

# Scientific Computing, Fall 2020

## Assignment II: Linear Systems

Aleksandar Donev

Courant Institute, NYU, [donev@courant.nyu.edu](mailto:donev@courant.nyu.edu)

September 18th, 2019

Due by Sunday **October 4th**

For the purpose of grading 100% will be taken to be 100 points, but 110 points are possible. Make sure to follow good programming practices in your MATLAB codes. For example, make sure that parameters, such as the number of variables  $n$ , are not hard-wired into the code and are thus easy to change. Use *fprintf* to format your output nicely for inclusion in your report.

### 1 [35 pts] Ill-Conditioned Systems: The Hilbert Matrix

Consider solving linear systems with the matrix of coefficients  $\mathbf{A}$  defined by

$$a_{ij} = \frac{1}{i + j - 1},$$

which is a well-known example of an ill-conditioned symmetric positive-definite matrix, see for example this Wikipedia article

[http://en.wikipedia.org/wiki/Hilbert\\_matrix](http://en.wikipedia.org/wiki/Hilbert_matrix)

#### 1.1 [10 pts] Conditioning numbers

[10pts] Form the Hilbert matrix in MATLAB and compute the conditioning number for increasing size of the matrix  $n$  for the  $L_1$ ,  $L_2$  and  $L_\infty$  (the column sum, row sum, and spectral matrix) norms based on the definition  $\kappa(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$  and using MATLAB's *norm* function. Note that the inverse of the Hilbert matrix can be computed analytically, and is available in MATLAB as *invhilb*. Compare to the answer with that returned by the built-in exact calculation *cond* and the estimate returned by the function *rcond* (check the help pages for details).

#### 1.2 [25 pts] Solving ill-conditioned systems

[5pts] Compute the right-hand side (rhs) vector  $\mathbf{b} = \mathbf{A}\mathbf{x}$  so that the exact solution is  $\mathbf{x} = \mathbf{1}$  (all unit entries). Solve the linear system using MATLAB's built-in solver (explain what method was used to solve the system in your report) and see how many digits of accuracy you get in the solution for several  $n$ , using, for example, the infinity norm. Also compute the relative norm of the residual  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\| / \|\mathbf{b}\|$  and explain how it changes with the conditioning number of the matrix.

Note: A method is called *backward stable* if it computes the *exact* solution to a nearby problem, i.e., if the residual is small.

[5pts] Do your results conform to the theoretical expectation discussed in class? After what  $n$  does it no longer make sense to even try solving the system due to severe ill-conditioning?

[5pts] Now do the same but solve the system using the Cholesky factorization of  $\mathbf{A}$ , and compare the results and report if anything has changed.

[5pts] Now try computing the solution by using the numerically-computed matrix inverse,  $\mathbf{x} = \text{inv}(\mathbf{A}) * \mathbf{b}$ , and compute the relative errors in the solution and in the residual and see how they behave. What difference do you see between using the LU or Cholesky factorization versus using matrix inverse? Which one is better and why?

[5pts] Finally, try  $\mathbf{x} = \text{invhilb}(\mathbf{A}) * \mathbf{b}$  and compare with using a numerical inverse. Report what you see and try to explain your observations.

## 2 [35 points] Least-Squares Fitting

Consider fitting a data series  $(x_i, y_i)$ ,  $i = 1, \dots, n$ , consisting of  $n = 100$  data points that approximately follow a polynomial relation,

$$y = f(x) = \sum_{k=0}^d c_k x^k,$$

where  $c_k$  are some unknown coefficients that we need to estimate from the data points, and  $d$  is the degree of the polynomial. Observe that we can rewrite the problem of least-squares fitting of the data in the form of an overdetermined linear system

$$[\mathbf{A}(\mathbf{x})] \mathbf{c} = \mathbf{y},$$

where the matrix  $\mathbf{A}$  will depend on the  $x$ -coordinates of the data points, and the right hand side is formed from the  $y$ -coordinates.

Let the correct solution for the unknown coefficients  $\mathbf{c}$  be given by  $c_k = k$ , and the degree be  $d = 9$ . Using the built-in function *rand* generate synthetic (artificial) data points by choosing  $n$  points  $0 \leq x_i \leq 1$  randomly, uniformly distributed from 0 to 1. Then calculate

$$\mathbf{y} = f(\mathbf{x}) + \epsilon \boldsymbol{\delta},$$

where  $\boldsymbol{\delta}$  is a random vector of normally-distributed perturbations (e.g., experimental measurement errors) with mean zero and unit variance, generated using the function *randn*. Here  $\epsilon$  is a parameter that measures the magnitude of the uncertainty in the data points. [Hint: *Plot your data for some small value of  $\epsilon$  to make sure the data points approximately follow  $y = f(x)$ .*]

### 2.1 [20pts] Different Methods

For several logarithmically-spaced perturbations (for example,  $\epsilon = 10^{-i}$  for  $i = 0, 1, \dots, 16$ ), estimate the coefficients  $\tilde{\mathbf{c}}$  from the least-squares fit to the synthetic data and report the error  $\|\mathbf{c} - \tilde{\mathbf{c}}\|$ . Do this using three different methods available in MATLAB to do the fitting:

- a) [5pts] The built-in function *polyfit*, which fits a polynomial of a given degree to data points [Hint: *Note that in MATLAB vectors are indexed from 1 and thus the order of the coefficients that polyfit returns is the opposite of the one we use here, namely,  $c_1$  is the coefficient of  $x^d$ .*]
- b) [5pts] Using the backslash operator to solve the overdetermined linear system  $\mathbf{A}\tilde{\mathbf{c}} = \mathbf{y}$ .
- c) [5pts] Forming the system of normal equations discussed in class,

$$(\mathbf{A}^T \mathbf{A}) \mathbf{c} = \mathbf{A}^T \mathbf{y},$$

and solving that system using the backslash operator.

[5pts] Report the results for different  $\epsilon$  from all three methods in one printout or plot, and explain what you observe.

### 2.2 [15pts] The Best Method

[10pts] If  $\epsilon = 0$  we should get the exact result from the fitting. What is the highest accuracy you can achieve with each of the three methods? Is one of the three methods clearly inferior to the others? Can you explain your results? *Hint: Theory suggests that the conditioning number of solving overdetermined linear systems is the square root of the conditioning number of the matrix in the normal system of equations,*  
 $\kappa(\mathbf{A}) = \sqrt{\kappa(\mathbf{A}^T \mathbf{A})}.$

[5pts] Test empirically whether the conditioning of the problem get better or worse as the polynomial degree  $d$  is increased.

### 3 [40 points] Rank-1 Matrix Updates

In a range of applications, such as for example machine learning, the linear system  $\mathbf{A}\mathbf{x} = \mathbf{b}$  needs to be re-solved after a rank-1 update of the matrix,

$$\mathbf{A} \rightarrow \tilde{\mathbf{A}} = \mathbf{A} + \mathbf{u}\mathbf{v}^T,$$

for some given vectors  $\mathbf{v}$  and  $\mathbf{u}$ . More generally, problems of *updating a matrix factorization* (linear solver) after small updates to the matrix appear very frequently and many algorithms have been developed for special forms of the updates. The rank-1 update is perhaps the simplest and best known, and we explore it in this problem. From now on, assume that  $\mathbf{A}$  is invertible and its inverse or  $\mathbf{LU}$  factorizations are known, and that we want to update the solution after a rank-1 update of the matrix. We will work with random dense matrices for simplicity.

#### 3.1 [10pts] Direct update

[5pts] In MATLAB, generate a random (use the built-in function *randn*)  $n \times n$  matrix  $\mathbf{A}$  for some given input  $n$  and compute its  $\mathbf{LU}$  factorization (you will need it for later parts of this problem). Also generate a right-hand-side (rhs) vector  $\mathbf{b}$  and solve  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

[5pts] Now generate random vectors  $\mathbf{v}$  and  $\mathbf{u}$  and obtain the updated solution  $\tilde{\mathbf{x}}$  of the system  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b}$ . Verify the new solution  $\tilde{\mathbf{x}}$  by directly verifying that the residual  $\mathbf{r} = \mathbf{b} - \tilde{\mathbf{A}}\tilde{\mathbf{x}}$  is small.

#### 3.2 [30pts] SMW Formula

It is not hard to show that  $\tilde{\mathbf{A}}$  is invertible if and only if  $\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq -1$ , and in that case

$$\tilde{\mathbf{A}}^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u} \mathbf{v}^T \mathbf{A}^{-1}}{1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}}. \quad (1)$$

This is the so-called Sherman-Morrison formula, a generalization of which is the Woodbury formula, as discussed on Wikipedia:

[http://en.wikipedia.org/wiki/Sherman-Morrison-Woodbury\\_formula](http://en.wikipedia.org/wiki/Sherman-Morrison-Woodbury_formula). The SMW formula (1) can be used to compute a new solution  $\tilde{\mathbf{x}} = \tilde{\mathbf{A}}^{-1} \mathbf{b}$  without actually computing  $\tilde{\mathbf{A}}^{-1}$  or factorizing  $\tilde{\mathbf{A}}$ .

[10pts] For some  $n$  (say  $n = 100$ ), compare the result from using the formula (1) versus solving the updated system  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b}$  directly as in part 3.1.

[20pts] Implement the SMW formula (1) carefully to compute  $\tilde{\mathbf{x}} = \tilde{\mathbf{A}}^{-1} \mathbf{b}$  as robustly and efficiently as you can, that, is, not actually calculating matrix inverses but rather (re)using the LU factorization of  $\mathbf{A}$  [Hint: You only need to solve two triangular systems to update the solution once you have the factorization of  $\mathbf{A}$ ]. Explain how you computed  $\tilde{\mathbf{x}}$  and estimate how expensive this is in terms of number of operations (FLOPS). Explain how much faster it is then solving  $\tilde{\mathbf{A}}\tilde{\mathbf{x}} = \mathbf{b}$  directly as a function of  $n$ .