A. (1973). <u>The Advanced Theory of Statistics</u>. NY: Hafner Publishing Co.

Kish, L. (1965). Survey Sampling. NY: Wiley.

Kish, L. and Frankel, M. (1974). Inference from complex samples (with discussion). *Journal of the Royal Statistical Society*, Series B, **36**, 1-37.

Koch, G.G., Freeman, D.H. and Freeman, J.L. (1975). Strategies in the multivariate analysis of data from complex surveys. *International Statistical Review* **43**, 59-78.

LaVange, L.M., Iannacchione, V.G., and Garfinkel, S. (1986). An application of logistic regression methods to survey data: predicting high cost users of medical care. <u>American</u> <u>Statistical Association, Proceedings of the Section on Survey Research Methods</u>, Washington, DC.

Rao, J. And Scott, A. (1987). On simple adjustments to chi-squared tests with sample survey data. <u>Annals of Statistics</u>, **15**, 385-397.

Sarndal, C., Swensson, B., and Wretman, J. (1992). <u>Model-Assisted Survey Sampling</u>. Springer-Verlag, NY.

Scott, A.J. and Holt, D. (1982). The effect of 2-stage sampling on ordinary least squares

methods. Journal of the American Statistical Association 77, 848-854.

Shah, B.V., Barnwell, B.G., and Bieler, G.S. (1996). <u>SUDAAN User's Manual, Release 7</u>, First Edition. Research Triangle Institute, RTP, NC.

Shah, B.V. and LaVange, L.M. (1994). Mixed models for survey data. American Statistical Association, <u>Proceedings of the Section on Survey Research Methods</u>.

Shah, B.V., Holt, M.M. and Folsom, R.E. (1977). Inference about regression models from sample survey data. <u>Bulletin of the International Statistical Institute</u> XLVII, **3**, 43-57.

Thomas, D.R. and Rao, J.N.K. (1987). Small-sample comparisons of level and power for simple goodness-of-fit statistics under cluster sampling. *JASA* **82**, 630-636.

Wolter, K.M. (1985). Introduction to Variance Estimation. NY: Springer-Verlag.

Woodruff, R. (1971). A simple method for approximating the variance of a complicated estimate. *Journal of the American Statistical Association* **66**, 411-414.

GEEs and Generalized Linear Models

Carr, G. and Portier, C. (1993). An evaluation of some methods for fitting dose-response models to quantal response data. *Biometrics* **49**, 779-791.

Diggle, P., Liang, K.Y., and Zeger, S.L. (1994). <u>Analysis of Longitudinal Data</u>. NY: Oxford University Press.

Dunlop, D.D. (1994). Regression for longitudinal data: a bridge from least squares regression. *American Statistician* **48**, 299-303.

Huber, P.J. (1967). The behavior of maximum likelihood estimates under non-standard conditions. In <u>Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and</u> <u>Probability</u> **1**, 221-233.

Karim, M.R. and Zeger, S.L. (1989). GEE: A SAS macro for longitudinal data analysis. Technical Report #674 from the Department of Biostatistics, The Johns Hopkins University.

Liang, K. and Zeger, S. (1986). Longitudinal data analysis using generalized linear models. *Biometrika* 73, 13-22.

Lipsitz, S.R., Kim, K. and Zhao, L. (1994a). Analysis of repeated categorical data using generalized stimating equations. *Statistics in Medicine* **13**, 1149-1163.

Lipsitz, S.R., Fitzmaurice, G.M., Orav, E.J. and Laird, N.M. (1994b). Performance of generalized estimating equations in practical situations. *Biometrics* **50**, 270-278.

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

Neuhaus, J.M. and Segal, M.R. (1993). Design effects for binary regression models fitted to dependent data. *Statistics in Medicine* **12**, 1259-1268.

Prentice, R.L. (1988). Correlated binary regression with covariates specific to each binary observation. *Biometrics* **44**, 1033-1048.

Rotnitzky, A. and Jewell, N.P. (1990). Hypothesis testing of regression parameters in semiparametric generalized linear models for cluster correlated data. *Biometrika* **77**, 485-497.

Royall, R.M. (1986). Model robust confidence intervals using maximum likelihood estimators. *International Statistical Review* **54**, 221-226.

Zeger, S. (1988). Commentary. Statistics in Medicine 7, 161-168.

Zeger, S. and Liang, K. (1986). Longitudinal data analysis for discrete and continuous outcomes. *Biometrics* **42**, 121-130.

Jackknife Variance Estimators

Carr, G. And Portier, C. (1993). An evaluation of some methods for fitting dose-response models to quantal-response developmental toxicity data. *Biometrics* **49**, 779-791.

Lipsitz, S. And Parzen, M. (1996). A jackknife estimator of variance for Cox regression for correlated survival data. *Biometrics* **52**, 291-298.

Lipsitz, S., Dear, K., and Zhao, L. (1994). Jackknife estimators of variance for parameter estimates from estimating equations with applications to clustered survival data. *Biometrics* **50**, 842-846.

Quenouille, M. (1956). Notes on bias in estimation. Biometrika 43, 353-360.

Tukey, J.W. (1958). Bias and confidence in not quite large samples. *Annals of Mathematical Statistics*, **29**, 614.

BRR Variance Estimators

McCarthy, PJ (1966). Replication: An Approach to the Analysis of Data from Complex Surveys. *Vital and Health Statistics*, Series 2, No. 14, National Center for Health Statistics, Public Health Service, Washington, DC.

McCarthy, PJ (1969). Pseudoreplication: Half-samples. *Review of the International Statistical Institute* **37**, 239-264.

Plackett, RL and Burman, JP (1946). The Design of Optimum Multifactorial Experiments. *Biometrika* **33**, 305-325.

Wolter, KM (1985). Introduction to Variance Estimation. NY: Springer-Verlag.