

18.335 Problem Set 4 Solutions

Problem 1: (10+10+10 points)

- (a) (Essentially the same recurrence is explained in equation 30.8 in the text.) You can derive this recurrence relation by recalling the expression for determinants in terms of cofactors or minors that you should have learned in a previous linear-algebra course. In $\det B$, the coefficient the (i, j) entry B_{ij} is $(-1)^{i+j}$ multiplied by the determinant of the submatrix (minor) of B with the i -th row and j -th column removed (this determinant, multiplied by $(-1)^{i+j}$, is the *cofactor* of B_{ij}). The cofactor formula tells us that the determinant is equal to any row (or column) of B multiplied by its cofactors. Take the last row of B , which has only two nonzero entries B_{mm} and $B_{m,m-1}$. The cofactor of B_{mm} is $+\det B_{1:m-1,1:m-1}$. The cofactor of $B_{m,m-1}$ is the -1 times the determinant of the matrix whose last column is all zeros except for $B_{m-1,m}$ in the lower-right. To get the determinant of that matrix we take everything in that column and multiply by their cofactors; everything is zero except for $B_{m-1,m}$, whose cofactor is $+\det B_{1:m-2,1:m-2}$. Putting these together, we have $\det B = B_{m,m} \det B_{1:m-1,1:m-1} - B_{m-1,m} B_{m,m-1} \det B_{1:m-2,1:m-2}$ as desired.

This gives a three-term recurrence relation, since the determinant of each submatrix is given in terms of the the determinants of the two smaller submatrices. At each iteration, we only need the determinants of the submatrices of size $n \times n$ and $(n-1) \times (n-1)$ to get the determinant of the submatrix of size $(n+1) \times (n+1)$ using three multiplications and one addition, or four flops. We have to apply the recurrence $m-1$ times to get the full $\det(H-zI)$, for $4(m-1) + m = O(m)$ flops (where the $+m$ flops was for subtracting z from the m diagonals). The code is given in part (b) below (including the derivative).

To test it, we computed a random H matrix as suggested in the problem, for $m=20$, and computed the relative error $(\text{evalpoly}(H,z) - \det(H-z*\text{eye}(m))) / \det(H-z*\text{eye}(m))$ for a few randomly chosen values of z . The relative error was always on the order of the machine precision $\approx 10^{-16}$.

- (b) If we simply take the derivative $\frac{d}{dz}$ of our recurrence relation for $\det B$, we get a three-term recurrence relation for $(\det B)'$:

$$(\det B)' = -\det B_{1:m-1,1:m-1} + B_{m,m}(\det B_{1:m-1,1:m-1})' - B_{m-1,m}B_{m,m-1}(\det B_{1:m-2,1:m-2}),$$

using the fact that for $B = H - zI$ we have $B' = -I$ and hence $B'_{mm} = -1$ while $B'_{m-1,m} = B'_{m,m-1} = 0$.

This is implemented in the following Matlab function:

```
function [p,pderiv] = evalpoly(H,z)
    m = size(H,1);
    d1 = (H(1,1) - z);
    d2 = 1;
    d1deriv = -1;
    d2deriv = 0;
    for i = 2:m
        d = (H(i,i) - z) * d1 - H(i-1,i)*H(i,i-1) * d2;
        dderiv = -d1 + (H(i,i) - z) * d1deriv - H(i-1,i)*H(i,i-1)* d2deriv;
        d2 = d1;
        d1 = d;
        d2deriv = d1deriv;
        d1deriv = dderiv;
    end
end
```

```
p = d1;
pderiv = d1deriv;
```

If we choose $\Delta z = 10^{-4}$, and compare $[p(z + \Delta z) - p(z - \Delta z)]/2\Delta z$ to `pderiv` computed by this function for a random H (as above) and for various random $z \in [0, 1]$, then the relative difference is usually on the order of 10^{-6} or so. Decreasing Δz to 10^{-6} , the relative difference is around 10^{-10} or so; exactly the convergence of $O(\Delta z^2)$ that one would expect: Taylor expanding $p(z + \Delta z) = p(z) + p'(z)\Delta z + p''(z)\Delta z^2/2 + O(\Delta z^3)$, one finds $[p(z + \Delta z) - p(z - \Delta z)]/2\Delta z = p'(z) + O(\Delta z^2)$. Of course, you can't make Δz too small or you run into the limitations of floating-point precision (e.g. $\Delta z = 10^{-20}$ will just give zero for $p(z + \Delta z) - p(z - \Delta z)$ if $z \sim 1$).

- (c) Newton's method is just the iteration $z \leftarrow z - p(z)/p'(z)$. You can do this “by hand” in Matlab, but I wrote a little Matlab script to repeat this until the answer stops changing to within $10\epsilon_{\text{machine}}$. (Better stopping criteria could be devised, but this works well enough for illustration purposes.)

```
function lam = newtpoly(H, lamguess)
    lam = lamguess;
    for i = 1:1000
        [p,pderiv] = evalpoly(H, lam);
        disp(lamnew = lam - p / pderiv);
        if (abs(lamnew - lam) < 10*eps*abs(lamnew))
            break;
        end
        lam = lamnew;
    end
end
```

Applying this to my random A and H from above, I run `eig(A)` to get the eigenvalues, and then execute the Newton solver with starting points that are somewhere nearby one of the eigenvalues (e.g. the eigenvalue rounded to two significant digits). You can run the command `format long` to have Matlab print out values to 15 decimal places, to better see what is going on. A typical output is:

```
>> newtpoly(H, 0.96)
0.961143436106079
0.961118817820461
0.961118806268620
0.961118806268617
0.961118806268617
ans = 0.961118806268617
```

In comparison, the corresponding eigenvalue returned by `eig(A)` is $\lambda \approx 0.961118806268619$, which differs only in the 15th decimal place from the Newton result. If you repeat this with several eigenvalues (i.e., different starting points), you find that it consistently finds λ to nearly machine precision. Notice that the number of accurate digits roughly doubles on each step, which is typical of Newton's method: this is “quadratic” convergence.

If you tried a test case with a bigger value of m , say $m = 1000$, you probably noticed a problem: the determinant grows so large that it overflows the range of floating-point arithmetic, leading to results of NaN or Inf! However, this is easily fixed. For example, as one is computing the determinant, one can periodically check the magnitude and rescale to prevent it from growing too large (or too small)—obviously, $p(z)$ multiplied by any constant still has the same roots. There are also slightly more clever recurrences

that rescale the result automatically; See Barth, Martin, and Wilkinson (1967).¹ The other key to finding eigenvalues in this way, described in lecture 30 of Trefethen, is that there is a simple technique to count the number of eigenvalues in any given interval, leading to the “bisection” algorithm I mentioned in class. I didn’t expect you to do any of that stuff however; I hope you all just picked a small enough m to avoid overflow.

Problem 2: Q’s ‘R us (10+15 points)

- (a) In finite precision, instead of $w = A^{-1}v$, we will get $\tilde{w} = w + \delta w$ where $\delta w = -(A + \delta A)^{-1} \delta A w$ (from the formula on page 95), where $\delta A = O(\epsilon_{\text{machine}}) \|A\|$ is the backwards error. [Note that we cannot use $\delta w \approx -A^{-1} \delta A w$, which neglects the $\delta A \delta w$ terms, because in this case δw is not small.] The key point, however, is to show that δw is mostly parallel to q_1 , the eigenvector corresponding to the smallest-magnitude eigenvalue λ_1 (it is given that all other eigenvalues have magnitude $\geq |\lambda_2| \gg |\lambda_1|$). Since w is also mostly parallel to q_1 , this will mean that $\tilde{w}/\|\tilde{w}\|_2 \approx q_1 \approx w/\|w\|_2$.

First, exactly as in our analysis of the power method, note that $w = A^{-1}v = \alpha_1 q_1 [1 + O(\lambda_1/\lambda_2)]$, since A^{-1} amplifies the q_1 component of v by $1/|\lambda_1|$ which is much bigger than the inverse of all the other eigenvalues. Thus, $w/\|w\|_2 = q_1 [1 + O(\lambda_1/\lambda_2)]$.

Second, if we Taylor-expand $(A + \delta A)^{-1}$ in powers of δA , i.e. in powers of $\epsilon_{\text{machine}}$, we obtain:² $(A + \delta A)^{-1} = A^{-1} - A^{-1} \delta A A^{-1} + O(\epsilon_{\text{machine}}^2)$. Since all of the terms in this expansion are multiplied on the *left* by A^{-1} , when multiplied by *any* vector they will again amplify the q_1 component much more than any other component. In particular, the vector $\delta A w$ is a vector in a random direction (since δA comes from roundoff and is essentially random) and hence will have some nonzero q_1 component. Thus, $\delta w = -(A + \delta A)^{-1} \delta A w = \beta_1 q_1 [1 + O(\lambda_1/\lambda_2)]$ for some constant β_1 .

Putting these things together, we see that $\tilde{w} = (\alpha_1 + \beta_1) q_1 [1 + O(\lambda_1/\lambda_2)]$, and hence $\tilde{w}/\|\tilde{w}\|_2 = q_1 [1 + O(\lambda_1/\lambda_2)] = \frac{w}{\|w\|_2} [1 + O(\lambda_1/\lambda_2)]$. Q.E.D.

- (b) Trefethen, problem 28.2:

- (i) In general, r_{ij} is nonzero (for $i < j$) if column i is non-orthogonal to column j . For a tridiagonal matrix A , only columns within two columns of one another are non-orthogonal (overlapping in the nonzero entries), so R should only be nonzero (in general) for the diagonals and for two entries above each diagonal; i.e. r_{ij} is nonzero only for $i = j$, $i = j - 1$, and $i = j - 2$.

Each column of the Q matrix involves a linear combination of all the previous columns, by induction (i.e. q_2 uses q_1 , q_3 uses q_2 and q_1 , q_4 uses q_3 and q_2 , q_5 uses q_4 and q_3 , and so on). This means that an entry (i, j) of Q is zero (in general) only if $a_{i,1:j} = 0$ (i.e., that entire row of A is zero up to the j -th column). For the case of tridiagonal A , this means that Q will have upper-Hessenberg form.

- (ii) **Note:** In the problem, you are told that A is symmetric and tridiagonal. You must also assume that A is real, or alternatively that A is Hermitian and tridiagonal. (This is what must be meant in the problem, since tridiagonal matrices only arise in the QR method if the starting point is Hermitian.) In contrast, if A is complex tridiagonal with $A^T = A$, the stated result is not true (RQ is not in general tridiagonal, as can easily be verified using a random tridiagonal complex A in Matlab).

¹W. Barth, R. S. Martin, and J. H. Wilkinson, “Calculation of the eigenvalues of a symmetric tridiagonal matrix by the bisection method,” *Numer. Math.* **9**, 386–393 (1967).

²Write $(A + \delta A)^{-1} = [A(I + A^{-1} \delta A)]^{-1} = (I + A^{-1} \delta A)^{-1} A^{-1} \approx (I - A^{-1} \delta A) A^{-1} = A^{-1} - A^{-1} \delta A A^{-1}$. Another approach is to let $B = (A + \delta A)^{-1} = B_0 + B_1 + \dots$ where B_k is the k -th order term in δA , collect terms order-by-order in $I = (B_0 + B_1 + \dots)(A + \delta A) = B_0 A + (B_0 \delta A + B_1 A) + \dots$, and you immediately find that $B_0 = A^{-1}$, $B_1 = -B_0 \delta A A^{-1} = -A^{-1} \delta A A^{-1}$, and so on.

