

This norm has the following advantages.

- (a) It allows us to evaluate partial derivatives of the norm in terms of polynomial and series manipulations. These can be used to express a sequence of least squares problems, whose solutions usually converge to a minimum perturbation $\|f - g\| = \|f - g - h\|$. The derivatives can also be used for Newton's method.
- (b) Minimizing $\|f\|$ gives a near-Chebyshev minimum on the unit disk [13].
- (c) It permits fast algorithms for the solution of subproblems at each iteration.

The expression of $\|f\|$ in the form (2) emphasizes the importance of the size of the values of $f(z)$ on the unit disk. This highlights the need for the following assumptions regarding the formulation of the problem:

- (a) The location of the origin has been chosen (thus making explicit an implied assumption in previous numerical polynomial algorithms),
- (b) The scale of $|z|$ has been chosen.

In particular, we assume that the problem context precludes a change of variable by an affine transformation $z = bz + a$.

Remark. There is also a purely computational reason for avoiding such transformations, as is set out in the next theorem.

Theorem 2.2. *Shifting from z to $z - a$ can amplify any uncertainties in the coefficients of f by an amount as much as $(1 + |a|)^n / (n + 1)$ in norm. This is exponential in n , for any $a \neq 0$. Moreover, the relative uncertainties in each coefficient can be amplified by arbitrarily large amounts.*

In other words, such shifts are ill-conditioned.

Proof. By examining the condition of the matrix that determines the Taylor coefficients of the shifted polynomial

$$f_a(z) = \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} (z - a)^k,$$

one quickly finds the worst case for perturbation in the 1-norm: choose $f(z) = z^n$. Then $\|f\|_1 = \|f\|_1 = 1$, but

$$f_a(z) = \sum_{k=0}^n \binom{n}{k} (-a)^k z^k$$

and hence $\|f_a\|_1 = (1 + |a|)^n$. Since the 1- and 2-norms are equivalent, with $\|f_a\|_2$

setting is that f is not monic, so

$$\begin{aligned} f(z) &= f_n + f_{n-1}z + f_{n-2}z^2 + \cdots + f_0z^n \\ &= h^{(0)}(z) = f_n^{1/m} + h_{d-1}z + \cdots + h_0z^d \end{aligned}$$

and for definiteness we choose the positive (real) root for h_d (recall that we scaled f so that $f = 1$ with leading coefficient greater than 0).

On the quality of the initial approximation

Let h_{min} and g_{min} be the functions that minimize f . Further, let $f + f_{min} = g_{min} h_{min}$. A bound showing the quality of $h^{(0)}$ is given by the next theorem.

Theorem 3.1. *There exists a constant K , depending on m and on the leading coefficient of*

Repeat steps 2 and 3 until sufficient accuracy is attained or your patience is exhausted. When this method converges, it converges linearly since it is just functional iteration (similar to that discussed in [5]).

Essentially, we ignore the interactions between the changes in h and the changes in g . By doing so, we for-

The normal equations (12) can be arranged to get

$$\sum_{k=0}^d T_k h = b_k, \quad 0 \leq k \leq d,$$

where

$$T_k = [z^{k-}] g(h(z)) \bar{g}(\bar{h}(1/z))$$

$$b_k = [z^k] (f(z) - g(h(z)) \bar{g}(\bar{h}(1/z))).$$

This derivation allows for a very fast computation of the entries in T through series manipulation. To solve stably and efficiently such a system it is also necessary to know that it is non-singular and positive definite as well as Hermitian and Toeplitz. To see this, we observe that T factors as $T = B B^*$, where B is an $(n + d) \times (d + 1)$ matrix with $B_k = [z^k] z g(h(z))$. This is a lower triangular Toeplitz matrix of full column rank when $g(h(z))$ is non-zero, whence $B B^*$ is non-singular and positive definite.

Once we have computed a h using this linear least-squares formulation we may then update $h := h + \delta h$. The entire process can then be iterated by re-linearizing around this new h and again approximating a h minimizing $f - g(h + \delta h)$. Since this nonlinear least-squares problem is only a component in the entire solver it is not clear that it is necessary to repeat this local process (for fixed f and g) until convergence occurs. However, we have seen examples in which substantial convergence is required in this sub-problem for a globally minimal decomposition to be obtained.

Computationally, each iteration of this nonlinear least-squares solver has very low cost. Each system T can be easily constructed using only series manipulations. T can be constructed with $O(n \log^2 n)$ operations using the series manipulation algorithms of Brent & Kung [4] and an FFT for polynomial multiplication. The solution to the system can be obtained via the stable Toeplitz solvers of Trench [14] using (d^2) operations, or the fast and stable methods (for positive definite matrices) [3, 12] which require $O(d \log^2 d)$ operations. In summary, each iteration requires $O(n \log^2 n)$ floating point operations.

5 NEWTON ITERATION

In this section we explore the direct use of Newton's method to solve the nonlinear minimization problem: find $g, h \in \mathbb{R}[z]$ minimizing $\|g h - f\|^2$. We give an effective method of computing the requisite derivatives analytically, and implement, test, and compare the method with the sequence of linear least-squares problems of the earlier section.

We consider

$$N_f(g + \delta g, h + \delta h) = f - (g + \delta g)(h + \delta h)^2$$

with

$$h(z) = \sum_{k=0}^d h_k z^k, \quad g(z) = \sum_{k=0}^{m-1} g_k z^k.$$

Assume for the purpose of exposition that $f, g, h, \delta f, \delta g$, and δh are all in $\mathbb{R}[z]$; however $z \in \mathbb{C}$. Lemma 2.1 is used to compute N_f . Denoting the integrand of the integral for

N_f by I_f , we expand to second order in $\delta g, \delta h$.

$$I_f(g + \delta g, h + \delta h) = (g(h(z)) - f(z))(g(h(\bar{z})) - f(\bar{z}))$$

$$+ g(h(\bar{z}))(g(h(z)) - f(z)) h(\bar{z})$$

$$+ (g(h(z)) - f(z)) g(h(\bar{z})) h(z)$$

$$+ g(h(z))g(h(\bar{z})) h(z) h(\bar{z})$$

$$+ (g(h(z)) g(h(\bar{z})))$$

$$+ g(h(\bar{z})) g(h(z)) h(\bar{z})$$

$$+ (g(h(z)) - f(z)) g(h(\bar{z})) h(z)$$

$$+ \frac{1}{2} g(h(\bar{z}))(g(h(z)) - f(z)) h^2(\bar{z})$$

$$+ \text{c. c.},$$

where c. c. indicates the complex conjugate of all non-real summands.

We write $\delta g = \sum_{k=0}^{m-1} \delta g_k z^k$, $\delta h = \sum_{k=0}^d \delta h_k z^k$.

small eigenvalues, of course, may be greatly perturbed in a relative sense, but, as we will see below, to stabilize the Newton step we will ignore eigenvalues that are too small.

We write $\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^t$, where $\mathbf{\Lambda}$ is the usual diagonal matrix of eigenvalues and \mathbf{Q} is orthonormal. Now substitute $\mathbf{x} = \mathbf{Q}\mathbf{y}$ in (13) to get

$$\begin{aligned} N_f + \mathbf{b}^t \mathbf{x} + \mathbf{x}^t \mathbf{A} \mathbf{x} &= N_f + \mathbf{b}^t \mathbf{Q} \mathbf{y} + \mathbf{y}^t \mathbf{Q}^t \mathbf{A} \mathbf{Q} \mathbf{y} \\ &= \mathbf{b}^t \mathbf{Q} \mathbf{y} + \mathbf{y}^t \mathbf{\Lambda} \mathbf{y}. \end{aligned}$$

The constant N_f can be dropped. Denoting $\mathbf{b}^t \mathbf{Q} = -2\mathbf{p}^t$, we have

$$\begin{aligned} \mathbf{b}^t \mathbf{Q} \mathbf{y} + \mathbf{y}^t \mathbf{\Lambda} \mathbf{y} &= -2p_1 y_1 - 2p_2 y_2 - \cdots - 2p_m y_m \\ &\quad + y_1^2 + y_2^2 + \cdots + y_m^2 \end{aligned}$$

Least Squares Iteration 2

$$g = 0.006265045950 + 0.$$