

# FUZZY-SET SAMPLING UNDER SEMI-DEFINITE CHAOS

Jimmy Thomas Efirid  
1430 Massachusetts Ave., Suite 306-51, Cambridge, MA 02138

## Introduction

Estimating the population prevalence of a binary response variable when sample information is incomplete or inconsistent, poses immense problems to the applied statistician. Traditionally, the analyst simply deletes these values and treats the remaining data as a representative sample from the population in question. Providing that stable covariate fields exist and between-variable correlation is zero, then multivariate ANOVA and Chi-square techniques will yield a quick check of this strategy. However, if sampling irregularities plague covariate fields or confounding presents a problem, then the analyst must seek other means of verification. One possible solution, given the availability of population information, is to construct a test statistic using a population projection technique known as Bayesian Logistic Regression.

## Bayesian Logistic Regression

Bayesian Logistic Regression is an estimation procedure that aims to eliminate within covariate group variance by utilizing the interrelational structure between the sample response variable and known population covariates. A new response vector with elements corresponding to the expected value of the underlying covariate patterns form the basis of Bayesian Logistic Regression. The expected values result from a Bayesian expansion of the outcome variable jointly conditional on the set of model covariates.

Sampling begins by randomly selecting variables from a population data base. Aiming for a computationally ideal model having three covariates, the next step involves reducing data set dimensionality via Stepwise Regression. Expressing in mathematical terms, let  $T_n = \{u, v, x, y, z\}$  represent a  $n$ -observation sample, where  $u$  denotes a binary outcome variable with three model covariates  $v, x$  and  $y$ . The  $z$ -covariate represents a model-free population field having a relatively stable multilevel element space. On the other hand, the situation assumes instabilities exist among one or more of the model covariates  $v, x$  and  $y$ . Indexing as follows,

- $i = 0, 1$
- $j, k, h, m = 1, 2, 3, \dots (v_n, x_n, y_n, z_n)$
- $z_n \rightarrow$  fixed population field
- $s \rightarrow$  stable vector field
- $c \rightarrow$  fuzzy-set/chaotic vector field

the expected response for a covariate pattern missing two elements (e.g. only the value for  $v$  is present) may be written as:

$$\begin{aligned} \Lambda_j &= \hat{P}(\mu=i | v_c=j) \\ &= \prod_{m=1}^{n_{z_s}} \hat{P}(\mu=i | v_c=j, z_s=m) \cdot \frac{\hat{P}(v_c=j | z_s=m) \cdot \hat{P}(z_s=m)}{\sum_{m=1}^{n_{z_s}} \hat{P}(v_c=j | z_s=m) \cdot \hat{P}(z_s=m)} \\ &= \prod_{m=1}^{n_{z_s}} \hat{P}(v_c=j | \mu=i, z_s=m) \cdot \hat{P}(z_s=m | \mu=i) \cdot \hat{P}(\mu=i) \cdot \frac{\left[ \frac{\hat{P}(z_s=m | v_c=j) \cdot \hat{P}(v_c=j)}{\sum_{j=1}^{n_{v_s}} \hat{P}(z_s=m | v_c=j) \cdot \hat{P}(v_c=j)} \cdot \hat{P}(z_s=m) \right]}{\prod_{m=1}^{n_{z_s}} \left[ \frac{\hat{P}(z_s=m | v_c=j) \cdot \hat{P}(v_c=j)}{\sum_{j=1}^{n_{v_s}} \hat{P}(z_s=m | v_c=j) \cdot \hat{P}(v_c=j)} \right]} \end{aligned}$$

Analogously,

$$\begin{aligned} \Lambda_{k,h} &= \hat{P}(\mu=i | x_c=k, y_c=h) \\ &= \prod_{m=1}^{n_{z_s}} \hat{P}(z_s=m | x_c=k, y_c=h) \cdot \hat{P}(\mu=i | x_c=k, y_c=h, z_s=m) \\ &= \prod_{m=1}^{n_{z_s}} \frac{\hat{P}(x_c=k | y_c=h, z_s=m) \cdot \hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)}{\sum_{h=1}^{n_{y_s}} \hat{P}(x_c=k | y_c=h) \cdot \hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)} \cdot \hat{P}(\mu=i | x_c=k, y_c=h, z_s=m) \\ &= \prod_{m=1}^{n_{z_s}} \frac{\hat{P}(x_c=k | y_c=h, z_s=m) \cdot \hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)}{\sum_{h=1}^{n_{y_s}} \hat{P}(x_c=k | y_c=h) \cdot \hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)} \cdot \hat{P}(\mu=i | x_c=k, y_c=h, z_s=m) \end{aligned}$$

$$= \sum_{m=1}^{n_{z_s}} \gamma_1 \cdot \gamma_2 \cdot \left[ \sum_{m=1}^{n_{z_s}} \gamma_1 \right]^{-1}$$

where

$$\begin{aligned} \gamma_1 &= \gamma_3 \cdot \gamma_4 \cdot \left[ \sum_{k=1}^{n_{x_s}} \gamma_3 \right]^{-1} \\ \gamma_2 &= \frac{\hat{P}(x_c=k | \mu=i, y_c=h, z_s=m) \cdot \hat{P}(y_c=h | \mu=i, z_s=m) \cdot \hat{P}(z_s=m | \mu=i) \cdot \hat{P}(\mu=i)}{\sum_{i=1}^{n_{\mu}} \hat{P}(x_c=k | \mu=i, y_c=h, z_s=m) \cdot \hat{P}(y_c=h | \mu=i, z_s=m) \cdot \hat{P}(z_s=m | \mu=i) \cdot \hat{P}(\mu=i)} \\ \gamma_3 &= \frac{\hat{P}(x_c=k | y_c=h, z_s=m) \cdot \hat{P}(y_c=h | z_s=m) \cdot \hat{P}(z_s=m)}{\sum_{k=1}^{n_{x_s}} \hat{P}(x_c=k | y_c=h, z_s=m) \cdot \hat{P}(y_c=h | z_s=m) \cdot \hat{P}(z_s=m)} \cdot \hat{P}(z_s=m | x_c=k) \cdot \hat{P}(x_c=k) \\ \gamma_4 &= \frac{\hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)}{\sum_{h=1}^{n_{y_s}} \hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)} \end{aligned}$$

and

$$\begin{aligned} \Lambda_{j,k,h} &= \hat{P}(\mu=i | v_c=j, x_c=k, y_c=h) \\ &= \prod_{m=1}^{n_{z_s}} \hat{P}(\mu=i | v_c=j, x_c=k, y_c=h, z_s=m) \cdot \frac{\hat{P}(v_c=j | x_c=k, y_c=h, z_s=m) \cdot \hat{P}(z_s=m | x_c=k, y_c=h)}{\prod_{m=1}^{n_{z_s}} \hat{P}(v_c=j | x_c=k, y_c=h, z_s=m) \cdot \hat{P}(z_s=m | x_c=k, y_c=h)} \end{aligned}$$

Simplifying the above expression,  $\Lambda_{j,k,h}$  becomes  $\prod_{m=1}^{n_{z_s}} \xi_1 \cdot \xi_2 \cdot \xi_3 \cdot \left[ \sum_{m=1}^{n_{z_s}} \xi_1 \right]^{-1}$  where  $\xi_1$  is given by

$$\begin{aligned} \xi_1 &= \frac{\hat{P}(v_c=j | \mu=i, x_c=k, y_c=h, z_s=m) \cdot \hat{P}(y_c=h | \mu=i, x_c=k, z_s=m) \cdot \hat{P}(z_s=m | \mu=i, x_c=k) \cdot \hat{P}(x_c=k | \mu=i) \cdot \hat{P}(\mu=i)}{\sum_{i=1}^{n_{\mu}} \hat{P}(v_c=j | \mu=i, x_c=k, y_c=h, z_s=m) \cdot \hat{P}(y_c=h | \mu=i, x_c=k, z_s=m) \cdot \hat{P}(z_s=m | \mu=i, x_c=k) \cdot \hat{P}(x_c=k | \mu=i) \cdot \hat{P}(\mu=i)} \\ \xi_2 &= \frac{\hat{P}(x_c=k | v_c=j, y_c=h, z_s=m) \cdot \hat{P}(y_c=h | v_c=j, z_s=m) \cdot \hat{P}(z_s=m | v_c=j) \cdot \hat{P}(v_c=j)}{\sum_{j=1}^{n_{v_s}} \hat{P}(x_c=k | v_c=j, y_c=h, z_s=m) \cdot \hat{P}(y_c=h | v_c=j, z_s=m) \cdot \hat{P}(z_s=m | v_c=j) \cdot \hat{P}(v_c=j)} \\ \xi_3 &= \xi_4 \cdot \left[ \sum_{k=1}^{n_{x_s}} \xi_4 \right]^{-1} \cdot \xi_5 \\ \xi_4 &= \frac{\hat{P}(x_c=k | y_c=h, z_s=m) \cdot \hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)}{\sum_{h=1}^{n_{y_s}} \hat{P}(x_c=k | y_c=h, z_s=m) \cdot \hat{P}(z_s=m | y_c=h) \cdot \hat{P}(y_c=h)} \cdot \hat{P}(z_s=m | x_c=k) \cdot \hat{P}(x_c=k) \\ \xi_5 &= \hat{P}(y_c=h | z_s=m) \cdot \hat{P}(z_s=m) \end{aligned}$$

Replacing  $(x_c=k, y_c=h)$  with  $(v_c=j, x_c=k)$  or  $(v_c=j, y_c=h)$  gives the equivalent forms  $\Lambda_{y,k}$  and  $\Lambda_{y,h}$ , while  $\Lambda_j$  is similarly expressed as  $\Lambda_k$  and  $\Lambda_h$  upon substituting  $(x_c=k)$  or  $(y_c=h)$  for  $(v_c=j)$ .

## A Homogeneity Test for Missing Response Elements

The statistical validity of discarding observations having missing or inconsistent response outcomes holds only if the true distribution of these items is indistinguishable from the set on non-missing elements. Unfortunately, routine testing of this assumption is often impossible due to covariate non-response and confounding.

An alternative approach involves partitioning the new response vector into two sets, depending upon whether or not the original response is missing. A test of homogeneity of the two groups, focusing on extreme deviation rather than on central tendency as a test criteria, is made possible by collapsing distinct covariate responses into single value blocks and subjecting these points to analysis using the parameter-free discrete range test. Designating non-missing element blocks as type-A and missing element blocks as type-B, combine values and rank accordingly, noting the relative position of each element block. Assuming element blocks to be numbered from 1 to  $N$ , define the range as the number on the highest type-B element block

minus the lowest type-B element. An unusually high or low range value raises suspicion that the two groups differ with respect to their tail regions. For illustrative purposes, consider the simple case where the number of type-A blocks ( $n_A$ ) equals type-B blocks ( $n_B$ ), e.g.  $N=n_A+n_B=2n$ . Letting  $\zeta$  designate the range of type-B blocks among type-A blocks and defining

$$\Omega_i = \left\{ 1 - \left[ 1 - \sum_{z=0}^{i-1} \frac{\xi \cdot \prod_{z=0}^{i-1} (n-z) \cdot \prod_{\beta=\min(\xi/2)+1}^{\xi-2} (n-\beta)}{2^{\min(\xi/2)+1} \cdot \prod_{\delta=1}^{\xi-2} (n-\delta-1) \cdot \prod_{v=0}^{\min(\frac{\xi+1}{2})-1} [2(n-v)-1]} \right] \cdot \left[ \frac{i \cdot \prod_{z=0}^{i-1} (n-z) \cdot \prod_{\beta=\min(i/2)+1}^{i-2} (n-\beta)}{2^{\min(i/2)+1} \cdot \prod_{\delta=1}^{i-2} (n-\delta-1) \cdot \prod_{v=0}^{\min(\frac{i+1}{2})-1} [2(n-v)-1]} \right] \right\}^{-1}$$

it follows that

$$\lim_{N \rightarrow \infty} P\{\zeta \leq N-i-1\} = \lim_{N \rightarrow \infty} P\{\zeta \leq N-i\} / \lim_{N \rightarrow \infty} \Omega_i.$$

Solving for  $i$  gives the probabilities and rejection regions listed in Tables 1 and 2. Proceeding in an analogous fashion will yield test values in the unequal group case. In Table 2, if  $\zeta$  exceed  $\zeta_0$  for the indicated sample size then reject  $H_0$  with Type-I error  $\leq \alpha$  (.05, .01).

**Simulation**

Preliminary simulation suggests that the above procedure performs well when both the number of elements within each distinct block and the overall number of blocks in each group remains high. This situation generally holds in large-scale sampling situation where the base population is greater than 100,000 with an initial sample greater than 10%.

**Table 1: Exact Asymptotic Probabilities for the Range Distribution**

$i$	$\lim_{N \rightarrow \infty} P\{\zeta \leq N-i\}$	$\lim_{N \rightarrow \infty} P\{\zeta = N-i\}$
1	1.00000000	0.25000000
2	0.75000000	0.25000000
3	0.50000000	0.18750000
4	0.31250000	0.12500000
5	0.18750000	0.07812500
6	0.10937500	0.04687500
7	0.06250000	0.02734380
8	0.03515630	0.01562500
9	0.01953130	0.00878906
10	0.01074220	0.00488281
11	0.00585938	*****

**Table 2: Rejection Regions for the Two-Sample ( $n_1=n_2$ ) Range Test**

Sup $\tau_0$ : $P\{\zeta \leq \tau_0\} \leq .05$	Sup $\tau_0$ : $P\{\zeta \leq \tau_0\} \leq .01$
$n \leq 4, \tau_0 = \{\emptyset\}$	$n \leq 5, \tau_0 = \{\emptyset\}$
$5 \leq n \leq 6, \tau_0 = N - 6$	$n = 6, \tau_0 = N - 7$
$7 \leq n \leq 26, \tau_0 = N - 7$	$7 \leq n \leq 10, \tau_0 = N - 8$
$n \geq 27, \tau_0 = N - 8$	$11 \leq n \leq 20, \tau_0 = N - 9$
	$21 \leq n \leq 204, \tau_0 = N - 10$
	$n \geq 205, \tau_0 = N - 11$

**REFERENCES**

Chernoff, H. and Savage, I.R. (1958). Asymptotic normality and efficiency of certain nonparametric test statistics. *Annals of Mathematical Statistics* 29: 972-994.

Csorgo, M. and Guttman, I. (1962). On the empty cell test. *Technometrics* 4: 235-247.

David, F.N. (1950). Two combinatorial tests of whether a sample has come from a given population. *Biometrika* 30: 7-111.

David, F.N. and Barton, D.E. (1962). *Combinatorial Chance*. Hafner Publishing Company.

Efrid, J.T. (1987). A Wald-Wolfowitz type test sensitive to outliers based on F.N. David's 'Ball and Cell' range statistic. *Joint Statistical Meeting, San Francisco, CA*.

Halperin, M. and Burrow, G.L. (1961). An asymptotic distribution for an occupancy problem with statistical applications. *Technometrics* 3: 79-86.

Jogdeo, K. (1968). Asymptotic normality in nonparametric methods. *Annals of Mathematical Statistics* 39: 905-922.

Le Cam, L. (1986). *Asymptotic Methods in Statistical Decision Theory*. Springer-Verlag.

Mood, A.M. (1954). On the asymptotic efficiency of certain nonparametric two-sample tests. *Annals of Mathematical Statistics* 25: 514-522.

Mood, A., Graybill, F.A. and Boes, D.C. (1974). *Introduction to the Theory of Statistics*. McGraw-Hill, Inc.

Nicholson, W.L. (1961). Occupancy probability distribution critical points. *Biometrika* 58: 175-180.

Wald, A. and Wolfowitz, J. (1940). On a test whether two samples are from the same population. *Annals of Mathematical Statistics* 2: 147-162.