

## On Certain Order Constrained Chebyshev Rational Approximations

BYRON L. EHLE

*Department of Mathematics, University of Victoria, Victoria, British Columbia, Canada*

*Communicated by Richard S. Varga*

Received March 22, 1974

Some rational approximations which share the properties of Padé and best uniform approximations are considered. The approximations are best in the Chebyshev sense, but the optimization is performed over subsets of the rational functions which have specified derivatives at one end point of the approximation interval. Explicit relationships between the Padé and uniform approximations are developed assuming the function being approximated satisfies easily verified constraints. The results are applied to the exponential function to determine the existence of best uniform  $A$ -acceptable approximations.

### I. INTRODUCTION

In this paper we consider rational approximations to a function  $f(x)$  which share the properties of both Padé and best uniform approximations. We shall require that the function  $f(x)$  being approximated satisfy the basic conditions:

- (1)  $f(x) \in C$  for  $x \in [0, b]$ ,  $0 < b < \infty$ ;
- (2)  $f(x) \in C^M$  at  $x = 0$ , for fixed  $M \geq 1$ ;
- (3)  $d^i f(x)/dx^i |_{x=0} = (i!) c_i$ ,  $i = 0, 1, \dots, M$ .

We then wish to study rational approximations to  $f(x)$  which are best in the Chebyshev sense, but where the optimization is done over subsets of the rationals which have specified derivatives at  $x = 0$ .

Let  $\Pi_m$  denote the collection of all real polynomials of degree at most  $m$  and let  $\Pi_{m,n}$  denote the collection of all real rational functions  $r_{m,n}(x)$  of the form

$$r_{m,n}(x) = q_m^{-1}(x) p_n(x), \quad (1.2)$$

where  $p_n \in \Pi_n$  and  $q_m \in \Pi_m$ . We normalize by requiring  $q(0) = 1$  and assume that  $q(x)$  does not vanish on the interval of approximation. In addition, let  $\Pi_{m,n,k}(f)$  be the subset of  $\Pi_{m,n}$  such that for  $0 \leq k \leq m + n$

$$r_{m,n,k}(x) \in \Pi_{m,n,k}(f) \leftrightarrow (d^i/dx^i) r_{m,n,k}(x)|_{x=0} = (i!) c_i, \quad i = 0, 1, \dots, k. \quad (1.3)$$

Now, consider the error  $\lambda_{m,n,k}$  associated with the best Chebyshev rational approximation of  $f(x)$  by members of  $\Pi_{m,n,k}(f)$  on  $[0, b]$ ,

$$\lambda_{m,n,k} = \inf_{r_{m,n,k} \in \Pi_{m,n,k}(f)} \{ \max_{0 \leq x \leq b} |r_{m,n,k}(x) - f(x)| \}. \quad (1.4)$$

It has recently been shown by Lawson [6] that there exists at least one member  $\bar{r}(x) \in \Pi_{m,n,k}(f)$  for which

$$\max_{0 \leq x \leq b} |\bar{r}(x) - f(x)| = \lambda_{m,n,k}$$

and that a rational function  $\bar{r}_{m-\mu,n-\nu,k}(x)$  is optimal in  $\Pi_{m,n,k}(f)$  in the Chebyshev sense if and only if there exists a set of points  $0 \leq x_1 < x_2 < \dots < x_N \leq b$ ,  $N = m + n + 1 - k - \min(\mu, \nu)$  and a constant  $\lambda$  for which

$$\bar{r}_{m-\mu,n-\nu,k}(x_i) - f(x_i) = (-1)^i \lambda, \quad i = 1, 2, \dots, N. \quad (1.5)$$

In this paper we wish to develop a further characterization of these approximations under the assumption that the function  $f(x)$  is normal of degree  $m + n$ . (The definition of normality, which depends only on the  $c_i$  of (1.1), is given in Section 2). In Section 2 we show that if  $f(x)$  is normal and if (1.3) is satisfied with  $k \geq m$  then  $r_{m,n,k}(x)$  can be written as an  $m + n - k$  parameter function constructed from Padé approximants to  $f(x)$ . As an example, Section 3 considers the problem of finding the best uniform order constrained approximations to the exponential function over the interval  $-\infty < x \leq 0$ . Section 4 is devoted to showing that if  $k \leq m + n - 3$  then the resulting best approximations are not  $A$ -acceptable, that is, they do not satisfy the condition  $|\bar{r}(z)| < 1$  for all  $z$  such that  $\text{Re}(z) < 0$ . Based on results shown previously in [2] and [3], it is shown that if  $k > m + n - 3$  then the best approximations are  $A$ -acceptable.

## 2. CHARACTERIZATION OF $r_{m,n,k}$ USING PADÉ APPROXIMANTS

To establish a connection between the elements of  $\Pi_{m,n,k}(f)$  and Padé approximants to  $f(x)$  we employ the following properties.

DEFINITION. Given a rational function  $r(x) = \frac{\sum_{i=0}^n \gamma_i x^i}{\sum_{i=0}^m \delta_i x^i}$ ,  $\delta_0 = 1$  and a set of constants  $c_j$ ,  $j = 0, 1, \dots, \eta$ ,  $\eta \geq m$ , then if the following system of equations is satisfied

$$\gamma_j - \sum_{i=0}^j \delta_i c_{j-i} = 0, \quad j = 0, 1, \dots, n;$$

$$\sum_{i=0}^m \delta_i c_{j-i} = 0, \quad j = n + 1, \dots, \eta$$

property  $A(\eta)$  is satisfied.

DEFINITION.. Given a set of constants  $c_j$ ,  $j = 0, 1, 2, \dots, \eta + \nu - 1$ ,  $\eta \geq 0$ ,  $\nu \geq 1$ , and  $c_j \equiv 0, j < 0$  if the Hankel determinant [4]

$$H_\nu^\eta = \begin{vmatrix} c_\eta & c_{\eta-1} & \cdots & c_{\eta-\nu+1} \\ c_{\eta+1} & c_\eta & & c_{\eta-\nu+2} \\ \vdots & \vdots & & \vdots \\ c_{\eta+\nu-1} & c_{\eta+\nu} & \cdots & c_\eta \end{vmatrix} \neq 0$$

then property  $B(\eta, \nu)$  is satisfied.

DEFINITION. For a given function  $f(x)$ , if the  $c_i, i = 0, 1, 2, \dots, m + n$  determined by (1.1) with  $M > m + n$  satisfy  $B(\eta, \nu)$  for all  $(\eta, \nu)$  such that  $\nu \leq m + 1$  and  $\eta \leq n + 1$  then  $f(x)$  is said to be normal of degree  $m + n$ .

LEMMA 2.1. If  $f(x)$  is normal of degree  $m + n$  then each entry  $R_{i,j}(x) \in \Pi_{i,j}$  of the Padé table of  $f(x)$  is uniquely determined in lowest terms and has numerator of exact degree  $j$  and denominator of exact degree  $i$ , when  $i \leq m, j \leq n$ .

*Proof.* Because  $f(x)$  is normal,  $R_{i,j}(x)$  satisfies property  $A(i + j)$  and properties  $B(j, i), B(j + 1, i)$ , and  $B(j, i + 1)$ . The uniqueness and nonzero value of the appropriate coefficients follows at once [4].

Assuming  $f(x)$  is normal of degree  $m + n$  we shall denote the unique Padé approximations with numerator of degree  $n - i$  and denominator of degree  $m$  by

$$R_{m,n-i}(x) = \frac{P_{m,n-i}(x)}{Q_{m,n-i}(x)}, \quad i = 0, 1, \dots, n.$$

For  $j \leq n$  we define a  $j$  parameter rational function based on these Padé approximations as follows:

$$\begin{aligned} \mathcal{R}_{m,n,j}(x; \mu_1, \mu_2, \dots, \mu_j) &= \frac{\mathcal{P}_{m,n,j}(x; \mu_1, \mu_2, \dots, \mu_j)}{\mathcal{Q}_{m,n,j}(x; \mu_1, \mu_2, \dots, \mu_j)} \\ &= \frac{(1 - \mu_1 - \mu_2 - \dots - \mu_j) P_{m,n-j}(x) + \sum_{i=1}^j \mu_i P_{m,n-i+1}(x)}{(1 - \mu_1 - \mu_2 - \dots - \mu_j) Q_{m,n-j}(x) + \sum_{i=1}^j \mu_i Q_{m,n-i+1}(x)}. \end{aligned} \quad (2.1)$$

In the next two Lemmas we shall establish conditions which guarantee that if  $r_{m,n,k}(x) \in \Pi_{m,n,k}(f)$  then

$$r_{m,n,k}(x) \equiv \mathcal{R}_{m,n,j}(x; \mu_1, \mu_2, \dots, \mu_j)$$

when  $j$  and  $(\mu_1, \mu_2, \dots, \mu_j)$  are suitably chosen.

LEMMA 2.2. *Let  $f(x)$  be normal of degree  $m + n$  and let  $r_{m,n,k}(x)$  satisfy condition (1.3), where  $k \geq m$ . Then there exists a unique set of constants  $(\mu_1^*, \mu_2^*, \dots, \mu_j^*)$  such that*

$$\begin{aligned} [p_{m,n,k}(x) - q_{m,n,k}(x)f(x)] - [\mathcal{P}_{m,n,j}(x; \mu_1^*, \dots, \mu_j^*) \\ - \mathcal{Q}_{m,n,j}(x; \mu_1^*, \dots, \mu_j^*)f(x)] = O(x^{m+n+1}) \end{aligned} \quad (2.2)$$

where  $j = m + n - k$ .

*Proof.* From the form of  $\mathcal{R}_{m,n,j}$  it is clear that for any  $(\mu_1, \mu_2, \dots, \mu_j)$  the difference in Eq. (2.2) is  $O(x^{k+1})$ .

Denote

$$\begin{aligned} p_{m,n,k}(x), q_{m,n,k}(x), P_{m,n-j}(x), \text{ and } Q_{m,n-j}(x) \quad \text{by} \\ p_{m,n,k}(x) = \sum_{i=0}^n a_i x^i, \quad q_{m,n,k}(x) = \sum_{i=0}^m b_i x^i \quad (2.3) \\ P_{m,n-j}(x) = \sum_{i=0}^{n-j} a_{m,n-j,i} x^i, \quad Q_{m,n-j}(x) = \sum_{i=0}^m b_{m,n-j,i} x^i. \end{aligned}$$

The left-hand side of (2.2) may then be written in the form

$$\begin{aligned} \sum_{i=k+1}^{m+n} d_i x^i - \sum_{i=k+1}^{m+n} \left[ (1 - \mu_1 - \mu_2 - \dots - \mu_j) d_{m,n-j,i} + \sum_{l=1}^j \mu_l d_{m,n-l+1,i} \right] x^i \\ + O(x^{m+n+1}), \end{aligned}$$

where

$$d_i = a_i - \sum_{l=0}^m b_l c_{i-l}, \quad i = k + 1, \dots, m + n; a_i \equiv 0 \text{ for } i > n;$$

and

$$d_{m,n-\nu,i} = a_{m,n-\nu,i} - \sum_{l=0}^m b_{m,n-\nu,l} c_{i-l}, \quad \nu = 0, 1, \dots, j; i = k + 1, \dots, m + n;$$

$$a_{m,n-\nu,i} \equiv 0 \text{ for } i > n - \nu.$$

From the form of the Padé approximants we observe that  $d_{m,n-\nu,i} \equiv 0$  for  $m + n - \nu \geq i$ . It follows that (2.2) is true if and only if a linear system of the following form is satisfied:

$$A\mu^T = \begin{bmatrix} \alpha_{1,1} & \alpha_{1,1} & \cdots & \alpha_{1,1} & \cdots & \alpha_{1,1} \\ \alpha_{2,1} & \alpha_{2,2} & \cdots & \alpha_{2,2} & \cdots & \alpha_{2,2} \\ \vdots & \vdots & & \vdots & & \vdots \\ \alpha_{i,1} & \alpha_{i,2} & \cdots & \alpha_{i,i} & \cdots & \alpha_{i,i} \\ \vdots & \vdots & & \vdots & & \vdots \\ \alpha_{j,1} & \alpha_{j,2} & \cdots & \alpha_{j,i} & \cdots & \alpha_{j,j} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_i \\ \vdots \\ \mu_j \end{bmatrix} = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_i \\ \vdots \\ e_j \end{bmatrix} \\ = e^T, \quad j = m + n - k. \tag{2.4}$$

In particular,  $\alpha_{i,i} = d_{m,n-j,k+i}$ ,  $i = 1, 2, \dots, m + n - k$  and  $\alpha_{i,i-1} = d_{m,n-j,k+i} - d_{m,n-j+i-1,k+i}$ ,  $i = 2, 3, \dots, m + n - k$ . By subtracting column  $j - 1$  from column  $j$ , column  $j - 2$  from column  $j - 1$ , etc., we observe that

$$\det(A) = \alpha_{1,1}(\alpha_{2,2} - \alpha_{2,1})(\alpha_{3,3} - \alpha_{3,2}) \cdots (\alpha_{j,j} - \alpha_{j,j-1}) \neq 0$$

since  $d_{m,n-j+i-1,k+i} \neq 0$  for  $i = 1, 2, 3, \dots, m + n - k$ . Consequently, (2.4) has a unique solution  $(\mu_1^*, \mu_2^*, \dots, \mu_j^*)$  which proves the theorem.

LEMMA 2.3. *Given two functions*

$$r_{m,n}(x) = \sum_{i=0}^n \psi_i x^i / \sum_{i=0}^m \omega_i x^i, \quad s_{m,n}(x) = \sum_{i=0}^n \alpha_i x^i / \sum_{i=0}^m \beta_i x^i$$

which both satisfy property  $A(k)$ ,  $k \geq m$ , and also satisfy the condition

$$\sum_{i=0}^m \omega_i c_{j-i} = d_i = \sum_{i=0}^m \beta_i c_{j-i}, \quad i = k + 1, \dots, m + n; \tag{2.5}$$

then  $r_{m,n}(x) = s_{m,n}(x)$  provided property  $B(n, m)$  is satisfied.

*Proof.* Property  $A(k)$  and Eqs. (2.5) specify two systems of  $m + n + 1$  linear equations in  $m + n + 1$  unknowns for determining the coefficients of  $r_{m,n}(x)$  and  $s_{m,n}(x)$ . The coefficients of these two systems are the same and property  $B(n, m)$  guarantees a unique solution to the system.

**THEOREM 2.1.** *Let  $f(x)$  be normal of degree  $m + n$  and let  $r_{m,n,k}(x)$  satisfy condition (1.3) with  $k \geq m$ . Then there exists a unique set of constants  $(\mu_1^*, \mu_2^*, \dots, \mu_j^*)$  such that*

$$r_{m,n,k}(x) \equiv \mathcal{R}_{m,n,m+n-k}(x; \mu_1^*, \mu_2^*, \dots, \mu_j^*).$$

*Proof.* The result follows at once from Lemmas 2.2 and 2.3

### 3. BEST-ORDER CONSTRAINED APPROXIMATIONS TO $e^x$

It is well known [1, 5] that the Padé approximations to the exponential have the form

$$R_{j,k}(x) = \frac{\sum_{m=0}^k ((j+k-m)! k! / (j+k)! m! (k-m)!) x^m}{\sum_{m=0}^j ((j+k-m)! j! / (j+k)! m! (j-m)! (-x)^m)} \quad (3.1)$$

for all  $j \geq 0, k \geq 0$ , and hence  $e^x$  is a normal function of any degree. It follows from the previous section that if  $r_{m,n,k}(x) \in \Pi_{m,n,k}(e^x)$ ,  $k \geq m$  then  $r_{m,n,k}(x) = \mathcal{R}_{m,n,j}(x; \mu_1, \dots, \mu_j)$  where  $j = m + n - k$ . Because of the continuity of  $e^x$  and the continuity assumption on  $r_{m,n,k}(x)$ , study of the points where  $e^x - r_{m,n,k}(x) = 0$  provides information about the possibility of  $r_{m,n,k}(x)$  satisfying (1.5). The following theorem characterizes the regions in the  $m + n - k$  dimension Euclidean space where exactly  $m + n - k$  exponential fittings occur. Excluded in this count is the exponential fit at  $x = 0$  and also the possible fit at  $-\infty$ .

**THEOREM 3.1.** *For any  $m, n > 0$  and any  $j \leq n$  let  $\mathcal{R}_{m,n,j}(x; \mu_1, \mu_2, \dots, \mu_j)$  denote the functions defined by Eq. (2.1) when  $f(x) = e^x$ . Then there exists a unique set of parameters  $(\mu_1^*, \mu_2^*, \dots, \mu_j^*)$  such that*

$$\mathcal{R}_{m,n,j}(x_i; \mu_1^*, \mu_2^*, \dots, \mu_j^*) - e^{x_i} = 0, \quad i = 1, 2, \dots, j$$

for arbitrary

$$-\infty < x_1 < x_2 < \dots < x_j < 0 \quad \text{where } \mu_i^* \geq 0, i = 1, 2, \dots, j, \text{ and} \\ \mu_1^* + \mu_2^* + \dots + \mu_j^* \leq 1.$$

*Proof.* For  $j = 1$  and  $n = m$  or  $n = m - 1$  the result is established in [2]. For  $j = 2$  and  $n = m$  the result is established in [3]. To establish the general result we shall use induction on  $k$ .

In general, for  $k = 1, n \geq 1$  and any  $x_1 < 0$  we have

$$\mathcal{R}_{m,n,1}(x_1; \mu_1) - e^{x_1} = \frac{(1 - \mu_1)P_{m,n-1}(x_1) + \mu_1 P_{m,n}(x_1)}{(1 - \mu_1)Q_{m,n-1}(x_1) + \mu_1 Q_{m,n}(x_1)} - e^{x_1}.$$

But by [8],  $R_{m,n-1}(x_1) - e^{x_1}$  and  $R_{m,n}(x_1) - e^{x_1}$  differ in sign. Hence, by the same argument used in [2], there exists a unique  $\mu_1^*, 0 \leq \mu_1^* \leq 1$  such that  $\mathcal{R}_{m,n,1}(x_1; \mu_1^*) = e^{x_1}$ .

Now assume that for  $k = \bar{k}$  the theorem is true for all  $n \geq \bar{k}$  and any set of  $\bar{k}$  distinct negative  $x$ 's. Let  $\mathcal{S}_{\bar{k}}$  denote any particular set of these  $x$ 's, and consider the function

$$\mathcal{R}_{m,n,\bar{k}+1}(x; \mu_1, \mu_2, \dots, \mu_{\bar{k}+1}). \tag{3.2}$$

Setting  $\mu_{\bar{k}-1} \equiv 0$ , (3.2) becomes  $\mathcal{R}_{m,n-1,\bar{k}}(x; \mu_1, \mu_2, \dots, \mu_{\bar{k}})$  while if  $\mu_1 \equiv 1 - \mu_2 - \dots - \mu_{\bar{k}+1}$ , (3.2) becomes  $\mathcal{R}_{m,n,\bar{k}}(x; \mu_2, \mu_3, \dots, \mu_{\bar{k}+1})$ . By the assumption, there is a set  $\{\mu_i^*\}_{i=1}^{\bar{k}}$  and a second set  $\{\bar{\mu}_i^*\}_{i=1}^{\bar{k}}$  such that

$$\begin{aligned} \mathcal{R}_{m,n-1,\bar{k}}(x; \mu_1^*, \dots, \mu_{\bar{k}}^*) - e^x &= 0 \\ \mathcal{R}_{m,n,\bar{k}}(x; \bar{\mu}_1^*, \dots, \bar{\mu}_{\bar{k}}^*) - e^x &= 0 \end{aligned} \tag{3.3}$$

for  $x \in \mathcal{S}_{\bar{k}}$ . Now consider any point on the line segment connecting the points  $(\mu_1^*, \mu_2^*, \dots, \mu_{\bar{k}}^*, 0)$  and  $(1 - \bar{\mu}_1^* - \bar{\mu}_2^* - \dots - \bar{\mu}_{\bar{k}}^*, \bar{\mu}_1^*, \bar{\mu}_2^*, \dots, \bar{\mu}_{\bar{k}}^*)$ . Denote the segment by

$$\begin{aligned} \mu_1 &= (1 - \theta)\mu_1^* + \theta(1 - \bar{\mu}_1^* - \bar{\mu}_2^* - \dots - \bar{\mu}_{\bar{k}}^*), \\ \mu_i &= (1 - \theta)\mu_i^* + \theta\bar{\mu}_{i-1}, \quad i = 2, 3, \dots, \bar{k} + 1, \theta \in [0, 1]. \end{aligned} \tag{3.4}$$

Then we have that

$$\mathcal{R}_{m,n,\bar{k}+1}(q; \mu_1, \mu_2, \dots, \mu_{\bar{k}+1}) = (\alpha + \beta\theta)/(\gamma + \delta\theta) \tag{3.5}$$

where  $\alpha, \beta, \gamma$ , and  $\delta$  are constants determined by  $q$ . An expression such as (3.5) is continuous if  $\gamma + \delta\theta \neq 0$  and is strictly increasing, strictly decreasing or constant. Since  $\mu_i \geq 0, i = 1, 2, \dots, \bar{k} + 1$  on the line segment determined by (3.4) we have from (3.1) that if  $q < 0, \gamma + \delta\theta \neq 0$ . Employing equations (3.3) it follows that

$$\mathcal{R}_{m,n,\bar{k}+1}(x; \mu_1, \mu_2, \dots, \mu_{\bar{k}+1}) \equiv e^x \tag{3.6}$$

for  $x \in \mathcal{S}_{\bar{k}}$ . Now consider some point  $x_{\bar{k}+1} < 0, x_{\bar{k}+1} \notin \mathcal{S}_{\bar{k}}$ . Observing that

as  $x \rightarrow -\infty$ ,  $\mathcal{R}_{m,n-1,\bar{k}}(x; \mu_1^*, \mu_2^*, \dots, \mu_k^*) - e^x$  and  $\mathcal{R}_{m,n,\bar{k}}(x; \bar{\mu}_1^*, \bar{\mu}_2^*, \dots, \bar{\mu}_k^*) - e^x$  differ in sign, these expressions must differ in sign for all  $x \notin \mathcal{S}_{\bar{k}}$ . Now

$$\mathcal{R}_{m,n,\bar{k}+1}(x_{\bar{k}+1}; \mu_1, \mu_2, \dots, \mu_{\bar{k}+1})$$

has the form of (3.5) but with different values of  $\alpha, \beta, \gamma$ , and  $\delta$ . Since continuity follows as before, there is a unique  $\bar{\theta}$ ,  $0 < \bar{\theta} < 1$  for which (3.6) is also satisfied with  $x = x_{\bar{k}+1}$ . This completes the proof since  $\bar{\theta}$  produces the required unique set of  $\mu_i$ ,  $i = 1, 2, \dots, \bar{k} + 1$ .

Now consider the problem of determining the best order constrained uniform approximation to  $e^x$  for  $-\infty < x \leq 0$ . Employing the transformation  $x = [-t/(1-t)]$ , we observe that this equivalent to finding the best approximation to

$$f(t) = \exp(-t/(1-t)), t \in [0, 1), \quad f(1) = 0 \quad (3.7)$$

which is a function matching the conditions of Section 1.

From Theorem 3.1 and the results of Lawson [6] summarized in Eq. (1.5) we have the following result.

**THEOREM 3.2.** *For  $k \geq m$ , the best order constrained Chebyshev approximation  $\bar{r}(t) \in \Pi_{m,n,k}(f(t))$ , where  $f(t)$  is defined by (3.7), has the form*

$$r(t) = \mathcal{R}_{m,n,j}(-t/(1-t); \mu_1, \mu_2, \dots, \mu_j), \quad j = m + n - k$$

where  $\mu_i \geq 0$ ,  $i = 1, 2, \dots, j$  and  $\mu_1 + \mu_2 + \dots + \mu_j \leq 1$ .

**COROLLARY 3.1.** *For  $k \geq m$ , the best order constrained uniform approximation  $\bar{r}(x) \in \Pi_{m,n,k}(e^x)$ ,  $-\infty < x \leq 0$ , has the form*

$$r(x) = \mathcal{R}_{m,n,j}(x; \mu_1, \mu_2, \dots, \mu_j), \quad j = m + n - k,$$

where  $\mu_i \geq 0$ ,  $i = 1, 2, \dots, j$ , and  $\mu_1 + \mu_2 + \dots + \mu_j \leq 1$ .

It has recently been shown by Saff and Varga [7] that certain sequences of these best approximations converge geometrically to  $e^x$  on  $-\infty < x \leq 0$ .

#### 4. $\mathcal{A}$ -ACCEPTABILITY OF ORDER CONSTRAINED APPROXIMATIONS TO $e^x$

Using Corollary 3.1 which provides a characterization of the form of best order constrained approximation to  $e^x$  along the negative real axis, we now investigate the  $\mathcal{A}$ -acceptability of these approximations. That is, we ask, which of these best approximations satisfy the condition  $|\bar{r}(z)| < 1$  for all  $z$  such that  $\text{Re}(z) < 0$ .



The  $A$ -acceptability of approximations to  $e^z$  of the form  $\mathcal{R}_{m,n,j}(z; \mu_1, \mu_2, \dots, \mu_j)$ , with  $j = 1$  and  $n = m$ ,  $n = m - 1$  and with  $j = 2$ ,  $n = m$  have been considered previously in [2, 3]. All members of these classes of approximation were shown to be  $A$ -acceptable for  $\mu_i \geq 0$ ,  $i = 1, 2, \dots, j$ ;  $\mu_1 + \mu_2 + \dots + \mu_j \leq 1$ .

Thus we have the following immediate result.

**THEOREM 4.1.** *Best order constrained, uniform approximations to  $e^x$  over  $-\infty < x \leq 0$  are  $A$ -acceptable approximations to  $e^z$  when  $n = m$  and  $k = 2m - 1$  or  $2m - 2$  when  $n = m - 1$  and  $k = 2m - 2$ .*

To illustrate that Theorem 4.1 cannot be generalized to all  $m, n \geq 0$ ,  $k \geq m$  consider

$$\begin{aligned} & \mathcal{R}_{3,1,1}(z; \mu_1) \\ &= \frac{(1 - \mu_1)1 + \mu_1(1 + (z/4))}{(1 - \mu_1)(1 - z + (z^2/2!) - (z^3/3!)) + \mu_1(1 - (3z/4) - (z^2/4) - (z^3/4!))}. \end{aligned}$$

For  $y$  real, consider the difference

$$\begin{aligned} & |\mathcal{Q}_{3,1,1}(iy; \mu_1)|^2 - |\mathcal{P}_{3,1,1}(iy; \mu_1)|^2 \\ &= y^4 \left[ \left[ (1 - \mu_1) \frac{(y - 3^{1/2})}{6} + \mu_1 \left( \frac{y}{24} \right) \right] \left[ (1 - \mu_1) \frac{(y + 3^{1/2})}{6} + \mu_1 \left( \frac{y}{24} \right) \right] \right. \\ & \quad \left. - \frac{(1 - \mu_1)\mu_1}{12} \right]. \end{aligned} \tag{4.1}$$

Using (4.1) it is easily verified that  $|\mathcal{R}_{3,1,1}(iy; \mu_1)| \leq 1$  for all  $y \in (-\infty, \infty)$  when  $0 < \mu_1 < 1$  and thus the best approximation for this case cannot be  $A$ -acceptable.

More generally, using results found in [1-3] it can be shown that

$$\begin{aligned} & |\mathcal{Q}_{m,m-2,1}(iy; \mu_1)|^2 - |\mathcal{P}_{m,m-2,1}(iy; \mu_1)|^2 \\ &= y^{2m-2} \left\{ \left[ (1 - \mu_1)(y - (m^2 - 2m)^{1/2}) \frac{(m - 3)!}{(2m - 3)!} + \mu_1 \frac{(m - 2)! y}{(2m - 2)!} \right] \right. \\ & \quad \times \left[ (1 - \mu_1)(y + (m^2 - 2m)^{1/2}) \frac{(m - 3)!}{(2m - 3)!} + \mu_1 \frac{(m - 2)! y}{(2m - 2)!} \right] \\ & \quad \left. - (1 - \mu_1)\mu_1(m)(m - 2) \left[ \frac{(m - 3)!}{(2m - 3)!} \right]^2 \right\} \end{aligned} \tag{4.2}$$

and we again see that  $|\mathcal{R}_{m,m-2,1}(iy; \mu_1)| \leq 1$  for all  $y \in (-\infty, \infty)$  for  $0 < \mu_1 < 1$ .

Based on the above result and fact that individual Padé approximants on or below the third subdiagonal are not  $A$ -acceptable, we state the following conjecture.

**THEOREM 4.2.** *Best, order constrained, uniform approximations to  $e^x$  over  $-\infty < x \leq 0$  are not  $A$ -acceptable approximations to  $e^z$  for any  $m \geq n \geq 0$  when  $m \leq k \leq 2m - 3$ .*

*Proof.* By producing expansions such as given in (4.2) the correctness of the theorem has been verified for the cases  $n = m$ ,  $n = m - 1$ , and  $n = m - 2$  all with  $k = m + n - 3$ .

In addition, the best approximations produced by Lawson [6] for  $m = n = k$ ,  $m = 2, 3, 4, 5$  were studied and only the approximation for  $m = 2$  was found to be  $A$ -acceptable as expected.

#### ACKNOWLEDGMENT

The results given above were obtained while the author was a visiting professor in the Department of Electrical and Computer Engineering of Syracuse University. Support by both the National Research Council of Canada (Grant No. A7637) and the University of Victoria was also received during the preparation of this paper.

#### REFERENCES

1. B. L. EHLE,  $A$ -stable methods and Padé approximation to the Exponential, *SIAM J. Math. Anal.* **4** (1973), 671–680.
2. B. L. EHLE, Some results on exponential approximation and stiff equations, *SIAM J. Numer. Anal.* to appear.
3. B. L. EHLE AND Z. PICEL, Two parameter, arbitrary order, exponential approximations for stiff equations, *Math. Comp.* **29** (1975), 501–511.
4. D. C. HANDSCOMB, "Methods of Numerical Approximation," Pergamon Press, Oxford, 1966.
5. P. M. HUMMEL AND C. L. SEEBECK, A generalization of Taylor's theorem, *Amer. Math. Monthly* **56** (1949), 243–247.
6. J. D. LAWSON, Order constrained Chebyshev rational approximation, *Math. Comp.* to appear.
7. E. B. SAFF AND R. S. VARGA, Convergence of Padé approximants to  $e^{-z}$  on unbounded sets, *J. Approximation Theory* **13** (1975), 470–488.
8. R. S. VARGA, On higher order stable implicit methods for solving parabolic partial differential equations, *J. Math. Phys.* **40** (1961), 220–231.