

# Numerical Optimization, 2023 Fall

## Homework 4

Name: Zhou Shouchen  
Student ID: 2021533042

Due 23:59 (CST), Jan. 4, 2024

Problem 1.  $f$  is a positive definite quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}, \quad \mathbf{A} \in \mathbb{S}_{++}^n, \mathbf{b} \in \mathbb{R}^n,$$

$\mathbf{x}^k$  is the current iteration point,  $\mathbf{d}^k$  is the descent direction. Derive the step size of exact linear search [20pts]

$$\alpha^k = \arg \min_{\alpha > 0} f(\mathbf{x}^k + \alpha \mathbf{d}^k).$$

$$\text{Let } g(\alpha) = f(\mathbf{x}^k + \alpha \mathbf{d}^k) = \frac{1}{2} (\mathbf{x}^k + \alpha \mathbf{d}^k)^T \mathbf{A} (\mathbf{x}^k + \alpha \mathbf{d}^k) + \mathbf{b}^T (\mathbf{x}^k + \alpha \mathbf{d}^k).$$

So

$$\frac{\partial g}{\partial \alpha} = \alpha (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k + \frac{1}{2} (\mathbf{d}^k)^T \mathbf{A} \mathbf{x}^k + \frac{1}{2} (\mathbf{x}^k)^T \mathbf{A} \mathbf{d}^k + \mathbf{b}^T \mathbf{d}^k$$

Since  $\mathbf{A} \in \mathbb{S}_{++}^n$ , i.e.  $\mathbf{A}$  is symmetric, so

$$(\mathbf{d}^k)^T \mathbf{A} \mathbf{x}^k = (\mathbf{x}^k)^T \mathbf{A} \mathbf{d}^k$$

So

$$\frac{\partial g}{\partial \alpha} = \alpha (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k + (\mathbf{d}^k)^T \mathbf{A} \mathbf{x}^k + \mathbf{b}^T \mathbf{d}^k$$

And

$$\frac{\partial^2 g}{\partial \alpha^2} = (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k$$

Since  $\mathbf{A} \in \mathbb{S}_{++}^n$ , which means that  $\mathbf{A}$  is a positive defined matrix, i.e.

$$\forall \mathbf{x} \in \mathbb{R}^n, \mathbf{x}^T \mathbf{A} \mathbf{x} > 0$$

So

$$\forall \mathbf{d}^k, \frac{\partial^2 g}{\partial \alpha^2} = (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k > 0$$

So  $g(\alpha)$  is a convex function. In order to find the minimum point  $\alpha^k = \arg \min_{\alpha > 0} g(\alpha)$ , we just need to let the gradient to be 0. i.e.

$$\frac{\partial g}{\partial \alpha}(\alpha^k) = 0$$

So

$$\alpha^k (\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k + (\mathbf{d}^k)^T \mathbf{A} \mathbf{x}^k + \mathbf{b}^T \mathbf{d}^k = 0$$

So

$$\alpha^k = - \frac{((\mathbf{d}^k)^T \mathbf{A} \mathbf{x}^k + \mathbf{b}^T \mathbf{d}^k)}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k}$$

Since we have known that  $(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k > 0$ , so the  $\alpha^k$  is a legal solution.

So above all, the step size of exact linear search is

$$\alpha^k = - \frac{((\mathbf{d}^k)^T \mathbf{A} \mathbf{x}^k + \mathbf{b}^T \mathbf{d}^k)}{(\mathbf{d}^k)^T \mathbf{A} \mathbf{d}^k}$$

Problem 2. Prove that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is affine if and only if  $f$  is both convex and concave. [20pts]

1. sufficiency:

If  $f$  is affine, then we can write  $f$  as  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ , where  $\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}$ .

Then we have

$$\nabla f(\mathbf{x}) = \mathbf{w}$$

and

$$\nabla^2 f(\mathbf{x}) = \mathbf{0}$$

So

$$\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}, \nabla^2 f(\mathbf{x}) \preceq \mathbf{0}$$

So  $f$  is both convex and concave.

2. necessity:

If  $f$  is convex and concave, then we have:

since  $f$  is concave, so

$$\forall \mathbf{x}_1, \mathbf{x}_2, \theta \in [0, 1], f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \geq \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2)$$

since  $f$  is convex, so

$$\forall \mathbf{x}_1, \mathbf{x}_2, \theta \in [0, 1], f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) \leq \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2)$$

So we have

$$\forall \mathbf{x}_1, \mathbf{x}_2, \theta \in [0, 1], f(\theta \mathbf{x}_1 + (1 - \theta) \mathbf{x}_2) = \theta f(\mathbf{x}_1) + (1 - \theta) f(\mathbf{x}_2)$$

Which is exactly the definition of affine function.

So  $f$  is affine.

So above all, we have proved that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is affine if and only if  $f$  is both convex and concave.

Problem 3. Solve the optimal solution of the Rosenbrock function

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2,$$

using MATLAB programming to implement three algorithms (each 20pts): gradient descent (GD) method, Newton method, and Quasi-Newton methods (either rank-1, DFP or BFGS). You are required to print iteration information of last 10 steps: including objective, step size, residual of gradient. Technical implementation: explain how to choose the step size, how to set the termination criteria, how to choose the initial point, the value of the required parameters, converge or not and convergence rate. (paste the code in the pdf to submit it, no need to submit the source code) [60pts]

1. Gradient descent method:

set step size to be constant:  $\alpha^k = 0.002$ .

direction:  $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$ .

The update policy: following the direction of the negative gradient, i.e.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha^k \nabla f(\mathbf{x}^k)$$

termination criteria: if the current point's gradient's norm  $\|\nabla f(\mathbf{x}^k)\|_2 < 10^{-4}$ , then we regard it converges.

initial point:  $\mathbf{x}^0 = (x, y)$ , where  $x$  and  $y$  are randomly uniformed generated in the range  $[-1, 1]$ .

There is no additional parameters.

The method converged, and from the variation of output, we can see that the convergence rate is linear convergence.

2. Newton method:

step size:  $\alpha^k = 1$ .

direction:  $\mathbf{d}^k = -\nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k)$ .

The update policy: find the zero point of the first order derivative, and then update the current point to the interaction point of the current point's tangent line and the axis, i.e.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \nabla^2 f(\mathbf{x}^k)^{-1} \nabla f(\mathbf{x}^k)$$

termination criteria: if the current point's gradient's norm  $\|\nabla f(\mathbf{x}^k)\|_2 < 10^{-4}$ , then we regard it converges.

initial point:  $\mathbf{x}^0 = (x, y)$ , where  $x$  and  $y$  are randomly uniformed generated in the range  $[-1, 1]$ .

There is no additional parameters.

The method converged, and from the variation of output, we can see that the convergence rate is quadratic convergence.

3. Quasi-Newton method:

step size: we can use Armijo linear search to find a suitable  $\alpha^k$ : setting  $\alpha^k = 1$ , and  $\gamma = 0.5$ , and let  $\alpha^k$  constantly times  $\gamma$  until the Armijo condition:  $f(\mathbf{x}^k + \alpha^k \mathbf{d}^k) \leq f(\mathbf{x}^k) + \alpha^k \gamma \nabla f(\mathbf{x}^k)^T \mathbf{d}^k$  is satisfied.

direction:  $\mathbf{d}^k = -H_k^{-1} \nabla f(\mathbf{x}^k)$ , where  $H_k$  is generated by BFGS method.

The update policy: similarly to Newton-method, but do not directly get the Hessian matrix but construct a similar  $H_k$ , and then update the current point to the interaction point of the current point and last point's secant line and the axis, i.e.

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha^k H_k^{-1} \nabla f(\mathbf{x}^k)$$

termination criteria: if the current point's gradient's norm  $\|\nabla f(\mathbf{x}^k)\|_2 < 10^{-4}$ , then we regard it converges.

initial point:  $\mathbf{x}^0 = (x, y)$ , where  $x$  and  $y$  are randomly uniformed generated in the range  $[-1, 1]$ .

The original  $H_0^{-1}$  is set to be  $I_2$ .

As for the update of  $H_k^{-1}$ , we set  $\mathbf{s}_k = \mathbf{x}^{k+1} - \mathbf{x}^k$  and  $\mathbf{y}^k = \nabla f(\mathbf{x}^{k+1}) - \nabla f(\mathbf{x}^k)$ , and we use the BFGS method to update  $H_k^{-1}$ .

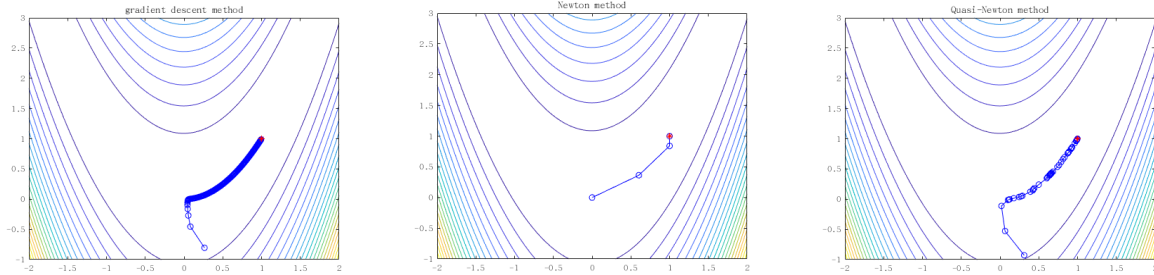
i.e. if  $\mathbf{s}_k^T \mathbf{y}_k > 0$ , then we update  $H_{k+1}^{-1} = (I_2 - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}) H_k^{-1} (I_2 - \frac{\mathbf{y}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}) + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}$ .

Otherwise it remains:  $H_{k+1}^{-1} = H_k^{-1}$ .

The method converged, and from the variation of output, we can see that the convergence rate is superlinear

convergence.

The output information and the trace of the three methods are as follows(it seems that the format of output is not aligned in latex, but it worked well in matlab file...):



output

```

1  =====
2  ===== start gradient descent method =====
3  k      x_k      y_k      f(x_k, y_k) step size  ||d||
4  =====
5  9873 9.9989e-01 9.9977e-01 1.2709e-08 2.0000e-03 1.0068e-04
6  9874 9.9989e-01 9.9977e-01 1.2689e-08 2.0000e-03 1.0060e-04
7  9875 9.9989e-01 9.9977e-01 1.2668e-08 2.0000e-03 1.0052e-04
8  9876 9.9989e-01 9.9977e-01 1.2648e-08 2.0000e-03 1.0044e-04
9  9877 9.9989e-01 9.9978e-01 1.2628e-08 2.0000e-03 1.0036e-04
10 9878 9.9989e-01 9.9978e-01 1.2608e-08 2.0000e-03 1.0028e-04
11 9879 9.9989e-01 9.9978e-01 1.2588e-08 2.0000e-03 1.0020e-04
12 9880 9.9989e-01 9.9978e-01 1.2567e-08 2.0000e-03 1.0012e-04
13 9881 9.9989e-01 9.9978e-01 1.2547e-08 2.0000e-03 1.0004e-04
14 9882 9.9989e-01 9.9978e-01 1.2527e-08 2.0000e-03 9.9958e-05
15 9883 9.9989e-01 9.9978e-01 1.2507e-08 -----
16 ===== finished gradient descent method =====
17 =====
18 =====
19 ===== start Newton method =====
20 k      x_k      y_k      f(x_k, y_k) step size  ||d||
21 =====
22 1 0.0000e+00 0.0000e+00 0.0000e+00 1.0000e+00 3.3601e+02
23 2 6.0241e-01 3.6289e-01 1.1613e+02 1.0000e+00 7.9437e-01
24 3 9.9973e-01 8.4159e-01 1.5808e-01 1.0000e+00 7.0584e+01
25 4 9.9974e-01 9.9948e-01 2.4921e+00 1.0000e+00 5.2493e-04
26 5 1.0000e+00 1.0000e+00 6.8896e-08 1.0000e+00 3.0811e-05
27 6 1.0000e+00 1.0000e+00 4.7467e-13 -----
28 ===== finished Newton method =====
29 =====
30 =====
31 ===== start Quasi-Newton method =====
32 k      x_k      y_k      f(x_k, y_k) step size  ||d||
33 =====
34 40 9.8100e-01 9.6149e-01 1.7181e-03 1.0000e+00 3.5169e-01
35 41 9.8457e-01 9.6871e-01 4.3737e-04 1.0000e+00 2.6875e-01
36 42 9.8695e-01 9.7348e-01 2.8300e-04 1.0000e+00 2.4288e-01
37 43 9.9422e-01 9.8801e-01 2.0630e-04 2.5000e-01 1.9886e-01
38 44 9.9809e-01 9.9578e-01 5.5454e-05 1.0000e+00 1.7670e-01
39 45 9.9890e-01 9.9761e-01 1.9908e-05 2.5000e-01 8.3814e-02
40 46 9.9897e-01 9.9790e-01 4.8890e-06 1.0000e+00 2.1574e-02
41 47 9.9914e-01 9.9826e-01 1.3252e-06 1.0000e+00 7.9361e-03
42 48 9.9954e-01 9.9908e-01 7.8646e-07 5.0000e-01 3.1274e-03
43 49 9.9990e-01 9.9995e-01 2.1516e-07 5.0000e-01 3.7022e-04
44 50 9.9999e-01 9.9998e-01 6.2918e-10 -----
45 ===== finished Quasi-Newton method =====
46 =====

```

And the codes are as follows:

#### Rosenbrock.m

```

1 % f(x, y) = (1-x)^2+100(y-x^2)^2
2 function [val, nabla, Hessian] = Rosenbrock(x,y)
3     val = (1 - x) ^ 2 + 100 * (y - x ^ 2) ^ 2;
4     nabla = [-2 * (1 - x) - 400 * x * (y - x ^ 2); 200 * (y - x ^ 2)];
5     Hessian = [2 - 400 * y + 1200 * x ^ 2, -400 * x; -400 * x, 200];
6 end

```

#### print\_info.m

```

1 function print_info(max_iter, obj, points, step_size, grad_res)
2     upper = max_iter - 1;
3     for i = max(1, max_iter - 10):upper
4         fprintf('%d %4e %4e %4e %4e %4e\n', i, points(i, 1), points(i, 2), obj(i), step_size(i), grad_res(i));
5     end
6     fprintf('%d %4e %4e %4e ---- ---- \n', max_iter, points(max_iter, 1), points(max_iter, 2),
7         obj(max_iter));
8 end

```

#### plot\_trace.m

```

1 function plot_trace(max_iter, points, id, name)
2     points = points(1 : max_iter - 1, :);
3
4     x = linspace(-2, 2, 100);
5     y = linspace(-1, 3, 100);
6     [X, Y] = meshgrid(x, y);
7     Z = 100 * (Y - X.^2).^2 + (1 - X).^2;
8     figure(id);
9     contour(X, Y, Z, 20);
10    hold on;
11
12    % plot the trace of the points
13    plot(points(:, 1), points(:, 2), 'b-o');
14
15    % mark the last point with a red star
16    plot(points(end, 1), points(end, 2), 'r*');
17
18    % set the axis
19    axis([-2, 2, -1, 3]);
20
21    % set name
22    title(name);
23
24 end

```

#### gradient\_descent.m

```

1 function gradient_descent()
2     disp('
3         =====
4         ')
5     disp('===== start gradient descent method =====')
6     disp(' k      x_k      y_k      f(x_k, y_k) step size ||d|| ')
7     disp('=====')
8
9     % (x, y) are random initial points
10    x = rand() * 2 - 1;
11    y = rand() * 2 - 1;

```

```

10
11 grad_res = zeros(1, 10000);
12 obj = zeros(1, 10000);
13 step_size = zeros(1, 10000);
14 points = zeros(10000, 2);
15
16 [val, grad, Hessian] = Rosenbrock(x, y);
17 obj(1) = val;
18 grad_res(1) = norm(grad);
19 points(1, :) = [x, y];
20
21 for iter = 2:10000
22     % update the value of x and y
23     [val, grad, Hessian] = Rosenbrock(x, y);
24     % [x, y] = [x, y] - 0.002 * nabla;
25
26     x = x - 0.002 * grad(1);
27     y = y - 0.002 * grad(2);
28
29     points(iter, :) = [x, y];
30     obj(iter) = val;
31     step_size(iter - 1) = 0.002;
32     grad_res(iter - 1) = norm(grad);
33
34     if (norm(grad) < 1e-4)
35         break;
36     end
37 end
38
39 print_info(iter, obj, points, step_size, grad_res);
40 plot_trace(iter, points, 1, 'gradient descent method');
41
42 disp('===== finished gradient descent method =====')
43 disp('
=====
')
44 end

```

### Newton\_method.m

```

1 function Newton_method()
2     disp('
=====
')
3     disp('===== start Newton method =====')
4     disp('k    x_k        y_k    f(x_k, y_k) step size    ||d||    ')
5     disp('=====
=====
=====')
6
7
8     % (x, y) are random initial points
9     x = rand() * 2 - 1;
10    y = rand() * 2 - 1;
11
12    grad_res = zeros(1, 100);
13    obj = zeros(1, 100);
14    step_size = zeros(1, 100);
15    points = zeros(100, 2);
16
17    for iter = 2:100
18        [val, grad, Hessian] = Rosenbrock(x, y);
19
20        % if Hessian is not invertible, then output error
21        if (abs(det(Hessian)) < 1e-7)
22            disp(' Error : Hessian matrix is invertible ');
23            break;
24        end
25

```

```

26     d = -inv(Hessian) * grad;
27     x = x + d(1);
28     y = y + d(2);
29
30     points(iter, :) = [x, y];
31     obj(iter) = val;
32     step_size(iter - 1) = 1;
33     grad_res(iter - 1) = norm(grad);
34
35     if (norm(d) < 1e-4)
36         break;
37     end
38 end
39
40
41 print_info(iter, obj, points, step_size, grad_res);
42 plot_trace(iter, points, 2, 'Newton method');
43
44 disp('===== finished Newton method =====')
45 disp('
=====
')
46 end

```

#### Quasi\_Newton\_method.m

```

1 function Quasi_Newton_method()
2     disp('
=====
')
3     disp('===== start Quasi-Newton method =====')
4     disp(' k    x_k    y_k    f(x_k, y_k) step size ||d|| ')
5     disp('=====')
6
7     % (x, y) are random initial points
8     x = rand() * 2 - 1;
9     y = rand() * 2 - 1;
10
11     grad_res = zeros(100, 1);
12     grad_record = zeros(100, 2);
13     obj = zeros(100, 1);
14     step_size = zeros(100, 1);
15     points = zeros(100, 2);
16
17     [val, grad, hess] = Rosenbrock(x, y);
18     obj(1) = val;
19     grad_res(1) = norm(grad);
20     grad_record(1, :) = [0, 0];
21     points(1, :) = [x, y];
22
23     H_inverse = eye(2);
24     for iter = 2:100
25         [val, grad, hess] = Rosenbrock(x, y);
26
27         d = -H_inverse * grad;
28
29         % use Armijo rule to determine the step size
30         alpha = 1;
31         gamma = 0.5;
32         while (Rosenbrock(x + alpha * d(1), y + alpha * d(2)) > val + gamma * alpha * grad' * d)
33             alpha = alpha * gamma;
34         end
35
36         x = x + alpha * d(1);
37         y = y + alpha * d(2);
38
39         s_k = [x, y] - points(iter - 1, :);

```



```

40     s_k = s_k';
41
42     y_k = (grad' - grad_record(iter - 1, :));
43
44     % the condition of update H_k:
45     if s_k' * y_k > 0
46         H_inverse = (eye(2) - s_k * y_k' / (y_k' * s_k)) * H_inverse * (eye(2) - y_k * s_k' / (y_k'
47             * s_k)) + s_k * s_k' / (y_k' * s_k);
48     end
49
50     points(iter, :) = [x, y];
51     obj(iter) = val;
52     step_size(iter - 1) = alpha;
53     grad_res(iter - 1) = norm(grad);
54
55     grad_record(iter, :) = grad;
56
57     if (norm(d) < 1e-4)
58         break;
59     end
60
61     print_info(iter, obj, points, step_size, grad_res);
62     plot_trace(iter, points, 3, 'Quasi-Newton method');
63
64     disp('===== finished Quasi-Newton method =====')
65     disp('
66         ')
67 end

```

implement.m

```

1 gradient_descent()
2 Newton_method()
3 Quasi_Newton_method()

```