

- Bayesian Methods for Multinomial Sampling with Missing Data Using Multiple Hypergeometric Functions<sup>1</sup>

*∞*`}

· · · ·

7 F

0

いい

9 a a

# рÀ

James M. Dickey, Jyh-Ming Jiang,

### and

### Joseph B. Kadane

#### U.S. Department of Justice National Institute of Justice

This document has been reproduced exactly as received from the person or organization originating it. Points of view or opinions stated in this document are those of the authors and do not neces. arily represent the official position or policies of the National Institute of Justice.

Permission to reproduce this copyrighted material has been granted by Public Domain/NIJ

U. S. Department of Justice

to the National Criminal Justice Reference Service (NCJRS).

Further reproduction outside of the NCJRS system requires permission of the convergent owner.

Research Report 6/83 - #15

State University of New York at Albany

Department of Mathematics and Statistics

5/30/83

<sup>1</sup>Work by Professor Dickey and Mr. Jiang sponsored by National Science Foundation, Grant MCS-8301335. Work by Professor Kadane sponsored by National Institutes of Justice, Grant 81-IJ-CX-0087, and the Office of Naval Research Contract, N00014-82-K-0622.



## National Criminal Justice Reference Service



This microfiche was produced from documents received for inclusion in the NCJRS data base. Since NCJRS cannot exercise control over the physical condition of the documents submitted, the individual frame quality will vary. The resolution chart on this frame may be used to evaluate the document quality.



MICROCOPY RESOLUTION TEST CHART NATIONAL BUREAU OF STANDARDS-1963-A



Microfilming procedures used to create this fiche comply with the standards set forth in 41CFR 101-11.504.

Points of view or opinions stated in this document are those of the author(s) and do not represent the official position or policies of the U. S. Department of Justice.

National Institute of Justice United States Department of Justice Washington, D. C. 20531





Bayesian Methods for Multinomial Sampling with Missing Data Using Multiple Hypergeometric Functions<sup>1</sup>

ж В

\*

1

**S S S** 

30.0

Ť

James M. Dickey, Jyh-Ming Jiang,

by

and

Joseph B. Kadane

U.S. Department of Justice National Institute of Justice

This document has been reproduced exactly as received from the person or organization originating it. Points of view or opinions stated in this document are those of the authors and do not necessarily represent the official position or policies of the National Institute of Justice.

Permission to reproduce this copyrighted material has been granted by

Public Domain/NIJ U. S. Department of Justice

to the National Criminal Justice Reference Service (NCJRS).

Further reproduction outside of the NCJRS system requires permission of the convright owner.

Research Report 6/83 - #15 State University of New York at Albany Department of Mathematics and Statistics

5/30/83

<sup>1</sup>Work by Professor Dickey and Mr. Jiang sponsored by National Science Foundation, Grant MCS-8301335. Work by Professor Kadane sponsored by National Institutes of Justice, Grant 81-IJ-CX-0087, and the Office of Naval Research Contract, N00014-82-K-0622.



### ABSTRACT

\$ ;

-. R.

0.2

ě

Distribution theory is given for Bayesian inference from multinomial (or multiple Bernoulli) sampling with missing category distinctions, such as a contingency table with supplemental purely marginal counts. A new conjugate family generalizes the usual Dirichlet prior distributions. The posterior moments and predictive probabilities are found to be proportional to ratics of Carlson's hypergeometric functions of matrix argument. Dimensionreducing integral identities and expansions are given for statistical use. Closed-form expressions are developed for cases of nested missing distinctions. Examples are given, together with a simple method for assessment of a Dirichlet subjective prior distribution. In his t popularized B sampling, in family of pri elegant. Den closed forms. and mixtures enough in man subjective un prior prejudi probabilities ing a sampled However, distinctions

Q

However, in the case of inference from data with some missing distinctions between categories this tractability seems to fade. The likelihood then contains factors which are powers of the sums of probabilities of the confused categories. The more data values that miss a particular distinction, the greater the power of the corresponding sum; and the greater the variety of such missed distinctions, the more such sums there are in the likelihood. In the least tractable cases, the sample contains moderate amounts of data of assorted kinds: complete data and data with various missing distinctions, without nesting or other logical constraints between the missed distinctions.

### 1. INTROUUCTION

In his two books (1950, 1965) and various articles, I.J. Good popularized Bayesian conjugate-prior inference for multinomial sampling, in which the Dirichlet distributions form the conjugate family of prior distributions. This mathematics is tractable and elegant. Densities, moments, and predictive probabilities take

closed forms. Furthermore, the family of Dirichlet distributions and mixtures of Dirichlets with a small number of terms is large enough in many situations to offer realistic models for predata subjective uncertainty. (An exception occurs when there is a prior prejudice for local smoothness between neighboring category probabilities, such as when the data frequencies result from grouping a sampled continuous variable (Dickey 1968b).)

1.1

Bayesian treatments of multinomial sampling with missing data were given a decade ago by Karson and Wrobleski (1970), Antelman (1972), and Kaufman and King (1973). These concerned two-way contingency tables containing data with missing information regarding row or column variables. The studies were restricted to the 2x2 case. Basu and Pereira (1982) extended consideration to Kx2 tables and summarized properties of the Dirichlet distribution under relevant changes of variable. For an account of frequentist results, see Chen and Fienberg (1974, 1976), Dempster et al (1977), and references cited therein (for example, Hartley 1958).

After formulating here the general problem of Bayesian multinomial inference for missing category distinctions, in Section 2, we shall introduce in Sections 3 and 4 a representation of posterior integrals from Dirichlet prior distributions as multiple hypergeometric functions, specifically, the elegant functions R and R of Carlson (1977). The posterior distributions are found to be identical to distributions developed in Dickey (1983). This is then extended in Section 4 to a new conjugate prior family for such sampling. A parallel development based on changes of variable in the Dirichlet distribution yields closed forms for the case of nested missing distinctions, Section 5. This leads to new simplications of R. Straightforward expansions of posterior distributions as probability mixtures of Dirichlet distributions are generalized in Section 6, and related expansions of R are given. Examples are briefly given in Section 7. An appendix outlines a simple method for assessment of a Dirichlet subjective prior distribution.

1.2

of  $\binom{M}{x}$  multiplied by (2.1).)

0

### 2. SAMPLING WITH MIUSING DISTINCTIONS

To begin, we shall need to establish notation for the nonmissing (uncensored) case. Denote by  $\underline{x} = (x_1, \dots, x_K)$ ' the vector of frequency counts for sampling from a distribution on a finite sample space having the unknown probability vector  $\underline{u} = (u_1, \dots, u_K)'$ . That is,  $x_i$  denotes the number of sample values falling in the i<sup>th</sup> category of probability  $u_i$ , i = 1,...,K. For Bayesian inference concerning <u>u</u> under <u>any noninformative-</u> stopping process, the vector  $\underline{x}$  can be treated as if it were the frequency count from a multiple-Bernoulli sequence  $\underline{\delta}_{x} = (\delta_{1}, \dots, \delta_{M})$ of prespecified length M. Then x = M, where x denotes  $x_1^{+} \dots + x_K^{-}$ . The corresponding likelihood function is proportional to the probability mass of such  $\underline{\delta}_{\mathbf{y}}$ ,

 $pr(\underline{\delta}_{x} | \underline{u}) = \prod_{i=1}^{K} u_{i}^{x_{i}}.$ 

(2.1)

Note that this function of  $\underline{u}$  is parameterized by the frequencies  $\underline{x}$ , a sufficient statistic. (The term "multiple-Bernoulli" refers to a fixed-length sampling distribution for  $\underline{\delta}_x$ , for which <u>x</u> has the "multinomial" distribution with mass equal to the product

The Bayesian predictive distribution is the subjective probability of an outcome sequence  $\underline{\delta}_x$  without conditioning on the unknown sampling parameters  $\underline{u}$ . That is,  $\underline{u}$  is averaged out of the sampling probability (2.1) according to a distribution of u. Whatever the prior distribution, the prior predictive probability takes the form of a prior moment,

$$pr(\underline{\delta}_{\mathbf{X}}) = E_{\mathbf{u}} pr(\underline{\delta}_{\mathbf{X}} | \underline{\mathbf{u}})$$
(2.2)

 $= E_{u} \Pi_{1}^{K} u_{1}^{x}$ 

This gives the subjective probability for the outcome sequence  $\delta_{\downarrow}$ conditional only on the total count x . A similar statement holds for posterior predictive probabilities of additional future outcomes in terms of posterior moments. Hence, to provide Bayesian predictive probabilities, we need merely develop the prior and posterior moments. (A similar statement will hold in the case of censored sampling.)

The conjugate prior family for uncensored sampling is the Dirichlet. The random vector u is said to have the Dirichlet distribution  $D(\underline{b})$  with parameter vector  $\underline{b} = (b_1, \dots, b_K)'$ , each  $b_1 > 0$ , if u has the density in any K-1 of its coordinates,

$$f(\underline{u};\underline{b}) \equiv B(\underline{b})^{-1} \Pi_{1}^{K} u_{1}^{b} i^{-1}, \qquad (2.3)$$
$$B(\underline{b}) = [\Pi_{1}^{K} \Gamma(b_{1})] / \Gamma(b_{.}),$$

for all  $\underline{u}$  in the probability simplex { $\underline{u}$ : each  $\underline{u} \ge 0$ ,  $\underline{u} = 1$ }.

·. 2.2

The corresponding prior moment is

g(<u>c;</u>

with updated parameters,

Consider the prior distribution,

 $\underline{u} \sim D(\underline{b})$ .

(2.4)

(2.5)

2.3

$$\underline{b} \equiv E_{u|b} \pi_{1}^{K} u_{1}^{c_{1}}$$
$$= B(\underline{b}+\underline{c})/B(\underline{b}).$$

The predictive distribution is then the Dirichlet-Bernoulli with mass  $pr(\underline{\delta}_x) = g(\underline{x};\underline{b})$ . Note that  $r(\underline{x})$  is the probability mass of the outcome sequence  $\frac{\delta}{x}$ , rather than the frequency count  $\underline{x}$ .

The posterior distribution corresponding to the conjugate prior (2.4) and likelihood function (2.1) is again Dirichlet,

 $\underline{u} | \underline{x} \sim D(b+x).$ 

(2.6)

The posterior Dirichlet density is then  $f(\underline{u};\underline{b}+\underline{x})$  and the posterior moment is  $g(\underline{c};\underline{b}+\underline{x})$ . Consider now the generalization to sampling from the distribution having probabilities  $\underline{u}$  when some observations are censored, that is, do not report unique categories, but rather mere sets of categories. For example, an observation may be reported as falling either in category i or in category j, for a particular pair i<j. Denote the frequency count of such observations by y<sub>i,j</sub>. Let  $y_{i,jk}$  be the frequency count of observations confusing a triple of categories, i<j<k. An analogous notation applies for any proper subset  $\sigma$  of at least two categories:  $y_{\sigma}$  for  $\sigma \in \{1, \dots, K\}$ . Denote a sequence of N such censored outcomes by  $\underline{\epsilon}_{v} = (\epsilon_{1}, \dots, \epsilon_{N})$  and

denote by y the corresponding vector counting the confusions of each possible type,

$$\underline{\mathbf{y}} = (\mathbf{y}_{12}, \dots, \mathbf{y}_{123}, \dots, \mathbf{y}_{23}, \dots, \mathbf{y}_{23})$$
(2.7)

2.4

The dimensionality of <u>y</u> is  $2^{K}-K-2$ , although in practice, <u>y</u> will be sparse with most coordinates zero.

Consider sampling under noninformative stopping and noninformative censoring (Dawid and Dickey 1977). The likelihood function is proportional to the likelihood function from an outcome sequence  $(\underline{\delta}_x, \underline{\epsilon}_y)$ ,  $\underline{\epsilon}_y = (\epsilon_1, \dots, \epsilon_N)$ , where  $\underline{\delta}_x$  and  $\underline{\epsilon}_y$  are independent and the coordinates  $\boldsymbol{\varepsilon}_i$  independently arise from prespecified partitions of the sample space. Then

$$\operatorname{pr}(\underline{\delta}_{x}, \underline{\varepsilon}_{y} | \underline{u}) = \operatorname{pr}(\underline{\delta}_{x} | \underline{u}) \operatorname{pr}(\underline{\varepsilon}_{y} | \underline{u}), \qquad (2.8)$$

where

$$pr(\underline{\varepsilon}_{y}|\underline{u}) = \Pi_{k=2}^{K-1} \Pi_{1} < \dots < i_{k} (u_{1} + \dots + u_{k})^{y_{1}} \cdots i_{k}$$
$$= \Pi_{\sigma} (\Sigma_{i\varepsilon\sigma} u_{1})^{y_{\sigma}}. \qquad (2.9)$$

Now, consider the effect of this likelihood function for inference from the Dirichlet prior density (2.3). The independence of the two factors in the likelihood (2.8) implies the following result.

Theorem 2.1. A data sequence  $(\underline{\delta}_x, \underline{\epsilon}_v)$  of prespecified length and censoring pattern has the predictive probability

$$pr(\underline{\delta}_{x}, \underline{\epsilon}_{y}) = E_{u|b}[pr(\underline{\delta}_{x}|\underline{u})pr(\underline{\epsilon}_{y}|\underline{u})]$$
$$= g(\underline{x}; \underline{b})h(\underline{y}; \underline{b} + \underline{x}), \qquad (2.10)$$

where the first factor g is given by (2.5) and the second factor is the expectation of (2.9) under the Dirichlet partial-posterior distribution,  $u \mid x \sim D(b+x)$ ,

## h(y;b+)

The posterior density of the unknown probability vector u has the expression in terms of this same quantity h;

## $f(\underline{u};\underline{b}+\underline{x})$

(The function f of three arguments generalizes the previous notation f of two arguments (2.3); the motivation for defining a negative third argument will follow later.) The corresponding posterior moment is proportional to a ratio of such quantities,

# <sup>E</sup>u|b+x,

where the proportionality factor  $g(\underline{c};\underline{b}+\underline{x})$  is the usual Dirichlet posterior moment (2.5), based on the uncensored part  $\underline{x}$  of the

data. 🛛

Attention focuses naturally, then, on the properties and cal-

culation of the quantities h (2.11).

$$\underline{x}) = E_{u|b+x} \operatorname{pr}(\underline{\varepsilon}_{y}|\underline{u}).$$
(2.11)

$$\underline{c};-\underline{y}) = f(\underline{u};\underline{b}+\underline{x})pr(\underline{\varepsilon}_{y}|\underline{u})/h(\underline{y};\underline{b}+\underline{x}).$$
(2.12)

$$y^{\left[\left(\Pi_{1}^{K}u_{1}^{c_{1}}\right)\Pi_{\sigma}\left(\Sigma_{i\varepsilon\sigma}u_{1}\right)^{e_{\sigma}}\right]}$$

(2.13)

= g(c;b+x)h(y+e;b+x+c)/h(y;b+x)

## 3. MULTIPLE HYPERGEOMETRIC FUNCTIONS

0

B. C. Carlson (1963,1971,1974,1977) developed the following class of multiple hypergeometric functions for an organizing and unifying role in the field of special functions. See also Dickey (1983) for probabilistic interpretations and statistical uses. We define the function R as a moment of a homogeneous linear form in the random vector  $u \sim D(b)$ ,

$$R_{a}(\underline{b},\underline{z}) = E_{u|b}^{(k-1)}(\underline{u}'\underline{z})^{a} , \qquad (3.1)$$

where  $\underline{z} = (z_1, \dots, z_K)'$  and  $\underline{u}'\underline{z} = u_1z_1 + \dots + u_Kz_K$ . (In this section we shall indicate the dimensionality of integrals in the notation we use for expectation.). This definition requires  $b_{i} \ge 0$  for all i, but more general definitions, including a function that generates R, appear in Carlson (1977). Dickey (1983) exhibits R as being, itself, the probability generating function of the Dirichlet-multinomial distribution, the conjugate Bayesian predictive distribution for multinomial sampling.

The following classical identity, attributed to Picard by Appell and Kampé de Fériet (1926), reduces the dimensionality of the integral representation, thus permitting computation of R by simple quadrature. Although this refers to a restricted range of parameters, a contour integral applies more generally (Carlson 1977, Theorem 6.8-2, p.155).

<u>Theorem 3.1</u> (Picard's identity). For  $-b_{a<0}$ ,  $w = (w_1, w_2)'$ , and d = (-a, b, +a)',

limiting form,

3.1

R\_-b

 $R_{a}(\underline{b})$ 

Z =

where

A simpler integral representation and a dimension-reducing integral identity are available for  $R_a$ , as follows, under the

$$R_{a}(\underline{b},\underline{z}) = E_{w|d}^{(1)} II_{i=1}^{K} (w_{1}z_{i}+w_{2})^{-b}i$$
  
=  $B(\underline{d})^{-1} \int_{0}^{1} v^{-a-1} (1-v)^{b} \cdot t^{a-1} [\Pi_{1}^{K} (vz_{i}+1-v)^{-b}i] dv.$ 

Under the parameter restriction  $a = -b_{.}$ , R takes the simple

$$(\underline{b},\underline{z}) = \Pi z_{\underline{i}}^{-\underline{b}_{\underline{i}}} .$$
(3.3)

A two-way multiple hypergeometric function generalizing R, a function of matrix argument, can be defined by considering a bilinear form  $\underline{u}' \underline{Zv}$  in independent random vectors,  $\underline{u} \sim D(\underline{b})$ (K coordinates) and  $\underline{v}\, \lor\, D(\underline{c})$  (L coordinates) for nonrandom matrix Z(K×L). Define the function R as a moment of the bilinear form,

$$F_{u|b}^{(K-1)} = E_{u|b}^{(K-1)} E_{v|c}^{(L-1)} (\underline{u}' Z \underline{v})^{a}$$

$$= E_{u|b}^{(K-1)} R_{a} (\underline{c}; \underline{u}' \underline{z}_{*1}, \dots, \underline{u}' \underline{z}_{*L}),$$
(3.4)

$$\begin{pmatrix} \frac{z}{1} \\ \frac{z}{2} \\ \frac{z}{2} \\ \frac{z}{K} \end{pmatrix} = (\underline{z}_{1}, \dots, \underline{z}_{K}).$$

3.2

(3.5)

3.3

parameter restriction  $a = -c_{1}$ .

Theorem 3.2.

$$R_{-c}(\underline{b}, \underline{z}, \underline{c}) = E_{u|b}^{(K-1)} \Pi_{1}^{L} (\underline{u}' \underline{z}_{*j})^{-c} j$$
  
=  $E_{w|d}^{(L)} \Pi_{1}^{K} (\underline{z}'_{*} \underline{v} + w_{L+1})^{-b} j$ , (3.6)

where

The second equality in (3.6), first given by Dickey (1968a), generalizes Picard's identity (L = 1).  $\Box$ 

Note that if b = c in (3.6),  $d_{L+1} = 0$  and the L-fold integral becomes an (L-1)-fold integral with  $w_{L+1} \equiv 0$ . More generally, as the following Corollary shows, the dimensionality of the integral can be reduced when the parameter vector c has zero coordinate values. This property is useful in missing data problems in which the high-dimensional vector y is sparse.

Corollary 3.3. If, without loss of generality, the vector <u>c</u> is taken in the form  $\underline{c} = (\underline{c}^{(1)'}, \underline{0}')'$  where the subvector  $\underline{c}^{(1)}$ has  $L^{(1)} \leq L$  coordinates, then

$$R_{-c} (\underline{b}, Z, \underline{c}) = R_{-c} (\underline{l}) (\underline{b}, Z^{(1)}, \underline{c}^{(1)})$$
(3.8)

where  $Z^{(1)}$  consists of the corresponding  $L^{(1)}$  columns of Z, Z =  $(Z^{(1)}, Z^{(2)})$ . A similar dimension-reduction occurs whenever there are some columns of Z which are proportional to each other, or proportional to the vector of unit entries l = (1, ..., l)'.

1

In principle, the integral identity (3.6) can be useful for computation of R by quadrature when L<<K (or  $L^{(1)}$ <<K). However, we have found this awkward in practice, because for interesting values of the parameters, the integrands tend to have poles in the range of integration. At present, series expansions of R appear more practical, as given later in the paper.

	•
4.1	
4. INFERENCE AND A NEW CONJUGATE FAMILY	and the correspond
	[B(b+x+c
Under the parameter restriction $a = -c$ , Carlson's two-way	
multiple hypergeometric function $R$ has the same form (3.6) as the	
mement quantities h (2.11) needed in Bayesian Dirichlet-prior	Note that our
inference for censored multinomial sampling. To express this	family as the prio
identification, specialize the matrix Z to the full-subsets in-	for the missing-da
dicator Z = $\tilde{Z}$ , consisting of zeros and ones whose L columns $\tilde{Z}_{*,i}$	development.
indicate all the proper subsets of two or more elements,	4.1 <u>A General Fam</u>
$\sigma_j \in \{1, \ldots, K\},$	
$\tilde{z} = \int l  if  i\varepsilon\sigma_j $ (4.3)	Dickey (1983)
ij (0 otherwise,	k t)
j=l,L, where	Λ, μ,
	$\underline{u} \sim \mathcal{D}(\underline{b}, \underline{b})$
L = 2 - K - 2. (4.2)	having density on
Theorem 4.1.	
$\mathbf{b}(\mathbf{d} \cdot \mathbf{b}) = \mathbf{F} \cdot \mathbf{u}^{\mathbf{L}} (\mathbf{\Sigma} \cdot \mathbf{u})^{\mathbf{J}} $ (4.3)	B( <u>b</u> ) <sup>-1</sup> (П.
$u b''j=1$ $u b''j=1$ $i\varepsilon\sigma_j$ $u'$	
$= R_{d} (\underline{b}, \overline{2}, -\underline{d}).  \Box$	
	(For the powers <u>d</u> :
This identification ties the problems of Bayesian statistical	in (4.6) is defined
inference for missing data to a mainstream segment of the theory	confusion in the lo
of special functions. The Dirichlet prior distribution $\underline{u} \sim D(\underline{b})$	The moments of
implies the predictive probability (2.10) with (4.3),	of R functions,

 $pr(\underline{\delta}_{x}, \underline{\varepsilon}_{y}) = [B(\underline{b}+\underline{x})/B(\underline{b})] \cdot R_{y}(\underline{b}+\underline{x}, \tilde{z}, -\underline{y}),$ (4.4)

Ģ

ling posterior moment (2.13) becomes

$$\underline{x+c}$$
 /B(b+x)]

(4.5)

4.2

PA.

$$R_{y,+e} \cdot (\underline{b} + \underline{x} + \underline{c}, \overline{z}, -\underline{y} - \underline{e}) / R_{y} \cdot (\underline{b} + \underline{x}, \overline{z}, -\underline{y}).$$

posterior distribution (2.12) is not in the same or, and thus the Dirichlet family is not conjugate ata likelihood. This motivates the following

## nily.

generalized the Dirichlet distributions to bitrary matrix parameter Z(K×L for arbitrary

$$\mathcal{D}(\underline{b}, \underline{Z}, -\underline{d}),$$
 (4.6)

the simplex of K probabilities,

$$\frac{1}{\left(\prod_{1}^{K} u_{1}^{b_{1}}\right)\left[\prod_{1}^{L} (\underline{u}' \underline{z})^{j_{1}}\right]}{\left(\frac{R_{d}}{2}, -\underline{d}\right)}$$

$$(4.7)$$

in (4.7), the sign of the third parameter  $-\underline{d}$ ed to match the usage of Carlson's R; this avoids ong run.)

f D are proportional, as in (4.5), to ratios

$$E_{u|b,Z,-d}^{(K-1)}[(\Pi_{1}^{K}u_{i}^{c_{i}})\Pi_{1}^{L}(\underline{u}'\underline{z}_{*j})^{e_{j}}] \qquad (4.8)$$

$$= [B(\underline{b}+\underline{c})/B(\underline{b})]$$

$$\cdot R_{d,+e} (\underline{b}+\underline{c},Z,-\underline{d}-\underline{e})/R_{d} (\underline{b},Z,-\underline{d}).$$

Zero parameter coordinates yield important reductions in complexity, as follows.

Lemma 4.2. If  $\underline{d} = (\underline{d}^{(1)'}, \underline{0}')'$  (without loss of generality)  $Z = (Z^{(1)}, Z^{(2)})$ , and  $\underline{u} \sim \mathcal{D}(\underline{b}, Z, -\underline{d})$ , then

$$\underline{\mathbf{u}} \sim \mathcal{D}(\underline{\mathbf{b}}, \mathbf{Z}^{(1)}, -\underline{\mathbf{d}}^{(1)}). \tag{4.9}$$

In particular,  $\mathcal{D}(\underline{b}, \underline{Z}, \underline{0}) \sim D(\underline{b})$ , regardless of the matrix Z.  $\Box$ 

<u>Lemma 4.3</u>. If  $\underline{b} = (\underline{b}_1', \underline{0}')'$ ,  $\underline{Z} = (\underline{Z}_1', \underline{Z}_2')'$  and  $\underline{u} = (\underline{u}_1', \underline{u}_2')'$ , then  $\underline{u} \sim \mathcal{D}(\underline{b}, \mathbb{Z}, -\underline{d})$  if and only if  $\underline{u}_2 = \underline{0}$  with probability one and

$$\underline{u}_{1} \sim \mathcal{D}(\underline{b}_{1}, Z_{1}, -\underline{d}). \quad \Box$$
(4.10)

Lemma 4.3 holds, in particular, for the Dirichlet distribution (d = 0). Compare to Corollary 3.3 for R.

4.2 Inference.

Ģ

Ó

To return to Dirichlet-prior inference for censored multinomial sampling, note that our posterior distribution (2.12) has the representation in the notation of (4.6),

$$\underline{u}|\underline{x},\underline{y} \sim \mathcal{D}(\underline{b}+\underline{x}, \tilde{Z}, -\underline{y}).$$
(4.11)

Hence, the posterior density f(u;b+x;-y) (2.12) is given by (4.7) with parameter values  $b+x, \tilde{z}, -y$ .

 $u \sim \mathcal{D}(b, \tilde{Z}, -d).$ 

by (4.8),

and the posterior distribution as an updating of parameters,

The corresponding posterior density  $f(\underline{u};\underline{b}+\underline{x};-(\underline{d}+\underline{y}))$  is given by (4.7) with parameter values  $\underline{b}+\underline{x}, Z, -(\underline{d}+\underline{y})$ . The Dirichlet-prior theory is the special case  $\underline{d} = 0$ .  $\Box$ The posterior moments for (4.14) are obtained from (4.8)under the posterior parameter values. For example, the posterior

mean is

where the coordinates  $\gamma_{1,i}$  of the vector  $\underline{\gamma}_i$  are defined to be zero except for unity in the ith coordinate. The posterior

Q

As stated earlier, the Dirichlet subfamily is not closed under censored multinomial sampling. However, the family of possible posterior distributions (4.11) does have this property, and thus we consider prior distributions of the following form.

Theorem 4.4. Consider the distributions

(4.12)

As a prior distribution, (4.12) yields the predictive probability,

 $pr(\underline{\delta}_{x}, \underline{\epsilon}_{v}) = [B(\underline{b}+\underline{x})/B(\underline{b})]$ (4.13)

 $R_{d,+y}(\underline{b}+\underline{x}, \tilde{Z}, -\underline{d}-y)/R_{d}(\underline{b}, \tilde{Z}, -\underline{d}),$ 

 $\underline{u} | \underline{x}, \underline{y} \sim \mathcal{D}(\underline{b} + \underline{x}, \widetilde{Z}, -(\underline{d} + \underline{y})).$ (4.14)

 $E(u_1 | \underline{x}, \underline{y}) = [(b_1 + x_1) / (b_1 + x_1)]$ (4.15) $\cdot R_{d_{1}+y_{1}}(\underline{b}+\underline{x}+\underline{\gamma}_{1}, \overline{2}, -\underline{d}-\underline{y})/R_{d_{1}+y_{1}}(\underline{b}+\underline{x}, \overline{2}, -\underline{d}-\underline{y})$ 

second moment is

 $E(u_{1}u_{j}|\underline{x},\underline{y}) = \frac{(b_{1}+x_{1}+\gamma_{1})(b_{j}+x_{j})}{(b_{1}+x_{1}+1)(b_{1}+x_{j})}$ 

• 
$$R_{d_{\cdot}+y_{\cdot}}(\underline{b}+\underline{x}+\underline{\gamma}_{1}+\underline{\gamma}_{1}, \overline{z}, -\underline{d}-\underline{y})/R_{d_{\cdot}+y_{\cdot}}(\underline{b}+\underline{x}, \overline{z}, -\underline{d}-\underline{y}).$$
 (4.16)

4.5

Dickey (1983) proposed the use of D as prior distributions for ordinary (noncensored) multinomial sampling. We suggest, however, that in practice, even for censored data, the usual Dirichlet distributions (or mixtures thereof with a small number of terms) may be preferred to D as prior distributions, thus yielding the new distributions D as posterior distributions. The Dirichlet has the advantage as a prior distribution of being convenient for subjective assessment (see Appendix), though more rigid in the subjective opinions it permits.

### 5. NESTED CENSORING

The posterior moments and distribution theory are considerably simplified when censoring is nested, that is, when for every two sets of confused categories, either one set is contained in the other or they are disjoint. For example, consider a multiplechoice questionnaire. If every person in the sample either answers all questions or is a complete nonrespondent then the censoring is nested. The Dirichlet-prior inference for such sampling was given by Basu and Pereira (1982).

The example can be extended to allow a particular question which the respondents can choose to ignore, or even a set of such questions that can be ignored as a whole. A series of sets of questions, linearly ordered by set-inclusion and the choice to answer all of each such set or only its subset is also permitted under nested censoring. For example, suppose the questions are presented in the same order to each respondent, who chooses some arbitrary place to stop answering questions, such as how far

he/she gets before time runs out. As long as the respondents answer all questions until the places they stop, the data is o nested. However, unconstrained decisions on whether to answer different questions would lead to non-nested censoring patterns in the data. Nested censoring is best treated in terms of nested partitions, which are defined in the following.

Define a matrix  $Z(K \times L)$  to be a <u>partition indicator</u> if its columns indicate disjoint and exhaustive subsets, that is, a partition of  $\{1, \ldots, K\}$ . Formally,

- 1. Each entry of Z is a zero or a one.
- 2. Orthogonal columns:  $\underline{z'} = 0$  for each  $j \neq K$ .
- 3. No row has all zero entries:  $\underline{z}_{i} \neq \underline{0}'$ , for each i = 1, ..., K.

Given a partition indicator Z and K-tuple u, define the corresponding partition-sum vector  $\underline{s}(\underline{u}) = (\underline{s}_1(\underline{u}), \dots, \underline{s}_L(\underline{u}))'$ ,

$$\underline{\mathbf{s}}(\underline{\mathbf{u}}) = \mathbf{Z}'\underline{\mathbf{u}} \quad . \tag{5.1}$$

Then each  $s_j(\underline{u}) = \underline{z}' \underline{u}$ . If  $\underline{u}$  is a probability vector, then the coordinates of  $\underline{s}(\underline{u})$  are the probabilities of the partition-element events.

Define the coordinatewise product of two vectors as the vector of products of corresponding coordinates,  $\underline{z} \times \underline{u} = (z_1 u_1, \dots, z_K u_K)'$ . Then given a partition indicator Z and K-tuple u, define the partition-inner-proportion vectors

$$\underline{\mathbf{r}}_{\mathbf{j}}(\underline{\mathbf{u}}) = \underline{\mathbf{z}}_{\mathbf{*}\mathbf{j}} \times \underline{\mathbf{u}} / \mathbf{s}_{\mathbf{j}}(\underline{\mathbf{u}}), \qquad (5.2)$$

j=1,...,L. Note that if  $\underline{u}$  is a probability vector, each  $\underline{r}_{j}$ consists of the conditional probabilities given the jth partitionelement event. As such; it has zero coordinates for all i not in the j<sup>th</sup> subset. We have the invertible mapping corresponding to (5.1) and (5.2):  $\underline{u} \leftrightarrow \underline{s}(\underline{u}), \underline{r}_{j}(\underline{u})(j=1,...,L)$ . The following theorem extends a result of Wilks (1962).

Theorem 5.1. Given a partition indicator Z, the vector u has the Dirichlet distribution  $u \sim D(b)$  if and only if the L+1 image vectors in (5.1), (5.2) are independently distributed and

r j = 1,...,L. N unlesg z<sub>ii</sub> = 1, probability one The follow literature. Corollary

. 5.2

Proof. By

R

S

Theorem 5.1 follows.

Theorem 5. if and only if for this same Z

Corollary

bution <u>u</u> ~  $\mathcal{D}(b, Z)$ same Z),

 $E[(\Pi u_1^{c_1})\Pi s_1(\underline{u})^{e_j}]$ (5.7)=  $[B(\underline{b}+\underline{c})R_{d}+e_{d},\underline{b}+\underline{c},\underline{z},-\underline{d}-\underline{e})]$  $/[B(\underline{b})R_{d}(\underline{b},Z,-\underline{d})]$ 

$$\underline{\underline{s}}(\underline{\underline{u}}) \sim D(\underline{\underline{s}}(\underline{\underline{b}}))$$
(5.3)  

$$\underline{\underline{r}}_{j}(\underline{\underline{u}}) \sim D(\underline{\underline{z}}_{*j} \times \underline{\underline{b}}),$$
(5.4)  
Note that  $\underline{\underline{z}}_{*j} \times \underline{\underline{b}} = \underline{s}_{j}(\underline{\underline{b}})\underline{r}_{j}(\underline{\underline{b}}),$  and that by Lemma 4.3,  
I, the i<sup>th</sup> coordinate of  $\underline{\underline{r}}_{j}(\underline{\underline{u}})$  equals zero with  
ne.  $\Box$   
by using simple form of R appears to be new in the  
 $\underline{\underline{r}}$  5.2. If Z is a partition indicator,  
 $\underline{R}_{d}.(\underline{\underline{b}}, \underline{Z}, -\underline{\underline{d}}) = B(\underline{\underline{s}}(\underline{\underline{b}}) + \underline{\underline{d}}) / B(\underline{\underline{s}}(\underline{\underline{b}})) . \Box$ (5.5)  
by (5.3),  
 $\underline{E}_{u|b}^{(K-1)} \Pi_{1}^{L}(\underline{\underline{z}}_{*j}, \underline{\underline{u}})^{dj} = E_{s(u)|s(b)}^{(L-1)} \Pi_{1}^{L}\underline{s}_{j}(\underline{\underline{u}})^{dj}.$ (5.6)  
5.1 generalizes immediately to distributions  $p$ , as  
 $\underline{\underline{r}}$ .  
Given a partition indicator Z,  $\underline{\underline{u}} \sim p(\underline{\underline{b}}, \underline{Z}, -\underline{\underline{d}})$   
the independent distributions (5.3), (5.4) hold  
Z, with  $\underline{\underline{s}}(\underline{\underline{b}})$  replaced by  $\underline{\underline{s}}(\underline{\underline{b}}) + \underline{\underline{d}}$ .  $\Box$   
 $\underline{\underline{5}}.\underline{\underline{4}}$ . Given a partition indicator Z, the distri-  
 $\underline{r}_{s}, \underline{Z}, -\underline{\underline{d}}$ ) has closed-form moments (in terms of the

=  $[B(\underline{b}+\underline{c})B(\underline{s}(\underline{b}+\underline{c})+\underline{d}+\underline{c})/B(\underline{s}(\underline{b}+\underline{c}))]$ / $[B(\underline{b})B(\underline{s}(\underline{b})+\underline{d})/B(\underline{s}(\underline{b}))]$ 

=	$\mathbb{B}(\underline{z}_{*j} \times (\underline{b} + \underline{c}))$	•	$\frac{B(\underline{s}(\underline{b}+\underline{c})+\underline{d}+\underline{e})}{B(\underline{s}(\underline{b})+\underline{d})}$	
	( <u>_</u> *j^ <u>_</u> )			ļ

By Lemma 4.2, this yields the moments of the prior distribution (4.12) and the posterior distributions (4.11), (4.14), under censoring by a partition, provided the prior distribution refers to the same partition (if to any).

The elements of a partition of a set are subsets, any one of which could then be partitioned further. Two partitions, of which one is a partition of an element of the other, are said to be <u>directly nested</u>. Then one is a partition of a set and the other is a partition of a subset. Two partitions are said to be <u>nested</u> if they belong to a sequence of successively directly nested partitions. Finally, define a collection of partitions to be <u>nested</u> if they form a tree; that is, each partition is nested with respect to a particular partition, the <u>root</u> of the collection. Then the root is the only partition in the collection that is a partition of the original set.

Define matrix  $Z(K \times L)$  to be a <u>two-level nested-partitions</u> <u>indicator</u> if it indicates a collection of partitions in which each nonroot partition is directly nested in the root partition. Namely,  $Z = (Z^{(1)}, Z^{(2)})$  where  $Z^{(1)}$  indicates the root partition and

 $z^{(2)} = (z_{j}^{(2)})_{j=1}^{L^{(1)}},$ 

of  $\underline{z}_{*i}^{(1)}$ .

5.4

Π

(5.8)

Such a two-level nested-partitions distribution  $\mathcal{D}$  can be expressed in terms of independent Dirichlet-distributed vectors, by application of Theorem 5.3 to (5.10). By iteration of Theorem 5.5, a nested-partitions distribution of any number of levels can be so expressed. Thus, the moments of any nested-partitions distribution can be obtained in closed form, as in Corollary 5.4. This yields further new closed forms for the function R for a nested-partitions indicator.

in which each submatrix  $Z_j^{(2)}$  indicates a partition of the subset indicated by the j<sup>th</sup> column  $\underline{z}_{*j}^{(1)}$  of  $Z^{(1)}$ , j=1,..., $L^{(1)}$ . The submatrices  $Z_j^{(2)}$  are filled out with zero rows for zero coordinates

<u>Theorem 5.5</u>. Given the two-level nested-partitions indicator Z,  $\underline{u} \sim \mathcal{D}(\underline{b}, Z, -\underline{d})$  if and only if independently

 $\underline{\mathbf{s}}^{(1)}(\underline{\mathbf{u}}) \sim D(\underline{\mathbf{s}}^{(1)}(\underline{\mathbf{b}}) + \underline{\mathbf{d}}^{(1)} + \underline{\overline{\mathbf{d}}}^{(2)})$ (5.9)

 $\underline{r}_{j}^{(1)}(\underline{u}) \sim \mathcal{D}(\underline{z}_{j}^{(1)} \times \underline{b}, \mathbb{Z}_{j}^{(2)}, -\underline{d}_{j}^{(2)}), \qquad (5.10)$ 

 $j=1,...,L^{(1)}, \text{ where } \underline{s}^{(1)} \text{ and } \underline{r}^{(1)}_{j}(j=1,...,L^{(1)}) \text{ refer to } Z^{(1)}.$ In conformity to  $(Z^{(1)},Z^{(2)}), \underline{d}' = (\underline{d}^{(1)'}, \underline{d}^{(2)'}) \text{ and}$  $\underline{d}^{(2)'} = (\underline{d}^{(2)'}_{j})_{\underline{j}=1}^{L^{(1)}}; \qquad (5.11)$ 

and  $\underline{d}^{(2)}$  has the j<sup>th</sup> coordinate  $\underline{d}_{1}^{(2)'} \underline{1}$ ; j=1,...,L<sup>(1)</sup>.

5.5

## 6. EXPANSION BY POSSIBLE DATA

A well known tautology in Bayesian statistics equates a posterior distribution based on censored data to a probability mixture of posterior distributions from the possible uncensored versions of the observed data. The following form of this tautology exhibits a D distribution as being, itself, a finite mixture of Dirichlet distributions. A useful expansion of R is obtained as a consequence.

Consider, for this purpose, matrices W of possible frequency counts refining the observed censored frequency vector y. Write

$$W = (w_{i\sigma}: i=1,...,K, \sigma \in \{1,...,K\})$$
 (6.1)

This K×L matrix has columns indexed by subsets  $\sigma(L=2^{K}-K-2)$ . For each subset  $\sigma < \{1, \ldots, K\}$ , the column  $\underline{w}_{*\sigma}$  is supported on  $\sigma$ ,

$$w_{i\sigma} = 0 \text{ for } i \not c \sigma,$$
 (6.2)

and these possible frequencies for the categories in  $\sigma$  sum to the observed count for  $\sigma$ ,

$$w_{\cdot\sigma} = \Sigma_{\sigma} \quad . \tag{6.3a}$$

This says, for the row-vector of column sums,

$$\underline{\mathbf{w}} = \underline{\mathbf{l}}' \mathbf{W} = \underline{\mathbf{y}}' \quad . \tag{6.3b}$$

Note that the column-vector of row sums,  $\underline{w}_{*} = \underline{Wl}$ , is a possible uncensored frequency vector for {1,...,K} harmonizing with the observed censored frequencies y.

Lemma 6.1. The sampling probability of a censored sequence is the weighted sum of probabilities over all possible refinements of

6.1

The notation W|y indicates that the summation is over nonnegative integer w<sub>in</sub> constrained by (6.2), (6.3). Proof.

In a more detailed notation, one can write,

e Theorem 6.2. For a Dirichlet prior distribution  $\underline{u} \sim D(\underline{b})$ , the posterior distribution  $\mathcal{D}$  (2.12), (4.11) based on <u>x</u>, <u>y</u> is a posteriorprobability mixture of posterior Dirichlet distributions based on

<u>X</u>,<u>W</u>,

where f with two and f with three arguments are defined by (2.3) and (4.11), (4.7), respectively,

6.2

(6.5)

the censored frequencies,

$$pr(\underline{\varepsilon}_{y}|\underline{u}) = \sum_{W|y} [\Pi_{\sigma} \begin{pmatrix} y_{\sigma} \\ \underline{w}_{*\sigma} \end{pmatrix}] \Pi_{i} u_{i}^{W} \cdot .$$
 (6.4)

$$pr(\underline{\varepsilon}_{y}|\underline{u}) = \Pi_{\sigma}(\sum_{i \in \sigma} u_{i})^{y_{\sigma}}$$
$$= \Pi_{\sigma}\sum_{\substack{w \neq \sigma}} \left[ y_{\sigma} \left[ \frac{y_{\sigma}}{\underline{w}_{\star\sigma}} \right] \Pi_{i} u_{i}^{w_{i\sigma}} \right]$$
$$= \sum_{\substack{w \neq y}} \Pi_{\sigma}[\left( \frac{y_{\sigma}}{\underline{w}_{\star\sigma}} \right] \Pi_{i} u_{i}^{w_{i\sigma}}] .$$

$$\Pi_{\sigma} \begin{pmatrix} y_{\sigma} \\ \frac{W}{*\sigma} \end{pmatrix} = \prod_{k=2}^{K-1} \prod_{\substack{1 \leq \cdots \leq i_{k} \\ w \neq 1 \leq \cdots \leq i_{k}}} \begin{pmatrix} y_{1} \cdots y_{k} \\ \frac{W}{*1} \cdots y_{k} \end{pmatrix} . \quad (6.6)$$

$$f(\underline{u};\underline{b}+\underline{x};-\underline{y})$$

$$= \sum_{W|Y} f(\underline{u};\underline{b}+\underline{x}+\underline{w}_{*}) pr(W|\underline{b}+\underline{x},\underline{y}),$$
(6.7)

6.3

(6.8)

(6.9)

pr(W | b+x, y)

 $= \left[ \Pi_{\sigma} \begin{pmatrix} y_{\sigma} \\ \underline{w} \end{pmatrix} \right] g(\underline{w}_{\star}; \underline{b} + \underline{x}) / h(\underline{y}; \underline{b} + \underline{x})$ 

and

h(y; b+x)

 $= \sum_{W \mid y} \left[ \Pi_{\sigma} \begin{pmatrix} y_{\sigma} \\ \underline{w}_{\star} \end{pmatrix} \right] g(\underline{w}_{\star}; \underline{b} + \underline{x}),$ 

where g and h are defined by (2.5), (2.11), (4.3).

Related results appear in Shefrin (1981). Theorem 6.2 can be generalized immediately for use of a prior distribution  $\mathcal{D}(\underline{b}, Z, -\underline{d})$ , in the spirit of Theorem 4.4, merely by substituting d+y for y.

The representation (6.7) can be reexpressed in a general notation for probability mixtures,

> $\underline{u} | \underline{x}, \underline{y} \sim D(\underline{b} + \underline{x} + \underline{w}) * W$ (6.10)

where W has the mass function (6.8), or its D-prior form. In this notation, our theorem generalizes simply to give a mixture-representation for any  $\mathcal{D}$  distribution, as follows.

Theorem 6.3. Consider arbitrary Z(K×L, for arbitrary K,L). Conformably partition  $Z = (Z^{(1)}, Z^{(2)}),$  $\underline{d} = (\underline{d}^{(1)'}, \underline{d}^{(2)'})' \text{ and } L = L^{(1)} + L^{(2)}. \text{ Then } \underline{u} \sim \mathcal{D}(\underline{b}, \mathbb{Z}, -\underline{d}) \text{ if and}$ only if

> $\underline{u} \sim \mathcal{D}(\underline{b}+\underline{w}, Z^{(1)}, -\underline{d}^{(1)}) * W,$ (6.11)

for  $\underline{w}_{} = \underline{d}^{(2)}$ .  $\Box$ 

Corollary 6.4.

Q

In our experience so far, the expansions of this section have been more convenient and economical for computations than the naive quadrature of analogs of the dimension-reduced integral (3.6).

6.4

(6.13)

where W(K×L<sup>(2)</sup>) has probability mass

$$pr(W) = \left[ \Pi_{j=1}^{L(2)} \begin{pmatrix} d_{j}^{(2)} \\ \vdots \\ w_{*j} \end{pmatrix} \Pi_{i=1}^{K} \overset{(2)^{W}ij}{ij} \right]$$

$$\frac{(\underline{b}+\underline{w})}{(\underline{b})} \cdot \frac{\underset{d}{R_{d}}(\underline{b},\underline{Z},-\underline{d})}{(\underline{b},\underline{Z},-\underline{d})}, \qquad (6.12)$$

Theorem 6.2 is the special case,  $L^{(1)} = 0$  with parameters  $\underline{b}+\underline{x}, \tilde{z}, -\underline{y}$ . Note that the generalization of restriction (6.2),  $w_{ij} = 0$ for  $z_{i,i}^{(2)} = 0$ , holds automatically for (6.12). Since (6.12) must sum to unity, we have an apparently new representation for R under the parameter restriction of (3.6).

$$R_{d} (\underline{b}, \underline{z}, -\underline{d}) = \sum_{w \mid d} (2) \begin{bmatrix} \pi^{L} \\ \mu \\ \mu \end{bmatrix} \begin{bmatrix} \pi^{L} \\ \mu \\ \mu \\ \mu \end{bmatrix} \begin{bmatrix} \pi^{K} z_{ij} \\ \mu \\ \mu \\ \mu \\ \mu \end{bmatrix}$$

- $[B(\underline{b}+\underline{w}_{*})/B(\underline{b})]R_{d}(1)^{(\underline{b}+\underline{w}_{*},Z^{(1)},-\underline{d}^{(1)})}$
- In case  $L^{(1)} = 0$  in (6.13),  $R_{d^{(1)}} = 1$  as in (6.9).

Such expansion methods can be easily combined with the partitionnesting methods of the preceding section to minimize the number of terms involved, as in Antelman (1972) for special cases. Recursion relations are easily set up for the multinomial coefficients.

Ŷ

6.5

Ŷ

An interesting special case of censored multinomial sampling is provided by contingency tables with supplemental purely marginal data. Consider a two-way table with three independent multinomial data sets, respectively for: the cross-classified table, itself; the row variable alone; and the column variable alone. By the likelihood principle, for suitable models, these data can be directly combined into the inferential equivalent of a single censored multinomial sample. Chen and Fienberg (1974) analyze such data, given here in Table 1, on M+N = 456 premature live births, classified by serum bilirium level (mg. per 100 ml, "Low" 0-1.0) and/or a composite health index (0-10, "Low" ("healthier") 0-6).

Serum bilirium

Low

High

Supplemental d on health inde

body of Table 1,

### 7. EXAMPLES

## Example 1. Contingency Tables with Supplemental Margins.

. Table 1. Data on premature live births (Chen and Fienberg 1974).

:	Health Low	Index High	÷.	Supp on s	lement erum b	al d ilir	lata ium
	35	75		a s	11		
	57	112			13		
ata x	117	36					-

For the cross-classified cell probabilities u (i serum bilirium, j health index), a uniform prior distribution,  $(u_{11}, u_{12}, u_{21}, u_{22}) \sim D(\underline{b}), \underline{b} = (1, 1, 1, 1)', yields a posterior$ distribution D(b+x,Z,-d), where x is given by the cross-classified 7.1

. 7.2

 $Z = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \end{pmatrix},$ (7.1)

and d = (117, 36, 11, 13)'. This posterior distribution can be expressed as a mixture of nested-partition distributions by Theorem 6.3, by expanding the two lower powers  $(u_{11}+u_{12})^{11}$  and  $(u_{21}+u_{22})^{13}$ . The resulting posterior means are given in Table 2 with the posterior standard deviations and correlations, and the maximum likelihood solution from Chen and Fienberg (1974). These Bayesian estimates are not very different from the maximum likelihood values, as should be expected from a uniform prior distribution and reasonable sample sizes. Hence, as is true of the likelihood inference, appreciable information is gained in the Bayesian analysis from the supplemental marginal data.

Table 2.	Estimates	of	cross-classified	cell'	probabilities	for

data	on	premature	live	births.

2				
Cell probabilities	ull	<sup>u</sup> 12	u <sub>21</sub>	<sup>u</sup> 22
Maxlikelihood est.	0.1880	0.2090	0.2960	0.3090
Posterior mean	0.1888	0.2096	0.2950	0.3065
Posterior SD	0.0255	0.0206	0.0277	0.0230
Posterior correla- tions <sup>u</sup> ll <sup>u</sup> l2	l	-0.21 1	-0.60 -0.21	-0.20 -0.40
<sup>u</sup> 21			1	-0.35

Kadane (1982) analyzed data from two sample surveys of attitudes on the death penalty. The primary categories are, for i = 1, 2, 3, 4:

nested (M+N=2338).

## Example 2. Combining Surveys with Different Questionnaires

- 1. Would not decide guilt versus innocence in a fair and impartial manner.
- 2. Fair and impartial on guilt versus innocence; and, on sentencing, would always vote for the death penalty, regardless of circumstances.
- 3. Fair and impartial; and would never vote for the death penalty.
- 4. Fair and impartial; and would sometimes and sometimes not vote for the death penalty.

A survey by the Field Research Corporation produced data,  $x_1 = 68$ ,  $x_3 = 97$ ,  $y_{2,4} = 674$  (M+N=839), and a Harris survey produced data,  $x_2$ =15,  $y_{1,3,4}$ =1484 (M+N=1499). By the likelihood principle, these two multinomial samples can be directly combined in the form of a single censored multinomial sample, for which the censoring is not

A genuine Bayesian analysis reports the coherent effect of data on prior distributions expressing actual expert opinion prior to knowledge of the data. To simulate aspects of such a process, we assessed a Dirichlet distribution by interactive elicitation of the opinion of a social psychologist with interests in legal

matters, a person who was not yet familiar with the two surveys. (Our assessment method is outlined in the Appendix.) The resulting prior and posterior distributions are, respectively,  $D(\underline{b})$  with  $\underline{b}$  as given in Table 3, and  $\mathcal{D}(\underline{b}+\underline{x}, Z, -\underline{d})$  where  $\underline{b}+\underline{x} = (70.8, 26.2, 118.0, 105.0)$ ,  $\begin{pmatrix} 0 \\ 1 \end{pmatrix}$ 

$$Z = \begin{bmatrix} 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix},$$
(7.2)

7.4

and d = (674, 1484)'.

Using the method of Theorem 6.3 to expand  $(u_2+u_4)^{674}$  in D, we obtain a probability mixture of nested-partition distributions, each having closed-form moments. The resulting posterior means, standard deviations, and correlations for the four cell probabilities are given in Table 3, together with the posterior means from a uniform prior distribution, and the estimates obtained by simply solving the following simultaneous linear equations. Setting each  $u_{i}$  and  $u_{\sigma}$  equal to the corresponding observed frequency ratio, we have the system,

$$u_1 = 68/839, u_3 = 97/839, u_2 + u_4 = 674/839;$$
 (7.3)

$$u_2 = 15/1499, u_1 + u_3 + u_4 = 1484/1499.$$
 (7.4)

two surveys anal;	yzed by Kadan	e (1982).		
Cell probabilities	ul	u <sub>2</sub>	u <sub>3</sub>	u <sub>4</sub>
Combined survey data				
<u>x</u>	68	15	97	0
<u>y</u>	y <sub>1,3,4</sub> =1484	$y_{2,4} = 674$		
(M+N = 2337)				
Estimate by observed				
frequency ratios	0.081	0.010	0.116	0.793
Posterior mean from				
uniform prior	0.082	0.011	0.116	0.791
Expert prior mean <u>b</u> /b	0.020	0.080	0.150	0.750
(b = 140)				
Expert prior SD	0.012	0.023	0.030	0.036
Posterior mean	0.073	0.016	0.122	0.789
Posterior SD	0.008	0.003	0.010	0.013
Posterior correlations				
ul	1	-0.01	-0.10	-0.55
u <sub>2</sub>		1	-0.01	-0.23
u <sub>n</sub>			1	-0.74

The three vector estimates are not very different. Note that the expert posterior mean vector is farther away from the

7.5

observed-frequency estimate than is the posterior mean yielded by the uniform prior distribution, and this is true in each coordinate separately. This seems reasonable in that the uniform prior expresses greater uncertainty than the expert prior. The effect of either prior is to move each coordinate of the observedfrequency estimate toward the corresponding prior mean. The uniform prior, however, moves the first coordinate estimate in the opposite direction than does the expert prior. Hence, the uniform prior distribution seems less reasonable than the expert prior in this case.

7.6

A subjective prior distribution in the Dirichlet family can be assessed, as follows, in the multiple-Bernoulli or multinomial sampling context and notation of Section 2. Our method is a "device of imaginary results", in the language of I.J. Good. That is, the expert whose opinion is being assessed imagines hypothetical data and the method elicits his pretended reaction to it. Of course, real data can be used in a similar way. The elicitations refer only to the Dirichlet-Bernoulli predictive distribution. Indeed the predictive distribution is primary and can be assessed by the method, regardless of the significance or not of assuming an imbedded multiple-Bernoulli process and Dirichlet mixing distribution. The assessment proceeds in three steps.

observation,

Q

pr

# APPENDIX. A Simple Method for Assessment of a Dirichlet Subjective Prior Distribution in Multinomial Sampling.

1. Elicit predictive probabilities for a single future

 $pr(\delta=i) = b_i/b_i, i=1,...,K.$ 

2. Condition opinion on an imaginary (or real) future sample with frequencies  $\underline{x} = (x_1, \dots, x_K)$ . Elicit conditional predictive probabilities for a further future observation,

$$(\delta=1|\underline{x}) = \frac{b}{b_{\cdot}+x_{\cdot}} \operatorname{pr}(\delta=1) + \frac{x_{\cdot}}{b_{\cdot}+x_{\cdot}} u_{1}, \qquad (A.2)$$

A.1

(A.1)

المتوج

where pr( $\delta$ =i) was obtained in Step 1 and  $\hat{u}_1$  denotes the usual relative-frequency estimate,

$$\hat{u}_{i} = x_{i}/x. \qquad (A.3)$$

Solve for b,

$$b_{\cdot} = x_{\cdot} \left[ \frac{\hat{u}_{1} - pr(\delta=i | \underline{x})}{pr(\delta=i | \underline{x}) - pr(\delta=i)} \right]$$
 (A.4)

3. Calculate b, ,

 $b_i = pr(\delta=i)b_i, i=1,\ldots,K$ . (A.5)

In practice, the numerator and denominator in (A.4) should have the same sign. Under the conjugate-prior model,  $pr(\delta=i|\underline{x})$ lies between  $pr(\delta=i)$  and  $\hat{u_1}$ . The ratio of its distance to  $\hat{u_1}$ to its distance to  $pr(\delta=i)$  is  $b_1/x_1$ , a (positive) constant in i. If  $pr(\delta=i \mid \underline{x})$  is elicited and (A.4) solved for b for various i values, then these b values can be averaged. Similarly, an average b can be obtained from consideration of several samples  $\underline{x}$ . Pooling of samples would permit nested conditioning and avoid the need for forgetting or unpretending.

The values given in Table A.l were elicited from a social psychologist for Example 2 of Section 7. The simple average of the assessed b values is  $\overline{b} = 140$ .

		•		
Tab	ole A.l.	Assessme	nt of an	expert
pri	or distr	ibution fo	or death	-penalty
att	itude su	rveys.		
1	1	2	3	4
Prior-predictive probabilities				
pr(\delta=i)				
[Elicited]	0.02	0.08	0.15	0.75
Imaginary sample data				
$x_{i}$	16	20	32	132
Relative-frequency estimate				
û	0.08	0.10	0.16	0.66
Posterior-predictive probabilities				
$pr(\delta=i   \underline{x})$ [Elic! - ed]	0.05	0.09	0.16	0.70
Solution				
$(\overline{b} = 140)$	200	200	с	160
$b_1 = \overline{b}.pr(\delta=1)$	2.8	11.2	21.0	105.0

-

Ŷ

A.2

A.3

### REFERENCES AND BIBLIOGRAPHY

Antelman, Gordon R. (1972) Interrelated Bernoulli processes.

J. Amer. Statist. Assoc. Vol. 67, 831-841. Appell, P., and J. Kampé de Fériet (1926). <u>Fonctions Hyper-</u><u>géométriques et Hypersphériques: Polynomes d'Hermite</u>.

Gauthier-Villars, Paris.

Basu, D., and Pereira, Carlos A. de B. (1982). On the Bayesian analysis of categorical data: The problem of nonresponse. J. of Statistical Planning and Inference, Vol. 6, 345-362.

Carlson, B.C. (1963). Lauricella's hypergeometric function  $F_D$ . J. of Mathl. Anal. and Appls., Vol. 7, 452-470.

Carlson, B.C. (1971). Appell functions and multiple averages.

S.I.A.M. J. of Math. Anal. Vol. 2, No. 3, (Aug'71),

420-430.

- Carlson, B.C. (1974). Inequalities for Jacobi polynomials and Diricnlet averages. <u>S.I.A.M. J. of Math. Anal</u>. Vol. 5, No. 4, (Aug'74), 586-596.
- Carlson, B.C. (1977). <u>Special Functions of Applied Mathematics</u>. Academic Press, New York.
- Chen, T. and Fienberg, S.E. (1974). Two-dimensional contingency tables with both completely and partially cross-classified data. <u>Biometrics</u>, Vol. 30, 629-642.

Chen, T. and Fienberg (1976). The analysis of contingency tables with incompletely classified data. <u>Biometrics</u>, Vol. 32, 133-144. Dawid, A.P. and Dickey, James M. (1977). Likelihood and Bayesian inference from selectively reported data. <u>Journal of the</u> <u>American Statistical Association</u>, Vol. 72, 845-850.



