

EE5121 Convex Optimization - Quiz 2

Anirudh R (EE18B103), Sreekar Sai R (EE18B154), Chandan Bhat(EE16D209)

September 29, 2021

DECLARATION

We, Anirudh R, Sreekar Sai R., and Chandan Bhat declare that the submitted work is entirely our original work.

Contributions:

- Anirudh coded up the basic conjugate gradient, PD matrix generation and generated all the plots for q1.
 - Sreekar developed the derivation for preconditioned CG, and modified the code accordingly.
 - Chandan coded up nonlinear CG and created all the visualisations.
- All edited and reviewed this report.

In this report, we analyse the conjugate gradient method and its variants. We introduce the problem in section 1, and then describe the conjugate gradient method and update steps in brief. In section 2.2, we describe a modification to CG incorporating preconditioning, and then compare the performance of vanilla and preconditioned conjugate gradient. We also implement a non-linear conjugate gradient method (Fletcher-Reeves and Polak-Ribiere) in section 3 and note some of our interesting observations.

1 Solving a linear system

We have a $n \times n$ symmetric positive definite matrix $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$. We seek to find the solution to the linear system of equations represented by

$$Ax = b \tag{1}$$

Performing Gaussian elimination is a straight-forward method to obtain the solution but its time complexity is $\mathcal{O}(n^3)$. Finding the inverse will also solve the linear system but the complexity of Gauss-Jordan method of inverse calculation is again $\mathcal{O}(n^3)$. Even other more efficient inverse calculation algorithms have time complexities **higher**[1] than $\mathcal{O}(n^2)$. In addition, when the size of the problem is large, we require large storage for the matrix A , which can become impossible in many practical scenarios.

Here, we pose this linear system of equations problem as an optimization problem by constructing the following loss function.

$$\phi(x) = \frac{x^T A x}{2} - b^T x \tag{2}$$

Here, we seek to find the global minima of this strictly convex loss function(as A (the hessian) is positive definite, $\phi(x)$ is strictly convex). At the minima, we know that the gradient will reach zero.

$$\nabla\phi(x) = Ax - b = 0 \quad (3)$$

Therefore, by finding the minima of $\phi(x)$, we will find x that satisfies $Ax = b$ which is what we require.

2 Conjugate Gradient Method

2.1 Conjugacy

The conjugate gradient method makes use of a conjugate vector of A as the direction of movement in each iteration. The conjugate vectors of A satisfy the following properties. If $\{p_0, p_1, \dots, p_{n-1}\}$ are conjugate vectors of A , then,

$$p_i^T A p_j = 0 \quad \forall i \neq j \quad (4)$$

$$p_i^T A p_i \neq 0 \quad \forall i \quad (5)$$

One important property of conjugate vectors is that they are linearly independent. There are multiple ways, but inefficient of generating the conjugate vectors.

- Eigen Decomposition of A : The eigen vectors of A are valid conjugate vectors of A . But, the time complexity of Eigen Value Decomposition is $\mathcal{O}(n^3)$. So, instead of doing this, we could just directly do Gaussian Elimination.
- Modified Gram-Schmidt Ortho-normalization: This is another method to generate conjugate vectors but is also $\mathcal{O}(n^3)$ in time complexity.

We use the following method of generating conjugate vectors that is pivotal to the efficiency of CG method. Residue $r_i = Ax_i - b$, Initial value: $p_0 = -r_0$.

$$p_k = -r_k + \beta_k p_{k-1} \quad \text{where } \beta_{k+1} = \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k} \text{ \& } r_{k+1} = r_k + \alpha_k A p_k \quad (6)$$

The overall update equation for x will be,

$$x_{k+1} = x_k + \alpha_k p_k \quad \text{where } \alpha_k = \frac{r_k^T r_k}{p_k^T A p_k} \quad (7)$$

2.2 Preconditioned CG

We know that the convergence of the conjugate gradient method, follows the following inequality

$$\|x_{k+1} - x^*\| \leq \frac{\lambda_{n-k} - \lambda_1}{(\lambda_{n-k} + \lambda_1)} \|x_0 - x^*\| \quad (8)$$

where $\lambda_n \geq \dots \geq \lambda_1$ are the eigenvalues of the matrix A . Notice that the lower bound quickly becomes close to zero if the larger eigenvalues are closer to λ_1 (i.e., if it has a lower condition number).

The key to accelerate CG is to modify the eigenvalue distribution of A . Instead of solving (1), we'd instead like to solve

$$P A x = P b \quad (9)$$

where P is a full rank matrix. We hope to choose a matrix P such that the condition number of PA is lesser than that of A (i.e., it has a better eigenvalue distribution).

$$L^T Ax = L^T b \quad (10)$$

$$(L^T AL)(L^{-1}x) = L^T b, \quad \text{or} \quad (11)$$

$$\hat{A}\hat{x} = \hat{b} \quad (12)$$

Ideally, we'd like to have our system to have all eigenvalues equal to each other, or simply the matrix $L^T AL$ to be equal to, or close to identity. One way of achieving such a form is by decomposing A using the cholesky decomposition:

$$A = CC^T \quad (13)$$

and using $L = (C^{-1})^T$, will lead to $\hat{A} = I$. However, the computational complexity of algorithms to compute the cholesky factorisation is $\mathcal{O}(n^3)$ in general. Instead, a much more cheaper alternative is to compute the incomplete cholesky decomposition, which is a *sparse approximation* of the cholesky factorization.

$$A \approx \tilde{C}\tilde{C}^T \quad (14)$$

Choosing $L = (\tilde{C}^{-1})^T$, we have a 'good enough' matrix L , such that the new system matrix $L^T AL$ is *close* to identity.

The naive approach to solving (12) is to simply solve for \hat{x} using this new system, and then compute $x = L\hat{x}$. Instead we can modify the update equations (section 2) to be in terms of our original variables x .

In our new system, the update directions \hat{p}_k must be conjugate wrt \hat{A} . In other words, they must satisfy,

$$\hat{p}_i^T (L^T AL) \hat{p}_j = 0 \quad \forall i \neq j \quad (15)$$

Defining $\hat{p}_k = L^{-1}p_k$ (p_k are conjugate with respect to A), we have

$$\hat{p}_i^T (L^T AL) \hat{p}_j = p_i^T A p_j = 0 \quad \forall i \neq j \quad (16)$$

therefore ensuring conjugacy of \hat{p}_k with respect to \hat{A} .

The new residue \hat{r}_k can be written as,

$$\hat{r}_k = \hat{A}\hat{x} - \hat{b} = L^T(Ax - b) = L^T r_k \quad (17)$$

Using the update equation in section 2 and equation 17, we have

$$\hat{\alpha}_k = \frac{\hat{r}_k^T \hat{r}_k}{\hat{p}_k^T \hat{A} \hat{p}_k} \quad (18)$$

$$= \frac{r_k^T L L^T r_k}{p_k^T A p_k} \quad (19)$$

For convenience, and to avoid repetitive calculations, we will define

$$y_k = M^{-1} r_k \quad (20)$$

where $M^{-1} = LL^T$, thus simplifying the expression for $\hat{\alpha}_k$. Similarly,

$$\hat{\beta}_{k+1} = \frac{r_{k+1}^T r_{k+1}}{\hat{r}_k^T \hat{r}_k} = \frac{r_{k+1}^T LL^T r_{k+1}}{r_k^T LL^T r_k} \quad (21)$$

$$= \frac{r_{k+1}^T y_{k+1}}{r_k^T y_k} \quad (22)$$

Now, using these, the update equations for \hat{x}_k and \hat{p}_k can be simplified to their corresponding equations in the original variables :

$$\hat{p}_k = -\hat{r}_k + \hat{\beta}_k \hat{p}_{k-1} \quad (23)$$

$$L^{-1} p_k = -L^T r_k + \hat{\beta}_k L^{-1} p_{k-1}, \quad \text{or} \quad (24)$$

$$p_k = -y_k + \hat{\beta}_k p_{k-1} \quad (25)$$

and

$$\hat{x}_{k+1} = \hat{x}_k + \hat{\alpha}_k \hat{p}_k \quad (26)$$

$$L^{-1} x_{k+1} = L^{-1} x_k + \hat{\alpha}_k L^{-1} p_k, \quad \text{or} \quad (27)$$

$$x_{k+1} = x_k + \hat{\alpha}_k p_k \quad (28)$$

Therefore, the update equations remain the same, except for α_k and β_k being replaced by $\hat{\alpha}_k$ and $\hat{\beta}_k$ respectively, and an additional variable y_k . The final update equations are summarized in Algorithm 1

Algorithm 1: Preconditioned CG

```

1 while  $\|r\| > \eta$  do
2    $\hat{\beta}_{k+1} = \frac{r_{k+1}^T y_{k+1}}{r_k^T y_k}$ 
3    $y_k = M^{-1} r_k$ 
4    $p_k = -y_k + \hat{\beta}_k p_{k-1}$ , (next conjugate direction)
5    $\hat{\alpha}_k = \frac{r_k^T y_k}{p_k^T A p_k}$ 
6    $x_{k+1} = x_k + \hat{\alpha}_k p_k$  (move)
7    $k \leftarrow k + 1$ 
8 end while
```

with the initial direction $p_0 = -y_0$. It is easy to see that M^{-1} is precomputed (before the start of iterations) to be the inverse of $M = \tilde{C}\tilde{C}^T$ from the incomplete cholesky decomposition.

The additional overhead we have here is storing an extra variable y_k at each step, and an upfront cost of computing M^{-1} .

2.3 Vanilla CG (vs) Preconditioned CG

In order to verify if our implementations are correct, we verified that the solutions provided are the same as the direct linear system solution provided by MATLAB `mldivide` ($A \setminus b$) for simple 2×2 and 3×3 systems of equations. For the case of $n = 60$, both methods converge to the optimum x^* solution calculated using `mldivide` as shown in Figure 2.

- **Eigen Values & Condition number:** We see (in Figure 1) that the eigenvalue spread of the preconditioned matrix, \hat{A} is much better with lower condition number(8.13 compared to 78.31 of A). The lower condition number and better eigenvalue spread of the preconditioned matrix, \hat{A} has led to convergence of the algorithm in far fewer iterations.(see Figure 2) The reduction in number of iterations is about **3-fold**. Since we deal with a better conditioned problem, the solution will also be more robust to noise and floating point errors in the solution.
- **Time taken:** Even though the number of iterations is lesser for the preconditioned CG, the amount of time taken by preconditioned CG(about 1.69 ms in one of the runs) is about the same or more compared to vanilla CG(about 1.59 ms in one of the runs). It is important to note that the time taken varies quite a bit during each run with the exact same conditions. So, this observation must be taken with a pinch of salt. So, preconditioned CG provides faster convergence in terms of number of iterations but is about the same or worse in terms of time taken.
- **Quality of Solution:** The quality of solution obtained by both methods at convergence is approximately the same. It is also evident from Figure 2 that the norm of error in the solution plateaus when introducing noise in b . (i.e., it converges to a different, 'wrong' minima)
- **Rate of Convergence:** The rate of convergence/rate of decrease in error for both methods **remained constant** throughout the execution of the respective algorithms in the log scale. But the decrease in norm of error per iteration was higher for the preconditioned CG method which led to **convergence of similar quality** in fewer iterations. Also, as the error drops linearly in the log scale, the loss function(in linear scale) drops drastically in the first few iterations and then stabilizes close to the solution.

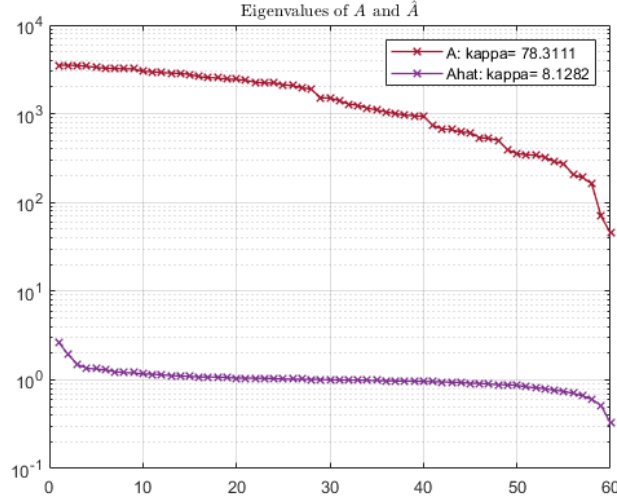


Figure 1: Eigenvalue spread of the original matrix, and the matrix corresponding to the preconditioned system $L^T A L$

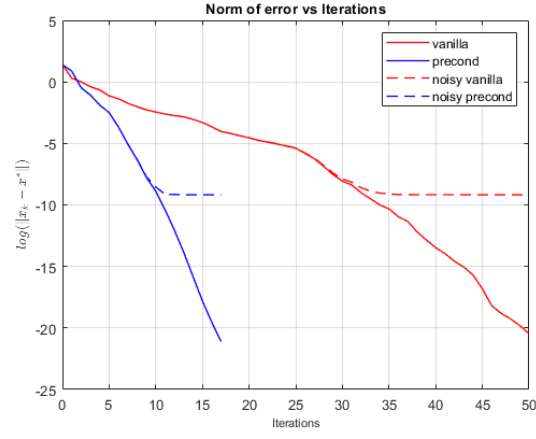


Figure 2: Evolution of the Error norm $\|x_k - x^*\|$

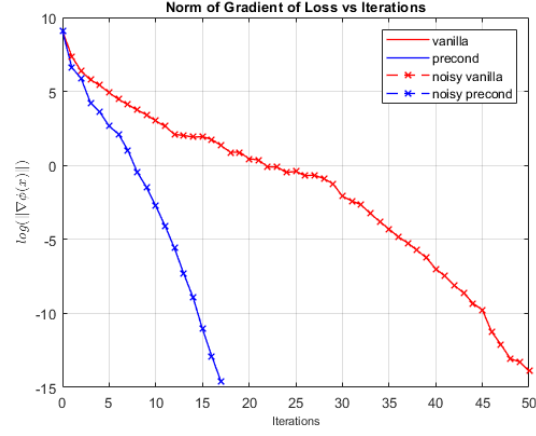


Figure 3: Norm of the gradient (=residue) r_k

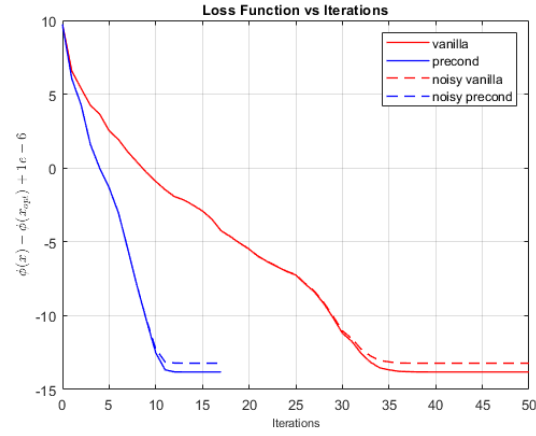


Figure 4: Objective Function values

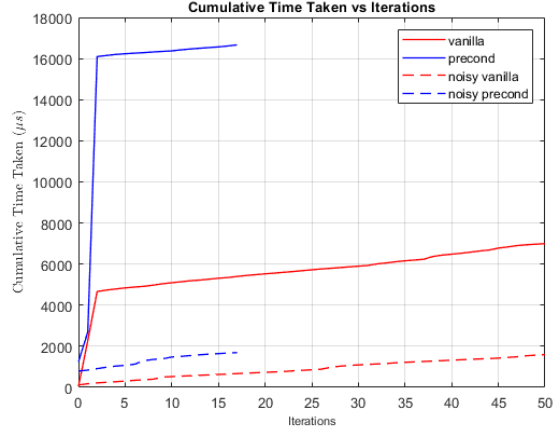


Figure 5: Cumulative time taken vs Iterations

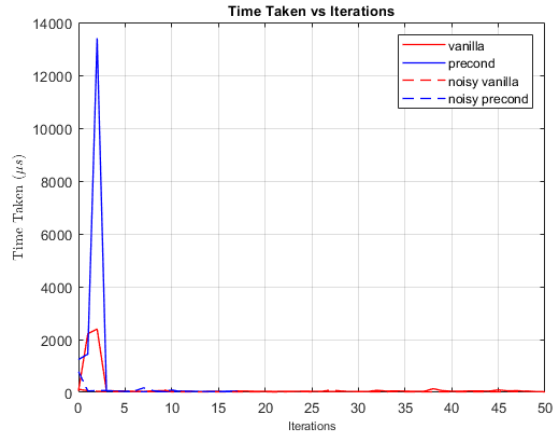


Figure 6: Time taken in each iteration

3 Nonlinear CG

Convex functions give a guarantee of achieving global minima, and preconditioning (section 2.2) efficiently solutions. Similar variants of the CG can be applied for minimizing general non linear functions. However, minimizing non - linear functions comes with the challenges of its own. To start with, lets consider a non-linear function $f(x_1, x_2)$ given by :

$$f(x_1, x_2) = -a \exp \left(-b \sqrt{\frac{1}{2}(x_1^2 + x_2^2)} \right) - \exp \left(\frac{1}{2}(\cos(cx_1) + \cos(cx_2)) \right) + a + \exp(1) \quad (29)$$

The function $f(x_1, x_2)$ is known as *Ackley function* whose behaviour is mostly flat outer region with a global minima at the center. Fig. 7 shows the distribution of f in the region $x_1 \in [-2, 2]$ and $x_2 \in [-2, 2]$.

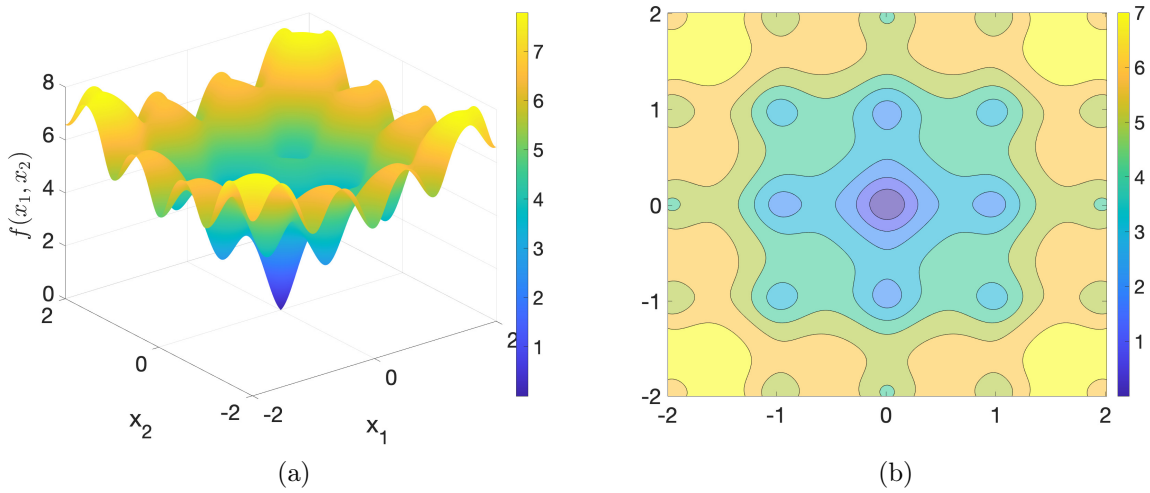


Figure 7: (a) Surface plot of Ackley function (b) Contour plots indicating local minima

Deviation from nice convex functions - Earlier we represented the function to be minimized using matrix relation and enforced the matrix to be positive definite, ensuring strict convexity. The gradients, step length α was calculated using the relation Eq. (3,19) which involved the linear operator A and the residue was just the gradient of ϕ . For non linear optimization we use similar steps but replace the gradient of ϕ with the gradient of non-linear function $f(x_1, x_2)$.

While there are various methods to solve the non-linear optimization, we implement two such algorithms:

1. Fletcher-Reeves method
2. Polak-Ribiere method

The non-linear methods are similar, in the sense that we determine a suitable (legitimate) search direction p_k and *then* the step length (α_k). One major difference is that we cannot clearly define any *conjugacy* wrt the nonlinear system, but we update the search direction with an expression identical to the previous case (which is different for each of these methods). The general flow of both the algorithms is as follows :

1. Set x_0 , the initial starting point and evaluate $f(x_0)$ and compute $\nabla f(x_0)$. Set descent direction $p_0 = -\nabla f(x_0)$
2. Compute α_k using line search (using Wolfe conditions)
3. Compute $x_{k+1} = x_k + \alpha_k p_k$
4. Compute β_{k+1} (different for each method)
5. Find descent direction $p_{k+1} = \nabla f_{k+1} + \beta_{k+1} p_k$
6. $k = k + 1$

Repeat step 2-6 until minima is reached. The ideal stopping condition would be $\nabla f_k = 0$, but is impractical. Instead, the stopping condition is set by thresholding the gradient, upper bounded by the number of iterations.

The major implementation difference in Fletcher-Reeves (FR) and Polak-Ribiere (PR) method is in choice of β_{k+1} .

$$\text{FR method : } \beta_{k+1} = \frac{\nabla f_{k+1}^T \nabla f_{k+1}}{\nabla f_k^T \nabla f_k} \quad (30)$$

$$\text{PR method : } \beta_{k+1} = \frac{\nabla f_{k+1}^T (\nabla f_{k+1} - \nabla f_k)}{\|\nabla f_k\|^2} \quad (31)$$

Line Search using strong Wolfe condition - It is important to get α_k correct, such that the *next* iteration has a legitimate descent direction. To illustrate this lets us consider the inner product with the update equation for descent direction at k^{th} step

$$\nabla f_k^T p_k = \|\nabla f_k\|^2 + \beta_k \nabla f_k^T p_{k-1} \quad (32)$$

For the inexact line search there is a possibility to get $\nabla f_k^T p_k > 0$, when the second term in the RHS of Eq. (32) dominates. This implies p_k is not the descent direction. This situation is avoided by considering the step length α_k that satisfies strong Wolfe conditions.[2]

$$f(x_k + \alpha_k p_k) \leq -f(x_k) + c_1 \alpha_k \nabla f_k^T p_k \quad (33)$$

$$|\nabla f(x_k + \alpha_k p_k)^T p_k| \leq -c_2 \nabla f_k^T p_k \quad (34)$$

where $0 < c_1 < c_2 < \frac{1}{2}$.

3.1 Implementation

The gradient of the Ackley function Eq. (29) is given by

$$\nabla f = \begin{bmatrix} c \exp(\cos cx_1/2 + \cos cx_2/2) \sin cx_1/2 + (abx_1 \exp(-b(x_1^2/2 + x_2^2/2)^{1/2}))/((2(x_1^2/2 + x_2^2/2)^{1/2})) \\ c \exp(\cos cx_1/2 + \cos cx_2/2) \sin cx_2/2 + (abx_2 \exp(-b(x_1^2/2 + x_2^2/2)^{1/2}))/((2(x_1^2/2 + x_2^2/2)^{1/2})) \end{bmatrix}$$

For better visualization we consider a region which has the global minima and few local minimas. We consider various different starting points, vary the parameters of sufficient decrease and curvature condition. For Fletcher Reeves to perform correctly the descent direction should always be kept legitimate, hence we use Wolfe conditions for line search. The different starting points (SP) considered are as follows

1. Start near to global minima - $P_1 = [0.03 \ 0.31]^T$
2. Start with the farther point to global minima - $P_2 = [0.57 \ 0.58]^T$
3. Start near to local minima $P_3 = [1 \ 0.2]^T$
4. Start farther to local minima $P_4 = [1.2 \ 0.27]^T$

The locations are chosen to generate all possible scenarios. All the figures corresponds to Fletcher-Reeves method with parameters $c_1 = 1e^{-4}$, $c_2 = 0.1$.

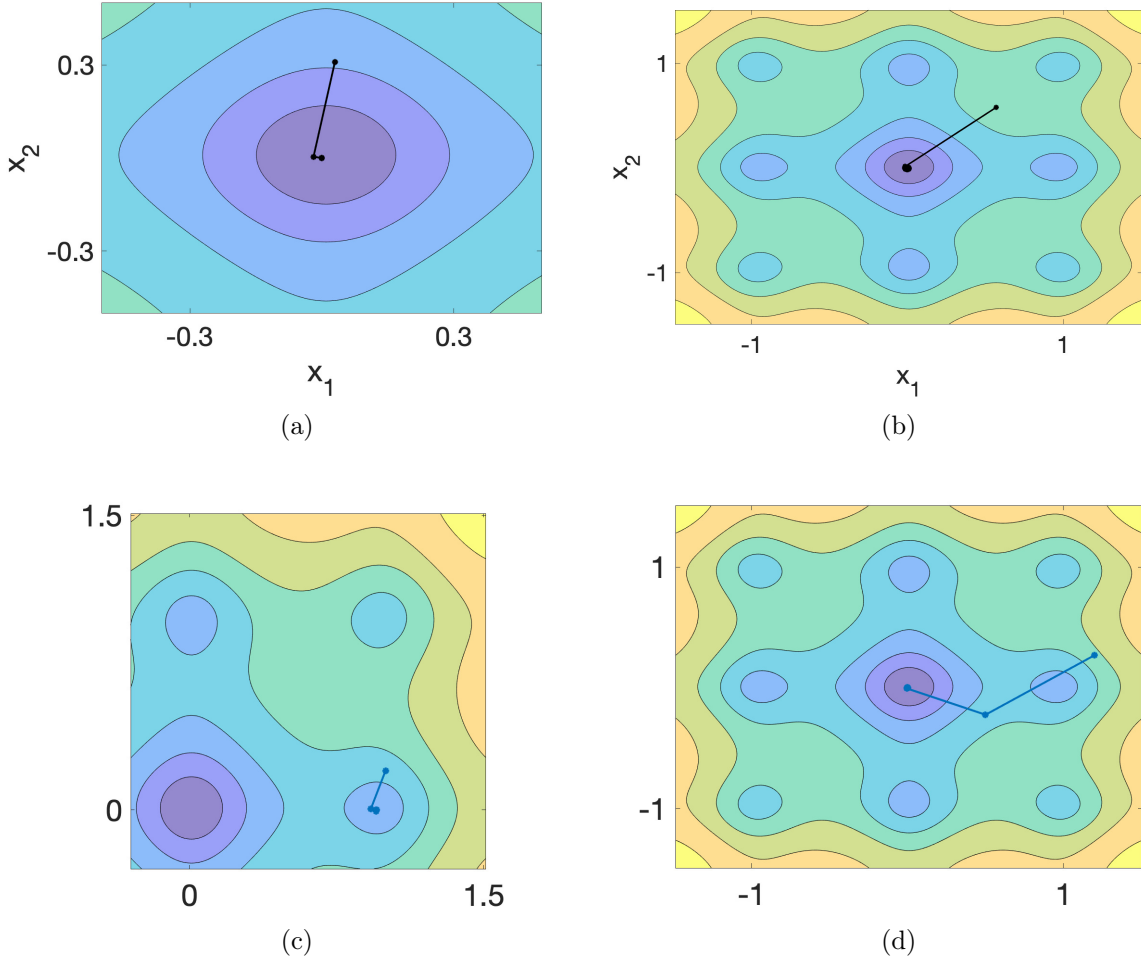


Figure 8: (a) $P_1 - [0.03 \ 0.31]^T$, SP is very close to the global minima, reaches the minimum easily (b) $P_2 - [0.57 \ 0.58]^T$ SP is far from global minima, but it is equally far from other local minima, the initial direction helps to fall to global minima (c) $P_3 - [1 \ 0.2]^T$, SP is very close to local minima and the algorithm stops at the local minima and (d) $P_4 = [1.2 \ 0.27]^T$, SP is near the same local minima but still the algorithm manages to reach the global minimum

Fig. 9 visualizes how the local minima is skipped. Similar results as seen in Fig. 8 is observed for PR method.

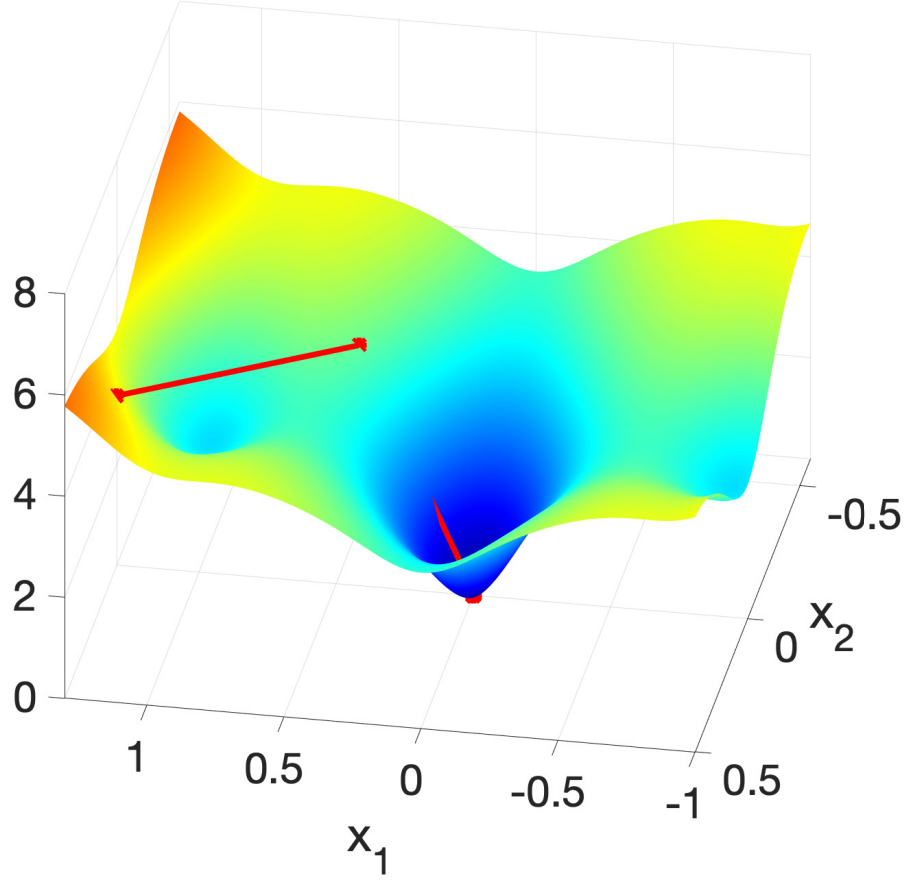


Figure 9: Surface plot showing the progression of x_k

3.2 Analysis

Let us start with observing the update equations x_{k+1} and p_{k+1} . We have

$$x_{k+1} = x_k + \alpha_k p_k$$

$$p_{k+1} = -\nabla f_k + \beta_k p_k$$

The above equation implies

$$p_{k-1} = \frac{x_k - x_{k-1}}{\alpha_{k-1}} \quad (35)$$

Substituting above two equations in x_{k+1} we get

$$x_{k+1} = x_k - \alpha_k \nabla f_k + \frac{\alpha_k \beta_{k-1}}{\alpha_{k-1}} (x_k - x_{k-1}) \quad (36)$$

Eq. 36 is similar to steepest descent update equation but with an additional momentum term. FR and PR methods are momentum based steepest descent methods.

Momentum problem : A general question is does the momentum based methods always perform better compared to SD approach. To understand this let us go back to evaluating α_k which was dependent on c_1 and c_2 . Let's change $c_1 = 1e^{-4}$ and $c_2 = 0.4$ and compare SD and PR for starting point $P = [1.2, 0.24]$.

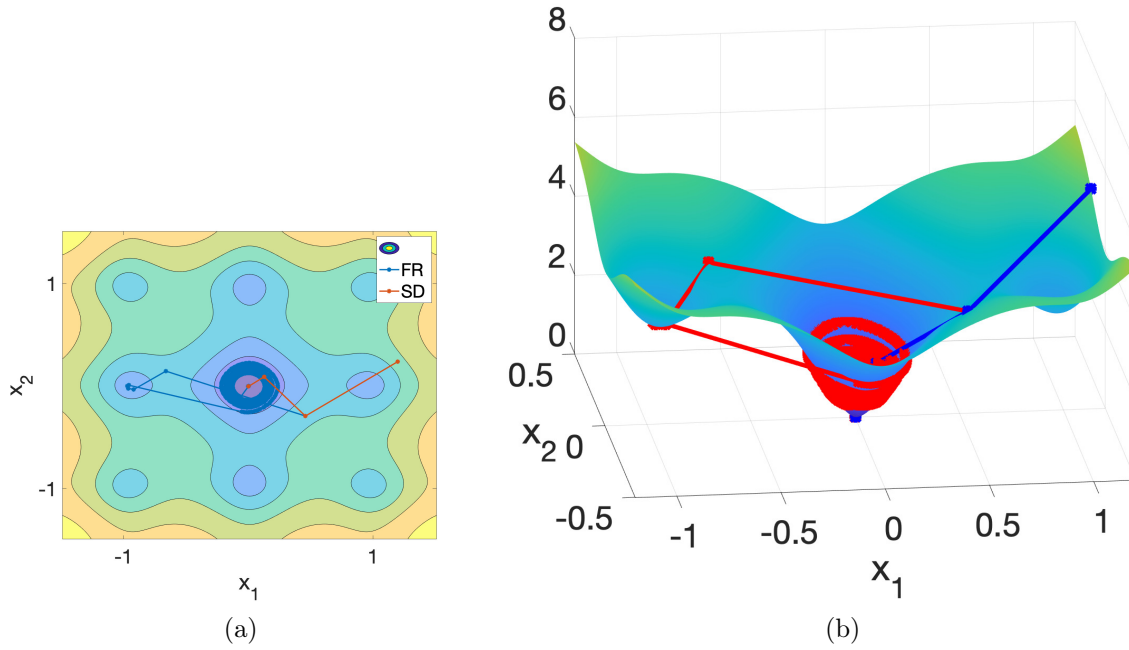


Figure 10: (a) Comparison of steepest descent and Fletcher-Reeves method. The sufficient and curvature parameters are changed. SD achieves minima with just 20 steps and FR takes 780 steps. For this run the threshold for minima is set at $1e^{-4}$, which means when function value reaches below the threshold the algorithm stops. (b) Due to the momentum created the Fletcher-Reeves method skips the global minima and reach a local minima and again jump to global minima from that point. The parameters c_1 and c_2 should be carefully chosen so that the momentum effect should not override the optimization problem.

Resetting the β_k : When the change in the α_k becomes very negligible than FR produces the descent direction which is approximately the same as in previous iteration. This makes the FR algorithm to get stuck and not improve any further. To overcome this problem the β_k is set to zero, which changes the descent direction to $-\nabla f_k$.

A MATLAB Code

A.1 Ordinary CG vs. Preconditioned CG

Listing 1: main script

```
1 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
2 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
3 % EE5121 — Convex Optimization
4 % Conjugate Gradient Method
5 %
6 % Authors: Anirudh R, Sreekar Sai R, Chandan Bhat
7 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
8 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
9
10 %Clearing and closing all figures
11 clear
12 close all;
13 %Clear screen
14 clc;
15 %Seeding the random number generator for repeatability
16 rng(5, 'twister');
17
18 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
19
20 %Generating a symmetric positive definite matrix
21 mat = generatePDMatrix(5, 25);
22 %Extracting eigen values of B
23 eig_values = eig(mat);
24 %Checking if eigen values are all >0
25 if eig_values > 0
26     disp(['Matrix generated is positive definite.']);
27 end
28
29 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
30
31 %'size' determines the size of A and b
32 N = 60;
33 %Amplitude of noise added to original b
34 noise = 1e-2;
35 %Factor multiplied with b_orig
36 mult = 1;
37 %Maximum eigen value we are willing to allow
38 max_eig = N^2;
39 %Generating matrix A and column vector b for calculations
40 A = generatePDMatrix(N, max_eig);
41 b_orig = mult * rand(N,1);
42 b = b_orig + noise * rand(N,1);
43
```

```

44 %Optimal solution to the linear equation  $Ax = b_{orig}$ 
45 x_opt = A \ b_orig;
46
47 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
48
49 %Parameters for running the iterative method
50 %Random initial point
51 x0 = rand(N, 1);
52 %Maximum number of iterations
53 max_iter = N;
54 %Tolerance for norm of gradient of loss at convergence
55 tolerance = 1e-6;
56
57 %Calling the conjugateGrad and preconditionedCG functions with pure b_orig
58 [x_hist, gf_hist, time_taken, num_iters] = conjugateGrad(A, b_orig, x0,
    max_iter, tolerance);
59 [x_hist_pre, gf_hist_pre, time_taken_pre, num_iters_pre, eig_values_pre,
    kappa] = preconditionedCG(A, b_orig, x0, max_iter, tolerance);
60
61 %Calling the functions with noisy b
62 [x_hist_noisy, gf_hist_noisy, time_taken_noisy, num_iters_noisy] =
    conjugateGrad(A, b, x0, max_iter, tolerance);
63 [x_hist_pre_noisy, gf_hist_pre_noisy, time_taken_pre_noisy,
    num_iters_pre_noisy, eig_values_pre_noisy, kappa_noisy] =
    preconditionedCG(A, b, x0, max_iter, tolerance);
64
65 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
66
67 %Plotting log of norm of gradient of loss function
68 figure;
69 % subplot(2,3,1);
70 plot(0:num_iters, log(gf_hist), 'r', 'LineWidth', 1, 'DisplayName', 'vanilla');
71 hold on;
72 plot(0:num_iters_pre, log(gf_hist_pre), 'b', 'LineWidth', 1, 'DisplayName', '
    precond');
73 hold on;
74 plot(0:num_iters_noisy, log(gf_hist_noisy), 'r-x', 'LineWidth', 1, '
    DisplayName', 'noisy vanilla');
75 hold on;
76 plot(0:num_iters_pre_noisy, log(gf_hist_pre_noisy), 'b-x', 'LineWidth', 1, '
    DisplayName', 'noisy precond');
77
78 grid;
79 title('Norm of Gradient of Loss vs Iterations', 'FontSize', 10);
80 ylabel('$\log( \|\nabla \phi(x)\| )$', 'interpreter', 'latex', 'FontSize',
    10);
81 xlabel('Iterations', 'FontSize', 8);
82 legend();

```

```

83 print('grad_norm','-dpdf');
84
85 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
86
87 %Plotting eigen values of A
88 figure;
89 % subplot(2,3,2);
90 semilogy(sort(eig(A),'descend'),'Marker','x','LineStyle','-', 'Color',
    [0.6350 0.0780 0.1840],'LineWidth',1,'DisplayName', ['A: kappa= '
    num2str(cond(A))]);
91 hold on;
92 semilogy(sort(eig_values_pre,'descend'),'Marker','x','LineStyle','-', '
    Color', [0.4940 0.1840 0.5560],'LineWidth',1,'DisplayName', ['Ahat:
    kappa= ' num2str(kappa)]);
93 grid;
94 title('Eigenvalues of $A$ and $\hat{A}$', 'interpreter','latex', '
    FontSize', 10);
95 legend();
96 print('eig','-dpdf');
97
98 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
99
100 %Plotting log of norm of error
101 figure;
102 % subplot(2,3,3);
103 plot(0:num_iters, log( vecnorm( x_hist- x_opt ) ),'r','LineWidth',1,'
    DisplayName', 'vanilla');
104 hold on;
105 plot(0:num_iters_pre, log( vecnorm( x_hist_pre- x_opt ) ),'b','LineWidth'
    ,1,'DisplayName', 'precond');
106 hold on;
107 plot(0:num_iters_noisy, log( vecnorm( x_hist_noisy- x_opt ) ),'r--','
    LineWidth',1,'DisplayName', 'noisy vanilla');
108 hold on;
109 plot(0:num_iters_pre_noisy, log( vecnorm( x_hist_pre_noisy- x_opt ) ),'b--
    ','LineWidth',1,'DisplayName', 'noisy precondition');
110
111 grid;
112 title('Norm of error vs Iterations', 'FontSize', 10);
113 ylabel('$\log(\|x_k-x^*\|)$', 'interpreter','latex', 'FontSize', 10);
114 xlabel('Iterations', 'FontSize', 8);
115 legend();
116 print('error_norm','-dpdf');
117
118 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
119
120 %Plotting time taken
121 figure;

```

```

122 % subplot(2,3,4);
123 plot(0:num_iters, 10^6*time_taken, 'r', 'LineWidth', 1, 'DisplayName', 'vanilla
    ');
124 hold on;
125 plot(0:num_iters_pre, 10^6*time_taken_pre, 'b', 'LineWidth', 1, 'DisplayName',
    'precond');
126 hold on;
127 plot(0:num_iters_noisy, 10^6*time_taken_noisy, 'r—', 'LineWidth', 1, '
    DisplayName', 'noisy vanilla');
128 hold on;
129 plot(0:num_iters_pre_noisy, 10^6*time_taken_pre_noisy, 'b—', 'LineWidth', 1,
    'DisplayName', 'noisy precond');
130
131 grid;
132 title('Time Taken vs Iterations', 'FontSize', 10);
133 ylabel('Time Taken ( $\mu$  s)', 'interpreter', 'latex', 'FontSize', 10);
134 xlabel('Iterations', 'FontSize', 8);
135 legend();
136 print('time', '-dpdf');
137
138 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
139
140 %Plotting cumulative time taken vs iterations
141 %Calculating cumulative times taken
142 time_taken_cum = zeros(num_iters+1, 1);
143 time_taken_cum(1) = time_taken(1);
144 for i=2:num_iters+1
145     time_taken_cum(i) = time_taken_cum(i-1) + time_taken(i);
146 end
147
148 time_taken_cum_pre = zeros(num_iters_pre+1, 1);
149 time_taken_cum_pre(1) = time_taken_pre(1);
150 for i=2:num_iters_pre+1
151     time_taken_cum_pre(i) = time_taken_cum_pre(i-1) + time_taken_pre(i);
152 end
153
154 time_taken_cum_noisy = zeros(num_iters_noisy+1, 1);
155 time_taken_cum_noisy(1) = time_taken_noisy(1);
156 for i=2:num_iters_noisy+1
157     time_taken_cum_noisy(i) = time_taken_cum_noisy(i-1) + time_taken_noisy
        (i);
158 end
159
160 time_taken_cum_pre_noisy = zeros(num_iters_pre_noisy+1, 1);
161 time_taken_cum_pre_noisy(1) = time_taken_pre_noisy(1);
162 for i=2:num_iters_pre_noisy+1
163     time_taken_cum_pre_noisy(i) = time_taken_cum_pre_noisy(i-1) +
        time_taken_pre_noisy(i);

```



```

164 end
165
166 figure;
167 % subplot(2,3,5);
168 plot(0:num_iters, 10^6*time_taken_cum, 'r', 'LineWidth', 1, 'DisplayName', '
    vanilla');
169 hold on;
170 plot(0:num_iters_pre, 10^6*time_taken_cum_pre, 'b', 'LineWidth', 1, '
    DisplayName', 'precond');
171 hold on;
172 plot(0:num_iters_noisy, 10^6*time_taken_cum_noisy, 'r—', 'LineWidth', 1, '
    DisplayName', 'noisy vanilla');
173 hold on;
174 plot(0:num_iters_pre_noisy, 10^6*time_taken_cum_pre_noisy, 'b—', 'LineWidth
    ', 1, 'DisplayName', 'noisy precondition');
175
176 grid;
177 title('Cumulative Time Taken vs Iterations', 'FontSize', 10);
178 ylabel('Cumulative Time Taken ( $\mu$  s)', 'interpreter', 'latex', 'FontSize'
    , 10);
179 xlabel('Iterations', 'FontSize', 8);
180 legend();
181 print('time_cum', '-dpdf');
182
183 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
184
185 %Plotting the values taken by the loss function for each x in x_hist and
    x_hist_pre
186 %Defining the loss function
187 f = @(x) 0.5*x'*A*x - b_orig'*x;
188
189 %Calculating loss value at columns in x_hist, x_hist_pre, x_hist_noisy and
    x_hist_pre_noisy
190 func_vals = zeros(num_iters+1, 1);
191 for i=1:(num_iters+1)
192     func_vals(i, 1) = f(x_hist(:, i));
193 end
194
195 func_vals_pre = zeros(num_iters_pre+1, 1);
196 for i=1:(num_iters_pre+1)
197     func_vals_pre(i, 1) = f(x_hist_pre(:, i));
198 end
199
200 func_vals_noisy = zeros(num_iters_noisy+1, 1);
201 for i=1:(num_iters_noisy+1)
202     func_vals_noisy(i, 1) = f(x_hist_noisy(:, i));
203 end
204

```

```

205 func_vals_pre_noisy = zeros(num_iters_pre_noisy+1, 1);
206 for i=1:(num_iters_pre_noisy+1)
207     func_vals_pre_noisy(i, 1) = f(x_hist_pre_noisy(:, i));
208 end
209
210 %Plotting log of loss wrt loss at minima(Loss at x_opt is subtracted from
    loss at every point)
211 figure;
212 % subplot(2,3,6);
213 plot(0:num_iters, log(func_vals - f(x_opt)+1e-6), 'r', 'LineWidth', 1, '
    DisplayName', 'vanilla');
214 hold on;
215 plot(0:num_iters_pre, log(func_vals_pre - f(x_opt)+1e-6), 'b', 'LineWidth'
    , 1, 'DisplayName', 'precond');
216 hold on;
217 plot(0:num_iters_noisy, log(func_vals_noisy - f(x_opt)+1e-6), 'r--', '
    LineWidth', 1, 'DisplayName', 'noisy vanilla');
218 hold on;
219 plot(0:num_iters_pre_noisy, log(func_vals_pre_noisy - f(x_opt)+1e-6), 'b--'
    , 'LineWidth', 1, 'DisplayName', 'noisy precondition');
220
221 grid;
222 title('Loss Function vs Iterations', 'FontSize', 10);
223 ylabel('$\phi(x) - \phi(x_{opt}) + 1e-6$', 'interpreter', 'latex', '
    FontSize', 10);
224 xlabel('Iterations', 'FontSize', 8);
225 legend();
226 print('phi', '-dpdf');
227
228 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

```

Listing 2: Code to implemented Vanilla CG

```

1 function [x_hist, gf_hist, time_taken, k] = conjugateGrad(A, b, x0,
    max_iter, tolerance)
2     %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
3     % Function to run the Conjugate Gradient Method(CGM).
4     % Argument:
5     %     A                Symmetric PD matrix describing the
        linear system
6     %     b                Target of the linear system
7     %     x0              Initial value of 'x' to start the
        iterative algorithm
8     %     max_iter        Maximum number of iterations allowed
9     %     tolerance       Tolerance for norm of gradient of loss
        function at convergence
10    % Return:
11    %     x_hist           History of values taken by x through the
        iterations

```

```

12 % gf_hist      History of norm of gradient of loss function
13 %             through the iterations
14 %             time_taken      Time taken for execution of each iteration
15 %             k                Number of iterations that actually
16 %             occurred
17 %%%%%%%%%%%%%%
18 %Starting measurement of time for initialization tasks
19 tic
20 %The gradient of loss function
21 gradf = @(x) A*x - b;
22
23 %x = x + alpha p, where alpha = r'*r/(p'*A*p)
24 x = x0;
25 r = gradf(x);
26 p = -r;
27 %Storing x and norm of gradient in the history
28 x_hist = [x];
29 gf_hist = [norm(r)];
30 k = 0;
31
32 %Logging the time taken before iterating
33 time_taken = toc;
34
35 %Displaying values before the iterations
36 fprintf('\n');
37 disp('Vanilla Conjugate Gradient Method:')
38 fprintf('\n');
39 disp(['CG: k=0, gf=' num2str(norm(r)) ' , time elapsed: ' num2str
40      (time_taken*10^6) ' micro seconds']);
41
42 while norm(r) > tolerance && k < max_iter
43     %Starting measurement of time for this iteration
44     tic
45
46     %Step size calculated according to CGM
47     step_size = r'*r/(p'*A*p);
48     %Update of x
49     x_new = x + step_size * p;
50
51     %x is updated. Updating parameters for next iteration
52     r_new = r + step_size * A * p;
53     beta_ = r_new'*r_new/(r'*r);
54     p = -r_new + beta_ *p;
55     r = r_new;
56     x = x_new;

```

```

57         %Storing the norm of grad and updated x
58         gf_hist = [gf_hist, norm(r_new)];
59         x_hist = [x_hist, x_new];
60         %Updating iterate for next iteration
61         k = k+1;
62
63         %Ending measurement of time and logging
64         time_taken = [time_taken, toc];
65
66         disp(['CG: k=' num2str(k) ' , gf=' num2str(norm(r_new)) '
67             , time elapsed: ' num2str(time_taken(end)*10^6) '
68             micro seconds']);
69     end
70     %Displaying relevant details of the CGM execution
71     fprintf('\n');
72     disp(['Total number of iterations: ' num2str(k) ' , Total time
73         taken: ' num2str(sum(time_taken)*10^6) ' micro seconds']);
74     fprintf('\n');
75 end

```

Listing 3: Function to generate random positive definite matrix

```

1 function [matrix] = generatePDMatrix(n, max_eig)
2     %%%%%%%%%%%
3     % Function to generate a nxn positive definite matrix.
4     % Argument:
5     %         n                Size of the square matrix
6     % Return:
7     %         matrix           nxn positive definite matrix
8     %%%%%%%%%%%
9
10    %Generating a random nxn matrix
11    temp = rand(n, n);
12    %Making a nxn symmetric matrix from random nxn matrix
13    temp = 0.5*(temp + temp');
14    %Eigen decomposition of nxn symmetric matrix
15    [V, D] = eig(temp);
16    %Choosing positive eigen values to construct our positive definite
17    %matrix
18    eig_vals = randi(max_eig, n, 1);
19    %Constructing the positive definite matrix using eig_vals and V
20    %corresponding to matrix temp
21    matrix = V*diag(eig_vals)*V';
22 end

```

Listing 4: Function to implemented Preconditioned CG

```

1 function [x_hist, gf_hist, time_taken, k, eig_values, kappa] =
2     preconditionedCG(A, b, x0, max_iter, tolerance)

```

```

2      %%%%%%%%%%%
3      % Function to run the Preconditioned Conjugate Gradient Method(CGM
4      % Argument:
5      %      A                      Symmetric PD matrix describing the
        linear system
6      %      b                      Target of the linear system
7      %      x0                     Initial value of 'x' to start the
        iterative algorithm
8      %      max_iter              Maximum number of iterations allowed
9      %      tolerance             Tolerance for norm of gradient of loss
        function at convergence
10     % Return:
11     %      x_hist                History of values taken by x through the
        iterations
12     %      gf_hist              History of norm of gradient of loss function
        through the iterations
13     %      time_taken           Time taken for execution of each iteration
14     %      k                    Number of iterations that actually
        occurred
15     %      eig_values           Eigenvalues of Ahat
16     %      kappa                Condition number of Ahat
17     %%%%%%%%%%%
18
19     %Starting measurement of time for initialization tasks
20     tic
21
22     %      Preconditioning
23     op.type = 'ict';
24     op.droptol = 1e-2;
25     C = ichol(sparse(A),op);
26     Minv = inv(C*C');
27
28     %The gradient of loss function
29     gradf = @(x) A*x - b;
30
31     %x = x + alpha p, where alpha = r'*r/(p'*A*p)
32     x = x0;
33     r = gradf(x);
34     y = Minv*r;
35     p = -y;
36     %Storing x and norm of gradient in the history
37     x_hist = [x];
38     gf_hist = [norm(r)];
39     k = 0;
40
41     %Logging the time taken before iterating
42     time_taken = toc;

```

```

43
44 %Displaying values before the iterations
45 fprintf('\n');
46 disp('Preconditioned Conjugate Gradient Method:')
47 fprintf('\n');
48 disp(['PCG: k=0, gf=' num2str(norm(r)) ' , time elapsed: '
      num2str(time_taken*10^6) ' micro seconds']);
49
50 while norm(r) > tolerance && k < max_iter
51     %Starting measurement of time for this iteration
52     tic
53
54     %Step size calculated according to CGM
55     step_size = r'*y/(p'*A*p);
56     %Update of x
57     x_new = x + step_size * p;
58
59     %x is updated. Updating parameters for next iteration
60     r_new = r + step_size * A * p;
61     y_new = Minv*r_new;
62
63     beta_ = r_new'*y_new/(r'*y);
64     p = -y_new + beta_ *p;
65     r = r_new;
66     x = x_new;
67     y = y_new;
68
69     %Storing the norm of grad and updated x
70     gf_hist = [gf_hist, norm(r_new)];
71     x_hist = [x_hist, x_new];
72     %Updating iterate for next iteration
73     k = k+1;
74
75     %Ending measurement of time and logging
76     time_taken = [time_taken, toc];
77
78     disp(['PCG: k=' num2str(k) ' , gf=' num2str(norm(r_new)) '
      , time elapsed: ' num2str(time_taken(end)*10^6) '
      micro seconds']);
79 end
80 %Displaying relevant details of the CGM execution
81 fprintf('\n');
82 disp(['Total number of iterations: ' num2str(k) ' , Total time
      taken: ' num2str(sum(time_taken)*10^6) ' micro seconds']);
83 fprintf('\n');
84
85 %Calculating the  $\hat{A}$  to study the eigenvalues and condition number
86 Cinv = inv(C);

```

```

87     Ahat = Cinv*A*Cinv';
88     eig_values = eig(Ahat);
89     kappa = cond(Ahat);
90 end

```

A.2 Nonlinear CG

Listing 5: lineSearch.m

```

1
2
3
4 function alpha = lineSearch(f,gradf,gradphi_0,c1,c2,x,phi_x,pk,rho)
5
6 % This code performs Line search for finding optimal alpha (step length)
7 % using strong Wolfe conditions. The algorithm for line search method is
8 % used from Numerical optimization by Jorge Nocedal Stephen J. Wright
9 % Chapter 3, Algorithm 3.5.
10
11 % The parameters passed to the search algorithm are
12 % f      — objective function
13 % gradf  — gradient of the function
14 % c1     — sufficient decrease constant
15 % c2     — curvature condition constant
16 % x      — solution at current iterate
17 % pk     — search direction
18 % phi(alpha) = f(x + alpha p_k)
19 % phi_x   — function value at current x
20 % gradphi_0 — gradient of function at x
21
22
23 % Parameters used in the algorithm
24 % xi      — the updated x (x_prev + alpha_i*pk)
25 % alpha_lo — alpha which supplied to zoom function
26 % alpha_hi — higher alpha
27
28
29 % Initialize the alpha by specifying the range it can take. alpha > 0
30 alpha_0 = 0; alpha_max = 200; alpha_1 = 1; % alpha_1 \in (0,alpha_max)
31 alpha_i = alpha_1; alpha_prev = alpha_0;
32 phi_prev = phi_x; phi_0 = phi_x;
33 x_prev = x; i = 0;
34
35 while true && (abs(alpha_prev - alpha_i) > 1e-4)
36
37     xi      = x_prev + alpha_i*pk;
38     phi_xi   = f(xi);
39     gradphi_xi = gradf(xi)'*pk;

```

```

40
41 % alpha_1 violates the sufficient condition
42 if (phi_xi > phi_0 + c1*alpha_i*gradphi_0) || ((phi_xi >= phi_prev )
    && (i > 0))
43     alpha = zoom(f,gradf,gradphi_0,c1,c2,x,phi_0,phi_prev,pk,
        alpha_prev,alpha_i);
44     break;
45 end
46
47 % strong Wolfe check
48 if(abs(gradphi_xi) <= -c2*gradphi_0)
49     alpha = alpha_i;
50     break;
51 end
52
53 % Check if we moved ahead
54 if (gradphi_xi >= 0)
55     alpha = zoom(f,gradf,gradphi_0,c1,c2,x,phi_0,phi_prev,pk,
        alpha_prev,alpha_i);
56     break;
57 end
58
59 alpha_prev = alpha_i;
60 alpha_i = min(rho*alpha_i,alpha_max);
61 phi_prev = phi_xi;
62 i = i + 1;
63
64 end
65
66 end

```

Listing 6: NLG1.m

```

1
2
3 function [x_k,val,f_hist,x_hist,k,alpha_hist] = NL_G1(f,gradf,x0,pk,c1,c2
    ,rho,method)
4
5 % This code performs conjugate gradient on Non — linear functions.
6 % The update equation for beta differs by the method of choice
7 % (a) Fletcher Reeves (b) Polak Riebere
8
9 % Initialization
10 tol = 1e-4;
11 gradf_x0 = gradf(x0);
12 gradphi_prev = gradf_x0'*pk;
13 x_prev = x0; phi_xk = f(x0); phi_prev = phi_xk + 10;
14 x_k = x0; gradf_prev = 1;
15 gradf_xk = gradf_x0; k = 0;

```



```

16 x_hist = x_k; f_hist = f(x_k); alpha_hist = [];
17
18 while norm(gradf_xk) > tol && norm(phi_xk - phi_prev) > tol
19
20     alpha_k = lineSearch(f,gradf,gradphi_prev,c1,c2,x_k,phi_xk,pk,rho);
21     x_k     = x_prev + alpha_k*pk;
22     gradf_prev = gradf_xk;
23     gradf_xk  = gradf(x_k);
24     phi_xk    = f(x_k);
25     phi_prev  = f(x_prev);
26     if strcmp(method,'fletcher')
27         beta_k = (gradf_xk'*gradf_xk)/(gradf_prev'*gradf_prev);
28     elseif strcmp(method,'polak')
29         beta_k = (gradf_xk'*(gradf_xk - gradf_prev))/(gradf_prev'*
30             gradf_prev);
31         if beta_k < 0
32             beta_k = 0;
33         end
34     elseif strcmp(method,'SD')
35         beta_k = 0;
36     end
37     alpha_hist = [alpha_hist,alpha_k];
38     if alpha_k < 1e-5
39         beta_k = 0;
40     end
41     pk = -gradf_xk + beta_k*pk;
42     k = k + 1;
43     x_prev = x_k;
44     gradphi_prev = gradf_xk'*pk;
45     x_hist = [x_hist,x_k];
46     f_hist = [f_hist,f(x_k)];
47 end
48 val = f(x_k);
49
50 end

```

Listing 7: TestNLCG.m

```

1 clear all; close all; clc
2
3 method = 'fletcher';
4 % Ackley Function, global minimum at (0,0)
5 syms x
6 a = 20; b = 0.2; c = 2*pi;
7 f1 = @(x) -a*exp(-b*sqrt(0.5*(x(1)^2 + x(2)^2))) -exp(0.5*(cos(c*x(1))
8     +cos(c*x(2)))) + a + exp(1);
9 gradf1 = @(x)[(c*exp(cos(c*x(1)))/2 + cos(c*x(2))/2)*sin(c*x(1)))/2 + (a*b*
10     x(1)*exp(-b*(x(1)^2/2 + x(2)^2/2)^(1/2)))/(2*(x(1)^2/2 + x(2)^2/2)

```

```

    ^ (1/2));
9    (c*exp(cos(c*x(1))/2 + cos(c*x(2))/2)*sin(c*x(2)))/2 + (a*b*x(2)*exp(-
    b*(x(1)^2/2 + x(2)^2/2)^(1/2)))/(2*(x(1)^2/2 + x(2)^2/2)^(1/2)));
10
11 % Parameters
12 c1 = 1e-4; c2 = 0.4; rho = 2;
13
14 % starting point for the non-linear CG
15 % Choose any one and comment the rest
16 x0 = [1.2;0.27];
17 %x0 = [-0.55;-0.03];
18 %x0 = [0.57;0.5];
19 %x0 = [0.57;0.58];
20 %x0 = [1;0.2];
21 x0 = [1.2;0.24];
22 %x0 = [0.03; 0.31]
23 % Test for a simple convex function, uncomment following lines
24 % f1 = @(x) 4*x(1)^2+x(2)^2;
25 % gradf1 =@(x) [8*x(1);2*x(2)];
26
27 pk = -gradf1(x0); % initial search direction
28 phi_xk = f1(x0); gradphi_prev = gradf1(x0)'*pk;
29 [~,~,f_hist,x_hist,k,alpha_hist] = NL_CG1(f1,gradf1,x0,pk,c1,c2,rho,method
    );
30
31
32 % Plotting
33 plotting = true;
34
35 if plotting
36     syms x1 y1
37     f2 = @(x1,y1) -a*exp(-b.*sqrt(0.5*(x1.^2 + y1.^2))) -exp(0.5.*(cos(c
        .*x1)+cos(c.*y1))) + a + exp(1);
38     sampling = 0.01;
39     xc = -1.5 : sampling : 1.5;
40     [Xc,Yc] = meshgrid(xc);
41     f_out = f2(Xc,Yc);
42     figure
43     surf(Xc,Yc,f_out)
44     shading interp
45     hold on
46     plot3(x_hist(1,:),x_hist(2,:),f_hist(:),'*-r','MarkerSize',10,'
        LineWidth',4); hold on
47
48     figure
49     [~, contourObj] = contourf(f_out);
50     % This is the secret that 'keeps' the transparency.
51     eventFcn = @(srcObj, e) updateTransparency(srcObj);

```

```

52     addlistener(contourObj, 'MarkedClean', eventFcn);
53     hold on
54     plot(x_hist(1,:)*(1/sampling)+floor(length(xc)/2),x_hist(2,:)*(1/
        sampling)+floor(length(xc)/2),'*-', 'MarkerSize',5,'LineWidth',2);
        hold on
55
56 end
57
58 display = true;
59 if display
60
61 fprintf('
        _____\n')
62
63 fprintf('Initial Point : ');
64 fprintf('%g, ', x0(1:end-1));
65 fprintf('%g]\n', x0(end));
66
67 fprintf('Minimum Point : ');
68 fprintf('%g, ', x_hist(1:end-1,end));
69 fprintf('%g]\n', x_hist(end,end));
70
71 fprintf('Minimum value : %g \n',f_hist(end));
72 fprintf('Method : %s \n',method);
73
74 fprintf('Steps taken %d \n',k);
75
76 fprintf('
        _____\n')
77
78 end
79
80 comparison = true;
81 if comparison
82
83 method = 'SD';
84 [~,~,f_hist_sd,x_hist_sd,k_sd,alpha_hist_sd] = NL_CG1(f1,gradf1,x0,pk,c1,
        c2,rho,method);
85 method = 'fletcher';
86 [~,~,f_hist_fr,x_hist_fr,k_fr,alpha_hist_fr] = NL_CG1(f1,gradf1,x0,pk,c1,
        c2,rho,method);
87 method = 'polak';
88 [~,~,f_hist_pr,x_hist_pr,k_pr,alpha_hist_pr] = NL_CG1(f1,gradf1,x0,pk,c1,
        c2,rho,method);
89
90 fprintf('
        _____\n')
91

```

```

92 fprintf('Minimum value : %g \n',f_hist_sd(end));
93 fprintf('Method : %s \n','SD');
94 fprintf('Minimum value : %g \n',f_hist_fr(end));
95 fprintf('Method : %s \n','fletcher');
96 fprintf('Minimum value : %g \n',f_hist_pr(end));
97 fprintf('Method : %s \n','polak');
98
99
100
101 fprintf('Steps taken SD %d \n',k_sd);
102 fprintf('Steps taken FR %d \n',k_fr);
103 fprintf('Steps taken PR %d \n',k_pr);
104
105 fprintf('
    _____\n')
106
107 end
108
109 plotting_cmp = true;
110 if plotting_cmp
111     figure
112     plot(f_hist_sd); hold on
113     plot(f_hist_fr); hold on
114     plot(f_hist_pr); hold on
115 end

```

Listing 8: zoom.m

```

1
2
3 function alpha = zoom(f,gradf,gradphi_0,c1,c2,x,phi_0,phi_prev,pk,alpha_lo
    ,alpha_hi)
4
5 % Zoom function — each iteration generates alpha between alpha_lo and
6 % alpha_hi, and then replaces one of these endpoints by alpha, in such a
7 % way that Wolfe condition is satisfied in the interval!
8
9 while true
10
11     % Choose an alpha in (alpha_lo, alpha_hi). There are many ways to
12     % select the alpha (a) interpolation or (b) just select a midpoint or
13     % (c)
14
15     alpha = (alpha_lo+alpha_hi)/2;
16     xi = x + alpha*pk;
17     phi_xi = f(xi);
18     gradphi_xi = gradf(xi)'*pk;
19

```

```

20     if phi_xi > phi_0 + c1*alpha*gradphi_0 || phi_xi >= phi_prev
21         alpha_hi = alpha;
22     else
23         gradphi_xi = gradf(xi)'*pk;
24         % Check strong Wolfe condition
25         if abs(gradphi_xi) <= -c2*gradphi_0
26             return;
27         end
28
29         if gradphi_xi*(alpha_hi-alpha_lo) >= 0
30             alpha_hi = alpha_lo;
31         end
32         alpha_lo = alpha;
33         phi_prev = phi_xi;
34
35     end
36
37     if alpha < 1e-5
38         return
39     end
40
41
42     end
43
44 end

```

References

- [1] https://en.wikipedia.org/wiki/Computational_complexity_of_mathematical_operations
- [2] Numerical Optimization by Jorge Nocedal and Stephen J. Wright, Springer, 2006.