

## A FENCHEL DUALITY ASPECT OF ITERATIVE I-PROJECTION PROCEDURES

BHASKAR BHATTACHARYA<sup>1</sup> AND RICHARD L. DYKSTRA<sup>2\*</sup>

<sup>1</sup>*Department of Mathematics, Southern Illinois University,  
Carbondale, IL 62901-4408, U.S.A.*

<sup>2</sup>*Department of Statistics and Actuarial Science, University of Iowa,  
Iowa City, IA 52242, U.S.A.*

(Received May 17, 1995; revised October 3, 1996)

**Abstract.** In this paper we interpret Dykstra's iterative procedure for finding an I projection onto the intersection of closed, convex sets in terms of its Fenchel dual. Seen in terms of its dual formulation, Dykstra's algorithm is intuitive and can be shown to converge monotonically to the correct solution. Moreover, we show that it is possible to sharply bound the location of the constrained optimal solution.

*Key words and phrases:* Algorithm, convex sets, Fenchel duality, I-projections, iterative, Kullback-Leibler.

### 1. Introduction

For finite probability vectors (PV's)  $\mathbf{p} = (p(1), \dots, p(m))'$  and  $\mathbf{q} = (q(1), \dots, q(m))'$  of length  $m$ , the *I-divergence* of  $\mathbf{p}$  with respect to  $\mathbf{q}$  (also known as the *Kullback-Leibler information number*) is given by

$$I(\mathbf{p} | \mathbf{q}) = \sum_{k=1}^m p(k) \ln \left( \frac{p(k)}{q(k)} \right).$$

Since  $I(\mathbf{p} | \mathbf{q}) \geq 0$ , and equals zero if and only if  $\mathbf{p} = \mathbf{q}$ ,  $I(\mathbf{p} | \mathbf{q})$  is often treated heuristically as a measure of distance or divergence between  $\mathbf{p}$  and  $\mathbf{q}$ . It is natural to consider the "closest" PV to  $\mathbf{q}$  which lies within a specified set of PV's  $\mathcal{C}$ . A PV  $\mathbf{u} \in \mathcal{C}$  that satisfies

$$(1.1) \quad I(\mathbf{u} | \mathbf{q}) = \min_{\mathbf{p} \in \mathcal{C}} I(\mathbf{p} | \mathbf{q}) < \infty$$

is said to be an *I-projection* of  $\mathbf{q}$  onto  $\mathcal{C}$ . The I-projection  $\mathbf{u}$  always exists uniquely if  $\mathcal{C}$  is closed and convex (Csiszar (1975)).

---

\* Partial support was provided by National Science Foundation Grant DMS 91-04673.

I-projections play a basic role in the information theoretic approach to statistics (Kullback (1959), Good (1963)). They are also important in the theory of large deviations (Sanov (1957)) and in statistical physics for the maximization of entropy (Rao (1965), Jaynes (1957)). For a duality approach to I-projections for general probability distributions, see Bhattacharya and Dykstra (1995).

Depending on the form of the set  $\mathcal{C}$ , it may be difficult to find a solution to the I-projection problem in (1.1). Csiszar (1975) has shown that if  $\mathcal{C}$  can be expressed as  $\bigcap_{i=1}^t \mathcal{C}_i$ , where each  $\mathcal{C}_i$  is a closed, linear set, then the sequence of cyclic iterated I-projections converges to the solution of (1.1). Dykstra (1985a) modified Csiszar's procedure to encompass the case where each  $\mathcal{C}_i$  is an arbitrary closed, convex set. He showed that the desired I-projection can be obtained as the limit of cyclic I-projections onto the  $\mathcal{C}_i$  if the projected vectors are appropriately modified and a mild condition holds. Winkler (1990) showed that the stated condition always holds, and hence the algorithm is always valid. In this paper, we use an extension to a Fenchel duality theorem for  $\mathbb{R}^m$  to prove Dykstra's result in a much shorter and more intuitive fashion.

## 2. Main results

For a closed, proper, convex function  $f$  on  $\mathbb{R}^m$ , the *convex conjugate*  $f^*$  defined on  $\mathbb{R}^m$  is given by

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} \left\{ \sum_{k=1}^m x(k)y(k) - f(\mathbf{x}) \right\}.$$

It is well known (Rockafellar (1970)) that  $f^*$  is also a closed, proper, convex function on  $\mathbb{R}^m$  and that  $f^{**} = f$ . For a convex cone  $K$  in  $\mathbb{R}^m$ , the *dual cone*  $K^*$  in  $\mathbb{R}^m$  is defined as

$$K^* = \left\{ \mathbf{y} \in \mathbb{R}^m : \sum_{k=1}^m x(k)y(k) \leq 0, \forall \mathbf{x} \in K \right\}.$$

Corollary 2.1, from Rockafellar ((1970), p. 335), will be used to identify the dual problem. By "ri" we mean relative interior as defined in Rockafellar (1970)

**COROLLARY 2.1.** *Let  $f$  be a closed, proper, convex function on  $\mathbb{R}^m$ , and let  $K$  be a nonempty, closed, convex cone in  $\mathbb{R}^m$ . Then*

$$(2.1) \quad \text{(I)} \quad \inf_{\mathbf{x} \in K} f(\mathbf{x}) = \sup_{\mathbf{y} \in K^*} -f^*(-\mathbf{y}) \quad \text{(II)}$$

*provided  $\text{ri}(\text{dom } f) \cap \text{ri}(K) \neq \emptyset$ , and moreover, the supremum on the right side is attained. In general,  $\hat{\mathbf{x}}$  and  $\hat{\mathbf{y}}$  are solutions to these problems which satisfy*

$$f(\hat{\mathbf{x}}) = \inf_{\mathbf{x} \in K} f(\mathbf{x}) = \sup_{\mathbf{y} \in K^*} -f^*(-\mathbf{y}) = -f^*(-\hat{\mathbf{y}})$$

if and only if

- (i)  $-\hat{\mathbf{y}}$  is a subgradient of  $f$  at  $\hat{\mathbf{x}}$ ,
- (ii)  $\hat{\mathbf{x}} \in K$ ,
- (iii)  $\hat{\mathbf{y}} \in K^*$ , and
- (iv)  $\sum_{k=1}^m \hat{x}(k)\hat{y}(k) = 0$ .

For a given convex set  $\mathcal{C}$  of PV's, the I-projection problem (1.1) can be expressed as the left side of (2.1), if we define

$$(2.2) \quad f(\mathbf{x}) = \begin{cases} \sum_{k=1}^m x(k) \ln \left( \frac{x(k)}{q(k)} \right), & \text{if } x(k) > 0, \forall k, \sum_{k=1}^m x(k) = 1 \\ +\infty, & \text{otherwise,} \end{cases}$$

and  $K$  by

$$K = \{\alpha \mathbf{x} : \mathbf{x} \in \mathcal{C}, \alpha \geq 0\}.$$

Dykstra (1985*b*) has shown that the convex conjugate of the function in (2.2) is given by

$$f^*(\mathbf{y}) = \ln \left( \sum_{k=1}^m q(k) \exp(y(k)) \right)$$

( $\mathbf{q}$  need not sum to 1). Note that  $f^*$  gives the cumulant generating function of  $\mathbf{q}$  (with probability on a particular set of values  $z(1), z(2), \dots, z(m)$ ) if  $f^*$  is evaluated along the one-dimensional path  $\mathbf{y} \cdot (z(1), z(2), \dots, z(m))'$ . Thus minimizing  $f^*$  over the (one-dimensional) region  $K^* = \{\mathbf{v} \in \mathbb{R}^m : \mathbf{v} = \alpha \mathbf{z}, \alpha \geq 0\}$  is equivalent to minimizing  $I(\mathbf{p} | \mathbf{q})$  over the set  $-K = \{\mathbf{x} \in \mathbb{R}^m : \sum_{i=1}^m x(i)z(i) \geq 0\}$ , since  $K^*$  is the dual of  $K$ . Restricting  $-K$  to the PV's, the duality theorem then shows that minimizing the cumulant generating function over nonnegative values of  $\mathbf{y}$  is equivalent to finding the I-projection of  $\mathbf{q}$  onto the set of distributions over  $(z(1), z(2), \dots, z(m))$  with a nonnegative mean.

Dykstra (1985*b*) has also shown that if  $\hat{\mathbf{y}}$  solves the dual problem (II) in (2.1), then the vector  $\hat{\mathbf{p}} = (\hat{p}(1), \hat{p}(2), \dots, \hat{p}(m))'$  given by

$$(2.3) \quad \hat{p}(k) = \frac{q(k) \exp(-\hat{y}(k))}{\sum_{s=1}^m q(s) \exp(-\hat{y}(s))}, \quad k = 1, 2, \dots, m,$$

solves the I-projection problem (1.1) if  $\mathcal{C}$  is closed and  $\text{ri}(\text{dom } f) \cap \text{ri}(K) \neq \emptyset$ .

In some problems  $\mathcal{C}$  may constrain some of the  $x(k)$  to always be zero, which causes the condition  $\text{ri}(\text{dom } f) \cap \text{ri}(K) \neq \emptyset$  to be violated. It can be shown that the dual formulation above is still valid if the domain of  $f^*$  is expanded to be  $\mathbb{R}_1^* \times \mathbb{R}_2^* \times \dots \times \mathbb{R}_m^*$ , where  $\mathbb{R}_i^* = \mathbb{R} \cup \{\pm\infty\}$ . The results hold as before with the interpretation that  $e^{-\infty} = 0$  and  $0 \cdot \infty = 0$ .

We now let  $\mathbf{q}$  be a fixed PV of length  $m$ . We wish to find the I-projection of  $\mathbf{q}$  onto the nonempty set  $\mathcal{C} = \bigcap_{i=1}^t \mathcal{C}_i$ , where each  $\mathcal{C}_i$  is a closed, convex set of PV's of length  $m$ . We only assume that there exists an  $\mathbf{r} \in \mathcal{C}$  such that  $I(\mathbf{r} \mid \mathbf{q}) < \infty$ . Winkler (1990) has shown that each of the I-projections required for the algorithm exists. We can use the Fenchel duality theorem in each step of the algorithm to identify a corresponding dual problem. Our point is that matters are much simpler and more intuitive in the dual formulation.

We present the algorithm below incorporating the primal and the dual formulations separately. Multiplication and division of vectors is done coordinatewise. To implement the algorithm for the primal problem at the  $i$ -th cycle and  $j$ -th step, we first form  $\mathbf{s}_{ij}$ , the vector to be projected; the I-projection obtained by projecting  $\mathbf{s}_{ij}$  onto  $\mathcal{C}_j$  is denoted by  $\mathbf{p}_{ij}$  (although  $\mathbf{s}_{ij}$  may not be a PV, we can define the I-projection  $\mathbf{p}_{ij}$  of  $\mathbf{s}_{ij}$  in the same way, for further details see Dykstra (1985a)); and the solution to the corresponding dual problem is  $\mathbf{y}_{ij}$  (used in the dual formulation of the algorithm). We note that it is straightforward to show that at the  $i$ -th cycle,  $j$ -th step,  $p_{ij}(k) = 0$  if and only if  $s_{ij}(k) = 0$  or there exists  $v$  such that  $p(k) = 0$ , for all  $\mathbf{p} \in \mathcal{C}_v$  (Winkler (1990)).

We first state algorithm in the primal form (Dykstra (1985a)) and then reformulate it in terms of the much more intuitive dual problem. The key point is that the intersection constraints in the primal problem translate to direct sum constraints in the dual problem. This is what allows the dual formulation of the algorithm to be phrased as a cyclic, descent algorithm which successively minimizes over one vector at a time while the rest are held fixed. We use the duality structure to give a short proof that the algorithm must work correctly. It will be convenient to let  $\pi_i(\mathbf{s})$  denote the I-projection of the vector  $\mathbf{s}$  onto the set  $\mathcal{C}_i$ .

*Primal formulation of the algorithm.*

- (1) *Initialization.* Set  $\mathbf{s}_{0,i} = \mathbf{p}_{0,i} = \mathbf{q}$ , and begin with  $n = 1$ ,  $i = 1$ .
- (2) *Implementation.*
  - (i) For  $i = 1$ , set  $\mathbf{s}_{n,1} = \mathbf{p}_{n-1,t}/(\mathbf{p}_{n-1,1}/\mathbf{s}_{n-1,1})$ ; for  $2 \leq i \leq t$ , set  $\mathbf{s}_{n,i} = \mathbf{p}_{n,i-1}/(\mathbf{p}_{n-1,i}/\mathbf{s}_{n-1,i})$ ; (we assume that  $0/0 = 0$ ).
  - (ii) Let  $\mathbf{p}_{n,i} = \pi_i(\mathbf{s}_{n,i})$ .
  - (iii) If  $i < t$ , increment  $i$  by 1 and repeat (2). If  $i = t$ , increment  $n$  by 1, set  $i = 1$  and repeat (2).

Of course, the key point is that  $\mathbf{p}_{n,i}$  must converge to the I-projection of  $\mathbf{q}$  onto  $\mathcal{C} = \bigcap_{i=1}^t \mathcal{C}_i$ . However, the proof (Dykstra (1985a)) is quite complicated and involved. It is much cleaner in the dual formulation.

The dual problem is equivalent to

$$\begin{aligned} \inf_{\mathbf{y} \in (\bigcap_{i=1}^t K_i)^*} \sum_{k=1}^m q(k) e^{-y(k)} &= \inf_{\mathbf{y} \in \text{cl}(K_1^* + \dots + K_t^*)} \sum_{k=1}^m q(k) e^{-y(k)} \\ &= \inf_{\mathbf{y}_i \in K_i^*, 1 \leq i \leq t} \sum_{k=1}^m q(k) e^{-y_1(k) - \dots - y_t(k)}. \end{aligned}$$

Typically the dual constraint region would be the direct sum  $K_1^* + \dots + K_t^*$ . However, the direct sum of closed sets need not be closed (Hestenes (1975)) and

hence our dual constraint region is the closure of the direct sum. This possible lack of closure of the direct sum complicates the proof substantially.

One of the drawbacks to iterative methods of optimization are the difficulties in specifying the distance of the current estimate from the actual solution. Small changes in the objective function between successive steps of the iterative procedure are no guarantee that the actual optimal value is close by.

However, consider the situation where  $\mathbf{p}^*$  is a feasible PV (i.e.  $\mathbf{p}^* \in \bigcap_1^t K_i$ ) which is close to  $\mathbf{p}_{n,t}$ . We let  $\hat{\mathbf{p}}$  denote the actual, constrained solution to our problem. We note that  $\mathbf{p}^*$  may be obtained by a least square projection onto a subset of  $\bigcap_{i=1}^t K_i$ , by a cyclic descent method if the constraint region can also be written as a direct sum of convex sets, or by some other procedure. However, since

$$(2.4) \quad \sum_{k=1}^m \hat{p}(k) \ln \left( \frac{\hat{p}(k)}{q(k)} \right) \geq \sum_{i=1}^t \sum_{k=1}^m p_{n,i}(k) \ln \left( \frac{p_{n,i}(k)}{s_{n,i}(k)} \right),$$

with the difference between the two sides of (2.4) converging monotonically to zero as  $n \rightarrow \infty$ , we can estimate  $I(\hat{\mathbf{p}} \mid \mathbf{q})$  to arbitrary accuracy by choosing  $n$  sufficiently large. Moreover, since

$$(2.5) \quad I(\mathbf{p}^* \mid \mathbf{q}) - \sum_{i=1}^t I(\mathbf{p}_{n,i} \mid \mathbf{s}_{n,i}) \geq I(\mathbf{p}^* \mid \mathbf{q}) - I(\hat{\mathbf{p}} \mid \mathbf{q}) \\ \geq I(\mathbf{p}^* \mid \hat{\mathbf{p}}) \geq \frac{1}{2} \left( \sum_{k=1}^m |p^*(k) - \hat{p}(k)| \right)^2$$

(Csiszar (1975)), we can specify an upper bound on both the I-divergence distance and the variation distance ( $L_1$ -norm) between the vector  $\mathbf{p}^*$  and the true solution  $\hat{\mathbf{p}}$ .

In similar fashion, we can expand  $I(\mathbf{p}^* \mid \hat{\mathbf{p}})$  using a Taylor series expansion to obtain the bound

$$(2.6) \quad I(\mathbf{p}^* \mid \hat{\mathbf{p}}) \geq \frac{1}{2} \sum_{k=1}^m \left( 1 - \frac{p^*(k) \wedge \hat{p}(k)}{p^*(k) \vee \hat{p}(k)} \right)^2 p^*(k)$$

where “ $\wedge$ ” (“ $\vee$ ”) indicates infimum (supremum). Thus, for example, if  $I(\mathbf{p}^* \mid \hat{\mathbf{p}})$  is bounded above by  $\epsilon$ , it is straightforward to show that for each  $i$ ,

$$a_i p^*(i) < \hat{p}(i) < \frac{p^*(i)}{a_i}$$

where  $a_i = 1 - \sqrt{2\epsilon/p^*(i)}$ . For  $\epsilon$  sufficiently small, this gives a bound on how far  $\hat{\mathbf{p}}$  can fall from the vector  $\mathbf{p}^*$ . Of course one can also find an approximation region for possible values of  $\hat{\mathbf{p}}$  by numerically checking whether candidate  $\mathbf{p}$ 's (plugged in for  $\hat{\mathbf{p}}$ ) satisfy (2.5) and (2.6). However this becomes rather intractable for large  $m$ .

The efficiency of the algorithm depends largely on the nature of the  $K_i$ 's. If all the  $K_i$ 's are orthogonal with each other (in terms of I-divergence) a single pass through each constraint will suffice.

The algorithm stated in terms of the individual dual problems is just the following simple, cyclic, descent procedure. In the dual formulation, the algorithm amounts to just sequentially minimizing the objective function with all  $\mathbf{y}$ 's but one held fixed, and then updating that  $\mathbf{y}$ .

*Dual formulation of the algorithm.*

- (1) *Initialization.* Set  $\mathbf{y}_{0,i} \equiv \mathbf{0}$ , and begin with  $n = 1$ ,  $i = 1$ .
- (2) *Implementation.*
  - (i) Let  $\mathbf{y}_{n,i}$  denote the solution to

$$\inf_{\mathbf{y} \in K_i^*} \sum_{k=1}^m q(k) e^{-y_{n,1}(k) - \cdots - y_{n,i-1}(k) - y(k) - y_{n-1,i+1}(k) - \cdots - y_{n-1,t}(k)}.$$

- (ii) If  $i < t$ , increment  $i$  by 1 and repeat (2). If  $i = t$ , increment  $n$  by 1, set  $i = 1$  and repeat (2).

Careful inspection of the two algorithms together with (2.3) will reveal that

$$p_{n,i}(k) = \frac{q(k) e^{-y_{n,1}(k) - \cdots - y_{n,i}(k) - y_{n-1,i+1}(k) - \cdots - y_{n-1,t}(k)}}{\sum_{s=1}^m q(s) e^{-y_{n,1}(s) - \cdots - y_{n,i}(s) - y_{n-1,i+1}(s) - \cdots - y_{n-1,t}(s)}}$$

for  $1 \leq k \leq m$ .

If it should happen that

$$\mathbf{y}_{n,1} + \cdots + \mathbf{y}_{n,i} + \mathbf{y}_{n-1,i+1} + \cdots + \mathbf{y}_{n-1,t} \rightarrow \hat{\mathbf{y}}$$

for some vector  $\hat{\mathbf{y}}$  as  $n \rightarrow \infty$ , then it would easily follow by continuity that

$$p_{n,i}(k) \rightarrow \frac{q(k) e^{-\hat{y}(k)}}{\sum_{s=1}^m q(s) e^{-\hat{y}(s)}} = \hat{p}(k), \quad k = 1, \dots, m,$$

as  $n \rightarrow \infty$ . Moreover, if  $\hat{\mathbf{p}}$  and  $\hat{\mathbf{y}}$  should satisfy conditions (i)–(iv) of Corollary 2.1, then  $\hat{\mathbf{p}}$  and  $\hat{\mathbf{y}}$  would have to solve (1.1) and its dual, respectively, which would establish the validity of both the primal and dual forms of the algorithm. That this is indeed the case is the gist of the following theorem.

**THEOREM 2.1.** *Let the sets  $C_i$ ,  $1 \leq i \leq t$ , be closed, convex sets of PV's of length  $m$  and let  $\mathbf{q}$  be a positive PV such that there exists  $\mathbf{r} \in C = \bigcap_{i=1}^t C_i$  with  $I(\mathbf{r} | \mathbf{q}) < \infty$ . Then  $\mathbf{p}_{n,i} \rightarrow \hat{\mathbf{p}}$  as  $n \rightarrow \infty$ , for  $1 < i < t$ , where  $\hat{\mathbf{p}}$  is the unique I-projection of  $\mathbf{q}$  onto  $C$  and the  $\mathbf{p}_{n,i}$  are obtained as in the algorithm.*

**PROOF.** It is easily shown that

$$\sum_{k=1}^m q(k) e^{-y_{n,1}(k) - \cdots - y_{n,i}(k) - y_{n-1,i+1}(k) - \cdots - y_{n-1,t}(k)}$$

is nonincreasing in  $n$  and  $i$ . Moreover, it easily follows that every subsequence of  $\{n\}_{n=1}^\infty$  has a further subsequence  $\{n_j\}_{j=1}^\infty$  such that

$$y_{n_j,1} + \dots + y_{n_j,t} \rightarrow \hat{y}$$

where  $\hat{y}$  may have coordinates that are  $+\infty$ . Clearly,  $\hat{y} \in cl(K_1^* + \dots + K_t^*) = (\bigcap_{i=1}^t K_i)^* = K^*$ .

It easily follows by continuity that

$$\begin{aligned} p_{n_j,t}(k) &= \frac{q(k)e^{-y_{n_j,1}(k) - \dots - y_{n_j,t}(k)}}{\sum_{s=1}^m q(s)e^{-y_{n_j,1}(s) - \dots - y_{n_j,t}(s)}} \\ &\rightarrow \frac{q(k)e^{-\hat{y}(k)}}{\sum_{s=1}^m q(s)e^{-\hat{y}(s)}} = \hat{p}(k), \quad k = 1, \dots, m. \end{aligned}$$

Dykstra (1985a) has shown that  $I(p_{n,i} \mid s_{n,i}) - I(p_{n-1,i} \mid s_{n-1,i}) \geq I(p_{n,i} \mid p_{n,i-1}) \rightarrow 0$ . Then using the inequality that

$$\sum_{s=1}^m |q(s) - r(s)| \leq [2I(q \mid r)]^{1/2}$$

for any two vectors  $q, r$ , it follows easily that  $p_{n_j,i}$  is arbitrarily close to  $p_{n_j,t}$  if  $j$  is sufficiently large, for any  $i$ . Then it must be the case that  $\hat{p} \in \bigcap_{i=1}^t K_i = K$ .

It follows from Dykstra ((1985b), Corollary 2.1) that  $-\hat{y}$  is a subgradient of the  $f$  from (2.2) at  $\hat{p}$ .

Thus if,

$$(2.7) \quad \sum_{k=1}^m \hat{y}(k)\hat{p}(k) = 0,$$

the four conditions of Corollary 2.1 will be met and  $\hat{p}$  and  $\hat{y}$  will be the solutions to the primal and dual problems. Since every subsequence of  $\{p_{n,i}\}_{n=1}^\infty$  has a sub-subsequence which converges to  $\hat{p}$ , the entire sequence must converge to  $\hat{p}$ .

To establish (2.7), recall that

$$s_{n,i} = p_{n,i-1} / (p_{n-1,i} / s_{n-1,i})$$

(for  $2 \leq i \leq t$ ). It follows that (suppressing the index of summation  $k$ ) for  $2 \leq i \leq t$ ,

$$\begin{aligned} I(p_{n_j,i} \mid s_{n_j,i}) &= \sum_k p_{n_j,i} \ln(p_{n_j,i} / s_{n_j,i}) \\ &= \sum_k p_{n_j,i} \ln(p_{n_j,i} / p_{n_j,i-1}) + \sum_k p_{n_j,i} \ln(p_{n_j-1,i} / s_{n_j-1,i}) \\ &= \sum_{m=n_h+1}^{n_j} \sum_k p_{n_j,i} \ln(p_{m,i} / p_{m,i-1}) \\ &\quad + \sum_k (p_{n_j,i} - p_{n_h,i}) \ln(p_{n_h,i} / s_{n_h,i}) \\ &\quad + \sum_k p_{n_h,i} \ln(p_{n_h,i} / s_{n_h,i}). \end{aligned}$$

A similar expression can be derived when  $i = 1$ .

Dykstra (1985a) has shown that the left side is nondecreasing in  $j$  and bounded above. Then, since the left side limit exists as  $j \rightarrow \infty$ , as does the limit of the last two summations on the right side, the limit of the first right side summation must also exist as  $j \rightarrow \infty$ .

Then summing over  $i$  and letting  $h \rightarrow \infty$ , we obtain

$$\begin{aligned} 0 &= \lim_{h \rightarrow \infty} \sum_{i=1}^t \sum_k (\mathbf{p}_{n_h, i} - \hat{\mathbf{p}}) \ln(\mathbf{p}_{n_h, i} / \mathbf{s}_{n_h, i}) \\ &= \lim_{h \rightarrow \infty} \sum_{i=1}^t \sum_k (\mathbf{p}_{n_h, i} - \hat{\mathbf{p}}) (-\mathbf{y}_{n_h, i}) \\ &= \lim_{h \rightarrow \infty} \sum_{i=1}^t \sum_k \hat{\mathbf{p}} \mathbf{y}_{n_h, i} \left( \text{since } \sum_k \mathbf{p}_{n, i} \mathbf{y}_{n, i} = 0, \forall n, i \right) \\ &= \lim_{h \rightarrow \infty} \sum_k \hat{\mathbf{p}} \sum_{i=1}^t \mathbf{y}_{n_h, i} \\ &= \sum_k \hat{\mathbf{p}} \hat{\mathbf{y}} \end{aligned}$$

which is condition (2.7).  $\square$

Csiszar (1975) showed that sequentially projecting the I-projections onto the individual constraint regions gives a sequence of probability vectors that converges to the true I-projection of  $\mathbf{q}$  if all the  $\mathcal{C}_i$  are linear sets (i.e.,  $p_1, p_2 \in \mathcal{C}_i \Rightarrow \alpha p_1 + (1 - \alpha)p_2 \in \mathcal{C}_i$ , for all  $\alpha$  for which  $\alpha p_1 + (1 - \alpha)p_2$  is a PV). However, this algorithm will typically not work if the  $\mathcal{C}_i$  are not linear sets and the reason is clear from the dual formulation. If the  $\mathcal{C}_i$  are linear sets, then each projection (in the dual formulation) allows each point in  $K_i^*$  as a possible value of  $\mathbf{y}_{ni}$  (thus all of  $K_i^*$  is feasible) since  $K_i^* = \{\mathbf{y} + \mathbf{z} : \mathbf{y} \in K_i^*\}$  for every  $\mathbf{z} \in K_i^*$ . This won't be true, however, if the  $\mathcal{C}_i$ 's are not linear sets. Although Csiszar's procedure is simpler in the primal formulation, it is more complex in the dual form (and, of course, only works if the constraint regions are linear sets).

### 3. Example

We consider the problem of finding the maximum likelihood estimates of the probabilities for a two-way classification with multinomial sampling under the constraints that the local odds ratios ( $\theta_{ij} = p_{ij}p_{i+1, j+1}/p_{i+1, j}p_{i, j+1}$ ) are all at least 1 without any additional model assumptions. As is well known  $\theta_{ij} = 1, \forall i, j$  implies the two classifications are stochastically independent and  $\theta_{ij} \geq 1, \forall i, j$  implies there is a positive association between the ordinal variables.

The problem we consider can be expressed as

$$(3.1) \quad \sup_{\mathbf{p} \in \prod_{i=1}^c \prod_{j=1}^r \mathcal{K}_{ij}} \prod_{j=1}^c \prod_{i=1}^r p_{ij}^{n_{ij}},$$



where  $\mathbf{p} \in P$  (the class of all PV's of length  $rc$ ), the  $n_{ij}$ 's are the observed frequencies and the closed convex cones  $K_{ij}$  are defined as

$$K_{ij} = \{\mathbf{x} : x_{ij} - x_{i,j+1} - x_{i+1,j} + x_{i+1,j+1} \geq 0\}$$

for  $i = 1, \dots, r - 1, j = 1, \dots, c - 1$ . Let  $\mathbf{g}$  denote the  $(r \times c)$  matrix of the empirical distribution. Dykstra and Lemke (1988) have shown that (3.1) has the same solution as the I-projection problem

$$(3.2) \quad \inf_{\mathbf{p} \in \bigcap_{i=1}^r \bigcap_{j=1}^c (g - H_{ij})} I(\mathbf{p} \mid \mathbf{q})$$

where

$$H_{ij} = \left\{ \mathbf{x} : \sum_{k=1}^i \sum_{m=1}^j x_{km} \leq 0 \right\}$$

when  $i = 1, \dots, r - 1, j = 1, \dots, c - 1$ , and

$$H_{ij} = \left\{ \mathbf{x} : \sum_{k=1}^i \sum_{m=1}^j x_{km} = 0 \right\}$$

when  $i = r$  or  $j = c$ , and  $\mathbf{q}$  is the uniform PV which has values  $q_{ij} = (rc)^{-1}$ . In words, this amounts to the sum of every upper-left corner of cell  $(i, j)$  of  $\mathbf{p}$  being as large as the same of  $\mathbf{g}$  and equality holds whenever  $i = r$  or  $j = c$ . The  $(i, j)$ -th constraint ( $1 \leq i \leq r - 1, 1 \leq j \leq c - 1$ ) can be expressed as

$$\mathbf{g} - H_{ij} = \{\mathbf{p} : \mathbf{p} \text{ is a PV and } \mathbf{a}'_{ij}\mathbf{p} \leq 0\}$$

where the  $(k, m)$ -th term of the vector  $\mathbf{a}_{ij}$  is given by

$$\sum_{s=1}^i \sum_{t=1}^j g_{st} - I(1 \leq k \leq i, 1 \leq m \leq j),$$

where  $I$  is the indicator function. The dual cone  $(\mathbf{g} - H_{ij})^*$  is then given by

$$(\mathbf{g} - H_{ij})^* = \{\alpha \mathbf{a}_{ij} : \alpha \geq 0\}$$

when  $i = 1, \dots, r - 1, j = 1, \dots, c - 1$ , and

$$(\mathbf{g} - H_{ij})^* = \{\alpha \mathbf{a}_{ij} : \alpha \in \mathbb{R}\}$$

when  $i = r$  or  $j = c$ .

Solving (3.1) is thus equivalent to solving the problem

$$\inf \sum_k \sum_m e^{-\sum_{i=1}^r \sum_{j=1}^c \alpha_{ij} a_{ij}(k,m)}$$

Table 1. Cross-classification of job satisfaction by income.

Income (US\$)	Job Satisfaction			
	VD	LD	MS	VS
<6000	20	24	80	82
6000–15,000	22	38	104	125
15,000–25,000	13	28	81	113
>25,000	7	18	54	92

(VD – Very Dissatisfied; LD – Little Dissatisfied; MS – Moderately Satisfied; VS = Very Satisfied).

Table 2. Maximum likelihood estimates of data in Table 1 subject to the constraints  $\theta_{ij} \geq 1$   $\forall i, j$  (values in first parentheses are the unrestricted male and values in second parentheses are the expected counts after smoothing).

Income (US\$)	Job Satisfaction			
	VD	LD	MS	VS
<6000	0.0222	0.0293	0.0862	0.0910
	(0.0222)	(0.0266)	(0.0888)	(0.0910)
	(20.00)	(26.37)	(77.63)	(82.00)
6000–15,000	0.0244	0.0400	0.1176	0.1387
	(0.0244)	(0.0422)	(0.1154)	(0.1387)
	(22.00)	(36.00)	(106.00)	(125.00)
15,000–25,000	0.0144	0.0307	0.0903	0.1254
	(0.0144)	(0.0311)	(0.0899)	(0.1254)
	(13.00)	(27.63)	(81.37)	(113.00)
>25,000	0.0078	0.0200	0.0599	0.1021
	(0.0078)	(0.0200)	(0.0599)	(0.1021)
	(7.00)	(18.00)	(54.00)	(92.00)

(VD – Very Dissatisfied; LD = Little Dissatisfied; MS = Moderately Satisfied; VS = Very Satisfied).

where the infimum is taken over the set

$$\{\alpha_{ij} \geq 0, 1 \leq i \leq r-1, 1 \leq j \leq c-1 \text{ and } \alpha_{ij} \in \mathbb{R}, i = r \text{ or } j = c\}.$$

However, equivalently, we can consider each I-projection problem at hand, find its dual problem, and the corresponding dual solution and then relate the primal solution to this dual solution.

We illustrate the procedure using the data described in Table 1 taken from Agresti ((1990), p. 21), also in 1984 General Social Survey (see Norušis (1988)). Thus we wish to solve the problem in (3.1) using the  $n_{ij}$  in Table 1. Since  $\theta_{12} = 0.82$ , some constraints will have to be active. When written in terms of (3.2) there

are essentially fifteen constraints (in addition to the always present probability vector constraint). For every constraint, the dual problem is solved by using a one-dimensional Newton-Raphson method. The maximum likelihood estimates subject to the constraints  $\theta_{ij} \geq 1 \forall i, j = 1, 2, 3$  are listed in Table 2.

We calculated the bound given in (2.5) for this example. We obtained  $\mathbf{p}^*$  by changing some elements of the solution given in Table 2 ( $p_{12}^* = 0.0291$ ,  $p_{13}^* = 0.0863$ ,  $p_{32}^* = 0.0308$ ,  $p_{33}^* = 0.0903$ ). It is easy to check that this  $\mathbf{p}^*$  is feasible. Then the left side of (2.5) is zero upto five decimal places (for large  $n$ ).

Patefield (1982) considered  $r \times c$  contingency tables and tests of hypotheses  $H_0 : \theta_{ij} = 1, \forall i, j$  versus  $H_1 : \theta_{ij} \geq 1, \forall i, j$ . The likelihood ratio statistic was calculated using numerical routines (for small values of  $r$  and  $c$ ). The restricted estimates as derived in this paper would prove to be useful for this purpose.

#### 4. Discussion

The iterative proportional fitting procedures, starting with Deming and Stephan (1940), are used widely in many facets of statistical applications. We have used the Fenchel duality theorem to show that Dykstra's iterative proportional fitting procedure amounts to a sequential coordinatewise minimization procedure in the dual space. This algorithm is simple and helps to explain why Dykstra's algorithm should work in the primal (I-projection) problem.

#### Acknowledgements

The authors thank the referees for helpful and knowledgeable comments.

#### REFERENCES

- Agresti, A. (1990). *Categorical Data Analysis*, Wiley, New York.
- Bhattacharya, B. and Dykstra, R. L. (1995). A general duality approach to I-projections, *J. Statist. Plann. Inference*, **3**, 146–159.
- Csiszar, I. (1975). I-divergence geometry of probability distributions and minimization problems, *Ann. Probab.*, **3**, 146–159.
- Deming, W. E. and Stephan, F. F. (1940). On a least square adjustment of a sampled frequency table when the expected marginal totals are known, *Ann. Math. Statist.*, **11**, 427–444.
- Dykstra, R. L. (1985a). An iterative procedure for obtaining I projections onto the intersection of convex sets, *Ann. Probab.*, **13**, 975–984.
- Dykstra, R. L. (1985b). Computational aspects of I-projections, *J. Statist. Comput. Simulation*, **21**, 265–274.
- Dykstra, R. L. and Lemke, J. H. (1988). Duality of I-projections and maximum likelihood estimation for log-linear models under cone constraints, *J. Amer. Statist. Assoc.*, **402**, 546–554.
- Good, I. J. (1963). Maximum entropy for hypothesis formulation, especially for multidimensional contingency tables, *Ann. Math. Statist.*, **34**, 911–934.
- Hestenes, M. R. (1975). *Optimization Theory*, Wiley, New York.
- Jaynes, E. T. (1957). Information theory and statistical mechanics, *Phys. Rev.*, **106**, 620–630.
- Kullback, S. (1959). *Information Theory and Statistics*, Wiley, New York.
- Norušis, M. J. (1988). *SPSSX Advanced Statistics Guide*, 2nd ed., McGraw-Hill, New York.
- Patefield, W. M. (1982). Exact tests for trends in ordered contingency tables, *J. Roy. Statist. Soc. Ser. C*, **31**, 32–43.

- Rao, C. R. (1965). *Linear Statistical Inference and Its Applications*, Wiley, New York.
- Rockafellar, R. T. (1970). *Convex Analysis*, Princeton University Press, New York.
- Sanov, I. N. (1957). On the probability of large deviations of random variables, *Mat. Sb.*, **42**, 11–44
- Winkler, W. (1990). On Dykstra's iterative proportional fitting procedure, *Ann. Probab.*, **18**, 1410–1415.