

Numerical Optimization

Re-exam Handin 2

Dmitry Serykh (qwl888)

June 14, 2020

1 The Setup

For this assignment, I used Python, and all the relevant functions from `numpy` and `scipy` packages. In particular, I have utilized the `minimize` function from the `scipy.optimize` library. I have chosen the BFGS optimizer, since it is the default choice for the unconstrained optimization in this library. BFGS stands for Broyden–Fletcher–Goldfarb–Shanno algorithm and belongs to the family of quasi-Newton optimizers that is a state-of-the-art alternative to Newton methods, thus it does not require the Hessian matrix. Moreover, it can be classified as a Line-Search algorithm, where we first determine the direction p_k and the step length in that direction α afterwards.

1.1 Parameters

The `scipy` implementation of the BFGS algorithm has following parameters:

- *maxiter* - Maximum number of iterations to perform.
- *gtol* - Gradient norm must be less than *gtol* before successful termination.
- *norm* - Order of norm
- *eps* - Step size used when the gradient is approximated.

I used the value of 1000 for the *maxiter* parameter, but it was never used, since the gradient magnitude was used as a primary stopping criterion. For *gtol*, I used the default value of 10^{-5} , and ∞ for the *norm*. *eps* was not used, since I have implemented a gradient for all five functions and we do not need to approximate it. For all the functions, I have used a dimensionality of two.

2 Testing

2.1 What constitutes a good testing protocol?

According to Chapter 1 in the book, good optimization algorithms must have following properties:

- **Robustness.** Ability of the algorithm to find minima of vast array of different functions. Furthermore, the algorithm should be indifferent to the starting point x_0 .
- **Efficiency.** A good optimization algorithm should not consume too much computational resources (and memory/storage).
- **Accuracy.** Good optimization algorithms should find the minima accurately to some tolerance ε .

In order to process and then visualize the gathered data from the experiments, I must use one of the existing aggregation methods. These include: minimum, maximum, average(mean), median, variance, standard deviation. I decided both to use the median and the mean. The main difference between these methods is how they manage the influence of the outliers. When using the mean, we maintain the influence of the outliers in the set, while it is excluded when using the median.

2.2 My testing protocol

In order to test the effectiveness of the BFGS implementation from the `scipy.optimize`, I came up with a testing protocol:

	f_1	f_2	f_3	f_4	f_5
efficiency	7	68	16	14	12
accuracy	$1.877e - 8$	$3.864e - 7$	$1.426e - 13$	$6.512e - 7$	$8.72e - 5$

Table 1: Median number of iterations until algorithm termination ($\|\nabla f\| < 1e - 5$) and final distance to the optimum for 100 random starting points in the interval $[-10, 10]$

- In order to measure the robustness of the algorithm, I have used the algorithm to find the minima of all five problem functions: The Ellipsoid function(f_1), the Rosenbrock Banana function (f_2), The Log-Ellipsoid function (f_3) and the Attractive Sector functions (f_4 and f_5). To test if the algorithm is indifferent to the starting point x_0 , I have used a random starting point taken from the uniform distribution in the interval between -10 to 10 for both dimensions, repeated all the experiments 100 times and took the median of the results.
- I decided to measure the efficiency of the algorithms in number of iterations until the gradient magnitude reaches 10^{-5} . It should however be noted that steps of some algorithms are more computationally expensive than others. For example the Newton step requires a computation of Hessian and therefore more expensive than a steepest-descent step. The results can be seen on Table 1.
- The accuracy of the algorithm was tested by measuring the Euclidean distance between the current point and the optima. The results can be seen on Table 1.
- The convergence plots with the number of iteration on the x-scale and the Euclidean distance between the current value of x and the optimum. The results can be seen on Figure 1.
- The convergence plots with the gradient magnitude can be seen on Figure 2.

3 Theoretical Part

3.1

3.1.1

Since $A \in \mathbb{R}^{d \times d}$ is invertible, $g(x) = Ax$ would constitute a linear map $\mathbb{R}^n \rightarrow \mathbb{R}^n$. Linear transformations preserve the ratios of lengths, therefore the inequality

$$f(\alpha(Ax) + (1 - \alpha)(Ay)) \leq \alpha f(Ax) + (1 - \alpha)f(Ay), \text{ for all } \alpha \in [0, 1]$$

would still hold, since we know that f is strictly convex and therefore:

$$f(\alpha(x) + (1 - \alpha)(y)) \leq \alpha f(x) + (1 - \alpha)f(y), \text{ for all } \alpha \in [0, 1]$$

3.1.2

We start by looking at an expression $x^T Q x$. We know that Q is symmetric positive definite, hence we can use Cholesky decomposition and obtain following: $x^T Q x = x^T L L^T x$, where L is triangular matrix with positive diagonal values, hence invertible.

Let $f(y) = y^T y$, then $\nabla f(y) = 2y$ and $\nabla^2 f(y) = 2I$. $2I$ is PD, therefore $f(y)$ is strictly convex. Then, if we set $y = L^T x$, we get:

$$f(y) = f(Lx) = x^T L L^T x = x^T Q x$$

$x^T Q x$ is therefore strictly convex, since $f(x)$ is convex and property from the previous subsection holds.

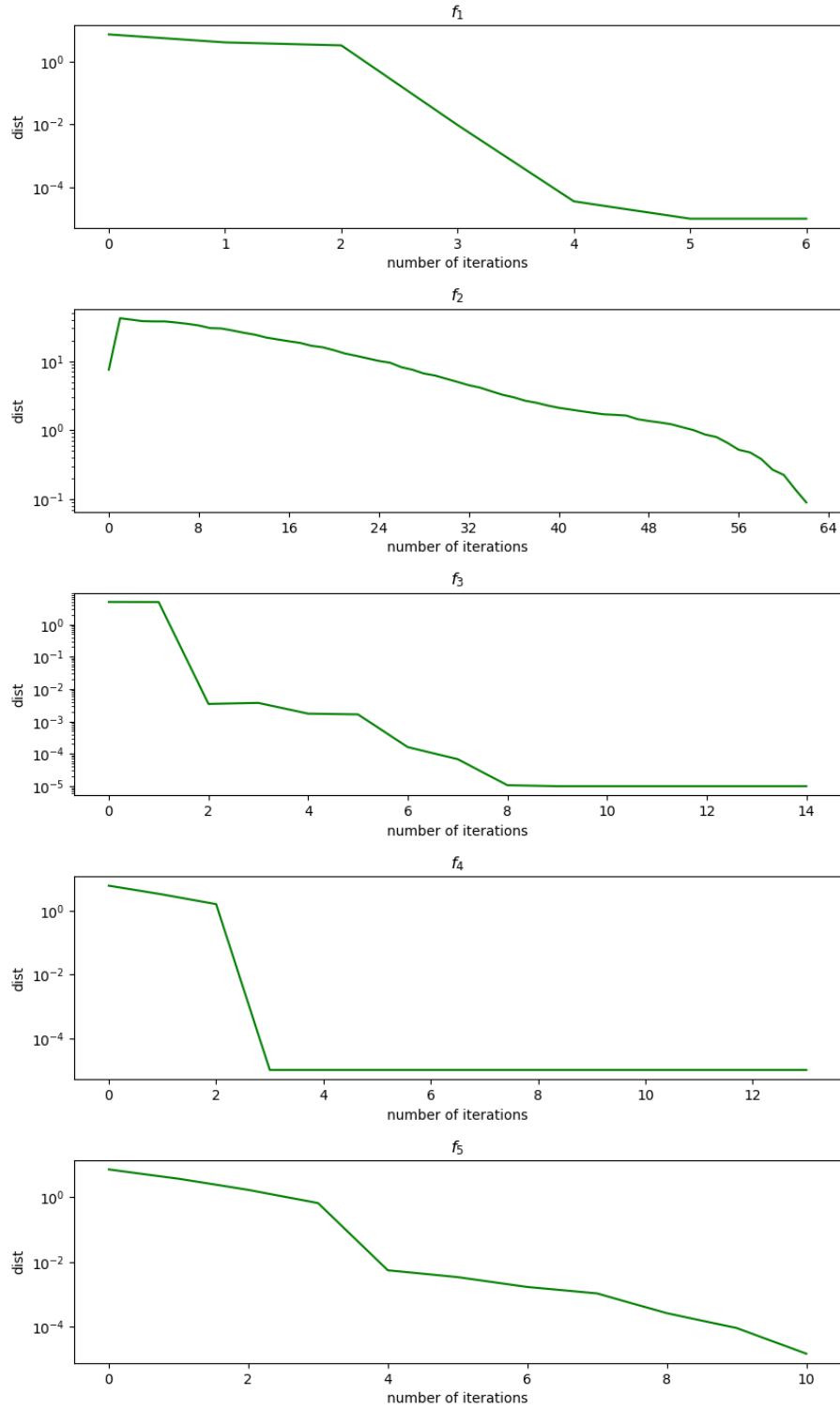


Figure 1: Convergence Plot for BFGS, median of 100 runs. Y-axis shows the Euclidean distance to the optimum of each f_1, \dots, f_7 with random starting point

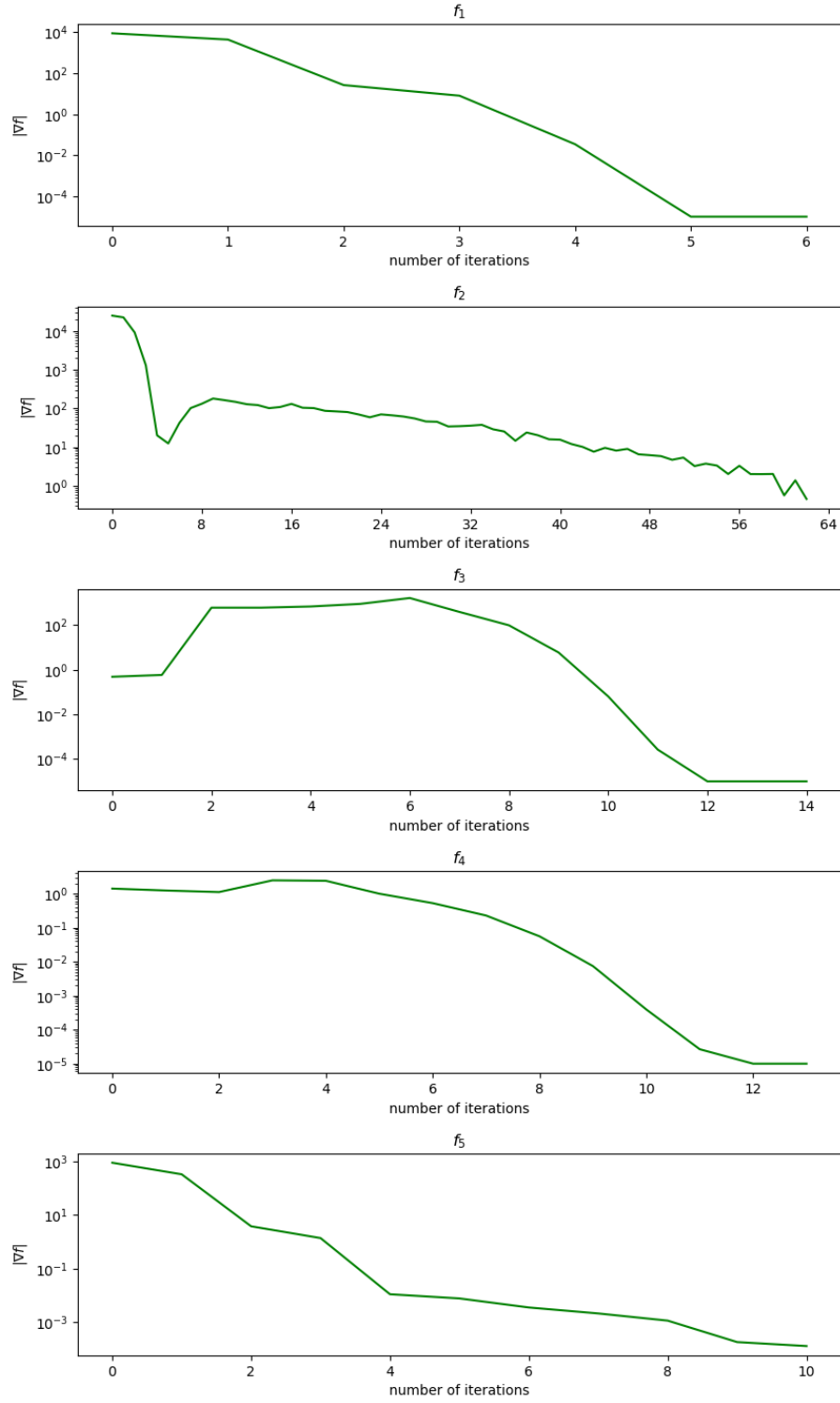


Figure 2: Convergence Plots for BFGS, median over 100 runs. Y-axis shows the norm of the gradient of each f_1, \dots, f_7 with random starting point

This point is relevant for quadratic-based functions such as f_1 and f_2 , where the convexity makes the optimizer converge on the optimum.

3.2

3.2.1

For $d = 1$, the Log-Ellipsoid function could be formulated as:

$$f_3(x) = \log(\epsilon + x^2)$$

Hence, the derivatives:

$$\begin{aligned} f_3'(x) &= \frac{2x}{\epsilon + x^2} \\ f_3''(x) &= \frac{2(\epsilon - x^2)}{(\epsilon + x^2)^2} \end{aligned}$$

For $|x| < \epsilon$, $(\epsilon - x^2) > 0$, hence $f_3''(x) > 0$, and f_3 would be strictly convex.

$$(-f_3)''(x) = \frac{2(x^2 - \epsilon)}{(\epsilon + x^2)^2}$$

Hence, for $|x| \geq \epsilon$, $-f_3''(x) \geq 0$, which means that $-f_3$ is convex, hence f_3 is concave for $|x| \geq \epsilon$.

This could potentially be problematic if the gradient-based optimizer gets its current value of $x > |\epsilon|$, then it would fail finding an optimum, because it could continue moving in that direction for an indefinite amount of time or until the number of iterations reaches the predefined limit.

4 Conclusion

Based on the results of application of the two performance metrics, the `scipy` implementation of the BFGS optimizer managed to find the minima of all case functions. It struggled most with the Rosenbrock function, due to its elongated and curved iso-level.