

CS 290H: Preconditioning iterative methods // Homework 1

Assigned October 3, 2004

Due by class time Monday, October 17

1. [20 points]

(a) Find a 2-by-2 matrix A that is symmetric and nonsingular, but for which neither A nor $-A$ is positive definite. What are the eigenvalues of A ? Find a 2-vector y such that $y^T A y < 0$.

(b) For A as above, find a 2-vector b such that the conjugate gradient algorithm, when started with the zero vector as an initial guess, does not converge to the solution of $Ax = b$. Show what happens on the first two iterations of CG, as described on Slide 3 of the September 26 class. How do you know it won't converge to the right answer?

2. [40 points] In this problem you'll actually prove that CG works in at most n steps, assuming that real numbers are represented exactly. (This is not a realistic assumption in floating-point arithmetic, or on any computer with a finite amount of hardware, but it gives a solid theoretical underpinning to CG.) Let A be an n -by- n symmetric, positive definite matrix, and let b be an n -vector.

We start with the idea of searching through n -dimensional space for the value of x that minimizes $f(x) = \frac{1}{2}x^T A x - b^T x$, which is the x that satisfies $Ax = b$. We begin by picking a set of n linearly independent search directions, called d_0, d_1, \dots, d_{n-1} . (Actually we don't know them in advance, but that's a detail.) At each iteration we proceed along the next direction until we are "lined up" with the final answer, the value of x at which $Ax = b$. In n -space, once we are lined up with the answer from n independent directions, we will be exactly on the answer.

The **first magic of CG** is that for the right kind of search directions, there is a way to define "lined up" for which we can actually compute how far to go along each search direction. The key definition uses *A-conjugate* vectors. Then "lined up" means that the error $e_i = x_i - x$ is exactly crossways to the search direction d_{i-1} , not in the sense of being perpendicular (which would mean $e_i^T d_{i-1} = 0$), but in the sense of being *A-conjugate*: $e_i^T A d_{i-1} = 0$.

An informal way to say that is, we proceed along the search direction until we are lined up with the solution as seen through *A-glasses*. The reason for lining up through *A-glasses* rather than bare eyes is that we can compute where to stop without knowing where the final answer is. We can't see and compute with x -space directly, but we can see the space where Ax and b live. And after lining up each of n independent directions in an n -dimensional space we are guaranteed to be sitting on top of the right answer, whether the independent directions are the conventional coordinate axes or the *A-conjugate* axes we see through our *A-glasses*.

To go along with this, we need to choose the search directions themselves to be mutually *A-conjugate*: we will require each d_i to be *A-conjugate* to all the earlier d_j 's, so $d_i^T A d_j = 0$ if $i \neq j$.

(a) Suppose we are given i mutually *A-conjugate* vectors d_0, \dots, d_{i-1} . Suppose $x_0 = 0$, and for each $j < i$ we have $x_j = x_{j-1} + \alpha_j d_{j-1}$. Write down and prove correct an expression for a scalar α_i such that, if we take $x_i = x_{i-1} + \alpha_i d_{i-1}$, then the error $e_i = x_i - x$ is *A-conjugate* to d_{i-1} .

Now, how do we get a sequence of A -conjugate directions to search along? In fact, we can start with any sequence of linearly independent directions, and convert them to A -conjugate directions by projecting out all the earlier search directions from each one, using Gram-Schmidt orthogonalization, as follows.

(b) Suppose we are given i mutually A -conjugate vectors d_0, \dots, d_{i-1} , and one more vector u_i that does not lie in their span. Write down and prove correct an expression for scalars $\beta_{i,j}$ such that, if we take

$$d_i = u_i + \sum_{j=0}^{i-1} \beta_{i,j} d_j,$$

then d_i is A -conjugate to all the earlier d_j .

Finally, the **second magic of CG** is that there is a way to choose a particular sequence of directions for which the Gram-Schmidt orthogonalization is really easy. If we choose the right directions to start with, we only need to project out *one* earlier direction, not all i of them. This is why the cost of one CG iteration is only $O(n)$, not $O(n^2)$.

(c) Suppose the vectors d_0, \dots, d_{i-1} , the vectors x_0, \dots, x_{i-1} , and the scalars α_j and $\beta_{i,j}$ are as above. Suppose in addition that at each stage we take $u_i = b - Ax_i$ (which is also known as r_i , the residual). First, prove that if this choice of u_i lies in the span of d_0, \dots, d_{i-1} , the CG iteration can stop with $x_i = x$. Second, show that this direction u_i is already A -conjugate to all of the d_j except d_{i-1} , and therefore we can take $\beta_{i,j} = 0$ for $j < i - 1$.

(d) One last detail: Prove that the CG code on the course slide does in fact compute the residual r_i correctly; that is, prove that $r_{i-1} - \alpha_i A d_{i-1}$ is in fact equal to $b - Ax_i$.

3. [40 points] You may do the programming assignments for this course in C or Fortran; I recommend C. In each case, you will set your code up with an interface so that it can be called from Matlab as a “mexFunction”. This will let you use Matlab to test and debug your code, and to plot results. This warmup assignment is just for you to learn how to write a Matlab interface using sparse matrices.

Write a C or Fortran mexFunction that can be called from Matlab as $y = \text{matvec}(A, x)$, which takes as input a sparse matrix A and a full column vector x , and returns a full column vector y whose value is the matrix-vector product Ax . Your routine can assume the matrix is real, but should not assume that it’s square. (You might want to check to make sure the sizes of A and x are compatible.)

Test your routine from Matlab with several sparse matrices you make up, verifying that it gives the same answer as Matlab’s $y = A*x$. (The norm of the difference, $\text{norm}(A*x - \text{matvec}(A, x))$, should be tiny. It may not be exactly zero because floating-point addition is not associative, and your routine may be doing arithmetic in a different order than Matlab’s.) See the Matlab functions `sprand` or `sprandn` to generate random sparse test matrices. Or download the Matlab interface to the UF Sparse Matrix Collection, at <http://www.cise.ufl.edu/research/sparse>.

For documentation on mexFunctions, open up a help window in Matlab, and look in the *External Interfaces/API* section, under the *Matlab* section. The C mexFunction syntax is quite a bit nicer than the Fortran. All mexFunctions have the routine name `mexFunction`.

You’ll need to use the following Matlab “mx” and “mex” routines:

<code>mxGetM</code>	returns the number of rows of a matrix
<code>mxGetN</code>	returns the number of columns of a matrix
<code>mxGetJc</code>	returns a pointer to the column pointer array (Ap)
<code>mxGetIr</code>	returns a pointer to the row indices (Ai)
<code>mxGetPr</code>	returns a pointer to the numerical values (Ax)

There is a more complex but very well-written example of mexFunctions for sparse matrices in Tim Davis's LDL code, which is at <http://www.cise.ufl.edu/research/sparse/ldl> and also on the CS290H reference page on the web. LDL computes a sparse Cholesky factorization; you don't need to understand how that works, but you can see from that code how a C program can get access to a Matlab sparse matrix.

Turn in all your code, and also a Matlab transcript of a session that tests your code and verifies that the output agrees with Matlab's. (Say "help diary" to Matlab to see how to record a transcript.)