

```

#define NFEW ..
#define NMANY ..

float *c,*d,*e,a,b;
Economize NMANY power series coefficients e[0..NMANY-1] in the range (a,b) into NFEW
coefficients d[0..NFEW-1].

c=vector(0,NMANY-1);
d=vector(0,NFEW-1);
e=vector(0,NMANY-1);
pcshft((-2.0-b-a)/(b-a),(2.0-b-a)/(b-a),e,NMANY);
pccheb(e,c,NMANY);
...
Here one would normally examine the Chebyshev coefficients c[0..NMANY-1] to decide
how small NFEW can be.
chebpc(c,d,NFEW);
pcshft(a,b,d,NFEW);

```

In our example, by the way, the 8th through 10th Chebyshev coefficients turn out to be on the order of  $-7 \times 10^{-6}$ ,  $3 \times 10^{-7}$ , and  $-9 \times 10^{-9}$ , so reasonable truncations (for single precision calculations) are somewhere in this range, yielding a polynomial with 8 – 10 terms instead of the original 13.

Replacing a 13-term polynomial with a (say) 10-term polynomial without any loss of accuracy — that does seem to be getting something for nothing. Is there some magic in this technique? Not really. The 13-term polynomial defined a function  $f(x)$ . Equivalent to economizing the series, we could instead have evaluated  $f(x)$  at enough points to construct its Chebyshev approximation in the interval of interest, by the methods of §5.8. We would have obtained just the same lower-order polynomial. The principal lesson is that the rate of convergence of Chebyshev coefficients has nothing to do with the rate of convergence of power series coefficients; and it is the *former* that dictates the number of terms needed in a polynomial approximation. A function might have a *divergent* power series in some region of interest, but if the function itself is well-behaved, it will have perfectly good polynomial approximations. These can be found by the methods of §5.8, but *not* by economization of series. There is slightly less to economization of series than meets the eye.

#### CITED REFERENCES AND FURTHER READING:

- Acton, F.S. 1970, *Numerical Methods That Work*, 1990, corrected edition (Washington: Mathematical Association of America), Chapter 12.
- Arfken, G. 1970, *Mathematical Methods for Physicists*, 2nd ed. (New York: Academic Press), p. 631. [1]

## 5.12 Padé Approximants

A *Padé approximant*, so called, is that rational function (of a specified order) whose power series expansion agrees with a given power series to the highest possible order. If the rational function is

$$R(x) \equiv \frac{\sum_{k=0}^M a_k x^k}{1 + \sum_{k=1}^N b_k x^k} \quad (5.12.1)$$

then  $R(x)$  is said to be a Padé approximant to the series

$$f(x) \equiv \sum_{k=0}^{\infty} c_k x^k \quad (5.12.2)$$

if

$$R(0) = f(0) \quad (5.12.3)$$

and also

$$\left. \frac{d^k}{dx^k} R(x) \right|_{x=0} = \left. \frac{d^k}{dx^k} f(x) \right|_{x=0}, \quad k = 1, 2, \dots, M + N \quad (5.12.4)$$

Equations (5.12.3) and (5.12.4) furnish  $M + N + 1$  equations for the unknowns  $a_0, \dots, a_M$  and  $b_1, \dots, b_N$ . The easiest way to see what these equations are is to equate (5.12.1) and (5.12.2), multiply both by the denominator of equation (5.12.1), and equate all powers of  $x$  that have either  $a$ 's or  $b$ 's in their coefficients. If we consider only the special case of a diagonal rational approximation,  $M = N$  (cf. §3.2), then we have  $a_0 = c_0$ , with the remaining  $a$ 's and  $b$ 's satisfying

$$\sum_{m=1}^N b_m c_{N-m+k} = -c_{N+k}, \quad k = 1, \dots, N \quad (5.12.5)$$

$$\sum_{m=0}^k b_m c_{k-m} = a_k, \quad k = 1, \dots, N \quad (5.12.6)$$

(note, in equation 5.12.1, that  $b_0 = 1$ ). To solve these, start with equations (5.12.5), which are a set of linear equations for all the unknown  $b$ 's. Although the set is in the form of a Toeplitz matrix (compare equation 2.8.8), experience shows that the equations are frequently close to singular, so that one should not solve them by the methods of §2.8, but rather by full  $LU$  decomposition. Additionally, it is a good idea to refine the solution by iterative improvement (routine `mprove` in §2.5) [1].

Once the  $b$ 's are known, then equation (5.12.6) gives an explicit formula for the unknown  $a$ 's, completing the solution.

Padé approximants are typically used when there is some unknown underlying function  $f(x)$ . We suppose that you are able somehow to compute, perhaps by laborious analytic expansions, the values of  $f(x)$  and a few of its derivatives at  $x = 0$ :  $f(0)$ ,  $f'(0)$ ,  $f''(0)$ , and so on. These are of course the first few coefficients in the power series expansion of  $f(x)$ ; but they are not necessarily getting small, and you have no idea where (or whether) the power series is convergent.

By contrast with techniques like Chebyshev approximation (§5.8) or economization of power series (§5.11) that only condense the information that you already know about a function, Padé approximants can give you genuinely new information about your function's values. It is sometimes quite mysterious how well this can work. (Like other mysteries in mathematics, it relates to *analyticity*.) An example will illustrate.

Imagine that, by extraordinary labors, you have ground out the first five terms in the power series expansion of an unknown function  $f(x)$ ,

$$f(x) \approx 2 + \frac{1}{9}x + \frac{1}{81}x^2 - \frac{49}{8748}x^3 + \frac{175}{78732}x^4 + \dots \quad (5.12.7)$$

(It is not really necessary that you know the coefficients in exact rational form — numerical values are just as good. We here write them as rationals to give you the impression that they derive from some side analytic calculation.) Equation (5.12.7) is plotted as the curve labeled “power series” in Figure 5.12.1. One sees that for  $x \gtrsim 4$  it is dominated by its largest, quartic, term.

We now take the five coefficients in equation (5.12.7) and run them through the routine `padé` listed below. It returns five rational coefficients, three  $a$ 's and two  $b$ 's, for use in equation (5.12.1) with  $M = N = 2$ . The curve in the figure labeled “Padé” plots the resulting rational function. Note that both solid curves derive from the *same* five original coefficient values.

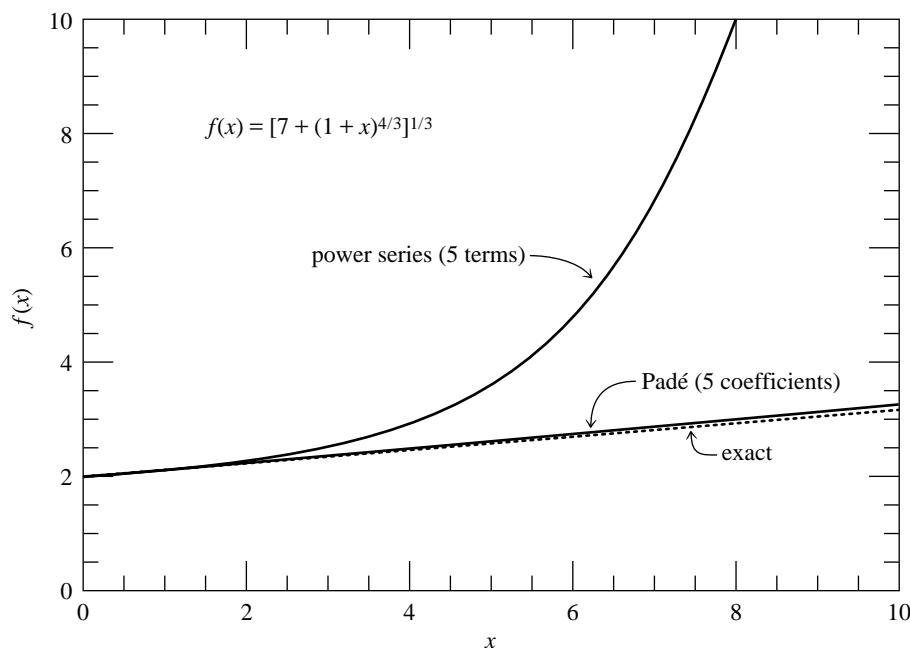


Figure 5.12.1. The five-term power series expansion and the derived five-coefficient Padé approximant for a sample function  $f(x)$ . The full power series converges only for  $x < 1$ . Note that the Padé approximant maintains accuracy far outside the radius of convergence of the series.

To evaluate the results, we need *Deus ex machina* (a useful fellow, when he is available) to tell us that equation (5.12.7) is in fact the power series expansion of the function

$$f(x) = [7 + (1+x)^{4/3}]^{1/3} \quad (5.12.8)$$

which is plotted as the dotted curve in the figure. This function has a branch point at  $x = -1$ , so its power series is convergent only in the range  $-1 < x < 1$ . In most of the range shown in the figure, the series is divergent, and the value of its truncation to five terms is rather meaningless. Nevertheless, those five terms, converted to a Padé approximant, give a remarkably good representation of the function up to at least  $x \sim 10$ .

Why does this work? Are there not other functions with the same first five terms in their power series, but completely different behavior in the range (say)  $2 < x < 10$ ? Indeed there are. Padé approximation has the uncanny knack of picking the function *you had in mind* from among all the possibilities. *Except when it doesn't!* That is the downside of Padé approximation: it is uncontrolled. There is, in general, no way to tell how accurate it is, or how far out in  $x$  it can usefully be extended. It is a powerful, but in the end still mysterious, technique.

Here is the routine that gets  $a$ 's and  $b$ 's from your  $c$ 's. Note that the routine is specialized to the case  $M = N$ , and also that, on output, the rational coefficients are arranged in a format for use with the evaluation routine `ratval` (§5.3). (Also for consistency with that routine, the array of  $c$ 's is passed in double precision.)

```
#include <math.h>
#include "nrutil.h"
#define BIG 1.0e30
```

```
void pade(double cof[], int n, float *resid)
```

Given `cof[0..2*n]`, the leading terms in the power series expansion of a function, solve the linear Padé equations to return the coefficients of a diagonal rational function approximation to the same function, namely  $(\text{cof}[0] + \text{cof}[1]x + \dots + \text{cof}[n]x^n) / (1 + \text{cof}[n+1]x + \dots +$

$\text{cof}[2*n]x^N$ ). The value `resid` is the norm of the residual vector; a small value indicates a well-converged solution. Note that `cof` is double precision for consistency with `ratval`.

```
{
void lubksb(float **a, int n, int *indx, float b[]);
void ludcmp(float **a, int n, int *indx, float *d);
void mprove(float **a, float **alud, int n, int indx[], float b[],
float x[]);
int j,k,*indx;
float d,rr,rrold,sum,**q,**qlu,*x,*y,*z;

indx=ivector(1,n);
q=matrix(1,n,1,n);
qlu=matrix(1,n,1,n);
x=vector(1,n);
y=vector(1,n);
z=vector(1,n);
for (j=1;j<=n;j++) {           Set up matrix for solving.
    y[j]=x[j]=cof[n+j];
    for (k=1;k<=n;k++) {
        q[j][k]=cof[j-k+n];
        qlu[j][k]=q[j][k];
    }
}
ludcmp(qlu,n,indx,&d);         Solve by LU decomposition and backsubstitu-
lubksb(qlu,n,indx,x);         tion.
rr=BIG;
do {                             Important to use iterative improvement, since
    rrold=rr;                     the Padé equations tend to be ill-conditioned.
    for (j=1;j<=n;j++) z[j]=x[j];
    mprove(q,qlu,n,indx,y,x);
    for (rr=0.0,j=1;j<=n;j++)     Calculate residual.
        rr += SQR(z[j]-x[j]);
} while (rr < rrold);          If it is no longer improving, call it quits.
*resid=sqrt(rr);
for (k=1;k<=n;k++) {           Calculate the remaining coefficients.
    for (sum=cof[k],j=1;j<=k;j++) sum -= x[j]*cof[k-j];
    y[k]=sum;
}
for (j=1;j<=n;j++) {           Copy answers to output.
    cof[j]=y[j];
    cof[j+n] = -x[j];
}
free_vector(z,1,n);
free_vector(y,1,n);
free_vector(x,1,n);
free_matrix(qlu,1,n,1,n);
free_matrix(q,1,n,1,n);
free_ivector(indx,1,n);
}
```

#### CITED REFERENCES AND FURTHER READING:

- Ralston, A. and Wilf, H.S. 1960, *Mathematical Methods for Digital Computers* (New York: Wiley), p. 14.
- Cuyt, A., and Wuytack, L. 1987, *Nonlinear Methods in Numerical Analysis* (Amsterdam: North-Holland), Chapter 2.
- Graves-Morris, P.R. 1979, in *Padé Approximation and Its Applications*, Lecture Notes in Mathematics, vol. 765, L. Wuytack, ed. (Berlin: Springer-Verlag). [1]

## 5.13 Rational Chebyshev Approximation

In §5.8 and §5.10 we learned how to find good polynomial approximations to a given function  $f(x)$  in a given interval  $a \leq x \leq b$ . Here, we want to generalize the task to find good approximations that are rational functions (see §5.3). The reason for doing so is that, for some functions and some intervals, the optimal rational function approximation is able to achieve substantially higher accuracy than the optimal polynomial approximation with the same number of coefficients. This must be weighed against the fact that finding a rational function approximation is not as straightforward as finding a polynomial approximation, which, as we saw, could be done elegantly via Chebyshev polynomials.

Let the desired rational function  $R(x)$  have numerator of degree  $m$  and denominator of degree  $k$ . Then we have

$$R(x) \equiv \frac{p_0 + p_1x + \cdots + p_mx^m}{1 + q_1x + \cdots + q_kx^k} \approx f(x) \quad \text{for } a \leq x \leq b \quad (5.13.1)$$

The unknown quantities that we need to find are  $p_0, \dots, p_m$  and  $q_1, \dots, q_k$ , that is,  $m + k + 1$  quantities in all. Let  $r(x)$  denote the deviation of  $R(x)$  from  $f(x)$ , and let  $r$  denote its maximum absolute value,

$$r(x) \equiv R(x) - f(x) \quad r \equiv \max_{a \leq x \leq b} |r(x)| \quad (5.13.2)$$

The ideal *minimax* solution would be that choice of  $p$ 's and  $q$ 's that minimizes  $r$ . Obviously there is *some* minimax solution, since  $r$  is bounded below by zero. How can we find it, or a reasonable approximation to it?

A first hint is furnished by the following fundamental theorem: If  $R(x)$  is nondegenerate (has no common polynomial factors in numerator and denominator), then there is a unique choice of  $p$ 's and  $q$ 's that minimizes  $r$ ; for this choice,  $r(x)$  has  $m + k + 2$  extrema in  $a \leq x \leq b$ , all of magnitude  $r$  and with alternating sign. (We have omitted some technical assumptions in this theorem. See Ralston [1] for a precise statement.) We thus learn that the situation with rational functions is quite analogous to that for minimax polynomials: In §5.8 we saw that the error term of an  $n$ th order approximation, with  $n + 1$  Chebyshev coefficients, was generally dominated by the first neglected Chebyshev term, namely  $T_{n+1}$ , which itself has  $n + 2$  extrema of equal magnitude and alternating sign. So, here, the number of rational coefficients,  $m + k + 1$ , plays the same role of the number of polynomial coefficients,  $n + 1$ .

A different way to see why  $r(x)$  should have  $m + k + 2$  extrema is to note that  $R(x)$  can be made exactly equal to  $f(x)$  at any  $m + k + 1$  points  $x_i$ . Multiplying equation (5.13.1) by its denominator gives the equations

$$p_0 + p_1x_i + \cdots + p_mx_i^m = f(x_i)(1 + q_1x_i + \cdots + q_kx_i^k) \quad (5.13.3)$$

$$i = 1, 2, \dots, m + k + 1$$

This is a set of  $m + k + 1$  linear equations for the unknown  $p$ 's and  $q$ 's, which can be solved by standard methods (e.g., *LU* decomposition). If we choose the  $x_i$ 's to all be in the interval  $(a, b)$ , then there will generically be an extremum between each chosen  $x_i$  and  $x_{i+1}$ , plus also extrema where the function goes out of the interval at  $a$  and  $b$ , for a total of  $m + k + 2$  extrema. For arbitrary  $x_i$ 's, the extrema will not have the same magnitude. The theorem says that, for one particular choice of  $x_i$ 's, the magnitudes can be beaten down to the identical, minimal, value of  $r$ .

Instead of making  $f(x_i)$  and  $R(x_i)$  equal at the points  $x_i$ , one can instead force the residual  $r(x_i)$  to any desired values  $y_i$  by solving the linear equations

$$p_0 + p_1x_i + \cdots + p_mx_i^m = [f(x_i) - y_i](1 + q_1x_i + \cdots + q_kx_i^k) \quad (5.13.4)$$

$$i = 1, 2, \dots, m + k + 1$$