

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

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 the convergence rate is **TODO**

2. Newton method:

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

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

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 the convergence rate is **TODO**

3. Quasi-Newton method:

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

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

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 is set to be I_2 .

As for the update of H_k , 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 .

i.e. $H_k =$

$H_k^{-1} =$

The method converged, and the convergence rate is **TODO**

And the codes and corresponding output are as follows:

Rosenbrock.m

```
1 % f(x, y) = (1-x)^2+100(y-x^2)^2
2 function [val, nabla] = 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 end
```

print_info.m

```
1 function print_info(obj, step_size, grad_res)
2     for i = max(1, length(obj) - 10):length(obj)
3         fprintf('Step %d: obj = %f, step size = %f, gradient residual = %f\n', i, obj(i), step_size(i),
4             grad_res(i));
5     end
6 end
```

plot_trace.m

```

1 function plot_trace(points, id, name)
2     % plot the contour of of Rosenbrock function
3     % then plot the trace of the points
4     % points is a matrix of size (n, 2)
5     % each row is a point
6     % the first column is x, the second column is y
7     % the last point is marked with a red star
8
9     % plot the contour of Rosenbrock function
10    x = linspace(-2, 2, 100);
11    y = linspace(-1, 3, 100);
12    [X, Y] = meshgrid(x, y);
13    Z = 100 * (Y - X.^2).^2 + (1 - X).^2;
14    figure(id);
15    contour(X, Y, Z, 20);
16    hold on;
17
18    % plot the trace of the points
19    plot(points(:, 1), points(:, 2), 'b-o');
20
21    % mark the last point with a red star
22    plot(points(end, 1), points(end, 2), 'r*');
23
24    % set the axis
25    axis([-2, 2, -1, 3]);
26
27    % set name
28    title(name);
29
30
31 end

```

gradient_descent.m

```

1 function gradient_descent(init_point)
2     disp('
3         =====
4         ')
5     disp('===== start gradient descent method =====')
6
7     % [x, y] = init_point;
8     x = init_point(1);
9     y = init_point(2);
10
11    points = [[x, y]];
12    obj = [];
13    step_size = [];
14    grad_res = [];
15    while (1)
16        % update the value of x and y
17        [val, nabla] = Rosenbrock(x, y);
18        % [x, y] = [x, y] - 0.001 * nabla;
19        x = x - 0.001 * nabla(1);
20        y = y - 0.001 * nabla(2);
21
22        points = [points; [x, y]];
23        obj = [obj, val];
24        step_size = [step_size, 0.001];
25        grad_res = [grad_res, norm(nabla)];
26
27        [newval, new_nabla] = Rosenbrock(x, y);
28        if (abs(newval - val) < 1e-15)
29            break;
30        end
31    end
32 end

```

```

30 print_info(obj, step_size, grad_res);
31 plot_trace(points, 1, 'gradient descent method');
32
33 disp('===== finished gradient descent method =====')
34 disp('
35     =====
36 end

```

Newton_method.m

```

1 function Newton_method(init_point)
2     disp('
3         =====
4     ')
5     disp('===== start Newton method =====')
6
7     % [x, y] = init_point;
8     x = init_point(1);
9     y = init_point(2);
10
11     points = [[x, y]];
12     obj = [];
13     step_size = [];
14     grad_res = [];
15
16     print_info(obj, step_size, grad_res);
17     plot_trace(points, 2, 'Newton method');
18
19     disp('===== finished Newton method =====')
20     disp('
21         =====
22 end

```

Quasi_Newton_method.m

```

1 function Quasi_Newton_method(init_point)
2     disp('
3         =====
4     ')
5     disp('===== start Quasi-Newton method =====')
6
7     % [x, y] = init_point;
8     x = init_point(1);
9     y = init_point(2);
10
11     points = [[x, y]];
12     obj = [];
13     step_size = [];
14     grad_res = [];
15
16     print_info(obj, step_size, grad_res);
17     plot_trace(points, 3, 'Quasi-Newton method');
18
19     disp('===== finished Quasi-Newton method =====')
20     disp('
21         =====
22 end

```

output

1 asfafsafafaas