

AmATH 584

Hw 5

Nate Whybra

12.1) We have

$$101 = \|A\|_F = \sqrt{\sigma_1^2 + \sum_{i=2}^{202} \sigma_i^2}$$

$$= \sqrt{\|A\|_2^2 + \sum_{i=2}^{202} \sigma_i^2}$$

$$= \sqrt{100^2 + \sum_{i=2}^{202} \sigma_i^2} \quad (1)$$

$$\Leftrightarrow \sum_{i=2}^{202} \sigma_i^2 = 101^2 - 100^2 = 201$$

suppose  $\sigma_i = \sigma_j \quad \forall i, j \geq 2$ , then

$$201 \sigma_i^2 = 201$$

$$\Leftrightarrow \sigma_i = 1 \quad (\text{since singular values are positive})$$

If all  $\sigma_i$  with  $i \geq 2$  were  $> 1$  the formula

(1) would not hold. So it must be that  $\exists i$  with  $\sigma_i \leq 1$ , so that the RHS of (1)  $\leq 101$ .

However, equality happens by setting  $\sigma_i = 1$  as demonstrated above, for all  $i \geq 2$ . So

$$K(A) = \frac{\sigma_1}{\sigma_{202}} \leq \frac{\sigma_1}{1} = \|A\|_2 = 100$$

and this is the tightest possible bound, because we have shown the smallest singular value of  $A$  is at least equal to 1.

21.1)

a) We have

$$\begin{bmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{bmatrix} = \begin{bmatrix} 1 & & & \\ 2 & 1 & & \\ 4 & 3 & 1 & \\ 3 & 4 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 & 0 \\ & 1 & 1 & 1 \\ & & 2 & 2 \\ & & & 2 \end{bmatrix}$$

So  $\det(A) = \det(L) \cdot \det(U)$ , but since both matrices are triangular, their determinants are the products of the diagonals, hence,

$$de + (A) = (1)^4 \cdot (2)^3 \cdot 1 = 8$$

b) we have,

$$P^{-1}A = L^{-1}U$$

we have  $\det(P) \cdot \det(A) = \det(L) \cdot \det(U)$ .  
 $\det(P) = -1$  since it is an odd number of row swaps from  $I$ .  
 As  $L$  and  $U$  are triangular, the determinants  
 are the products of the diagonals. So  
 $\det(L) = 16$  and  $\det(U) = -8 \cdot (-1) = +8$

$$\det(A) = -(1)^4 \cdot 8 \cdot \frac{\pi}{24} \cdot \frac{-\sqrt{6}^2}{\pi} \cdot \frac{1}{3} = -8 \cdot (-1) = +8$$

c) To compute  $\det(A)$  from  $PA = LU$  we can compute the product of the products of the diagonals of  $L$  and  $U$ , and multiply it by  $(-1)^{\# \text{ of row swaps for } P}$  however,  $L$  has 1's on the diagonals which simplifies things further, in general we can write

$$\det(A) = (-1)^{m - \sum_{i=1}^m P_{ii}} \cdot \prod_{i=1}^m \text{diag}(U)_i$$

where  $m - \sum_{i=1}^m P_{ii}$  counts how many columns of  $P$  are "not in the right place" considering if there were no permutations,  $P$  would be the identity  $I$ .

A1) If we make  $\beta = 7$  with finite precision  $t$ ,  $\frac{8}{7}$  can be represented as  $1.1$ , which can be verified by noticing that  $1 \cdot 7^0 + 1 \cdot 7^{-1} = 1 + \frac{1}{7} = \frac{8}{7}$ . There is no base  $\beta$  in which we can represent  $\pi$  with finite precision to exactly since it can be shown that irrational numbers have infinitely long non-repeating decimal representations in any <sup>integer</sup> base  $\beta$ . I proved this in combinatorics class once.



A2)

$$a) K = \frac{\|f'(1.2)\|}{\|f(1.2)\|} = (1.2) \frac{|\cos(1.2)|}{|\sin(1.2)|}$$

$$= (1.2) |\cot(1.2)| \approx 0.467 < 10^2$$

Since the <sup>relative</sup> condition number of this problem is small, it is well conditioned.

b) The plot does not behave as I expect. The error is large when  $h$  is large (as expected), but also large when  $h$  is small (not expected), and there seems to be a minimal error on the order of  $10^{-10}$  when  $10^{-10} < h < 10^{-5}$ .

c) (see next page)

c) Assuming these rounding errors we have that our derivative approximation is

$$\approx \underbrace{\frac{f(x_0+h) - f(x_0)}{h}}_{\text{(1) derivative approximation}} + \underbrace{O\left(\frac{\epsilon_{\text{machine}}}{h}\right)}_{\text{(2) rounding error}}$$

when  $h$  is large our error from term (1) is large but the error from term (2) is small, and when  $h$  is small the error in term (1) is small but the error from term (2) is large. This implies there should be some value of  $h$  where the error is minimized (a value of  $h$  where the contribution of errors from both terms is <sup>jointly</sup> minimized), which explains our observations from (b).

d) Using this alternate formula,  
the results are more like what I would  
expect. The error is small when  $h$   
is small and large when  $h$  is large,  
it is a linear function on this  $\log \log$   
scale.



A3)

a) The condition number is  $44.8023 < 10^2$ , so yeah, the matrix is well-conditioned.

b) we have  $\|x - x_{ge}\|_2 \approx 8.05$ , this is saying the GE solution is noticeably different from our exact solution, so I don't trust this solution.

c) we have  $\|x - x_{qr}\|_2 \approx 2.7e-14$  which is very very small, so I trust this solution.

d) we have  $\|x - x_{geop}\|_2 \approx 3.4e-15$  which is very very small, so I trust this solution.

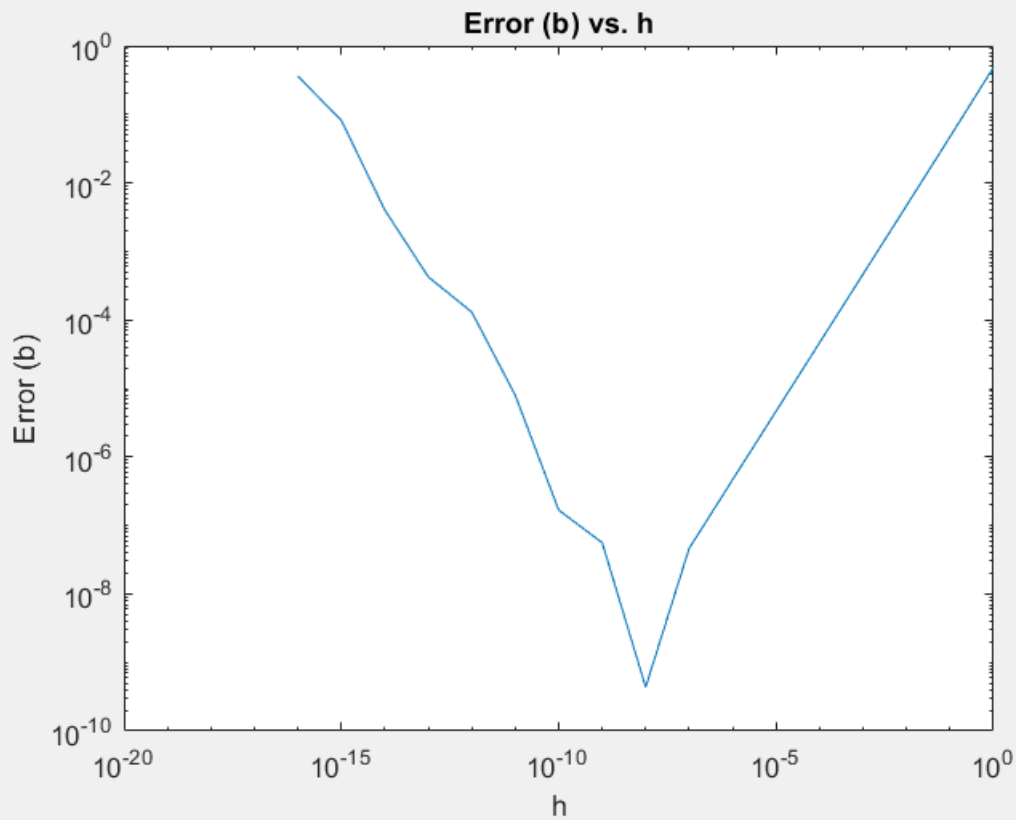


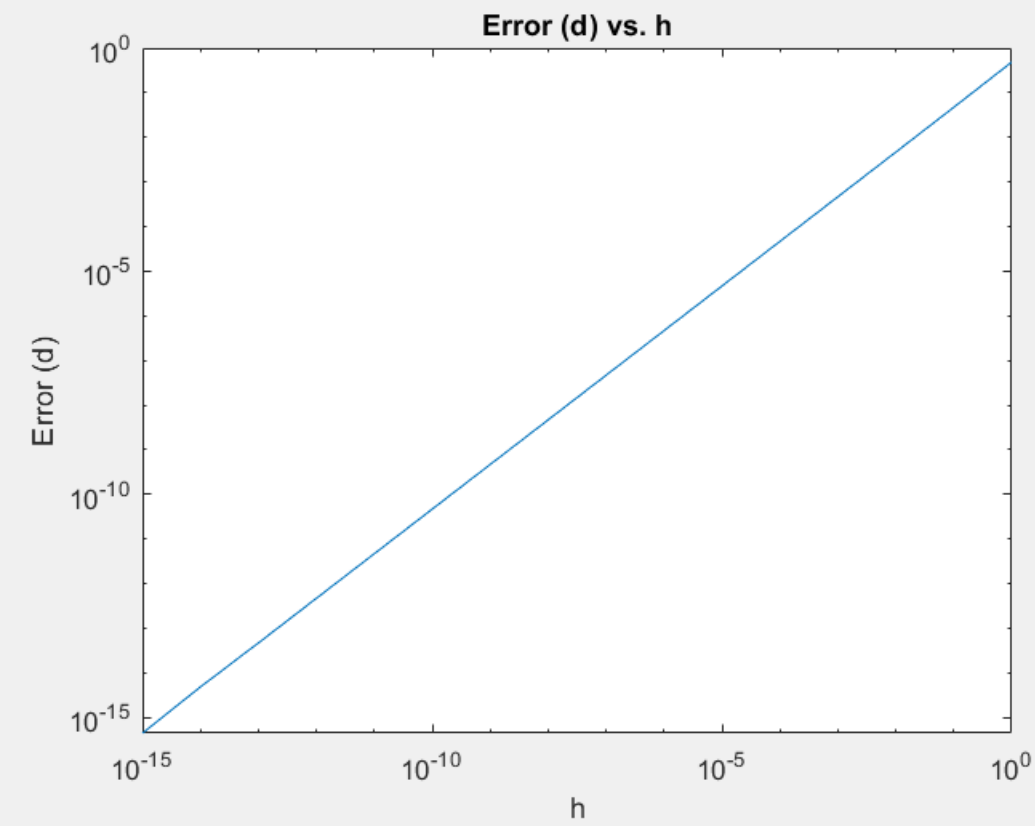
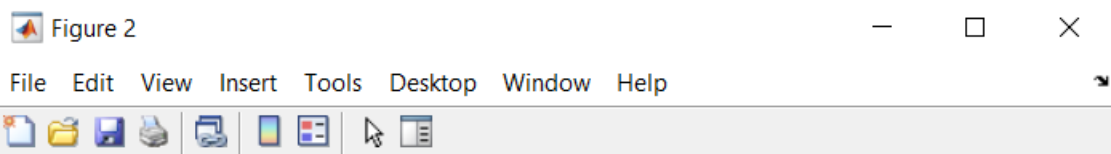
problem\_A2.m ✕ problem\_A3.m ✕ gecp.m ✕ +

```
1      x0 = 1.2;
2      f_prime_exact = cos(x0);
3      h = 10.^(0:-1:-16);
4      f_prime_approx = (sin(x0 + h) - sin(x0)) ./ h;
5      better_approx = (2 * cos(x0 + h/2) .* sin(h/2)) ./ h;
6      error_1 = abs(f_prime_exact - f_prime_approx);
7      error_2 = abs(f_prime_exact - better_approx);
8
9      figure;
10     loglog(h, error_1);
11     xlabel("h")
12     ylabel("Error (b)")
13     title("Error (b) vs. h")
14
15     figure;
16     loglog(h, error_2);
17     xlabel("h")
18     ylabel('Error (d)')
19     title('Error (d) vs. h')
```

Figure 1

File Edit View Insert Tools Desktop Window Help





problem\_A2.m x problem\_A3.m x gecp.m x +

```
1 % Make the matrix A as defined in the problem.
```

```
2 n = 100;
```

```
3 A = -1 * tril(ones(n), -1) + eye(n);
```

```
4 A(:, end) = 1;
```

```
6 % Make the random vector.
```

```
7 x = randn(100, 1);
```

```
9 % Calculate b.
```

```
10 b = A * x;
```

```
12 % Part (a).
```

```
13 k = cond(A);
```

```
15 % Part (b).
```

```
16 x_ge = A \ b;
```

```
17 norm_error_1 = norm(x - x_ge, 2);
```

```
19 % Part (c).
```

```
20 [Q,R] = qr(A, 0);
```

```
21 x_qr = R \ (Q' * b);
```

```
22 norm_error_2 = norm(x - x_qr, 2);
```

```
24 % Part (d).
```

```
25 % P*A*Q = L*U...
```

```
26 [L, U, P, Q] = gecp(A);
```

```
27 x_gecp = Q * inv(U) * inv(L) * P * b;
```

```
28 norm_error_3 = norm(x - x_gecp, 2);
```