LECTURE II: Solving linear Equations

This lecture is about solving linear systems of equations,

$$Ax = b, (1)$$

where A is a $n \times n$ matrix, and x, b are n-dimensional vectors.

LU-Decomposition

Possibly the first method that one learns for solving Eq. (1) is the Gaussian elimination. In computer algebra a slightly modified version of this method is the LU-decomposion, in which one tries to decompose the matrix A according to A = L U with L being a lower-triangular matrix and U and upper triangular matrix. Knowing the LU decomposition for a matrix A allows us to solve the linear system very easily:

$$Ax = b$$
$$LUx = b$$
$$Ux = L^{-1}b$$
$$x = U^{-1}(L^{-1}b),$$

where $L^{-1}b$ is computed using forward substitution, and $U^{-1}(L^{-1}b)$ by backward substitution. Note that sometimes an additional 'Pivoting step', is needed in which either rows (partial pivoting) or rows and columns (full pivoting) are reordered, e.g., if you would get zeros on the diagonal.

The question arises of how to obtain the LU decomposition? One way uses the recursive leading-row column LU algorithm, where we separate out the first row/column of A, U, L:

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & A_{22} \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{0} \\ l_{21} & L_{22} \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} \\ \mathbf{0} & U_{22} \end{pmatrix}$$

Here it is worth pointing out, that A_{22} is a $(n-1) \times (n-1)$ -matrix and a_{12}, a_{21} ... are vectors of length (n-1), either in column-form (subscript 21) or row-form (subscript 12). The elements a_{11} and u_{11} are real numbers.

Multiplying out, this matrix equation can be rewritten as:

$$a_{11} = u_{11}$$

$$u_{12} = a_{12}$$

$$l_{21} = \frac{1}{a_{11}}a_{21}$$

$$L_{22}U_{22} = A_{22} - \frac{1}{a_{11}}a_{21} \otimes a_{12}$$

The $(n-1) \times (n-1)$ matrix $A_{22} - \frac{1}{a_{11}}a_{21} \otimes a_{12}$ is the Schur complement and defines a new system of size $(n-1) \times (n-1)$ to solve.

Overall, the LU-decomposition has costs proportional to n^3 . Therefore, this method can only be used for 'small' matrices, $n \leq 10,000$.

Iterative Solvers

This subsection is based on the book 'Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods' by Berrett et al.

Iterative methods use successive approximations to obtain more accurate solutions to a linear system

$$x^{(k)} = Bx^{(k-1)} + c$$

at each step. Overall, one can distinguish 'Stationary methods' (e.g., Jacobi method, Gauss-Seidel method or Successive Overrelaxation) in which neither B nor c depend on the iteration k and 'Nonstationary methods' (e.g., Conjugate Gradient, Minimal Residual, Generalized Minimal Residual, Biconjugate Gradient Stabilized).

Preconditioning: The rate at which an iterative method converges depends greatly on the spectrum of the coefficient matrix. A '**preconditioning**' P is employed to transform the coefficient matrix into one with a more favorable spectrum. A good preconditioner improves the convergence of the iterative method, but sometimes is even required to ensure that the iterative method does not fail to converge. One multiples Eq. (1) by P,

$$PAx = Px,$$
(2)

to obtain a new linear system with the coefficient matrix being the product (PA). The trick is now to find a preconditioning matrix, such that (i) it is approximately the inverse of A and (ii) it is easy to calculate. There is an immense amount of knowledge and tools around preconditioning, although we don't have time to go into them.

The Jacobi Method

Considering a linear system consisting of n equations

$$Ax = b$$
,

the *i*-th equation can be written as

$$x_i = \frac{1}{a_{ii}} \Big(b_i - \sum_{j \neq i} a_{i,j} x_j \Big).$$

Then the simplest way to solve for x_i and keeping all other equations fixed, i.e., we consider all equations as being independent from each other:

$$x_i^{(k)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j \neq i} a_{i,j} x_j^{(k-1)} \right).$$

In matrix terms this can be rewritten as:

$$x^{(k)} = D^{-1}(-L - U)x^{(k-1)} + D^{-1}b$$

with D being the diagonal, U the strictly upper, L the strictly lower part of A.

#include<stdio.h>
#include<math.h>

3

```
4 /* We are solving
     3x + 20y - z = -18
5
      2x - 3y + 20z = 25
6
      20x + y - 2z = 17
7
8 */
9 /* Bring system in diagonally dominant form:
      20x + y - 2z = 17
10
     3x + 20y - z = -18
11
     2x - 3y + 20z = 25
12
13 */
14 /* Equations:
     x = (17. - y + 2z)/20.
15
     y = (-18.-3x+z)/20.
16
17
     z = (25.-2x+3y)/20.
18 */
19 /* Defining function */
20 #define f1(x, y, z) (17. - y + 2. * z)/20.
<sup>21</sup> #define f^{2}(x, y, z) = (-18, -3, *x+z)/20.
<sup>22</sup> #define f3(x,y,z) (25. -2.*x+3.*y)/20.
23
24 /* Main function */
25 int main()
26 {
   float x0=0, y0=0, z0=0, x1, y1, z1, e1, e2, e3;
27
   int count=1;
28
   float e = 1e-6;
29
30
   printf("\nCount\tx\ty\tz\n");
31
32
   do
   {
33
    /* Calculation */
34
    x1 = f1(x0, y0, z0);
35
36
    y1 = f2(x0, y0, z0);
    z1 = f3(x0, y0, z0);
37
     printf("%d\t%0.7f\t%0.7f\t%0.7f\n",count, x1,y1,z1);
38
39
    /* Error */
40
    e1 = fabs(x0-x1);
41
    e2 = fabs(y0-y1);
42
43
    e3 = fabs(z0-z1);
44
    count++;
45
46
    /* Set value for next iteration */
47
    x0 = x1;
48
    y0 = y1;
49
    z0 = z1;
50
   }while(e1>e && e2>e && e3>e);
51
52
   printf("\nSolution: x=\%7.6f, y=\%7.6f and z = \%7.6f \setminus n", x1, y1, z1);
53
54
55
   return 0;
56
57 }
```

Gauss-Seidel Method

A simple extension of the Jacobi method is the Gauss-Seidel method. Now we assume that the equations are examined one at a time in sequence and we will use previously computed results. This results in

$$x_i^{(k)} = \frac{1}{a_{ii}} \Big(b_i - \sum_{j < i} a_{i,j} x_j^{(k)} - \sum_{j > i} a_{i,j} x_j^{(k-1)} \Big),$$

or in matrix form:

$$x^{(k)} = (D+L)^{-1}(-Ux^{(k-1)}+b).$$

Successive Overrelaxation (SOR)

The Successive Overrelaxation Method, or SOR, is devised by applying extrapolation to the Gauss-Seidel method. This extrapolation takes the form of a weighted average between the previous iterate $x^{(k)}$ and the Gauss-Seidel iterate $\tilde{x}_i^{(k-1)}$, successively for each component:

$$x_i^{(k)} = \omega \tilde{x}_i^{(k)} + (1 - \omega) x_i^{(k-1)}, \qquad i = 1, \dots, n.$$

In matrix-form, this gives

$$x^{(k)} = (D + \omega L)^{-1} (-\omega U + (1 - \omega)D) x^{(k-1)} + \omega (D + \omega L)^{-1} b.$$

Usually, we have to pick $\omega \in]0,2[$. For $\omega = 1$ we obtain the normal Gauss-Seidel method, for $\omega < 1$ we are using an underrelaxation and for $\omega > 1$ this method gives overrelaxation.

Conjugate Gradient Method (CG)

Nonstationary methods differ from stationary methods in that the computations involve information that changes at each iteration.

The Conjugate Gradient method can be used for symmetric positive definite systems. (Note: positive definite: An $n \times n$ symmetric real matrix M is positive-definite if $x^T M x > 0$ for all non-zero x in \mathbb{R}^n .)

CG generates vector sequences of iterates (i.e., successive approximations to the solution), residuals corresponding to the iterates, and searches for directions used in updating the iterates and residuals.

The fundamental idea is to rewrite our system of equations as a minimization problem, i.e., Eq. (1) is rewritten as

$$E(x) := \frac{1}{2}x^T A x - b^T x \quad \to \quad \min$$

The gradient of E(x) is grad E = Ax - b, so that the minimum of E corresponds to the solution of Eq. (1).

We denote the residual of the k-th step as $r^{(k)} := b - Ax^{(k)}$, and in CG, one determines the next step along a search direction $p^{(k)}$, i.e.

$$x^{(k)} = x^{(k-1)} + \alpha_k p^{(k)},$$

with α_k chosen to yield the minimum of $E(x^{(k)})$ along this line, which is given by

$$\alpha_k = \frac{r^{(k-1)^T} r^{(k-1)}}{p^{(k)^T} A p^{(k)}}.$$

This α_k also minimizes $r^{(k)T}A^{(-1)}r^{(k)}$.

The trick in Conjugate Gradients is the choice of $p^{(k)}$. Just going as steeply downhill as possible (i.e. $p^{(k)} = -gradE(x^{(k-1)}) = r^{(k-1)}$) is not optimal. This "gradient method" is prone to going 'zig-zag', needing many steps.



Figure 1: Sketch of the gradient and conjugate gradient methods (plot from 'A gradientbased algorithm competitive with variational Bayesian EM for mixture of Gaussians' by Kuusela et al.)

Instead one chooses the search direction by

$$p^{(k)} = r^{(k)} + \beta_{(k)} p^{(k-1)}$$

with

$$\beta_{(k)} = \frac{r^{(k)T}r^{(k)}}{r^{(k-1)T}r^{(k-1)}}.$$

This choice of β ensures that $r^{(k)}$ and $r^{(k-1)}$ are orthogonal (and also orthogonal to all previous choices). This orthogonality, in turn prevents the 'zig-zag' of the gradient method, and leads to faster convergence.

In summary, the iteration procedure can be summarized by:

- 1. Compute r
- **2.** Compute β
- 3. Compute *p*
- 4. Compute α
- 5. Update *x*
- 6. Go back to (1)

In infinite precision, CG will arrive at the exact solution in at most n steps. In practical application, CG typically needs $N_{\text{its}} \ll n$. The number of iterations needed to reduce the

residual by a certain factor depends on the conditioning number of the matrix A.

$$N_{
m its} \propto \sqrt{rac{\lambda_{
m max}}{\lambda_{
m min}}},$$
 (3)

where $\lambda_{\text{max/min}}$ is the largest and smallest eigenvalue of *A*. (Preconditioning precisely aims to bring the preconditioned eigenvalues close to unity)

In the following, we want to consider one example, to get a better understanding of the method. Consider

$$\mathbf{A}\mathbf{x} = \begin{pmatrix} 4 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

We start with the initial guess

$$\mathbf{x_0} = \begin{pmatrix} 2\\1 \end{pmatrix}.$$

We start by computing the residual vector:

$$\mathbf{r_0} = \begin{pmatrix} 1\\2 \end{pmatrix} - \begin{pmatrix} 4&1\\1&3 \end{pmatrix} \begin{pmatrix} 2\\1 \end{pmatrix} = \begin{pmatrix} -8\\-3 \end{pmatrix} = \mathbf{p_0}$$

Now, we have to compute α :

$$\alpha_0 = \frac{r^{(0)}{}^T r^{(0)}}{p^{(0)}{}^T A p^{(0)}} = \frac{(-8,3) \cdot \begin{pmatrix} -8\\ -3 \end{pmatrix}}{(-8,3) \begin{pmatrix} 4 & 1\\ 1 & 3 \end{pmatrix} \begin{pmatrix} -8\\ -3 \end{pmatrix}} = \frac{73}{331}$$

This gives us our next solution:

$$\mathbf{x_1} = \mathbf{x_0} + \alpha_0 \mathbf{p_0} = \begin{pmatrix} 2\\1 \end{pmatrix} + \frac{73}{331} \begin{pmatrix} -8\\-3 \end{pmatrix} = \begin{pmatrix} 0.2356\\0.3384 \end{pmatrix}$$

We now need to move to the second iteration:

$$\mathbf{r_1} = \mathbf{r_0} - \alpha_0 \mathbf{A} \mathbf{p_0} = \begin{pmatrix} -8 \\ -3 \end{pmatrix} - \frac{73}{331} \begin{pmatrix} 4 & 1 \\ 1 & 3 \end{pmatrix} \begin{pmatrix} -8 \\ -3 \end{pmatrix} = \begin{pmatrix} -0.2810 \\ 0.7492 \end{pmatrix}$$

Then, we get for β :

$$\beta_{(1)} = \frac{r^{(1)} r^{(1)}}{r^{(0)} r^{(0)}} = \frac{(-0.2810, 0.7492) \cdot (-0.2810, 0.7492)^T}{(-8, -3) \cdot (-8, -3)^T} = 0.0088.$$

With β we can compute

$$\mathbf{p_1} = \mathbf{r_1} + \beta_1 \mathbf{p_0} = \begin{pmatrix} -0.3511\\ 0.7229 \end{pmatrix},$$

plugging this into the equation for α leads to

$$\alpha_1 = \frac{r^{(1)T}r^{(1)}}{p^{(1)T}Ap^{(1)}} = 0.4122$$

This leads to the final solution:

$$\mathbf{x_2} = \mathbf{x_1} + \alpha_1 \mathbf{p_1} = \begin{pmatrix} 0.0909\\ 0.6346 \end{pmatrix}$$

This solution is up to round up errors exact. In fact, it is possible to show that for exact arithmetic, the method converges to the correct solution within m steps where m determines the size of the $m \times m$ matrix A.

Below you also find a python code, that actually uses the CG and also Gradient method to solve our problem from above:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 ##setup plot
5 fig,(ax1,ax2)=plt.subplots(1,2,figsize=(10,5))
  ##define the problem
7
8 A = np.array([[4.,1.],[1.,3.]]) # change 3->23 for bad 'gradient'
     convergence
b = np.array([1.,2.])
10
11 #create grid for contour plot
xp = np.arange(-.6, 2.4, 0.01)
<sup>13</sup> yp = np.arange(-1, 2, 0.01)
14
15
_{16} X, Y = np.meshgrid(xp, yp)
17 ##Setup quadratic E for minimization
^{18} R = 0.5 * (A[0][0] * X + A[0][1] * Y) * X + 0.5 * (A[1][0] * X + A[1][1] * Y) * Y - b[0] * X - b
      [1]*Y
<sup>19</sup> ax1.contour(X,Y,R,25)
20
21
22 #initial guess
x = np.array([2.,1.])
_{24} beta = 0
       = np.array([0, 0])
25 p
       = np.array([1, 1])
26 r
27
28 #plot initial guess
29 ax1.scatter(x[0],x[1],color='red')
30
31
32 ### CG
33 #perform 2 steps, since we know that we are then at the final location
34 CGresiduals = []
35 for i in range(1,3):
    r0
         = r
36
         = b - np.dot(A, x)
    r
37
    if i >1:
38
        beta = np.dot(r.T,r)/np.dot(r0.T,r0)
39
         = r + beta * p
40
    р
         = np.dot(r.T,r)/np.dot(np.dot(p.T,A),p)
41
    а
    #storing it shortly in xnew for plotting
42
    xnew = x + a * p
43
    ax1.scatter(xnew[0],xnew[1],color='red')
44
    ax1.arrow(x[0], x[1], xnew[0]-x[0], xnew[1]-x[1], color='red')
45
    x = xnew
46
```

```
CGresiduals.append(np.dot(r.T,r)**0.5)
47
48
49 # append last residual
r = b - np.dot(A, x)
  CGresiduals.append(np.dot(r.T,r)**0.5)
51
52
53 ### Gradient method
54 #reset initial guess
x = np.array([2.,1.])
56
57 GradientResiduals = []
<sup>58</sup> for i in range(1,20):
         = b - np.dot(A, x)
59
   r
         = np.dot(r.T,r)/np.dot(np.dot(r.T,A),r)
    а
60
    xnew = x + a * r \# step in direction of 'r
61
62
    ax1.scatter(xnew[0],xnew[1],color='blue',alpha=0.5)
63
    ax1.arrow(x[0], x[1], xnew[0]-x[0], xnew[1]-x[1], color='blue', alpha=0.5)
64
    x = xnew
65
    GradientResiduals.append(np.dot(r.T,r)**0.5)
66
67
ax2.plot(CGresiduals, '-o', label='CG')
69 ax2.plot(GradientResiduals, '-+', label='gradient')
70 ax2.set_yscale('log');
71 ax2.legend();
```

Generalized Minimal Residuals (GMRES)

GMRES is important, because this iterative method also works for non-symmetric A and/or for A without positivity properties. This generality comes at a higher cost (both computations and storage) than for conjugate gradients. In the following, we will only sketch the ideas behind this method.

At iteration K, define the Krylov-subspace

$$K_k := \operatorname{span}\left(r^{(0)}, Ar^{(0)}, A^2 r^{(0)}, \dots A^{k-1} r^{(0)}\right)$$
(4)

At each iteration k, GMRES will determine $x^{(k)} = x^{(0)} + c^{(k)}$, with $c^{(k)} \in K_k$, such that the residal $||b - Ax^{(k)}||$ is minimized over all possible $c^{(k)}$.

To do this, GMRES needs to build and store a orthonormal basis $\{v^{(k)}\}$ of K_k , and needs to evaluate $x^{(k)} = x^{(0)} + y_1 v^{(1)} + y_2 v^{(2)} + \ldots y_k v^{(k)}$ for suitable coefficients $y_1, \ldots y_k$. It turns out the residual r^k can be computed from known data without forming $x^{(k)}$

It turns out the residual r^{κ} can be computed from known data without forming $x^{(\kappa)}$ explicitly, and therefore, the method can be summarized as follows:

- calculate the new orthonormal basis-vector v_k using the Arnoldi method
- find y_k that minimizes $||r^{(k)}||$
- Calculate residual $r^{(k)}$
- repeat until residual is small enough
- compute final $x^{(k)}$

The storage space for the orthonormal basis and the sums over the orthogonal basis, entails a cost that grows linearly with k. To mitigate this growing cost, it is typical to perform 'm' iterations, and then to re-start the entire algorithm with the then best solution.

This is called 'GMRES(m)'.

GMRES and Conjugate Gradients are actually simular in that they construct their approximate solutions using the Krylov subspace K_k , and in that residuals obey orthogonality conditions.

Conjugate Gradients takes advantage of the extra knowledge that A must be symmetric and positive-definite to express the scheme without the need to retain information from all previous iterations.