

Math 128A, Summer 2019

Lecture 28, 8/8/2019

1 Review: Homework 6

1.1 Problem 5

Consider the Fredholm equation:

$$u(x) + \int_0^1 \underbrace{K(x, y)}_{\text{causal}} u(y) dy = g(x)$$

which is a rank 1 perturbation of the identity.
versus the Volterra equation:

$$y(t) = y(0) + \int_0^t \underbrace{f(s, y(s))}_{\text{causal}} ds$$

and $x_i + \sum_{j=1}^n A_{ij}x_j = b_i$ which is $(I + A)x = b$ so $x = b - Ax$.
Fredholm are generally harder than Volterra, but in our problem we have a special kernel (K for kernel),

$$K(x, y) = \cos(x) \sin(y)$$

with

$$\int_0^1 \cos(x) \sin(y) u(y) dy = \cos(x) \int_0^1 \sin(y) u(y) dy$$

We need to solve forward and backward simultaneously. Recall that any rank-1 matrix, say A , is a row vector left-multiplied by a column vector:

$$A = ab^T = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} \begin{bmatrix} b_1 & \cdots & b_n \end{bmatrix} = \begin{bmatrix} a_1 b_1 & \cdots & a_1 b_n \\ \vdots & & \vdots \\ a_n b_1 & \cdots & a_n b_n \end{bmatrix}$$

Then let $(I + A)x = v$, so that:

$$\begin{aligned} x + a(b^T x) &= v \\ x &= v - \underbrace{(b^T x) a}_{\text{scalar}} \end{aligned}$$

We have an idea: let's take the dot product across the equation with b .

$$\begin{aligned} b^T x &= b^T v - (b^T x)(b^T a) \\ (1 + b^T a)b^T x &= b^T v \\ b^T x &= \frac{b^T v}{1 + b^T a}, \end{aligned}$$

if $b^T a \neq -1$. This is how we solve this rank-1 perturbation problem. In this particular problem, we have a rank-1 kernel, and thus our problem has a rank-1 perturbation.

In our problem, once we discretize the equation via gaussian quadrature, we take on the linear system to solve via Gaussian elimination with partial pivoting (GEPP).

2 Review

Yesterday, we took A symmetric positive definite, with $x^T A x > 0$ for $x \neq 0$ and $A = A^T$. We computed the Cholesky factorization,

$$A = R^T R = L L^T,$$

where R is (right) upper-triangular with $r_{ii} > 0$, we check:

$$\begin{aligned} x^T A x &= (R x)^T (R x) \\ &= \|R x\|^2 > 0, \end{aligned}$$

if $x \neq 0$, because R is invertible.

One of the checks if A is symmetric positive definite is for us to compute the cholesky factorization, and if multiplying them together doesn't return A , then A is not symmetric positive definite (we have nothing to lose).

Recall that we have three ways to compute the Cholesky factorization:

(1) $LDL^T = \underbrace{(\sqrt{DL^T})^T}_{R^T} \underbrace{(\sqrt{DL^T})}_R$

(2) Direct calculation, solving for one element at a time (and using symmetry). We specify an algorithm to do this. First we write the underlying math:

$$i \geq j: \quad a_{ij} = \sum_{k=1}^i r_{ki} \underbrace{r_{kj}}$$

where $r_{kj} = r_{1j}r_{1j} + r_{2j}r_{2j} + \dots + r_{ij} \underbrace{r_{ij}}$ where r_{ij} is what we solve for first, then the rest one at a time.

Additionally,

$$a_{ii} = \sum_{k=1}^i r_{ki}^2$$

so we find a_{11} first: $r_{11} = \sqrt{a_{11}}$.

and note $r_{kj} = 0$ when $k > j$ and $r_{ki} = 0$ when $k > i$. Our pseudocode is as follows: Find column j of \mathbb{R} .

```

1 for j = 1:n
2   for i = 1:j-1
3     r_{ij} = ( a_{ij} - \sum_{k=1}^{i-1} r_{ki} r_{kj} ) / r_{ii}
4   end
5   r_{jj} = \sqrt{ a_{jj} - \sum_{k=1}^{j-1} r_{kj}^2 }
6 end

```

(3) Use Newton's Method for the function we are trying to solve:

$$f(R) = R^T R - A = 0.$$

and of course, R is an upper-triangular matrix with nonnegative diagonal entries. We can say we want to find a correction E such that:

$$(R + E)^T (R + E) - A = E^T E,$$

where we usually set this equal to 0, but instead we cancel out the quadratic term. We check the case that if this is exact, then we would have to deal with a quadratic term, so we insert (recall from yesterday):

$$E := \text{uph}(R^T A R^{-1} - I) R$$

On the other hand, if F is the exact correction to the Cholesky factorization (and E is only what we can get),

$$(R + F)^T(R + F) - A = 0$$

Comparing these,

$$\begin{aligned} R^T R - A + E^T R + R^T E + E^T E &= E^T E \\ R^T R - A + F^T R + R^T F + F^T F &= 0 \end{aligned}$$

where subtracting gives

$$(E - F)^T R + R^T \underbrace{(E - F)}_{(R+E)-(R+F)} = F^T F$$

and rewriting,

$$\underbrace{R^T}_{low.tri} \underbrace{(E - F)^T}_{low.tri} + \underbrace{(E - F)}_{up.tri} \underbrace{R^{-1}}_{up.tri} = R^{-T} F^T F R^{-1},$$

where a lower-triangular times a lower-triangular matrix gives a lower triangular matrix and accordingly for upper-triangular matrices. Hence this equation takes the form of a lower-triangular matrix plus an upper-triangular matrix equals a full matrix.

Hence the error in the correction, the error in $(R + E)$ since $(R + F)$ is exactly Cholesky factor.

We have, for the error relative to the solution we started with,

$$\begin{aligned} (E - F)R^{-1} &= \text{uph}((FR^{-1})^T(FR^{-1})), \\ E - F &= \text{uph} \left(R^{-T} \underbrace{F^T F}_{(\text{err in } R)^T \cdot (\text{err in } R)} R^{-1} \right) R \end{aligned}$$

which is in terms of the relative error of the approximate old solution, so is **quadratic convergence** a la **Newton's**.

3 Floating Point Errors in Cholesky Factorization

When we compute the Cholesky factorizations,

$$fl(r_{11}) = fl(\sqrt{a_{11}}) = \sqrt{a_{11}}(1 + \delta_{11}),$$

where $|\delta_{11}| \leq \varepsilon$ is the exact result, correctly rounded. This gives:

$$\begin{aligned} r_{11}^2 &= a_{11}(1 + \delta_{11})^2 \\ &= a_{11}(1 + 2\delta_{11}) + \cancel{O(\varepsilon^2)} \\ &= \hat{a}_{11}, \end{aligned}$$

where $|\hat{a}_{11} - a_{11}| \leq 2\varepsilon|a_{11}|$. Our hope is that

$$\hat{R} = fl(R)$$

is the **exact** Cholesky factor of

$$|\hat{A} - A| \leq O(n)\varepsilon|A| = [a_{ij}].$$

Because we don't yet have the notion of a norm, we check each entry value and check that each can be bounded.

Remember from our algorithm earlier for directly computing the Cholesky factorization,

$$r_{ij} = \left(a_{ij} - \sum_{k=1}^{i-1} r_{ki} r_{kj} \right) / r_{ii}$$

where we can only work with what we know ('the guys with hats'). Recalling $fl(a - b) = a(1 + \delta) - b(1 + \delta)$, what actually happens is:

$$\hat{r}_{ij} = \left(a_{ij}(1 + \delta_{ij}) - \sum_{k=1}^{i-1} \hat{r}_{ki} \hat{r}_{kj}(1 + i\delta_{ij}) \right) / (r_{ii}(1 + \delta_{ii}))$$

We tend to not like forward error analysis (as the above). We prefer backwards error analysis as the computed results almost satisfying the exact equation it is trying to solve.

Multiplying across by the denominator gives:

$$\hat{r}_{ii} \hat{r}_{ij}(1 + \delta''_{ii}) = a_{ij}(1 + \delta_{ij}) - \sum_{k=1}^{i-1} \hat{r}_{ki} \hat{r}_{kj} - \sum_{k=1}^{i-1} \hat{r}_{ki} \hat{r}_{kj} i \delta_{ij}$$

What we have is doing exactly what we want, if we interpret everything a little differently.

$$\hat{r}_{ii} \hat{r}_{ij} = \underbrace{\left[a_{ij}(1 + \delta_{ij}) - \sum_{k=1}^i \hat{r}_{ki} \hat{r}_{kj} i \delta_{ij} \right]}_{\text{scapegoat}} - \sum_{k=1}^{i-1} \hat{r}_{ki} \hat{r}_{kj}$$

and this is the exact Cholesky factor of our matrix, where we labeled the bracketed quantity as the 'scapegoat' as the source of our error.

So we conclude:

\hat{R} is the exact Cholesky factor of $\hat{A} = \left[a_{ij}(1 + \delta_{ij}) - \sum_{k=1}^i \hat{r}_{ki} \hat{r}_{kj} i \delta_{ij} \right]$, and we write the error as a good quantity plus an uncertain quantity:

$$|\hat{A} - A| \leq \varepsilon |A| + (n\varepsilon) |\hat{R}^T| |\hat{R}|,$$

where $|M|$ denotes taking the absolute value of every element.

This means that all we need to do is take the absolute value of all entries, place them into a matrix R , multiply it by its transpose, then we get close to the exact Cholesky factors.

We may also want to check how close is \hat{R} to **the exact** cholesky factor? Then, what can we say about how close \hat{R} is to R . This gives a measure of **how sensitive** the Cholesky factors are to **perturbations** in A .

4 Sensitivity to Perturbations

Generally we look at the sensitivity of the inverse A^{-1} by perturbations in A , say $A \mapsto A + E$. We write per element, via geometric series:

$$\frac{1}{a + e} = \frac{1}{a \left(1 + \frac{e}{a}\right)} = \frac{1}{a} \left(1 - \frac{e}{a} + O\left(\frac{e}{a}\right)^2\right)$$

Now how do we do this for matrices?

$$(A + E)^{-1} = (A(I + A^{-1}E))^{-1}$$

We need a geometric series for matrices. Recall the derivation for the 1-dimensional case:

$$\begin{aligned} \frac{1 - x^{n+1}}{1 - x} &= 1 + x + x^2 + \cdots + x^n \\ 1 - x^{n+1} &= (1 - x)(1 + x + \cdots + x^n) \\ &= (1 + x + \cdots + x^n) - (x + x^2 + \cdots + x^{n+1}) \\ &= 1 - x^{n+1}. \end{aligned}$$

and in the matrix case,

$$\begin{aligned} (I - A^{n+1})(I - A)^{-1} &= I + A + \cdots + A^n \\ (I - A^{n+1}) &= (I + A + \cdots + A^n)(I - A) \\ &= I + A + \cdots + A^n - (A + \cdots + A^{n+1}) \end{aligned}$$

The first condition for an infinite geometric series of matrices is that $(I - A)^{-1}$ must **exist**, and the second is that $A^n \xrightarrow{n \rightarrow \infty} 0$. It turns out this is equivalent (\iff) all eigenvalues of A are < 1 in absolute values: $|\lambda_j(A)| < 1, \forall j$.

Lecture ends here.

Next time, we'll think of a way to prove this **without** using Jordan Normal Form. We'll talk about **perturbation theory** and perhaps **norms**.