

# Uniform Approximation by Rational Functions Having Restricted Denominators

E. H. KAUFMAN, JR.

*Department of Mathematics,  
Central Michigan University, Mount Pleasant, Michigan 48859*

AND

G. D. TAYLOR\*

*Department of Mathematics,  
Colorado State University, Fort Collins, Colorado 80523*

*Communicated by Richard S. Varga*

Received June 25, 1979

This paper considers approximation of continuous functions on a compact metric space by generalized rational functions for which the denominators have bounded coefficients and are bounded below by a fixed positive function. This lower bound alleviates numerical difficulties, and in some applications (e.g., digital filter design) has a useful physical interpretation. A “zero in the convex hull” characterization of best approximations is developed and used to prove uniqueness and de la Vallée Poussin results. Examples are given to illustrate this theory and its differences with the standard theory, where the denominators are merely required to be positive. A modified differential correction algorithm is presented and is proved to always converge at least linearly, and often quadratically.

## 1. INTRODUCTION

In this paper we consider approximation of continuous functions by generalized rational functions whose denominators are required to be bounded away from zero. This is in contrast to the standard theory, where the denominators are only required to be positive. There are at least four reasons for having this stronger requirement:

- (1) Best approximations always exist.

\* Supported in part by the Air Force Office of Scientific Research, Air Force Systems Command, USAF, under Grant F-49620-79-C-0124, and by the National Science Foundation under Grant MCS-78-05847.

(2) Iterative procedures for computing rational approximations may experience numerical difficulties if the denominators of the rational functions become very small.

(3) Even if a good rational approximation is found, it may not be very useful if its denominator is too small at some point.

(4) There may be some physical advantage in being able to control the denominator. For example, McCallig [6] used a version of the algorithm described in this paper to compute approximations to the desired magnitude-squared response of a digital filter; control of the denominator of the rational function amounts to control of the feedback gain of the resulting filter, which allows one to progress smoothly from "fully recursive" filter designs to nonrecursive designs, and to reduce sensitivity problems and hardware requirements.

Formally, the situation we are considering is as follows.  $X$  is a compact metric space,  $m$  and  $n$  are fixed positive integers,  $\mathcal{P}$  and  $\mathcal{Q}$  are subspaces of  $C[X]$  with bases  $\{\theta_1, \dots, \theta_m\}$  and  $\{\psi_1, \dots, \psi_n\}$ , respectively, and  $L$  is a strictly positive continuous function on  $X$  (which is often a constant in practice). Our family of approximating functions is then defined to be

$$\mathcal{R}_L = \{P/Q: P = p_1\theta_1 + \dots + p_m\theta_m \in \mathcal{P}, Q = q_1\psi_1 + \dots + q_n\psi_n \in \mathcal{Q}, \\ Q \geq L \text{ on } X, |q_j| \leq 1 \text{ for } j = 1, \dots, n\}.$$

We note in passing that without the restrictions  $|q_j| \leq 1$ , the restriction  $Q \geq L$  would be no stronger than the more usual restriction  $Q > 0$ , since it could always be satisfied by multiplying  $P$  and  $Q$  by a sufficiently large positive constant.

In order that  $\mathcal{R}_L$  be nonempty,  $\mathcal{Q}$  must contain at least one strictly positive function, so without loss of generality we will assume  $\psi_1 > 0$  on  $X$ . To insure  $\mathcal{R}_L \neq \emptyset$ , and for other reasons which will be clearer later, we will also assume  $\max_{x \in X} L(x) < \min_{x \in X} \psi_1(x)$ ; this requirement is no restriction in practice, since it can always be obtained by multiplying  $\psi_1$  by a suitable positive constant.

Given  $f \in C[X]$ , a best approximation to  $f$  is defined to be a function  $R^* \in \mathcal{R}_L$  such that  $\|f - R^*\| \leq \|f - R\|$  for all  $R \in \mathcal{R}_L$ , where for  $g \in C[X]$ ,  $\|g\| = \max_{x \in X} |g(x)|$ . The following theorem can be proved by standard techniques.

**THEOREM 1.** *If  $f \in C[X]$ , then there exists a best approximation to  $f$  from  $\mathcal{R}_L$ .*

In the remaining sections we will consider characterization of best approximations, uniqueness, de La Vallée Poussin results, and computation of

best approximations by a modified differential correction algorithm. The theory differs in some respects from the restricted range situation, where the entire rational function rather than just the denominator is restricted.

Suppose  $f \in C[X]$  and  $R^* \in \mathcal{R}_L$ , where  $R^* = P^*/Q^* = (p_1^*\theta_1 + \dots + p_m^*\theta_m)/(q_1^*\psi_1 + \dots + q_n^*\psi_n)$ . The following notation will be useful.

$$\sigma(x) = \text{sgn}(f(x) - R^*(x)) \quad \forall x \in X;$$

$$\mathcal{P} + R^*\mathcal{Q} = \{P + R^*Q: P \in \mathcal{P}, Q \in \mathcal{Q}\};$$

$$X_0 = \{x \in X: |f(x) - R^*(x)| = \|f - R^*\| \};$$

$$Y_0 = \{y \in X: Q^*(y) = L(y)\};$$

$$I_0 = \{j \in \{1, \dots, n\}: |q_j^*| = 1\};$$

$$S = \{\sigma(x) \hat{x}: x \in X_0\} \cup \{\psi(y): y \in Y_0\} \cup \{q_j^* \mathbf{e}_{m+j}: j \in I_0\},$$

where

$$\hat{x} = (\theta_1(x), \dots, \theta_m(x), R^*(x) \psi_1(x), \dots, R^*(x) \psi_n(x))^T,$$

$$\psi(y) = (0, \dots, 0, -\psi_1(y), \dots, -\psi_n(y))^T$$

and

$$\mathbf{e}_k = (\delta_{1k}, \dots, \delta_{m+n,k})^T,$$

$$\delta_{ij} = \text{Kronecker delta};$$

$$\mathcal{H}(S) = \text{the convex hull of } S = \left\{ \sum_{i=1}^k \lambda_i s_i: \right.$$

$$\left. k \text{ is a positive integer, } s_i \in S \forall i, \lambda_i \geq 0 \forall i, \sum_{i=1}^k \lambda_i = 1 \right\};$$

$$\text{int } \mathcal{H}(S) = \text{the interior of } \mathcal{H}(S).$$

## 2. CHARACTERIZATION

We first prove a Kolmogorov-type characterization theorem.

**THEOREM 2.** *Suppose  $f \in C[X] - \mathcal{R}_L$ . Then  $R^* = P^*/Q^* \in \mathcal{R}_L$  is a best approximation to  $f$  iff there is no  $\bar{P} = \bar{p}_1\theta_1 + \dots + \bar{p}_m\theta_m \in \mathcal{P}$ ,  $\bar{Q} = \bar{q}_1\psi_1 + \dots + \bar{q}_n\psi_n \in \mathcal{Q}$  satisfying*

$$(i) \quad \text{sgn}(\bar{P} + R^*\bar{Q})(x) = \text{sgn}(f(x) - R^*(x)), \quad \forall x \in X_0;$$

$$(ii) \quad \bar{Q}(y) < 0, \quad \forall y \in Y_0;$$

- (iii)  $\bar{q}_j > 0$  if  $q_j^* = 1$ ;  
 (iv)  $\bar{q}_j < 0$  if  $q_j^* = -1$ .

*Proof.* ( $\Leftarrow$ ) Suppose  $R^*$  is not a best approximation. Then there is a better approximation  $R = P/Q \in \mathcal{R}_L$ . By our assumptions in the previous section, we have  $(\max_{x \in X} L(x))/(\min_{x \in X} \psi_1(x)) < 1$ ; letting  $a$  be any number satisfying  $(\max_{x \in X} L(x))/(\min_{x \in X} \psi_1(x)) < a < 1$ , we define  $\bar{Q} = \bar{q}_1 \psi_1 + \dots + \bar{q}_n \psi_n \in \mathcal{Q}$  by  $\bar{Q} = Q^* - a\psi_1$ . Then  $\bar{Q}$  satisfies  $\bar{Q}(y) < 0 \quad \forall y \in Y_0$ ,  $\bar{q}_j > 0$  if  $q_j^* = 1$ , and  $\bar{q}_j < 0$  if  $q_j^* = -1$ . Now define  $\bar{P} \equiv P - P^*$ ,  $\bar{Q} \equiv Q^* - Q + \eta \bar{Q}$ , where  $\eta$  is a positive number. For  $x \in X_0$ , we have

$$\begin{aligned} \operatorname{sgn}(\bar{P} + R^*\bar{Q})(x) &= \operatorname{sgn}(P - P^* + R^*(Q^* - Q + \eta \bar{Q}))(x) \\ &= \operatorname{sgn}(P - R^*Q + \eta R^*\bar{Q})(x) \\ &= \operatorname{sgn} \left\{ Q(x) \left[ \left( R(x) - R^*(x) + \eta \frac{R^*(x) \bar{Q}(x)}{Q(x)} \right) \right] \right\} \\ &= \operatorname{sgn} \left[ f(x) - R^*(x) - (f(x) - R(x)) + \eta \frac{R^*(x) \bar{Q}(x)}{Q(x)} \right]. \end{aligned}$$

Now choosing  $\eta$  so small that  $|\eta(R^*(x) \bar{Q}(x)/Q(x))| < \|f - R^*\| - \|f - R\| \quad \forall x \in X_0$ , we have

$$\operatorname{sgn}(\bar{P} + R^*\bar{Q})(x) = \operatorname{sgn}(f(x) - R^*(x)) \quad \forall x \in X_0,$$

so (i) holds. For  $y \in Y_0$ , we have

$$\bar{Q}(y) = Q^*(y) - Q(y) + \eta \bar{Q}(y) \leq L(y) - L(y) + \eta \bar{Q}(y) < 0,$$

so (ii) holds. If  $q_j^* = 1$ , we have

$$\bar{q}_j = q_j^* - q_j + \eta \bar{q}_j \geq 1 - 1 + \eta \bar{q}_j > 0,$$

so (iii) holds. If  $q_j^* = -1$ , we have

$$\bar{q}_j = q_j^* - q_j + \eta \bar{q}_j \leq -1 - (-1) + \eta \bar{q}_j < 0,$$

so (iv) holds.

( $\Rightarrow$ ) Suppose there exists  $\bar{P} \in \mathcal{P}$ ,  $\bar{Q} \in \mathcal{Q}$  satisfying (i)–(iv) above. Let  $\lambda$  be a small positive number. Then

$$\frac{P^* + \lambda \bar{P}}{Q^* - \lambda \bar{Q}} - \frac{P^*}{Q^*} = \lambda \frac{\bar{P}Q^* + P^*\bar{Q}}{(Q^* - \lambda \bar{Q})Q^*} = \lambda \frac{\bar{P} + R^*\bar{Q}}{Q^* - \lambda \bar{Q}}.$$

Thus

$$f - \frac{P^* + \lambda \bar{P}}{Q^* - \lambda \bar{Q}} = f - \frac{P^*}{Q^*} - \lambda \frac{\bar{P} + R^* \bar{Q}}{Q^* - \lambda \bar{Q}},$$

and using the arguments of [2, pp. 159, 160] it can be shown that

$$\left\| f - \frac{P^* + \lambda \bar{P}}{Q^* - \lambda \bar{Q}} \right\| < \left\| f - \frac{P^*}{Q^*} \right\|$$

for all  $\lambda$  sufficiently small. It remains only to show that

$$\frac{P^* + \lambda \bar{P}}{Q^* - \lambda \bar{Q}} \in \mathcal{R}_L$$

for  $\lambda$  sufficiently small.

For all  $y \in Y_0$ , we have

$$Q^*(y) - \lambda \bar{Q}(y) > Q^*(y) = L(y).$$

Since  $Y_0$  is compact, there exists  $\zeta > 0$  such that  $\bar{Q}(y) \leq -\zeta \forall y \in Y_0$ . Let  $Y_1 = \{x \in X: \bar{Q}(x) < -\zeta/2\}$ ,  $Y_2 = X - Y_1$ . Then  $Y_2$  is compact, with  $Y_2 \cap Y_0 = \emptyset$ . Let  $\mu = \min\{Q^*(x) - L(x): x \in Y_2\} > 0$ . Choose  $\lambda$  to satisfy  $0 < \lambda < \mu/\max(1, \|\bar{Q}\|)$ . Then if  $x \in Y_1$ , we have  $Q^*(x) - \lambda \bar{Q}(x) > Q^*(x) \geq L(x)$ ; if  $x \in Y_2$ , we have  $Q^*(x) - \lambda \bar{Q}(x) \geq L(x) + \mu - \lambda \bar{Q}(x) \geq L(x)$ , so  $Q^*(x) - \lambda \bar{Q}(x) \geq L(x) \forall x \in X$ . Finally, if  $q_j^* = 1$  we have

$$q_j^* - \lambda \bar{q}_j = 1 - \lambda \bar{q}_j < 1,$$

and if  $q_j^* = -1$  we have

$$q_j^* - \lambda \bar{q}_j = -1 - \lambda \bar{q}_j > -1,$$

so choosing  $\lambda$  sufficiently small will insure that  $|q_j^* - \lambda \bar{q}_j| \leq 1, j = 1, \dots, n$ . Thus for  $\lambda$  sufficiently small we have

$$\frac{P^* + \lambda \bar{Q}}{Q^* - \lambda \bar{Q}} \in \mathcal{R}_L \quad \text{and} \quad \left\| f - \frac{P^* + \lambda \bar{Q}}{Q^* - \lambda \bar{Q}} \right\| < \left\| f - \frac{P^*}{Q^*} \right\|,$$

so  $P^*/Q^*$  is not a best approximation.

Q.E.D.

We can now prove a “zero in the convex hull” characterization of best approximations.

**THEOREM 3.** *Suppose  $f \in C[X] - \mathcal{R}_L$ . Then  $R^* = P^*/Q^* \in \mathcal{R}_L$  is a best approximation to  $f$  iff the origin of  $(m+n)$ -space lies in the convex hull of  $S$  (where  $S$  is the set defined in the Introduction).*

*Proof.* By the theorem on linear inequalities [2, p. 19],  $\mathbf{0} \notin \mathcal{H}(S)$  iff the system of inequalities

$$\langle z, s \rangle > 0, \quad s \in S,$$

is consistent (here  $\langle \cdot, \cdot \rangle$  denotes inner product). But this is true iff  $\exists$  a vector  $z = [z_1, \dots, z_{m+n}]^T$  satisfying

$$\begin{aligned} \sigma(x)[z_1 \theta_1(x) + \dots + z_m \theta_m(x) + R^*(x)(z_{m+1} \psi_1(x) + \dots + z_{m+n} \psi_n(x))] \\ > 0 \quad \forall x \in X_0; \\ z_{m+1} \psi_1(y) + \dots + z_{m+n} \psi_n(y) < 0 \quad \forall y \in Y_0; \\ z_{m+j} > 0 \quad \text{if } q_j^* = 1; \end{aligned}$$

and

$$z_{m+j} < 0 \quad \text{if } q_j^* = -1.$$

Letting  $\bar{P} = z_1 \theta_1 + \dots + z_m \theta_m$  and  $\bar{Q} = z_{m+1} \psi_1 + \dots + z_{m+n} \psi_n$ , we see by Theorem 2 that this is true iff  $R^*$  is not a best approximation. Thus  $\mathbf{0} \in \mathcal{H}(S)$  iff  $R^*$  is a best approximation. Q.E.D.

We illustrate the application of this theorem with the following example.

EXAMPLE 1. Let  $X = \{0, 1\}$ ,  $f(x) = x$ ,  $\mathcal{P} = \Pi_0$  = the set of all polynomials of degree  $\leq 0$ ,  $\mathcal{S} = \Pi_1$ ,  $L(x) \equiv 0.1$ ,  $R^*(x) = (1/11)/(1 - 0.9x)$ . We have  $X_0 = \{0, 1\}$ ,  $Y_0 = \{1\}$ ,  $I_0 = \{1\}$ ,  $\sigma(0) = -1$ ,  $\sigma(1) = 1$ . Thus

$$\sigma(x) \begin{bmatrix} \theta_1(x) \\ (R^* \psi_1)(x) \\ (R^* \psi_2)(x) \end{bmatrix} = \sigma(x) \begin{bmatrix} 1 \\ \frac{1}{11} / (1 - 0.9x) \\ \frac{x}{11} / (1 - 0.9x) \end{bmatrix},$$

so

$$S = \left\{ \begin{bmatrix} -1 \\ -\frac{1}{11} \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ \frac{10}{11} \\ \frac{10}{11} \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \right\}.$$

Solving the linear system

$$\begin{aligned} -\lambda_1 + \lambda_2 &= 0, \\ -\frac{1}{11}\lambda_1 + \frac{10}{11}\lambda_2 - \lambda_3 + \lambda_4 &= 0, \\ \frac{10}{11}\lambda_2 - \lambda_3 &= 0, \\ \lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 &= 1 \end{aligned}$$

yields the solution  $\lambda_1 = 1/3, \lambda_2 = 1/3, \lambda_3 = 10/33, \lambda_4 = 1/33$ . Since we have  $\lambda_i \geq 0 \forall i$ , we have  $\mathbf{0} \in \mathcal{H}(S)$ , so  $R^*$  is a best approximation.

We may observe that in this example  $\mathbf{0}$  is a positive convex combination of exactly  $m + n + 1 = 4$  vectors in  $S$ , and the coefficient matrix used is nonsingular. Thus by Cramer's rule  $(\delta_1, \delta_2, \delta_3)^T \in \mathcal{H}(S)$  if  $|\delta_1|, |\delta_2|, |\delta_3|$  are sufficiently small, so  $\mathbf{0}$  is actually in the interior of  $\mathcal{H}(S)$ ; this distinction will be important in the next two sections.

The next example illustrates what can happen if our assumption that  $\max_{x \in X} L(x) < \min_{x \in X} \psi_1(x)$  is violated.

EXAMPLE 2. Let  $X = [0, 1], \mathcal{P} = \Pi_0, \mathcal{Q} = \Pi_1, L(x) \equiv 1, R^*(x) = 1/(1 + 0.5x)$ . We have

$$\begin{bmatrix} 0 \\ -1 \\ 0 \end{bmatrix} \in S \quad \text{and} \quad \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \in S,$$

so  $\mathbf{0} \in \mathcal{H}(S)$  regardless of whether  $R^*$  is a best approximation to  $f$  or not. Intuitively, the trouble is that we are imposing a double restraint on the denominator at  $x = 0$  which ties it down completely there.

The next example shows that the standard alternation characterization of best approximations does not hold in our setting.

EXAMPLE 3. Let  $X = [0, 3]$ ,

$$\begin{aligned} f(x) &= 3.5x, & 0 \leq x \leq 1, \\ &= 2 + 1.5x, & 1 \leq x \leq 2, \\ &= 8 - 1.5x, & 2 \leq x \leq 3, \end{aligned}$$

$\mathcal{P} = \Pi_0, \mathcal{Q} = \Pi_2, L(x) \equiv 0.2, R^*(x) = 1/(1 - 0.8x + 0.2x^2)$ . We have  $X_0 = \{0, 1, 3\}, Y_0 = \{2\}, I_0 = \{1\}, \sigma(0) = -1, \sigma(1) = 1, \sigma(3) = 1$ . The fact that  $R^*$  is a best approximation to  $f$  is shown by the equality

$$\begin{aligned} & \frac{4}{17} \begin{bmatrix} -1 \\ -1 \\ 0 \\ 0 \end{bmatrix} + \frac{3}{17} \begin{bmatrix} 1 \\ 2.5 \\ 2.5 \\ 2.5 \end{bmatrix} + \frac{1}{17} \begin{bmatrix} 1 \\ 2.5 \\ 7.5 \\ 22.5 \end{bmatrix} + \frac{15}{34} \begin{bmatrix} 0 \\ -1 \\ -2 \\ -4 \end{bmatrix} \\ & + \frac{3}{34} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}. \end{aligned}$$

Without denominator restrictions we would expect  $(m-1) + (n-1) + 2 - \min(m-1 - \text{degree of } P^*, n-1 - \text{degree of } Q^*) = \text{four alternating points for } f - R^*$ ; since there is only one denominator restriction we might have hoped for three alternating extreme points, but there are only two. It is tempting to conjecture that denominator constraints act as extreme points with "sign of error" opposite that of the previous extreme point, but there are examples which show that there may still be fewer than expected alternating extreme points. Thus there does not appear to be a simple alternation theorem in this setting. Some partial results under suitable Haar assumptions are possible; however, they do not seem to add much insight.

### 3. UNIQUENESS AND DE LA VALLÉE POUSSIN RESULTS

Best approximations from  $\mathcal{R}_L$  need not be unique, as shown by the following example.

EXAMPLE 4. Let  $X = \{0, 1\}$ ,  $f(x) = x$ ,  $\mathcal{P} = \Pi_0$ ,  $\mathcal{Q} = \Pi_2$ ,  $L(x) \equiv 0.1$ ,  $R^*(x) = (1/11)/(1 - 0.9x)$ . We have  $X_0 = \{0, 1\}$ ,  $Y_0 = \{1\}$ ,  $I_0 = \{1\}$ ,  $\sigma(0) = -1$ ,  $\sigma(1) = 1$ . Thus

$$\sigma(x) \begin{bmatrix} \theta_1(x) \\ (R^*\psi_1)(x) \\ (R^*\psi_2)(x) \\ (R^*\psi_3)(x) \end{bmatrix} = \sigma(x) \begin{bmatrix} 1 \\ \frac{1}{11}/(1 - 0.9x) \\ \frac{x}{11}/(1 - 0.9x) \\ \frac{x^2}{11}/(1 - 0.9x) \end{bmatrix}$$



so

$$S = \left\{ \begin{bmatrix} -1 \\ \frac{1}{11} \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ \frac{10}{11} \\ \frac{10}{11} \\ \frac{10}{11} \end{bmatrix}, \begin{bmatrix} 0 \\ -1 \\ -1 \\ -1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \right\}$$

and  $R^*$  is shown to be a best approximation by the equality

$$\frac{1}{3} \begin{bmatrix} -1 \\ \frac{1}{11} \\ 0 \\ 0 \end{bmatrix} + \frac{1}{3} \begin{bmatrix} 1 \\ \frac{10}{11} \\ \frac{10}{11} \\ \frac{10}{11} \end{bmatrix} + \frac{10}{33} \begin{bmatrix} 0 \\ -1 \\ -1 \\ -1 \end{bmatrix} + \frac{1}{33} \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}.$$

But  $R(x) = (1/11)/(1 + q_2x + q_3x^2)$  is also a best approximation for any  $q_2, q_3$  satisfying  $|q_2| \leq 1, |q_3| \leq 1, q_2 + q_3 = -0.9$  since  $\|f - R\| = 1/11 = \|f - R^*\|$ .

It turns out uniqueness is assured if, unlike the situation in this example,  $\mathbf{0}$  is in the interior of  $\mathcal{H}(S)$ . To prove this, we need the following lemma.

LEMMA 1. Suppose  $f \in C[X]$  and  $R^* \in R_L$ . Suppose  $\tilde{X}_0$  is an arbitrary compact subset of  $X$ , and  $\tilde{S}$  is the  $S$  of earlier theorems with  $X_0$  replaced by  $\tilde{X}_0$ . If  $\mathbf{0} \in \text{int } \mathcal{H}(\tilde{S})$ , then  $\bar{P} \equiv 0, \bar{Q} \equiv 0$  is the only solution in  $\mathcal{P}, \mathcal{Q}$  to the inequalities

- (a)  $\sigma(x)(\bar{P} + R^*\bar{Q})(x) \geq 0, \forall x \in \tilde{X}_0;$
- (b)  $\bar{Q}(y) \leq 0, \forall y \in Y_0;$
- (c)  $\bar{q}_j \geq 0$  if  $q_j^* = 1;$
- (d)  $\bar{q}_j \leq 0$  if  $q_j^* = -1.$

*Proof.* Suppose  $\bar{P} \in \mathcal{P}, \bar{Q} \in \mathcal{Q}$  satisfy (a)–(d). Letting  $\mathbf{z} = [\bar{p}_1, \dots, \bar{p}_m, \bar{q}_1, \dots, \bar{q}_n]^T$ , the system (a)–(d) may be rewritten as  $\langle \mathbf{z}, \mathbf{s} \rangle \geq 0 \forall \mathbf{s} \in \tilde{S}$ . Suppose that there is a  $\mathbf{z} \neq \mathbf{0}$  satisfying these inequalities. Since  $\mathbf{0} \in \text{int } \mathcal{H}(\tilde{S})$ , for

$\delta > 0$  sufficiently small we have  $-\delta \mathbf{z} \in \mathcal{H}(\bar{S})$ . By Caratheodory's theorem [2, p. 17] for some integer  $k \leq m + n + 1$ ,  $\exists \mathbf{s}_1, \dots, \mathbf{s}_k \in \bar{S}$ ,  $\lambda_1, \dots, \lambda_k \geq 0$  with  $\sum_{i=1}^k \lambda_i = 1$  such that  $-\delta \mathbf{z} = \sum_{i=1}^k \lambda_i \mathbf{s}_i$ . Thus  $-\delta \langle \mathbf{z}, \mathbf{z} \rangle = \langle \mathbf{z}, -\delta \mathbf{z} \rangle = \langle \mathbf{z}, \sum_{i=1}^k \lambda_i \mathbf{s}_i \rangle = \sum_{i=1}^k \lambda_i \langle \mathbf{z}, \mathbf{s}_i \rangle \geq 0$ , which is contradiction. Thus  $\mathbf{z} = \mathbf{0}$ , so  $\bar{P} \equiv 0, \bar{Q} \equiv 0$ . Q.E.D.

Although we will not use it in this paper, the converse of this lemma can also be shown to be true.

We can now prove uniqueness of  $\mathbf{0} \in \text{int } \mathcal{H}(S)$ .

**THEOREM 4.** *Suppose  $f \in C[X] - \mathcal{R}_L$  and  $R^* = P^*/Q^* \in \mathcal{R}_L$ . If  $\mathbf{0} \in \text{int } \mathcal{H}(S)$ , then  $R^*$  is the unique best approximation to  $f$  from  $\mathcal{R}_L$ .*

*Proof.* By Theorem 3,  $R^*$  is a best approximation. Suppose  $R = P/Q \in \mathcal{R}_L$  were another best approximation. Let  $\bar{P} = P - P^*$ ,  $\bar{Q} = Q^* - Q$ . For  $x \in X_0$ , we have

$$\begin{aligned} \sigma(x)(\bar{P} + R^*\bar{Q})(x) &= \sigma(x)(P - R^*Q)(x) \\ &= \sigma(x)Q(x)[f(x) - R^*(x) - (f(x) - R(x))] \geq 0. \end{aligned}$$

For  $y \in Y_0$ , we have

$$\bar{Q}(y) = Q^*(y) - Q(y) = L(y) - Q(y) \leq 0.$$

Finally, if  $q_j^* = 1$ , we have

$$\bar{q}_j = q_j^* - q_j = 1 - q_j \geq 0$$

and if  $q_j^* = -1$ , we have

$$\bar{q}_j = q_j^* - q_j = -1 - q_j \leq 0.$$

Thus by Lemma 1,  $\bar{P} \equiv 0$  and  $\bar{Q} \equiv 0$ . Thus  $P \equiv P^*$  and  $Q \equiv Q^*$ , so  $R^*$  is unique. Q.E.D.

The next example shows that the converse of this theorem is false.

**EXAMPLE 5.** Let  $X = [0, 1]$ ,  $f(x) = 2 - 2x$ ,  $\mathcal{P} = \mathcal{Q} = \Pi_0$ ,  $L(x) \equiv 0.1$ ,  $R^*(x) \equiv 1/1$ . We have  $X_0 = \{0, 1\}$ ,  $Y_0 = \emptyset$ ,  $I_0 = (1)$ ,  $S = \{\begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \end{bmatrix}\}$ . Here  $R^*$  is the unique best approximation, but  $\mathbf{0} \notin \text{int } \mathcal{H}(S)$ .

In general, we always have  $\mathbf{0} \notin \text{int } \mathcal{H}(S)$  if  $R^*$  is a best approximation with  $Y_0 = \emptyset$ ; the reason is that  $\mathbf{0} \in \text{int } \mathcal{H}(S)$  implies that the coefficients of the best approximation are unique, but  $Y_0 = \emptyset$  implies that for some  $\alpha$  with  $0 < \alpha < 1$ ,  $(\alpha P^*)/(\alpha Q^*)$  is another best approximation in  $\mathcal{R}_L$  with different coefficients.

The hypotheses of Theorem 4 are actually sufficient to prove strong uniqueness.

**THEOREM 5.** *Suppose  $R^*$  is a best approximation to  $f \in C[X]$  with  $0 \in \text{int } \mathcal{H}(S)$ . Then there is a constant  $\gamma > 0$  such that for any  $R \in \mathcal{R}_L$ ,  $\|f - R\| \geq \|f - R^*\| + \gamma \|R - R^*\|$ .*

*Proof* (sketch). The proof follows the same general lines as the proof of strong uniqueness in the setting where the denominator is merely required to be positive [2, p. 165], with  $P - P^*$  playing the role of  $P$  in the standard proof and  $Q^* - Q$  playing the role of  $Q$ , and Lemma 1 of this paper used in place of Lemmas 1 and 2 in [2]. The other major change is that the definition of  $c$  on p. 166 of [2] is replaced by

$$c = \inf \{ \max [ -\sigma(x)(\bar{P} + R^*\bar{Q})(x) ] : \bar{p} \in \mathcal{P}, \bar{Q} \in \mathcal{Q}, \bar{Q}(y) \leq 0 \forall y \in Y_0, \\ \bar{q}_j \geq 0 \text{ if } q_j^* = 1, \bar{q}_j \leq 0 \text{ if } q_j^* = -1, \\ \max \{ \|\bar{P} + R^*\bar{Q}\|, \|\bar{P}\|, \|\bar{Q}\| \} = 1 \};$$

the extra complication in the last equality of this definition is needed since otherwise we could have  $\|\bar{P} + R^*\bar{Q}\| = 1$  with  $\|\bar{Q}\|$  and  $\|\bar{P}\|$  arbitrarily large. Q.E.D.

We finish this section with two de la Vallée Poussin estimates, which give lower bounds on error norms.

**THEOREM 6.** *Suppose  $f \in C[X] - \mathcal{R}_L$ ,  $R^* \in \mathcal{R}_L$  (not necessarily a best approximation), and  $\tilde{X}_0 \neq \emptyset$  is some compact subset of  $X$ . Suppose  $0 \in \mathcal{H}(S)$ , where  $S$  is the  $S$  of earlier theorems with  $X_0$  replaced by  $\tilde{X}_0$ . Then  $\inf \{ \|f - R\| : R \in \mathcal{R}_L \} \geq \min \{ |f(x) - R^*(x)| : x \in \tilde{X}_0 \}$ .*

*Proof.* Without loss of generality, we may assume  $f(x) \neq R^*(x) \forall x \in \tilde{X}_0$ . Suppose the conclusion of the theorem is false. Then  $\exists R = P/Q \in \mathcal{R}_L$  with  $\|f - R\| < \min \{ |f(x) - R^*(x)| : x \in \tilde{X}_0 \}$ . Let  $\bar{P} \equiv P - P^*$ ,  $\bar{Q} \equiv Q^* - Q$ . Proceeding as in the proof of Theorem 4, we get

- (a')  $\sigma(x)(\bar{P} + R^*\bar{Q})(x) > 0 \forall x \in \tilde{X}_0$ ;
- (b)  $\bar{Q}(y) \leq 0 \forall y \in Y_0$ ;
- (c)  $\bar{q}_j \geq 0$  if  $q_j^* = 1$ ;
- (d)  $\bar{q}_j \leq 0$  if  $q_j^* = -1$ .

Since  $0 \in \mathcal{H}(S)$ , Caratheodory's theorem [2, p. 17] implies that for some integer  $k \leq m + n + 1$ ,  $\exists s_1, \dots, s_k \in S$ ,  $\lambda_1, \dots, \lambda_k \geq 0$  with  $\sum_{l=1}^k \lambda_l = 1$  such that  $\sum_{l=1}^k \lambda_l s_l = 0$ . Thus there are nonnegative induces  $u, v, w$  with  $u + v + w = k$  such that  $\sum_{l=1}^u \lambda_l \sigma(x_l) \hat{x}_l + \sum_{l=1}^v \lambda_{u+l} \psi(y_l) + \sum_{l=1}^w \lambda_{u+v+l} q_l^* e_{m+l} = 0$ ,

where  $x_1, \dots, x_u \in \bar{X}_0$ ;  $y_1, \dots, y_v \in Y_0$ ;  $i_1, \dots, i_w \in I_0$  and  $\hat{x}_l = (\theta_1(x_l), \dots, \theta_m(x_l), R^*(x_l) \psi_1(x_l), \dots, R^*(x_l) \psi_n(x_l))^T$ . We will next show that

$$\sum_{l=1}^u \lambda_l \sigma(x_l) (\bar{P} + R^* \bar{Q})(x_l) \leq 0.$$

We have

$$\begin{aligned} \sum_{l=1}^u \lambda_l \sigma(x_l) \theta_i(x_l) &= 0 && \text{for } i = 1, \dots, m, \\ \sum_{l=1}^u \lambda_l \sigma(x_l) R^*(x_l) \psi_j(x_l) &= \sum_{l=1}^v \lambda_{u+l} \psi_j(y_l) - \hat{\lambda}_j && \text{for } j = 1, \dots, n, \end{aligned}$$

where

$$\begin{aligned} \hat{\lambda}_j &= \lambda_{u+v+l} q_j^* && \text{if } j = i_l \in I_0 \\ &= 0 && \text{otherwise.} \end{aligned}$$

So we get

$$\begin{aligned} &\sum_{l=1}^u \lambda_l \sigma(x_l) (\bar{P} + R^* \bar{Q})(x_l) \\ &= \sum_{l=1}^u \lambda_l \sigma(x_l) \left[ \sum_{i=1}^m \bar{p}_i \theta_i(x_l) + R^*(x_l) \sum_{j=1}^n \bar{q}_j \psi_j(x_l) \right] \\ &= \sum_{i=1}^m \bar{p}_i \left[ \sum_{l=1}^u \lambda_l \sigma(x_l) \theta_i(x_l) \right] + \sum_{j=1}^n \bar{q}_j \left[ \sum_{l=1}^u \lambda_l \sigma(x_l) R^*(x_l) \psi_j(x_l) \right] \\ &= 0 + \sum_{j=1}^n \bar{q}_j \left[ \sum_{l=1}^v \lambda_{u+l} \psi_j(y_l) - \hat{\lambda}_j \right] \\ &= \sum_{l=1}^v \lambda_{u+l} \left[ \sum_{j=1}^n \bar{q}_j \psi_j(y_l) \right] - \sum_{j=1}^n \hat{\lambda}_j \bar{q}_j \\ &= \sum_{l=1}^v \lambda_{u+l} \bar{Q}(y_l) - \sum_{l=1}^w \lambda_{u+v+l} q_l^* \bar{q}_l \leq 0 \end{aligned}$$

by properties (b)–(d), as claimed. But by properties (a'), this implies  $\lambda_1 = \dots = \lambda_u = 0$ . But this in turn implies that  $\mathbf{0} \in \mathcal{S}(S')$ , where  $S' = \{\psi(y): y \in Y_0\} \cup \{q_j^* e_{m+j}: j \in I_0\}$ . Thus, if  $X_0 = \{x \in X: |f(x) - R^*(x)| = \|f - R^*\|\}$  as before, we have  $\mathbf{0} \in \mathcal{S}(\{\sigma(x) \hat{x}: x \in X_0\} \cup S')$ . Thus, by Theorem 3,  $R^*$  is a best approximation to  $f$ . Thus  $\inf\{\|f - R\|: R \in \mathcal{R}_L\} = \|f - R^*\| = \max_{x \in X} |f(x) - R^*(x)| \geq \min\{|f(x) - R^*(x)|: x \in \bar{X}_0\}$ , contradicting the assumption that the conclusion of the theorem is false. Q.E.D.

If we assume that  $\mathbf{0} \in \text{int } \mathcal{R}(\tilde{S})$ , we can prove the following stronger result.

**THEOREM 7.** *Suppose  $f \in C[X]$ ,  $R^* = P^*/Q^* \in \mathcal{R}_L$  (not necessarily a best approximation), and  $\tilde{X}_0 \neq \emptyset$  is some compact subset of  $X$ . Suppose  $\mathbf{0} \in \text{int } \mathcal{R}(\tilde{S})$ , where  $\tilde{S}$  is the  $S$  of earlier theorems with  $X_0$  replaced by  $\tilde{X}_0$ . Then for every  $R = P/Q \in \mathcal{R}_L$  with  $R \neq R^*$  we have  $\max\{|f(x) - R(x)|: x \in \tilde{X}_0\} > \min\{|f(x) - R^*(x)|: x \in \tilde{X}_0\}$ .*

*Proof.* Suppose the conclusion is false. Let  $\bar{P} \equiv P - P^*$ ,  $\bar{Q} \equiv Q^* - Q$ . Then  $\exists R \in \mathcal{R}_L$  with  $R \neq R^*$  and  $\max\{|f(x) - R(x)|: x \in \tilde{X}_0\} \leq \min\{|f(x) - R^*(x)|: x \in \tilde{X}_0\}$ . Proceeding as in the proof of Theorem 4, we get

- (a'')  $\sigma(x)(\bar{P} + R^*\bar{Q})(x) \geq 0, \forall x \in \tilde{X}_0;$
- (b)  $\bar{Q}(y) \leq 0, \forall y \in Y_0;$
- (c)  $\bar{q}_j \geq 0$  if  $q_j^* = 1;$
- (d)  $\bar{q}_j \leq 0$  if  $q_j^* = -1.$

Thus Lemma 1 implies  $\bar{P} \equiv 0, \bar{Q} \equiv 0$ . Thus  $R \equiv R^*$ , contrary to assumption. Q.E.D.

We observe that Example 4 with  $X$  replaced by  $\{0, 0.1, 1\}$  and  $X_0 = \{0, 1\}$  shows that the conclusion of the theorem may fail if  $\mathbf{0} \notin \text{int } \mathcal{R}(\tilde{S})$ .

#### 4. COMPUTATION OF BEST APPROXIMATIONS

The differential correction algorithm introduced by Cheney and Loeb [3] and discussed further by Barrodale *et al.* [1] can be modified to compute approximations from  $\mathcal{R}_L$  by inserting extra constraints to force  $Q(x) \geq L(x)$ . We have

**ALGORITHM (Restricted-denominator differential correction—RDDC).**

- (i) Choose  $P_0/Q_0 \in \mathcal{R}_L;$
- (ii) Having found  $P_k/Q_k \in \mathcal{R}_L$  with  $\|f - R_k\| = \Delta_k$ , choose  $P_{k+1}, Q_{k+1}$  as a solution to the problem

$$\text{minimize: } \max_{x \in X} \frac{|f(x)Q(x) - P(x)| - \Delta_k Q(x)}{Q_k(x)}$$

subject to:  $|q_j| \leq 1, j = 1, \dots, n,$  and  $Q(x) \geq L(x), \forall x \in X;$

- (iii) continue until some stopping criterion is met.

One common stopping criterion is to stop when  $(\Delta_k - \Delta_{k+1})/\Delta_k < \varepsilon$  for some prescribed  $\varepsilon > 0$ , selecting  $R_{k+1}$  as the approximation returned by the algorithm if  $\Delta_{k+1} < \Delta_k$ , and selecting  $R_k$  otherwise. A convenient way of choosing  $P_0/Q_0$ , which often is considerably more efficient than making some arbitrary choice such as  $P_0/Q_0 \equiv 1/1$  (see Lee and Roberts [5] for numerical evidence in the unrestricted-denominator case), is to minimize  $\max_{x \in X} |f(x) Q(x) - P(x)|$  subject to  $Q(x) \geq L(x)$ ,  $\forall x \in X$  and  $|q_j| \leq 1$  for  $j = 1, \dots, n$ .

Using the techniques of Barrodale *et al.* [1], we prove

**THEOREM 8.** *The RDDC algorithm converges monotonically and at least linearly.*

*Proof.* Let  $M = \max_{x \in X} \sum_{j=1}^n |\psi_j(x)|$  (thus  $\|Q\| \leq M$  for all  $Q \in \mathcal{Q}$  with  $|q_j| \leq 1$ ). Suppose  $R_k$  is not a best approximation. Let  $R^*$  be a best approximation with  $\Delta^* = \|f - R^*\|$ . Let  $\delta = \min_{x \in X} L(x) > 0$ . We have

$$\begin{aligned} & \max_{x \in X} \frac{|f(x) Q_{k+1}(x) - P_{k+1}(x)| - \Delta_k Q_{k+1}(x)}{Q_k(x)} \\ & \leq \max_{x \in X} \frac{|f(x) Q^*(x) - P^*(x)| - \Delta_k Q^*(x)}{Q_k(x)} \\ & = \max_{x \in X} \left\{ [|f(x) - R^*(x)| - \Delta_k] \cdot \frac{Q^*(x)}{Q_k(x)} \right\} \leq [\Delta^* - \Delta_k] \cdot \frac{\delta}{M} < 0, \end{aligned}$$

$\therefore \forall x \in X$ ,

$$\left[ \left| f(x) - \frac{P_{k+1}(x)}{Q_{k+1}(x)} \right| - \Delta_k \right] \cdot \frac{Q_{k+1}(x)}{Q_k(x)} \leq [\Delta^* - \Delta_k] \cdot \frac{\delta}{M} < 0,$$

$$\therefore \forall x \in X, \quad \left| f(x) - \frac{P_{k+1}(x)}{Q_{k+1}(x)} \right| - \Delta_k < 0,$$

$\therefore \Delta_{k+1} < \Delta_k$ , so the convergence is monotonic.

$\forall x \in X$ , we have

$$\begin{aligned} \left| f(x) - \frac{P_{k+1}(x)}{Q_{k+1}(x)} \right| - \Delta_k & \leq [\Delta^* - \Delta_k] \cdot \frac{\delta}{M} \cdot \frac{Q_k(x)}{Q_{k+1}(x)} \\ & \leq [\Delta^* - \Delta_k] \cdot \left( \frac{\delta}{M} \cdot \frac{\delta}{M} \right), \end{aligned}$$

$$\therefore \Delta_{k+1} - \Delta_k \leq [\Delta^* - \Delta_k] \cdot \frac{\delta^2}{M^2},$$

$$\therefore \Delta_{k+1} - \Delta^* - (\Delta_k - \Delta^*) \leq -\frac{\delta^2}{M^2} (\Delta_k - \Delta^*),$$

$$\therefore \Delta_{k+1} - \Delta^* \leq \left(1 - \frac{\delta^2}{M^2}\right) (\Delta_k - \Delta^*),$$

$\therefore \Delta_k$  converges at least linearly to  $\Delta^*$ . Q.E.D.

We observe that this theorem is stronger than the corresponding theorem in the unrestricted-denominator case [1, Theorems 1, 2] in that it gives information on the rate of convergence, and finiteness of  $X$  is not required to prove convergence. Finiteness of  $X$  is required, however, in order to run the algorithm in the usual way.

The following lemma was proved and used by Barrodale *et al.* [1] for the case  $\mathcal{P} = \Pi_{m-1}$  and  $\mathcal{Q} = \Pi_{n-1}$ . The Haar subspace assumption of our lemma is equivalent to their assumption that  $\min(m-1 - \text{degree of } P^*, n-1 - \text{degree of } Q^*) = 0$ , with  $P^*$  and  $Q^*$  having no common nonconstant factors in their setting.

**LEMMA 2.** *Suppose  $X$  contains at least  $m+n+1$  distinct points and  $R^* = P^*/Q^* \in \mathcal{R}_0 \equiv \{P/Q : P = p_1\theta_1 + \dots + p_m\theta_m \in \mathcal{P}, Q = q_1\psi_1 + \dots + q_n\psi_n \in \mathcal{Q}, Q > 0 \text{ on } X, \max_{1 \leq j \leq n} |q_j| = 1\}$ . Suppose that the space spanned by  $\{\theta_1, \dots, \theta_m, R^*\psi_1, \dots, R^*\psi_n\}$  is a Haar subspace of dimension  $m+n-1$ ; that is, the space has dimension  $m+n-1$ , and no nontrivial element of it has more than  $m+n-2$  distinct zeros in  $X$ . Then  $\exists \theta > 0$  such that for all  $R = P/Q \in \mathcal{R}_0$  we have  $\|Q - Q^*\| \leq \theta \|R - R^*\|$ .*

*Proof.* Dua and Loeb [4] prove this lemma in the case where  $X = [0, 1]$ ,  $\mathcal{P} = \Pi_{m-1}$ , and  $\mathcal{Q} = \Pi_{n-1}$ , but their proof requires these extra conditions only in proving that if  $Q \geq 0$  on  $X$  and  $P \equiv R^*Q$  on  $X$ , then  $P \equiv P^*$  on  $X$  and  $Q \equiv Q^*$  on  $X$ . This fact, however, follows from an argument of the type given by Cheney [2, p. 165]. Q.E.D.

We can now prove quadratic convergence of the RDDC in some circumstances.

**THEOREM 9.** *Suppose  $X$  contains at least  $m+n+1$  distinct points.  $R^* = P^*/Q^* \in \mathcal{R}_L$  is a best approximation to  $f \in C[X] - \mathcal{R}_L$ , and the space spanned by  $\{\theta_1, \dots, \theta_m, R^*\psi_1, \dots, R^*\psi_n\}$  is a Haar subspace of dimension  $m+n-1$ . Suppose that either of the following two conditions holds:*

- (A)  $Y_0 = \emptyset$  or
- (B)  $0 \in \text{int } \mathcal{H}(S)$ .

*Then the rate of convergence of the restricted-denominator differential correction algorithm is at least quadratic.*

*Proof.* If some  $R_k$  is a best approximation the conclusion of the theorem is true, so we assume this is not the case.

Then the approximations produced by the RDDC algorithm (except possibly  $P_0/Q_0$ ) satisfy the normalization  $\max_{1 \leq j \leq n} |q_j| = 1$ , since otherwise the (negative) minimum computed in step (ii) of the algorithm could be decreased by renormalization. Since we may also assume  $R^*$  satisfies this normalization, the hypotheses of Lemma 2 hold for  $R^*$  and  $R_k$ ,  $k \geq 1$ . We now claim that if condition (A) holds, then  $R^*$  is a best approximation to  $f$  from  $\mathcal{R}_0$ . To see this, we note that by Theorem 3 we have  $\mathbf{0} \in \mathcal{H}(\{\sigma(x) \hat{\mathbf{x}}: x \in X_0\} \cup \{q_j^* \mathbf{e}_{m+j}: j \in I_0\})$  (see Introduction). If there were  $R = P/Q \in \mathcal{R}_0$  satisfying  $\|f - R\| < \|f - R^*\|$ , letting  $\tilde{L}(x) \equiv \frac{1}{2} \min(\min_{x \in X} L(x), \min_{x \in X} Q(x))$  we have  $R^* \in \mathcal{R}_{\tilde{L}}$  and  $R \in \mathcal{R}_{\tilde{L}}$ , with  $Q^* > \tilde{L}$  and  $Q > \tilde{L}$ . But the convex hull statement above and Theorem 3 now imply that  $R^*$  is a best approximation from  $\mathcal{R}_{\tilde{L}}$ , which is a contradiction. Thus, if condition (A) holds, the strong uniqueness of  $R^*$  holds by a theorem in Cheney [2, p. 165], while if condition (B) holds, strong uniqueness holds by Theorem 5. With the strong uniqueness of  $R^*$  and the conclusion of Lemma 2 available, the rest of the proof is as given by Barrodale *et al.* [1, Theorem 3]. Q.E.D.

As noted earlier, condition (A) and (B) of this theorem are mutually exclusive; Example 4 shows that they are not exhaustive. Under the hypotheses of Theorem 9 the absolute values of the differences of the coefficients of  $R_k$  and  $R^*$  can be shown to be bounded by sequences which converge quadratically to zero.

It is sometimes desirable to ignore the function  $f$  at some points of  $X$ , but still apply the denominator restrictions on all of  $X$ ; the theory of this paper goes through unchanged if the subset of  $X$  on which  $f$  is to be approximated is compact. This situation is illustrated in the following example, where  $f$  is the desired magnitude squared response of a digital filter, and we do not wish to approximate  $f$  in the "transition band" (0.1, 0.11).

EXAMPLE 6. Let  $X = \{0, 0.005, 0.01, 0.015, \dots, 0.5\}$ ,

$$\begin{aligned} f(x) &= 1, & 1 \leq x \leq 0.1, \\ &= \text{undefined}, & 0.1 < x < 0.11, \\ &= 0.14, & 0.11 \leq x \leq 0.5, \end{aligned}$$

$$\mathcal{P} = \mathcal{Q} = \left\{ \sum_{i=1}^6 a_i \cos(2\pi(i-1)x) \right\},$$

$L(x) \equiv 0.125$ . Applying the RDDC algorithm on a CDC CYBER 172, which has roughly 14 digits of accuracy, we get (rounded to five places)



$$\begin{aligned}
 R^*(x) = & (0.34408 + 0.35185 \cos 2\pi x + 0.08270 \cos 4\pi x + 0.15552 \cos 6\pi x \\
 & + 0.31629 \cos 8\pi x + 0.23541 \cos 10\pi x)/(1.00000 \\
 & + 0.27473 \cos 2\pi x - 0.56073 \cos 4\pi x - 0.04879 \cos 6\pi x \\
 & + 0.68555 \cos 8\pi x + 0.37588 \cos 10\pi x),
 \end{aligned}$$

with  $\Delta^* = 0.13945$ ,  $X_0 = \{0^+, 0.085^-, 0.1^+, 0.11^-, 0.125^+, 0.215^-, 0.34^+, 0.345^+, 0.405^-, 0.41^-, 0.5^+\}$  (where the sign indicates the sign of  $f - R^*$ ),  $Y_0 = \{0.105\}$ ,  $I_0 = \{1\}$ . Nine iterations were required, with the quantities  $\Delta_k - \Delta^*$  (for  $k = 0, 1, \dots, 8$ ) being approximately  $1 \times 10^{-1}$ ,  $8 \times 10^{-2}$ ,  $2 \times 10^{-2}$ ,  $1 \times 10^{-2}$ ,  $4 \times 10^{-4}$ ,  $1 \times 10^{-6}$ ,  $8 \times 10^{-12}$ ,  $2 \times 10^{-14}$ ,  $8 \times 10^{-15}$ ; this sequence indicates the eventual quadratic nature of the convergence, up to machine accuracy.

We finally remark that the theory of this paper can also be extended to include restricted range conditions and a positive continuous multiplicative weight function  $w$ ; that is, we may further require the functions  $R$  in  $\mathcal{R}_L$  to satisfy  $R \leq u$  on  $X_1$  and  $R \geq l$  on  $X_2$ , where  $u$  and  $l$  are given continuous functions defined on compact subsets  $X_1$  and  $X_2$  of  $X$ , respectively, and we wish to minimize  $\|w \cdot (f - R)\|$  instead of  $\|f - R\|$ . The set  $S$  must be expanded to include the vectors of the form  $\tilde{\sigma}(x) \hat{x}$  ( $x \in Z_0$ ), where  $Z_0 = \{x \in X: R^*(x) = u(x) \text{ or } R^*(x) = l(x)\}$  and

$$\begin{aligned}
 \tilde{\sigma}(x) = -1 & \quad \text{if } R^*(x) = u(x) \\
 = 1 & \quad \text{if } R^*(x) = l(x).
 \end{aligned}$$

In this setting, the assumption  $\max_{x \in X} L(x) < \min_{x \in X} \psi_1(x)$  must be replaced by the assumption that there exists a member of  $\mathcal{R}_L$  satisfying all the restrictions (denominator and restricted range) strictly. The restrictions  $P \leq Q \cdot u$  on  $X_1$  and  $P \geq Q \cdot l$  on  $X_2$  must be added to the RDDC algorithm, and in the expression to be minimized in steps (i) and (ii)  $P(x)$  must be replaced by  $w(x)P(x)$ . The results of this paper still hold essentially as stated (with a few minor changes in the proofs) except for Theorems 2 and 9 and Lemma 1; in Theorem 2 we must add the condition  $\text{sgn}(\bar{P} + R^*\bar{Q})(x) = \tilde{\sigma}(x) \forall x \in Z_0$ ; in Lemma 1 we must add the condition  $\tilde{\sigma}(x)(\bar{P} + R^*\bar{Q}) \geq 0 \forall x \in Z_0$ ; and in Theorem 9 we must add to condition (A) either  $Z_0 = \emptyset$ , or  $f \leq u$  on  $X_1$  and  $f \geq l$  on  $X_2$ . For numerical examples in this extended setting, see McCallig [6].

REFERENCES

1. I. BARRODALE, M. J. D. POWELL, AND F. D. K. ROBERTS, The differential correction algorithm for rational  $l_\infty$  approximation, *SIAM J. Numer. Anal.* **9** (1972), 493–504.
2. E. W. CHENEY, "Introduction to Approximation Theory," McGraw-Hill, New York, 1966.

3. E. W. CHENEY AND H. L. LOEB, Two new algorithms for rational approximation, *Numer. Math.* **3** (1961), 72–75.
4. S. N. DUA AND H. L. LOEB, Further remarks on the differential correction algorithm, *SIAM J. Numer. Anal.* **10** (1973), 123–126.
5. C. M. LEE AND F. D. K. ROBERTS, A comparison of algorithms for rational  $l_\infty$  approximation, *Math. Comp.* **27** (1973), 111–121.
6. M. T. MCCALLIG, R. KURTH, AND B. STEEL, Recursive digital filters with low coefficient sensitivity, in “Proceedings, 1979 IEEE International Conference on Accoustics, Speech, and Signal Processing, Washington D.C., April 1979.”