

A generalization of the Padé-Approximation to e^{-x} on $[0, \infty)$.

by H.U. Opitz and K. Scherer

0. In a recent paper [4] the authors have introduced a generalization of the Padé-Approximation to e^{-x} to obtain sharp upper bounds of the number

$$R := \limsup_{n \rightarrow \infty} (\lambda_{n,n})^{1/n}$$

where for any integers m, n with $m \leq n$

$$\lambda_{m,n} := \inf_{P_m, Q_n} \|e^{-x} - P_m(x)/Q_n(x)\|_{\infty, [0, \infty)}$$

and P_m, Q_n are polynomials of degree m and n , respectively. The problem of determining R has been posed by Cody-Meinardus-Varga [1] who obtained already the upper bound $R \leq 1/2.298 \dots$. There is strong numerical evidence that R should be equal to $1/9.289 \dots$ according to [1] and a recent paper [6] of Trefethen-Gutknecht who used quite a different method. However these results are limited to finite n .

Rahman-Schmeisser [5] improved the upper bound to $R \leq 1/4.0982 \dots$ and Nemeth [2] claimed to have shown $R \leq 1/6.475 \dots$ by means of Laguerre-Padé-Approximation. However his proof is incomplete. On the other hand the method in [4] allows the construction of rational approximations which yield the bound $R \leq 1/9.033$. In the following we give an outline of these results. In addition we discuss here the connection to other rational approximations of e^{-x} .

1. The generalization of the Padé-Approximation mentioned above is based on the following simple idea:

Given two polynomials $p_m(x), q_n(x)$ of degree m and n , respectively, set

$$(1.1) \quad P_m(x) := \int_a^{\infty} e^{-u} q_n(u) p_m(u+x) du$$

$$(1.2) \quad Q_n(x) := \int_b^{\infty} e^{-u} p_m(u) q_n(u-x) du$$

where a, b are real numbers. An easy computation yields

$$(1.3) \quad e^x \frac{p_m(x)}{q_n(x)} = \frac{\int_b^{a+x} e^{x-u} q_n(u-x) p_m(u) du}{\int_b^\infty e^{-u} q_n(u-x) p_m(u) du}$$

This idea is well known from Padé-Approximation. In fact it constitutes the special case $a=b=0$, $q_n(u) = u^n$ and $p_m = u^m$. However the usefulness of (1.3) for the error analysis of other rational approximations so far does not seem to have been recognized.

We now show first how Laguerre-Padé-Approximation can be covered in this way. To this end take $a=b=0$ and consider the coefficients I_j of the Laguerre-expansion of the numerator in (1.3) (x being replaced by $-x$), thus

$$I_j := \int_0^\infty dx e^{-x} L_j(x) \int_{-x}^0 e^{-(x+u)} q_n(u+x) p_m(u) du \quad j = 0, 1, \dots$$

where $L_j(x)$ denotes the j -th Laguerre-polynomial. Interchanging the variables gives

$$\begin{aligned} I_j &= \int_0^\infty dx e^{-x} L_j(x) \int_0^x e^{-v} q_n(v) p_m(v-x) dv \\ &= \int_0^\infty dv e^{-v} q_n(v) \int_v^\infty e^{-x} L_j(x) p_m(v-x) dx = \int_0^\infty e^{-2v} q_n(v) A_{j,m}(v) dv \end{aligned}$$

where $A_{j,m}(v) := \int_0^\infty e^{-t} L_j(t+v) p_m(-t) dt$. Then the choice

$$(1.4) \quad p_m(-t) = L_m(t)$$

yields $A_{j,m}(v) \equiv 0$ for $j < m$ and hence $I_j = 0$ for $j < m$. Moreover for $j \geq m$ we obtain by the Rodrigues formula for L_j

$$A_{j,m}(v) = \frac{1}{j!} \sum_{l=0}^j \binom{j}{l} v^l \int_0^\infty (e^{-y} y^{j-l})^{(j)} L_m(y) dy = \sum_{l=0}^{j-m} c_{j,m}(l) v^l$$

with constants $c_{j,m}(l)$. It follows for $j \geq m$ that $I_j = \sum_{l=0}^{j-m} c_{j,m}(l) \int_0^\infty e^{-2v} v^l q_n(v) dv$. Hence we have $I_j = 0$ also for $j < m+n$ if we take

$$(1.5) \quad q_n(v) = L_n(2v)$$

Substituting (1.4), (1.5) into (1.1), (1.2) yields explicit expressions for this (m,n) rational approximation of e^{-x} which is essentially equivalent to Laguerre-Padé-Approximation. The only difference is that instead of having $I_{m+n} = 0$ the error vanishes at zero (However (1.3) yields an error expression involving only real integrals in contrast to [2] where a complex integral is needed).

More generally one can ask whether any rational approximation to e^{-x} does admit an error analysis in this way. To this end we investigate first under what circumstances relation (1.1) is invertible. Writing $P_m(x) := \sum_{j=0}^m c_j x^j$, $p_m(u) = \sum_{i=0}^m b_i u^i$ we obtain from (1.1)

$$\sum_{j=0}^m c_j x^j = \sum_{i=0}^m b_i \sum_{j=0}^i \binom{i}{j} x^j \int_a^\infty e^{-u} q_n(u) u^{i-j} du$$

Denoting

$$\gamma_i := \int_a^\infty e^{-u} q_n(u) u^i du$$

a comparison of the coefficients of x^j yields the systems of equations

$$(1.6) \quad c_j = \sum_{i=j}^m \binom{i}{j} \gamma_{i-j} b_i$$

Since it is a triangular system it is uniquely solvable for the b_i provided $\gamma_0 \neq 0$. In view of the first equation $c_m = \gamma_0 b_m$ this is guaranteed if $P_m(x)$ is a polynomial of proper degree m .

So let $P_m(x)$, $Q_n(x)$ be of proper degree m and n , respectively. Then the following iterative procedure to find the corresponding $p_m(x)$, $q_n(x)$ in (1.1), (1.2) makes sense: Given Polynomials $p_m^{(v)}$, $q_n^{(v)}$ of degree m and n , respectively, form

a) Set $\tilde{q}_n^{(v)} := q_n^{(v)} / \|q_n^{(v)}\|$

b) Solve $P_m(x) = \int_a^\infty e^{-u} \tilde{q}_n^{(v)}(u) p_m^{(v+1)}(u+x) du$ for $p_m^{(v+1)}$

c) Set $\tilde{p}_{m+1}^{(v+1)} = p_m^{(v+1)} / \|p_m^{(v+1)}\|$

d) Solve $Q_n(x) = \int_b^\infty e^{-u} \tilde{p}_{m+1}^{(v+1)}(u) q_n^{(v+1)}(u-x) du$ for $q_n^{(v+1)}$

Here $\| \cdot \|$ denotes some norm chosen in advance. Since the steps b), d) always can be carried out, the sequence $(\tilde{p}_m^{(v)}, \tilde{q}_n^{(v)})$ must contain a convergent subsequence with a limit \tilde{p}_m^* , \tilde{q}_n^* , say. So there must exist polynomials p_m^* , q_n^* such that

$$P_m(x) = \int_a^\infty e^{-u} \tilde{q}_n^*(u) p_m^*(u+x) du,$$

$$Q_n(x) = \int_b^\infty e^{-u} \tilde{p}_m^*(u) q_n^*(u-x) du,$$

and $p_m^*/\|p_m^*\| = \tilde{p}_m$, $q_n^*/\|q_n^*\| = \tilde{q}_n$. It follows that the polynomials $\tilde{p}_m := P_m/\|p_m^*\|$ and $\tilde{q}_n := Q_n/\|q_n^*\|$ can be written in the form (1.1), (1.2) with polynomials \tilde{q}_n , \tilde{p}_m . Thus any rational function $P_m(x)/Q_n(x)$ has the representation

$$\frac{P_m(x)}{Q_n(x)} = e^\alpha \frac{\tilde{P}_m(x)}{\tilde{Q}_n(x)}$$

with some constant α and \tilde{P}_m, \tilde{Q}_n being of the form (1.1), (1.2). Setting $e^\alpha \tilde{P}_m(x) := \tilde{P}_m(x-\alpha)$ and $\tilde{Q}_n(x) := \tilde{Q}_n(x-\alpha)$ one easily verifies that \tilde{P}_m, \tilde{Q}_n are of type (1.1), (1.2) with a $-\alpha$ instead of a and representing polynomials $\tilde{p}_m(u)$ and $\tilde{q}_n(u+\alpha)$, respectively. Thus we conclude that any rational function can be written in the form (1.1), (1.2) provided we admit still a translation. In fact, this latter idea was applied in [5] to Padé-Approximation. Optimal choices of the translation and the quotient $r=m/n$ there yielded the above mentioned bound for R .

2. The motivation for our use of (1.3) came from [3] where it is shown that the error of the Padé-Approximation has exactly one extremum on $(-\infty, 0)$ for any fixed n . The idea was then to construct rational approximations of the form (1.1), (1.2) for which the error in (1.3) has several extrema in order to mimic the alternation property of the error of best approximation.

To this end we choose $k+1$ fixed real numbers satisfying

$$(2.1) \quad -\infty = b_0 < b_1 < \dots < b_k = 0 = a_1 < \dots < a_k < a_{k+1} = +\infty$$

and define

$$(2.2) \quad p_m^*(v) = \left[\prod_{j=1}^k (v + na_j) \right]^{m/k}$$

$$(2.3) \quad q_n^*(v) = \left[\prod_{j=1}^k (v + nb_j) \right]^{n/l}$$

where m, n are such that m/k and n/l are even integers. Then for $a = b = 0$ the rational approximations in (1.1), (1.2) take the form

$$P_m^*(-nx) = \int_0^\infty [e^{-u} \prod_{j=1}^k |u + b_j|^{1/l} \prod_{j=1}^k |u - x + a_j|^{r/k}]^n du$$

$$Q_n^*(-nx) = \int_0^\infty [e^{-u} \prod_{j=1}^k |u + a_j|^{r/k} \prod_{j=1}^k |u + x + b_j|^{1/l}]^n du$$

where $r=m/n$. The error in (1.3) becomes

$$(2.4) \quad E_{m,n}(x) = e^{-nx} - \frac{P_m^*(-nx)}{Q_n^*(-nx)} = e^{-nx} \frac{\int_0^{-x} [H(v, v+x)]^n dv}{\int_0^\infty [H(v, v+x)]^n dv}$$

with

$$(2.5) \quad H(v, z) = e^{-v} \left(\prod_{s=1}^k |v + a_s| \right)^{r/k} \left(\prod_{s=1}^l |z + b_s| \right)^{1/l}$$

By this construction we can prove

THEOREM 1: Under certain additional restrictions on the parameters a_s, b_s the holds for $r = m/n \in (0, 1)$ and the above conditions on m, n

$$R \leq \limsup_{n \rightarrow \infty} \left(\| e^{-x} - P_m^*(-x)/Q_n^*(-x) \|_{[0, \infty)} \right)^{1/n} \leq \max_{\substack{0 \leq i \leq l-1 \\ 1 \leq j \leq k}} \exp \| \phi_{ij} \|$$

Each of the $k \cdot l$ numbers $\| \phi_{ij} \|$ can be computed as follows:

Set $I_{ij}(x) := [-x - b_{i+1}, -x - b_i] \cap [-a_{j+1}, -a_j] \cap [-x, 0]$ and define for x with $I_{ij}(x) \neq \emptyset$ the functions

$$(2.6) \quad \phi_{ij}(x) = -x + \log H(v_{ij}(x), x + v_{ij}(x)) - \log H(v_0(x), x + v_0(x))$$

where $v_{ij}(x), v_0(x)$ are given as the unique solutions of

$$(2.7) \quad 0 = -1 + \frac{r}{k} \sum_{s=1}^k \frac{1}{v(x) + a_s} + \frac{1}{l} \sum_{s=1}^l \frac{1}{x + v(x) + b_s}$$

under the constraints $v_{ij}(x) \in I_{ij}(x)$ and $v_0(x) \in \max(0, -x - b_1)$.

Then we have

$$(2.8) \quad \| \phi_{ij} \| = \max_{I_{ij}(x) \neq \emptyset} \phi_{ij}(x) = \phi_{ij}(x_{ij})$$

where x_{ij} is uniquely determined via the equation

$$(2.9) \quad 0 = \phi'_{ij}(x) = -1 + \frac{1}{l} \sum_{s=1}^l \frac{1}{v_{ij}(x) + x + b_s} - \frac{1}{l} \sum_{s=1}^l \frac{1}{v_0(x) + x + b_s}$$

So this theorem says that each set of parameters $\{a_s\}_{s=1}^k, \{b_s\}_{s=1}^l$ yields an upper bound for R which can be computed by solving $k \cdot l$ equations in one variable x where each function evaluation again requires the solution of two non-linear equations. The additional conditions on a_s, b_s we have mentioned above will ensure that this procedure yields a unique solution for each i, j which can be reliably computed by any convergent iteration scheme. We will come back to this point later on.

Next we describe the way leading to the equations (2.7), (2.9). Since the interval $[-x, 0]$ is covered by the union of all $I_{ij}(x)$ (some of them may be empty for fixed x) we obtain from (2.4) for each fixed $x > 0$ the estimate

$$(2.10) \quad |E_{m,n}(x)| \leq k \cdot l e^{-nx} \frac{\max_{ij} \int_{I_{ij}(x)} H(v, v+x)^n dv}{\int_{I_0(x)} H(v, v+x)^n dv}$$

Now from the definition of $I_{ij}(x)$ it is almost clear that $H(v, v+x)$ has exactly one maximum for $v \in I_{ij}(x)$. Since the values of $H(v, v+x)$ are taken to the n -th power it is plain to estimate the integrals in (2.10) by the so-called "methode du col" or Laplace-method.

THEOREM 2: a) For $(i, j) \neq (0, k)$ there holds

$$(2.11) \quad \int_{I_{ij}(x)} H(v, v+x)^n dv \leq \sqrt{\frac{2\pi}{nA_{ij}}} H(v_{ij}(x), v_{ij}(x)+x)^n$$

with

$$A_{ij} := \frac{r}{k} \left[\sum_{s=1}^j (a_{j+1} - a_s)^{-2} + \sum_{s=j+1}^k (a_s - a_j)^2 \right] + \frac{1}{T} \left[\sum_{s=1}^i (b_{i+1} - b_s)^{-2} + \sum_{s=i+1}^1 (b_s - b_i)^2 \right]$$

b) In case $(i, j) = (0, k)$ there holds

$$(2.12) \quad \int_{I_{0,k}(x)} H(v, v+x)^n dv \leq 2(\beta + z_{0k}^+) H(v_{0k}(x), v_{0k}(x)+x)^n$$

where $\beta > 0$ is any number with $e^{-\beta} (1 + \beta/|b_1|) < 1/2$ and $z_{0,k}^+$ is the unique solution in $(-b_1, \infty)$ of

$$(2.13) \quad 1 = \frac{1}{T} \sum_{s=1}^1 \frac{1}{z_{0k}^+ + b_s}$$

c) There holds with $K = \max(k/r, 1)$

$$(2.14) \quad \int_{I_0(x)} H(v, v+x)^n dx \geq \sqrt{\frac{\pi}{2nK}} H(v_0(x), v_0(x)+x)^n$$

Concerning a detailed proof we refer the reader to our paper [4].

COROLLARY: There holds (the lim sup is taken only over even multiples of k and l !)

$$\limsup_{n \rightarrow \infty} \left(\| e^{-x - P_m^*(-x)/Q_n^*(-x)} \|_{\infty, [0, \infty)} \right)^{1/n} \leq \max_{ij} \sup_{I_{ij}(x) \neq \emptyset} \exp \phi_{ij}(x)$$

This follows immediately from the estimates (2.10)–(2.12), (2.14) since they are uniform in x and – except for a factor \sqrt{n} – also uniform in n .

As it is not hard to show that the above lim sup gives also an upper bound for R it remains to compute the maxima of the functions $\phi_{ij}(x)$. Implicit differentiation in (2.6) yields the formulae (2.9) for $\phi'_{ij}(x)$ and furthermore

$$(2.15) \quad \phi''_{ij}(x) = \frac{1}{\alpha_0(x) + \beta_0(x)} - \frac{1}{\alpha_{ij}(x) + \beta_{ij}(x)}$$

where we have set

$$\alpha_0(x) := \left[\frac{r}{k} \sum_{s=1}^k (v_0 + a_s)^{-2} \right]^{-1}; \quad \beta_0(x) = \left[\frac{1}{T} \sum_{s=1}^1 (v_0 + x + b_s)^{-2} \right]^{-1},$$

and $\alpha_{ij}(x), \beta_{ij}(x)$ are defined similarly. It is clear that $\phi_{ij}(x)$ has exactly one maximum (observe $r < 1!$) if we can show that for each x satisfying (2.9) there holds $\phi'_{ij}(x) < 0$. The next lemma gives a sufficient criterion for this to be true. It is a simplified version of corresponding results in [4, Section 6] and we restrict ourselves to the case $k = 1$.

LEMMA: Let $k = 1$, and define $\delta := (v_0^+ + b_1 - b_{i+1} + a_j/k)(v_0^+ + x + b_1)$ with $v_0^+ := v_0(a_{j+1} - b_j)$. Then there holds for each x satisfying (2.9)

$$(2.16) \quad |\phi'_{ij}(x)| \geq \frac{(\alpha_0 + \delta - \bar{\beta}_i)}{(B_0 + \bar{\beta}_i - \delta)(B_0 + \alpha_0)}, \quad 0 \leq i \leq l-1; 1 \leq j \leq k$$

where

$$\begin{aligned} \bar{\beta}_i &:= \left[\frac{1}{T} \sum_{s=1}^i (b_{i+1} - b_s)^{-2} + \frac{1}{T} \sum_{s=i+1}^l (b_s - b_i)^{-2} \right]^{-1} \\ \bar{\beta}_0 &:= \left[\frac{1}{T} \sum_{s=1}^l (z_{0k}^+ + b_s)^{-1} \right]^{-1}, \quad 1 = \frac{1}{T} \sum_{s=1}^l (z_{0k}^+ + b_s)^{-1} \\ \alpha_0 &:= \left[\frac{1}{T} \sum_{s=1}^l (v_0^+ + a_s)^{-2} \right]^{-1}, \quad B_0 := (v_0(a_j - b_{i+1}) + x)^2 \end{aligned}$$

Proof: We start from the equation ($v_0 = v_0(x), v_{ij} = v_{ij}(x)$)

$$(2.17) \quad \frac{1}{k} \sum_{s=1}^k \frac{1}{v_0 + x + b_s} = -\frac{r}{k} \sum_{s=1}^k \frac{1}{v_{ij} + a_s}$$

which follows by substitution of (2.7) into (2.9). Elementary operations then yield

$$\begin{aligned} \beta_0(x)^{-1} &\leq \frac{1}{k(v_0 + x + b_1)} \sum_{s=1}^k \frac{1}{v_0 + x + b_s} = \frac{-r}{k} \sum_{s=1}^k \frac{1}{(v_{ij} + a_s)(v_0 + x + b_1)} \\ &= \alpha_{ij}(x)^{-1} - \frac{r}{k} \sum_{s=1}^k \frac{x + v_{ij} + a_s + v_0 + b_1}{(v_{ij} + a_s)(v_0 + x + b_1)} = \alpha_{ij}(x)^{-1} - S \end{aligned}$$

In view of $x + v_{ij} \in [-b_{i+1}, -b_i]$ the sum S can be estimated as

$$\begin{aligned} (v_0 + x + b_1)S &\geq \frac{r}{k} \sum_{s=1}^k (a_s + v_0 + b_1 - b_{i+1})(v_{ij} + a_s)^{-2} \\ &\geq (v_0 + b_1 - b_{i+1}) \alpha_{ij}(x)^{-1} + r a_j (v_{ij} + a_j)^{-2/k} \\ &\geq [(v_0 + b_1 - b_{i+1}) + a_j/k] \alpha_{ij}(x)^{-1} \end{aligned}$$

Substitution of this relation into the above inequality leads to

$$(2.18) \quad \alpha_{ij}(x) \leq \left(1 - \frac{v_0 + b_1 - b_{i+1} + a_j/k}{v_0 + x + b_1} \right) \beta_0(x) \leq \beta_0(x) - \delta$$

since $\beta_0(x) \geq (v_0 + x + b_1)^2$. This proves the lemma because $\underline{\alpha}_0$ and $\underline{\beta}_i$ are easily seen to be lower and upper bounds of $\alpha_0(x)$ and $\beta_{ij}(x)$, respectively (cf. [4]) so that by (2.18)

$$|\phi''_{ij}(x)| \geq \frac{1}{\beta_0(x) - \delta + \underline{\beta}_i} - \frac{1}{\beta_0(x) + \underline{\alpha}_0}$$

Inequality (2.16) is then an immediate consequence.

Condition (2.16) represents one possible restriction on the parameters a_s, b_s which is sufficient for the validity of Theorem 1. It guarantees that the maxima x_{ij} for the functions $\phi_{ij}(x)$ are uniquely determined as solution of the equation (2.9). For any choice of parameters a_s, b_s (2.16) has to be checked numerically. In particular v_0^+ has to be calculated from (2.7) with $x = a_{j+1} - b_i$ which is an upper bound for x_k with $I_{ij}(x) \neq \emptyset$. In case $j=k$ we calculate v_0^+ as the positive solution of $k = r \sum_{s=1}^k (v+a_s)^{-1}$ which corresponds to $x = +\infty = a_{k+1}$.

To make use of Theorem 1 for the computation of close upper bounds for R one has to develop a strategy for finding optimal parameters a_s, b_s . This may be done by setting up a system of equations for a_s, b_s requiring that $k+1$ of the numbers $\|\phi_{ij}\|$ are equal. One can apply Newton's method to solve the system since there are simple explicit formula for the partial derivatives of $\|\phi_{ij}\|$ with respect to the a_s, b_s available (see [4]).

In this way we constructed rational approximations up to $k=1=60$ parameters in which case we also obtained our best bound (for further results see [4]):

k = 1	r	a_k	b_1	$\max_{i,j} \exp \ \phi_{ij}\ $
60	0.98	876.03939	-1.154006	0.110659
60	0.99	1831.5232	-1.154049	0.110649
60	1	8547.7862	-1.154055	0.1106602

For $r=1$ Theorem 1 has to be modified for $j=k$ since then in general $x_{ik} = +\infty$ (cf. [4]). Certain numerical difficulties only appeared when computing the solutions of (2.9) for large j near k . Indeed, in case $a_j - b_{i+1} > 100$ we determined $\|\phi_{ij}\|$ directly by bisection since even a "damped version" of Newton's method had difficulties in solving (2.9). Then this problem is ill-conditioned described by a bound M for $|\phi''_{ij}(x)|$ for all x near the solution as condition number. If all equations were satisfied up to certain working accuracy tol , the (absolute) error in determining x_{ij} would be $M \cdot \text{tol}$. To get an idea of the size of M let us consider the

