

# AN ITERATIVE PROCEDURE FOR OBTAINING I-PROJECTIONS ONTO THE INTERSECTION OF CONVEX SETS<sup>1</sup>

BY RICHARD L. DYKSTRA

*The University of Iowa*

A frequently occurring problem is to find a probability distribution lying within a set  $\mathcal{E}$  which minimizes the  $I$ -divergence between it and a given distribution  $R$ . This is referred to as the  $I$ -projection of  $R$  onto  $\mathcal{E}$ . Csiszar (1975) has shown that when  $\mathcal{E} = \cap_{i=1}^m \mathcal{E}_i$  is a finite intersection of closed, linear sets, a cyclic, iterative procedure which projects onto the individual  $\mathcal{E}_i$  must converge to the desired  $I$ -projection on  $\mathcal{E}$ , provided the sample space is finite.

Here we propose an iterative procedure, which requires only that the  $\mathcal{E}_i$  be convex (and not necessarily linear), which under general conditions will converge to the desired  $I$ -projection of  $R$  onto  $\cap_{i=1}^m \mathcal{E}_i$ .

**1. Introduction.** Suppose  $p$  and  $q$  are probability measures defined on subsets of the finite set  $\mathcal{X}$ , which without loss of generality we take to be the first  $m$  positive integers. The  $I$ -divergence of  $p$  with respect to  $q$ , also called the Kullback-Liebler information number, cross-entropy between  $p$  and  $q$ , information for discrimination, entropy of  $p$  relative to  $q$ , etc., is given by

$$(1.1) \quad I(p \parallel q) = \sum_{k=1}^m p(k) \ln(p(k)/q(k))$$

where  $p(k)$  is the mass placed by the probability measure  $p$  at the point  $k$ . We let  $P$  denote the set of all probability measures on  $\mathcal{X}$ , and use the convention that products (or quotients) are to be interpreted as pointwise multiplication (or division). For example,  $s_{12} = rp_{11}/s_{11}$  means that the measure  $s_{12}$  puts mass  $r(k)p_{11}(k)/s_{11}(k)$  at  $k$ .

We mention that  $I(p \parallel q)$  is defined analogously for general probability measures on infinite spaces, but for simplicity, we will only consider finite sets. (See Kullback, 1959, or Csiszar, 1975, for the general definition.)

It is well known that  $I(p \parallel q) \geq 0$ , and that  $I(p \parallel q) = 0$  iff  $p \equiv q$ . Thus it is heuristically reasonable to think of  $I(p \parallel q)$  as representing a "distance" between  $p$  and  $q$ . However,  $I(\cdot \parallel \cdot)$  is not a metric, nor is the symmetrized version

$$J(p \parallel q) = \frac{I(p \parallel q) + I(q \parallel p)}{2}$$

used by Jeffreys (1948). Nevertheless,

$$(1.2) \quad \sum_{k=1}^m |p(k) - q(k)| \leq [2I(p \parallel q)]^{1/2},$$

---

Received November 1983; revised November 1984.

<sup>1</sup> This work was partially supported by ONR Contract N00014-83-K-0249.

AMS 1970 subject classifications. Primary 90C99; secondary 49D99.

Key words and phrases.  $I$ -divergence,  $I$ -projections, convexity, Kullback-Liebler information number, cross-entropy, iterative projections, iterative proportional fitting procedure.

as shown independently by Kullback (1967), Kemperman (1967), and Csiszar (1967), so that we have some idea of what small values of  $I(p \parallel q)$  imply.

If we interpret  $I(p \parallel q)$  as distance, then it seems natural to define the  $I$ -projection of the probability distribution  $r$  onto a set  $\mathcal{E}$  of probability distributions as being a  $q \in \mathcal{E}$  such that  $I(q \parallel r) < \infty$  and

$$(1.3) \quad I(q \parallel r) = \min_{p \in \mathcal{E}} I(p \parallel r).$$

In some sense,  $q$  is a measure in  $\mathcal{E}$  that lies as close as possible to  $r$ .

Minimization problems of the form (1.3) play a key role in the information-theoretic approach to statistics (e.g., Kullback, 1959; Good, 1963; etc.) and also occur in other areas such as the theory of large deviations (Sanov, 1957) and maximization of entropy (Rao, 1965; and Jaynes, 1957).

However, in statistical circles,  $I$ -projections are probably most important for being dual problems to certain log-linear model maximum likelihood estimation (MLE) problems. In particular, it is known that the multinomial MLE problem:

$$(1.4) \quad \text{Maximize } \prod_{k=1}^m p(k)^{n(k)} \quad \text{subject to } p \in P, \ln p \in \mathcal{M}$$

(where  $\mathcal{M}$  is some subspace of  $\bar{R}^m$  containing the constant vectors) has precisely the same solution as the  $I$ -projection problem:

$$(1.5) \quad \text{Minimize } I(p \parallel u) \quad \text{over } P \quad \text{such that } s - p \in \mathcal{M}^\perp$$

(where  $s(k) = (\sum_1^m n(i))^{-1}n(k)$ , and  $u$  is the uniform distribution on  $\mathcal{X}$ ). Note that for a subspace spanned by the vectors  $a_1, \dots, a_t$ , we have

$$\mathcal{M}^\perp = \{a_1, \dots, a_t\}^\perp = \cap_1^t \{a_i\}^\perp,$$

and hence (1.4) is equivalent to minimizing  $I(p \parallel u)$  for  $p$  in the set  $\cap_1^t (s - \{a_i\}^\perp)$ . Since the  $s - \{a_i\}^\perp$  are linear spaces, Csiszar's algorithm of cyclic, iterated  $I$ -projections is appropriate here. In this format, it is easy to see the connection between Csiszar's procedure and the IPFP (iterative proportional fitting procedure) which has received so much attention in the general area of categorical data. Meyer (1980) has an extensive discussion and several examples where he relates Csiszar's procedure and the general IPFP.

Suppose now that  $\mathcal{M} = K_1 + \dots + K_t$  is a closed, convex cone expressible as a direct sum of closed, convex cones containing the constant vectors, rather than a direct sum of subspaces. Such a configuration would arise naturally if one were considering order constraints in a log-linear model. Then it is well known that the dual (polar) cone of  $\mathcal{M}$ , defined as

$$\mathcal{M}^* = \{y; \sum_1^m y(i)x(i) \leq 0 \text{ for all } x \in \mathcal{M}\}$$

is expressible as

$$\mathcal{M}^* = (K_1 + K_2 + \dots + K_t)^* = K_1^* \cap K_2^* \cap \dots \cap K_t^*.$$

Lemke and Dykstra (1984) have generalized the (1.4)–(1.5) duality results to the case where  $\mathcal{M}$  is a closed, convex cone (with  $\mathcal{M}^*$  replacing  $\mathcal{M}^\perp$  in (1.5)). This means that many MLE problems involving partial orders in log-linear models

are equivalent to *I*-projection problems of the form:

$$\text{Minimize } I(p \| u) \text{ for } p \in \cap_1^t (s - K_i^*).$$

If the  $K_i^*$ 's are not subspaces, Csiszar's cyclic, iterated scheme need not work. However, the procedure described in this paper will work since the sets  $s - K_i^*$  will be closed, convex sets of probability distributions.

Of course, we would really like to be able to identify structure in these log-linear model situations, which leads to the area of inference for various competing models. While these are important questions, we shall only be concerned with the MLE problem in this paper.

Csiszar (1975) discusses *I*-projections in great detail, and has a geometric development for *I*-projections which is quite appealing. (Cencov, 1972, also has a geometric development of *I*-projections, but with the arguments interchanged.) Csiszar also discusses the existence of *I*-projections, and shows that if  $\mathcal{E}$  is a convex set of probability distributions (PDs) which is variation closed, the unique existence of the *I*-projection of  $r$  onto  $\mathcal{E}$  is guaranteed provided there exists a  $p \in \mathcal{E}$  such that  $I(p \| r) < \infty$ . This result is clearly applicable for closed, convex sets of PDs on the finite set  $\mathcal{X}$ . We shall make repeated use of the elegant characterization of *I*-projections given in the following theorem.

**THEOREM 1.1 (Csiszar).** *A probability  $q \in \mathcal{E}$  with  $I(q \| r)$  finite is the *I*-projection of  $r$  onto the convex set  $\mathcal{E}$  of probability distributions iff*

$$(1.6) \quad I(p \| r) \geq I(p \| q) + I(q \| r) \quad \forall p \in \mathcal{E}.$$

*Note that in our setting of finite  $\mathcal{X}$ , it follows from (1.1) that (1.6) is equivalent to*

$$(1.7) \quad \sum_{k=1}^m (p(k) - q(k)) \ln(q(k)/r(k)) \geq 0 \quad \forall p \in \mathcal{E}.$$

If in fact  $q$  is an algebraic inner point of  $\mathcal{E}$ , i.e., for every  $p$  in  $\mathcal{E} - \{q\}$ , there exists  $0 < \alpha < 1$  and  $p' \in \mathcal{E}$  such that  $q = \alpha p + (1 - \alpha)p'$ , equality must hold in (1.6) and (1.7).

This situation is roughly akin to projecting onto subspaces in least squares theory. In particular, Csiszar defines an  $\mathcal{E}$  to be a linear set if  $p, p' \in \mathcal{E}$  implies  $\alpha p + (1 - \alpha)p' \in \mathcal{E}$  for every  $\alpha$  for which it is a probability. If  $\mathcal{E}$  is a linear set, then the inequality sign in (1.6) and (1.7) may be replaced by an equality sign as long as  $\mathcal{X}$  is finite. Based upon this characterization, Csiszar is able to prove that if  $\mathcal{X}$  is finite,  $\mathcal{E} = \cap_1^t \mathcal{E}_i$  is a finite intersection of arbitrary linear sets, and there exists a  $p \in \mathcal{E}$  such that  $I(p \| r) < \infty$ , then successive, cyclic, iterated *I*-projections onto the individual sets must converge to the *I*-projection of  $r$  onto  $\mathcal{E}$ . Thus if  $q_0 = r$ , and  $q_n$  denotes the *I*-projection of  $q_{n-1}$  onto  $\mathcal{E}_n$  (where  $\mathcal{E}_n = \mathcal{E}_i$  if  $n = t + i, 1 \leq i \leq t$ ), then  $q_n$  must converge to  $q \in \mathcal{E}$  as  $n \rightarrow \infty$  where  $I(q \| r) = \min_{p \in \mathcal{E}} I(p \| r)$ .

This result is very much dependent upon the assumption that the  $\mathcal{E}_i$  be linear sets (in fact it is not true in general) and the accompanying fact that equality holds in (1.6).

**2. The procedure.** We now propose a procedure which will enable one to obtain *I*-projections onto a finite intersection of arbitrary closed, convex sets of

probabilities by iteratively finding  $I$ -projections onto the individual sets. We will prove that that under a mild restriction, the procedure must give the correct solution, and then we examine an example.

First let us note that we can also define an  $I$ -projection onto  $\mathcal{E}$  for any nonzero, finite measure  $r$  on  $\mathcal{X}$  as being a probability in  $\mathcal{E}$  which minimizes

$$\sum_i^n q(k)\ln(q(k)/r(k))$$

over  $\mathcal{E}$  and has finite  $I$ -divergence with respect to  $r$ .

Of course, since  $\mathcal{E}$  contains only probabilities, the  $I$ -projection of  $r$  onto a convex set  $\mathcal{E}$  is identical to the  $I$ -projection onto  $\mathcal{E}$  of the normalized measure  $r'$  defined by  $r'(k) = r(k)/r(\mathcal{X})$ . It is easily shown that the characterization of  $I$ -projections given in (1.6) and (1.7) is still valid, even though  $I(p \parallel r) \geq 0$  need no longer be true.

Let us now state our algorithm. We assume that we wish to find the  $I$ -projection of  $r$  onto  $\mathcal{E} = \cap_i^t \mathcal{E}_i$ , where each  $\mathcal{E}_i$  is a closed, convex set of probabilities. We assume that we can project onto each  $\mathcal{E}_i$  individually and shall denote the  $I$ -projection of  $s$  onto  $\mathcal{E}_i$  by  $\pi_i(s)$  and the  $I$ -projection of  $s$  onto  $\mathcal{E}$  by  $\pi(s)$ . We also assume there exists a  $t \in \mathcal{E}$  such that  $I(t \parallel r) < \infty$ . (Remember that multiplication and division of vectors will refer to the operations being performed coordinate-wise.)

**ALGORITHM.**

1. Let  $s_{11} = r$ , and let  $p_{11} = \pi_1(s_{11})$ . We then set  $s_{12} = p_{11} = r(p_{11}/s_{11})$ . (We note that if  $s_{11}(k) = 0$ , then so is  $p_{11}(k)$ . We take  $0/0$  to be 1.)
2. Let  $p_{12} = \pi_2(s_{12})$ . Set  $s_{13} = p_{12} = r(p_{11}/s_{11})(p_{12}/s_{12})$ .
3. Continue, until  $p_{1t} = \pi_t(s_{1t})$  where  $s_{1t} = p_{1,t-1} = r(p_{11}/s_{11}) \cdots (p_{1,t-1}/s_{1,t-1})$ . Set  $s_{21} = r(p_{12}/s_{12}) \cdots (p_{1t}/s_{1t})$ . Note that  $s_{21} = p_{1t}/(p_{11}/s_{11})$ .
4. Set  $p_{21} = \pi_1(s_{21})$ , and then set  $s_{22} = r(p_{21}/s_{21})(p_{13}/s_{13}) \cdots (p_{1t}/s_{1t})$ , or equivalently,  $s_{22} = p_{21}/(p_{12}/s_{12})$ .
5. Continue. In general, set  $s_{ni} = p_{n,i-1}/(p_{n-1,i}/s_{n-1,i})$ ,  $2 \leq i \leq t$  and set  $s_{n,1} = p_{n-1,t}/(p_{n-1,1}/s_{n-1,1})$ . We then let  $p_{ni} = \pi_i(s_{ni})$ , and define  $s_{n,i+1}$ , or  $s_{n+1,1}$  if  $i = t$ , in similar fashion.

Suppose now that the  $\mathcal{E}_i$  are actually linear sets so that equality holds in (1.7). Noting that for any  $p \in \mathcal{E}_i$

$$(2.1) \quad \begin{aligned} I(p \parallel s_{n,i}) &= \sum_k p(k)\ln[p(k)/s_{ni}(k)] \\ &= \sum_k p(k)\ln[p(k)/(p_{n,i-1}(k)/(p_{n-1,i}(k)/s_{n-1,i}(k)))] \end{aligned}$$

$$(2.2) \quad = \sum_k p(k)\ln[p(k)/p_{n,i-1}(k)] + \sum_k p(k)\ln[p_{n-1,i}(k)/s_{n-1,i}(k)],$$

we observe that the last term must be equal to

$$\sum_k p_{n-1,i}(k)\ln[p_{n-1,i}(k)/s_{n-1,i}(k)]$$

and hence free of  $p$ . Thus the  $p \in \mathcal{E}_i$  which minimizes (2.1), is also the one which minimizes the first part of (2.2), i.e., the  $I$ -projection of  $p_{n,i-1}$  onto  $\mathcal{E}_i$ . It easily follows that our procedure reduces to the cyclic, iterative procedure given by Csiszar when the  $\mathcal{E}_i$  are closed, linear sets.

For a simple example to show that Csiszar's procedure does not work for general convex sets, consider the following:

$$\mathcal{E}_1 = \left\{ \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}; p_{11} \geq p_{12}, p_{21} \geq p_{22}, p_{ij} \geq 0, \sum_1^2 \sum_1^2 p_{ij} = 1 \right\},$$

$$\mathcal{E}_2 = \left\{ \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}; p_{11} \geq p_{21}, p_{12} \geq p_{22}, p_{ij} \geq 0, \sum_1^2 \sum_1^2 p_{ij} = 1 \right\},$$

and

$$r = \begin{pmatrix} 1/16 & 3/16 \\ 7/16 & 5/16 \end{pmatrix}.$$

Csiszar's procedure yields  $(9/32 \ 7/32, \ 9/32 \ 7/32)$ , whereas  $(1/4 \ 1/4, \ 1/4 \ 1/4)$  is the correct solution.

Everything hinges on the following theorem, which retains the notation used in describing the algorithm.

**THEOREM 2.1.** *Assume  $\mathcal{E} = \cap_1^t \mathcal{E}_i$ , where the  $\mathcal{E}_i$  are closed, convex sets of probability distributions and  $r \neq 0$  is a nonnegative vector such that there exists a  $t \in \mathcal{E}$  where  $I(t \| r) < \infty$ . If there exists a convergent subsequence  $p_{n,i} \rightarrow p$  for some  $i$  such that*

$$(2.3) \quad \liminf_j \sum_k (p_{n_j,i}(k) - p(k)) \ln \left( \frac{p_{n_j,i}(k)}{s_{n_j,i}(k)} \right) \geq 0$$

for every  $i$ , then  $p_{n,i} \rightarrow p$  as  $n \rightarrow \infty$  and  $p = \pi(r)$ .

**PROOF.** Recall that  $p_{ni} = \pi_i(s_{ni})$  where

$$(2.4) \quad \begin{aligned} s_{ni} &= r \frac{p_{n1}}{s_{n1}} \dots \frac{p_{n,i-1}}{s_{n,i-1}} \frac{p_{n-1,i+1}}{s_{n-1,i+1}} \dots \frac{p_{n-1,t}}{s_{n-1,t}} \\ &= \begin{cases} p_{n-1,t} \left( \frac{p_{n-1,1}}{s_{n-1,1}} \right)^{-1} & \text{if } i = 1 \\ p_{n,i-1} \left( \frac{p_{n-1,i}}{s_{n-1,i}} \right)^{-1} & \text{if } 2 \leq i \leq t. \end{cases} \end{aligned}$$

Thus

$$\begin{aligned}
 (2.5) \quad & I(p_{ni} \parallel s_{ni}) - I(p_{n-1,i} \parallel s_{n-1,i}) \\
 &= \sum p_{ni} \ln \left( \frac{p_{ni}}{s_{ni}} \right) - \sum p_{n-1,i} \ln \left( \frac{p_{n-1,i}}{s_{n-1,i}} \right) \\
 &= \sum p_{ni} \ln \left( \frac{p_{ni}}{p_{n,i-1}} \right) + \sum p_{ni} \ln \left( \frac{p_{n-1,i}}{s_{n-1,i}} \right) - \sum p_{n-1,i} \ln \left( \frac{p_{n-1,i}}{s_{n-1,i}} \right) \\
 &= \sum p_{ni} \ln \left( \frac{p_{ni}}{p_{n,i-1}} \right) + \sum (p_{ni} - p_{n-1,i}) \ln \left( \frac{p_{n-1,i}}{s_{n-1,i}} \right) \\
 &\geq I(p_{ni} \parallel p_{n,i-1}) \quad \text{if } 2 \leq i \leq t, \quad \text{and similarly if } i = 1,
 \end{aligned}$$

since (by 1.7) the last term must be nonnegative because  $p_{ni} \in \mathcal{E}_i$ . Noting that  $I(p_{ni} \parallel p_{n,i-1}) \geq 0$  since  $p_{ni}$  and  $p_{n,i-1}$  are probabilities, we have that  $I(p_{ni} \parallel s_{ni})$  is nondecreasing in  $n$  for each  $i$ . Let us now show that these sequences are bounded above.

For  $v \in \cap_1^t \mathcal{E}_i$ ,

$$\begin{aligned}
 (2.6) \quad & I(v \parallel p_{ni}) = \sum v \ln \left( \frac{v}{p_{ni}} \right) \\
 &= \sum v \ln v - \sum v \ln r \left[ \frac{p_{n1}}{s_{n1}} \dots \frac{p_{ni}}{s_{ni}} \frac{p_{n-1,i+1}}{s_{n-1,i+1}} \dots \frac{p_{nt}}{s_{nt}} \right] \\
 &= \sum v \ln v - \sum v \ln r - \sum_{j=1}^t \left[ \sum v \ln \frac{p_{aj}}{s_{aj}} - I(p_{aj} \parallel s_{aj}) \right] \\
 &\quad - \sum_{j=1}^t I(p_{aj} \parallel s_{aj}) \quad \text{where } a = \begin{cases} n, & j \leq i \\ n-1, & j > i \end{cases} \\
 &\leq I(v \parallel r) - \sum_{j=1}^t I(p_{aj} \parallel s_{aj})
 \end{aligned}$$

by (1.6) and the fact that  $v$  belongs to every  $\mathcal{E}_i$ . Thus, choosing  $v$  such that  $I(v \parallel r) < \infty$ , we have a uniform upper bound on  $I(p_{ni} \parallel s_{ni})$ . Hence the  $\lim_{n \rightarrow \infty} I(p_{ni} \parallel s_{ni})$  exists finite for every  $i$ , and by (2.5)

$$(2.7) \quad I(p_{ni} \parallel p_{n,i-1}) \rightarrow 0 \quad \text{as } n \rightarrow \infty \quad \text{for } 2 \leq i \leq t.$$

(Similarly  $I(p_{n1} \parallel p_{n-1,t}) \rightarrow 0$  as  $n \rightarrow \infty$ .)

Note that (1.2) and (2.7) imply that if  $p_{n_j,i} \rightarrow p$  for some  $i$ , then  $p_{n_j,i} \rightarrow p$  for every  $i$ . Thus  $p \in \cap_1^t \mathcal{E}_i$  since  $p_{n_j,i} \in \mathcal{E}_i$  and the  $\mathcal{E}_i$  are closed.

Using (1.7) and (2.3), for any  $v \in \cap_1^t \mathcal{E}_i$ ,

$$\begin{aligned}
 & 0 \leq \sum_i \liminf_j \sum_k (p_{n_j,i} - p) \ln(p_{n_j,i}/s_{n_j,i}) + \sum_i \sum_k (v - p_{n_j,i}) \ln(p_{n_j,i}/s_{n_j,i}) \\
 & \leq \liminf_j \sum_i \sum_k (v - p) \ln(p_{n_j,i}/s_{n_j,i}) \\
 & = \liminf_j \sum_k (v - p) \ln r \prod_i (p_{n_j,i}/s_{n_j,i})/r \\
 & = \lim_j \sum_k (v - p) \ln(p_{n_j}, t/r) = \sum_k (v - p) \ln(p/r).
 \end{aligned}$$

Thus  $p = \pi(r)$  by (1.7). Setting  $v = p$  in (2.6) and using (2.3), it follows that for the subsequence  $\{n_j\}$ ,

$$\sum_{i=1}^t I(p_{n_j,i} \parallel s_{n_j,i}) \rightarrow I(p \parallel r).$$

However, by the monotonicity in  $n$  of  $I(p_{ni} \parallel s_{ni})$ , we have that

$$\sum_{i=1}^t I(p_{n,i} \parallel s_{n,i}) \rightarrow I(p \parallel r).$$

Thus  $I(p \parallel p_{nt}) \rightarrow 0$ , so that by (1.2) and (2.5),  $p_{ni} \rightarrow p$  as  $n \rightarrow \infty$  for all  $i$ .

We remark that it would take rather surprising behavior of the  $p_{ni}$  for condition (2.3) to not hold, and we strongly conjecture that this condition is always true. As we note in the following corollary, if the  $p_{n,i} \nrightarrow \pi(r)$  we must have  $\sup_{n,i} \sum_k s_{ni}(k) \rightarrow \infty$ , and we have been unable to construct examples where this happens. We point out that when one uses the algorithm, a step can be put in to check the value of  $\sup_{n,i} \sum_k s_{ni}(k)$ . If the algorithm is not going to converge correctly, then this value must become excessively large. Otherwise, the algorithm must converge to the correct solution.

**COROLLARY 2.1.** *If the algorithm should not converge correctly, then  $\sup_{n,i} \sum_k s_{ni}(k) \rightarrow \infty$  as  $n \rightarrow \infty$ .*

**PROOF.** Suppose  $s_{ni}(k)$  is uniformly bounded above, and  $p_{n_j,i} \rightarrow p$  is a convergent subsequence (which must exist since  $0 \leq p_{ni}(k) \leq 1$ ). If there exists  $0 < m \leq s_{n_j,i}(k)$  for all  $n_j, k$ , then (2.3) follows by standard continuity properties. If not, there exists at least one  $k$  such that  $s_{n_j,i}(k)$  takes on arbitrarily small values. Then either  $p_{n_j,i}(k) \rightarrow 0$  for all such  $k$  (which cause no problems), or there exists some  $k$  such that  $p_{n_j,i}(k) \rightarrow p(k) > 0$ , while  $s_{n_j,i}(k)$  takes on arbitrarily small values. This contradicts the uniform upper bound of  $I(p_{ni} \parallel s_{ni})$ .

We recommend that when one uses the algorithm, one should compute the average value over an entire cycle  $((1/t) \sum_{i=1}^t p_{n,i})$  rather than a single projection  $p_{ni}$  to estimate the  $I$ -projection. Convergence is still guaranteed, and this value is much more stable and seems to converge much more quickly to the correct solution.

**3. An example.** We consider the case where  $\mathcal{X}$  is an  $n \times n$  table. A probability on  $\mathcal{X}$  is represented by an  $n \times n$  matrix  $(m_{ij})$  of nonnegative numbers which sum to one. We denote the corresponding marginals by

$$m_{k.} = \sum_{j=1}^n m_{kj} \quad \text{and} \quad m_{.j} = \sum_{k=1}^n m_{kj}, \quad k, j = 1, \dots, n.$$

We now consider the problem of finding the  $I$ -projection of a fixed array  $(r_{kj})$  subject to the marginal PD's being stochastically ordered, i.e.,

$$\sum_{\ell=1}^i m_{\ell.} \geq \sum_{\ell=1}^i m_{. \ell}, \quad \text{for all } i.$$

(Kullback, 1971, has given an iterative procedure for  $I$ -projections where equality is forced to hold for all  $i$ , also known as marginal homogeneity.)

Equivalently, we want to find the  $I$ -projection of  $(r_{kj})$  onto

$$(3.1) \quad \mathcal{E} = \bigcap_{i=1}^{n-1} \mathcal{E}_i \quad \text{where} \quad \mathcal{E}_i = \{(m_{kj}); \sum_{\ell=1}^i \sum_{h=1}^n m_{\ell h} \geq \sum_{h=1}^i \sum_{\ell=1}^n m_{\ell h}\}.$$

Note that the  $\mathcal{E}_i$  are closed, convex sets of probabilities which are *not* linear sets.  $I$ -projections of  $(v_{kj})$  onto the  $\mathcal{E}_i$  can be found by forcing equality in the constraint if the  $(v_{kj})$  violate the constraint. (See Theorem 2.11 of Barlow et al., 1972, which can be modified to apply to arbitrary, convex functions).

To express  $\pi_i(v)$ , we let

$$A_i = \{(\ell, m); 1 \leq \ell \leq i, i + 1 \leq m \leq n\},$$

$$B_i = \{(\ell, m); 1 \leq m \leq i, i + 1 \leq \ell \leq n\},$$

$$C_i = \{(\ell, m); \ell, m = 1, \dots, i\} \cup \{(\ell, m); \ell, m = i + 1, \dots, n\}$$

and

$$\hat{v}_i = (\sum_{A_i} v_{\ell m} \sum_{B_i} v_{\ell m})^{1/2}.$$

If  $(v_{kj})$  satisfies the constraint of  $\mathcal{E}_i$ , namely  $\sum_{A_i} v_{\ell m} \geq \sum_{B_i} v_{\ell m}$ , then  $\pi_i(v)$  has  $(k, j)$ th component

$$\pi_i(v)_{kj} = v_{kj} / \sum_1^n \sum_1^n v_{\ell m} \quad \text{for all } k, j.$$

If  $(v_{kj})$  does not satisfy the constraint imposed by  $\mathcal{E}_i$ , that is  $\sum_{A_i} v_{\ell m} < \sum_{B_i} v_{\ell m}$ , then

$$\pi_i(v)_{kj} = \begin{cases} v_{kj} \left( \frac{\sum_{B_i} v_{\ell m}}{\sum_{A_i} v_{\ell m}} \right)^{1/2} [2\hat{v}_i + \sum_{C_i} v_{\ell m}]^{-1}, & (k, j) \in A_i, \\ v_{kj} \left( \frac{\sum_{A_i} v_{\ell m}}{\sum_{B_i} v_{\ell m}} \right)^{1/2} [2\hat{v}_i + \sum_{C_i} v_{\ell m}]^{-1}, & (k, j) \in B_i, \\ v_{kj} [2\hat{v}_i + \sum_{C_i} v_{\ell m}]^{-1}, & (k, j) \in C_i. \end{cases}$$

The key point is that finding the  $I$ -projection onto  $\mathcal{E}_i$  is quite easy (and easily programmed), while finding the  $I$ -projection onto  $\mathcal{E} = \bigcap^{n-1} \mathcal{E}_i$  is very difficult. However our algorithm enables one to find the latter  $I$ -projection using only the ability to handle the  $I$ -projections onto the individual  $\mathcal{E}_i$ .

To illustrate our example with some numbers, we consider some rather famous data from Stuart (1953) concerning grades of unaided distance vision for left and

TABLE 1  
Unaided distance vision (From Kendall, 1974)

Right eye	Left eye				Totals
	Highest grade	Second grade	Third grade	Lowest grade	
Highest grade	821	112	85	35	1053
Second grade	116	494	145	27	782
Third grade	72	151	583	87	893
Lowest grade	43	34	106	331	514
Totals	1052	791	919	480	3242



TABLE 2  
*I*-Projection of data in Table 1  
 (Values in parentheses are Table 1 values normed to sum to unity)

Right eye	Left eye				Totals
	Highest grade	Second grade	Third grade	Lowest grade	
Highest grade	.2534 (.2532)	.0344 (.0345)	.0262 (.0262)	.0120 (.0108)	.3260 (.3247)
Second grade	.0358 (.0358)	.1525 (.1524)	.0447 (.0447)	.0092 (.0083)	.2422 (.2412)
Third grade	.0222 (.0222)	.0466 (.0466)	.1799 (.1798)	.0298 (.0268)	.2785 (.2754)
Lowest grade	.0120 (.0133)	.0095 (.0105)	.0295 (.0327)	.1022 (.1021)	.1532 (.1586)
Totals	.3234 (.3245)	.2430 (.2440)	.2803 (.2834)	.1532 (.1480)	.9999 (.9999)

right eyes. If one wished to estimate the probabilities of falling into the various categories, subject to the provision that right eye vision is at least as good as left eye vision, one might find the *I*-projection of the data in Table 1 onto the  $\mathcal{E}$  given in (3.1). Using this algorithm, we have essentially obtained convergence to the true *I*-projection by 3 cycles. These values are listed in Table 2 (with the unrestricted MLEs given in parentheses).

This estimate might prove useful for constructing a likelihood ratio type test for testing whether right eye vision is better than left eye vision. This data is treated by Plackett (1981) who uses it for tests involving marginal homogeneity and quasi-symmetry. It also appears in Kendall (1974).

**4. Acknowledgement.** The author would like to thank a referee for helpful, knowledgeable comments and Peter Wollan for doing the calculations in Table 2.

## REFERENCES

- BARLOW, R. E, BARTHOLOMEW, D. J., BREMNER, J. M. and BRUNK, H. D. (1972). *Statistical Inference Under Order Restrictions*. Wiley, New York.
- CENCOV, N. N. (1972). *Statistical Decision Rules and Optimal Inference*. Translations of Mathematical Monographs, American Mathematical Society, Vol. 53, 115–125.
- CSISZAR, I. (1967). Information-type measures of difference of probability distributions and indirect observations. *Studia Sci. Math. Hungar.* **2** 299–318.
- CSISZAR, I. (1975). *I*-divergence geometry of probability distributions and minimization problems. *Ann. Probab.* **3** 146–159.
- GOOD, I. J. (1963). Maximum entropy for hypothesis formulation, especially for multidimensional contingency tables. *Ann. Math. Statist.* **34** 911–934.
- JAYNES, E. T. (1957). Information theory and statistical mechanics. *Phys. Rev.* **106** 620–630.
- JEFFREYS, H. (1948). *Theory of Probability*. 2nd ed. Clarendon Press, Oxford.
- KEMPERMAN, J. H. B. (1967). On the optimum rate of transmitting information. Probability and Information Theory. *Lecture Notes in Mathematics*, 126–169. Springer-Verlag, New York.

- KENDALL, M. G. (1974). *Rank Correlation Methods*. Griffin, London.
- KULLBACK, S. (1959). *Information Theory and Statistics*. Wiley, New York.
- KULLBACK, S. (1967). A lower bound for discrimination information in terms of variation. *IEEE Trans. Informa. Theory* **IT-13** 126-127.
- KULLBACK, S. (1971). Marginal homogeneity of multidimensional contingency tables. *Ann. Math. Statist.* **42** 594-606.
- LEMKE, J. H. and DYKSTRA, R. L. (1984). Multinomial maximum likelihood estimation with multiple cone restrictions (submitted, *Ann. Statist.*).
- MEYER, M. M. (1980). Generalizing the iterative proportional fitting procedure. Dept. of Applied Statistics, Tech. Report No. 371, Univ. of Minnesota.
- PLACKETT, R. L. (1981). *The Analysis of Categorical Data*. Macmillan, New York.
- RAO, C. R. (1965). *Linear Statistical Inference and its Applications*. Wiley, New York.
- SANOV, I. N. (1957). On the probability of large deviations of random variables. *Mat. Sb.* **42** 11-44.
- STUART, A. (1953). The estimation and comparison of strengths of association in contingency tables. *Biometrika* **40** 105-110.

DEPARTMENT OF STATISTICS  
AND ACTUARIAL SCIENCES  
THE UNIVERSITY OF IOWA  
IOWA CITY, IOWA 52242