

# Numerical Optimization for Large Scale Problems

## Assignment Report

Giacomo Maino s338682, Mattia Molinari s337194, Alessandro Perlo s337131

February 21<sup>st</sup>, 2025

### Abstract

In this report, we explore optimization techniques for large-scale unconstrained problems, focusing on variations of Newton method. Specifically, we implement and analyze the Modified Newton Method and the Truncated Newton Method, comparing their performance on several test functions. Both exact and finite difference-based Hessian and gradient computations are considered. Our experiments evaluate convergence rates, computational efficiency, and the impact of preconditioning. Additionally, we discuss the challenges posed by finite difference approximations and the sensitivity of each method to different problem structures.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Modified Newton Method	2
1.2	Truncated Newton method	2
1.3	Finite differences	3
<b>2</b>	<b>Rosenbrock function</b>	<b>3</b>
<b>3</b>	<b>Extended Rosenbrock function</b>	<b>4</b>
3.1	Exact gradient and Hessian	4
3.2	Finite differences gradient and Hessian	5
<b>4</b>	<b>Generalized Broyden tridiagonal function</b>	<b>8</b>
4.1	Exact gradient and Hessian	8
4.2	Finite differences gradient and Hessian	10
<b>5</b>	<b>Banded trigonometric function</b>	<b>13</b>
5.1	Exact gradient and Hessian	14
5.2	Finite differences gradient and Hessian	15
<b>6</b>	<b>Conclusions</b>	<b>18</b>
<b>A</b>	<b>Code snippets</b>	<b>18</b>
A.1	Code for method implementations	18
A.1.1	Modified Newton method	18
A.1.2	Truncated Newton method	21
A.2	Code for objective functions	26
A.2.1	Rosenbrock function	26
A.2.2	Extended Rosenbrock function	27
A.2.3	Generalized Broyden tridiagonal function	28
A.2.4	Banded trigonometric problem	30
A.3	Utility code for running experiments	31

## 1 Introduction

In this section we will describe the implementation details of the algorithms used to solve the optimization problems, namely modified Newton method and truncated Newton method, focusing on the differences with respect to the standard Newton method. These methods will be tested against the Rosenbrock function and

three test problems from [1]. The chosen test problems are the extended Rosenbrock function (problem 25), the generalized Broyden tridiagonal function (problem 32) and the banded trigonometric function (problem 16) and results are contained in sections 3, 4 and 5 respectively.

The experiments were conducted using 11 points: a predefined starting point and 10 additional randomly generated points uniformly distributed in a hypercube around the initial guess. For each test function, we performed optimizations at problem dimensions  $n = 10^3, 10^4, 10^5$ . We implemented backtracking line search with a sufficient decrease condition, using standard parameters  $\rho = 0.5$  and  $c = 10^{-4}$ , but further tuning has been necessary in the case of the banded trigonometric function. Each method was evaluated in terms of success rate, number of iterations to convergence, execution time, and experimental convergence rate. The experimental convergence rate was computed using the formula:

$$q = \frac{\log(\|e_{k+1}\|/\|e_k\|)}{\log(\|e_k\|/\|e_{k-1}\|)} \quad (1)$$

where  $e_k$  denotes the error at iteration  $k$ . Error at iteration  $k$  is approximated as the norm of the difference between the point at current iteration and the point at previous iteration, i.e.  $\hat{e}_k = x_k - x_{k-1}$ . For each experiment, we report average metrics over the successful runs, where a run is considered successful if the method converges (gradient norm is less than  $tol = 10^{-6}$ ) within a maximum of  $k_{max} = 10^3$  iterations and Armijo condition is satisfied within  $bt_{max} = 50$  backtracking attempts.

## 1.1 Modified Newton Method

The modified Newton method aims to enhance robustness of the standard Newton method by ensuring positive-definiteness of the Hessian matrix. At iteration  $k$ , it is necessary to check whether the Hessian matrix  $H_k$  is positive definite: in case it is not, the matrix is modified by adding a matrix  $B_k$  in order to ensure positive definiteness. A common choice for  $B_k$  is a multiple of the identity matrix, i.e.  $B_k = \tau_k I$  so that the whole spectrum of  $H_k$  is shifted by  $\tau_k$ . Then we want to find the smallest  $\tau_k$  such that  $H_k + \tau_k I$  is positive definite, which is  $-\lambda_{k,min} + \beta$  where  $\lambda_{k,min}$  is the negative eigenvalue with the largest module.

To avoid to have to compute  $\lambda_{k,min}$ , we adopted the *Cholesky with Added Multiple of the Identity* algorithm outlined in [2] that consists in building a sequence of  $\tau_k$  until the modified matrix is positive definite. The sequence is built starting from  $\tau_k = \min_i h_{ii} + \beta$  where  $\min_i h_{ii}$  is the smallest diagonal element of  $H_k$ . Then, at each iteration:

1. positive-definiteness is assessed trying to perform a Cholesky factorization of  $H_k + \tau_k I$ ;
2. if the factorization is not successful,  $\tau_k$  is increased by a factor  $c$  and the process is repeated for a limited number of times  $k_{chol,max}$ .

In all the experiments we choose  $\beta = 10^{-3}$ ,  $k_{chol,max} = 100$ . A good value for the constant factor is  $c = 2$ , but as we will discuss in section 3 for the extended Rosenbrock function a larger value  $c = 5$  is beneficial. The method is endowed with a line search strategy with backtracking. We carry on experiments both with and without preconditioning, using the incomplete Cholesky factorization as preconditioner.

## 1.2 Truncated Newton method

The truncated Newton method aims to reduce the computational cost of the Newton method by adopting the following strategies:

- the newton system  $H_k p_k = -\nabla f(x_k)$  is solved approximately by means of an iterative method (i.e. conjugate gradient method), with a tolerance that depends on  $\|\nabla f(x_k)\|$ ;
- whenever a direction of negative curvature is found in the execution of the iterative method, the method is stopped and the direction is used as the search direction to prevent a non-negative curvature direction to be chosen in case of a non-positive definite  $H_k$ .

In all the experiments, we choose the relative tolerance for the iterative method at iteration  $k$  to be

$$\eta_k = \min\{0.5, \sqrt{\|\nabla f(x_k)\|}\}$$

that is a forcing term that is proven to yield a superlinear convergence rate. The method is endowed with a line search strategy with backtracking. We carry on experiments both with and without preconditioning, using the incomplete Cholesky factorization as preconditioner for the Newton system whenever the Hessian matrix is positive definite.

Starting Point	Algorithm	Iterations	$\ \nabla f(x_k)\ $	$f(x_k)$	Time (s)	$x_k$
$x_0^{(1)} = (1.2, 1.2)$	Modified Newton	8	1.436e-11	1.0883e-25	0.009753	(1.0000, 1.0000)
	Truncated Newton	9	1.0471e-07	5.5531e-18	0.010447	(0.9999, 0.9999)
$x_0^{(2)} = (-1.2, 1)$	Modified Newton	21	4.4733e-10	3.744e-21	0.001906	(0.9999, 0.9999)
	Truncated Newton	64	9.1038e-15	2.072e-29	0.00128	(0.9999, 0.9999)

Table 1: Results of optimization algorithms on the Rosenbrock function with different starting points.

### 1.3 Finite differences

Experiments in subsequent sections will adopt both exact and finite differences gradient and Hessian to perform the optimization. When finite differences are adopted, the gradient will be estimated using centered finite differences

$$\frac{\partial f}{\partial x_k} \approx \frac{f(x + he_k) - f(x - he_k)}{2h} \quad (2)$$

while the Hessian will be estimated using forward finite differences, using the following formula

$$\frac{\partial^2 f}{\partial x_k \partial x_j} \approx \frac{f(x + he_k + he_j) - f(x + he_k) - f(x + he_j) + f(x)}{h^2} \quad (3)$$

where  $e_k$  and  $e_j$  are the  $k$ -th and the  $j$ -th canonical basis vectors respectively. Moreover, two different approaches will be adopted to choose the step size  $h$ : the first one will use a fixed step size while the second one will use a step size that depends on the current point  $x$  and that is different for each component, defined as follows

$$h_{k,i} = h|x_{k,i}|$$

where  $h_{k,i}$  is the increment for component  $i$  at step  $k$ ,  $h$  is a relative step size and  $x_{k,i}$  is the  $i$ -th component of the point at step  $k$ . Due to the large scale nature of the problems, the finite differences method is expected to be slower than the exact method, so ad-hoc implementations that will exploit the sparsity of the Hessian matrix and the separability of the specific functions will be used.

## 2 Rosenbrock function

We first test our implementation of the optimization algorithms on the Rosenbrock function, defined as follows.

$$f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2 \quad (4)$$

The Rosenbrock function is a non-convex function that is commonly used to test optimization algorithms. It has a global minimum at  $x^* = (1, 1)$ , where  $f(x^*) = 0$ . The function is characterized by a narrow, curved valley, which makes it difficult for optimization algorithms to converge to the global minimum. We test the algorithms on the Rosenbrock function with two different starting points:  $x_0^{(1)} = (1.2, 1.2)$  and  $x_0^{(2)} = (-1.2, 1)$ . The results are shown in Figure 1. Both methods converge to the global minimum in both cases, but the Modified Newton method converges faster than the Truncated Newton method. The results are summarized in Table 1.

- For the starting point  $x_0^{(1)} = (1.2, 1.2)$ , the Modified Newton method converges in a number of iterations that is comparable to the Truncated Newton method: this is expected since the starting point is close to the minimum, so the Newton system is solved accurately.
- For the starting point  $x_0^{(2)} = (-1.2, 1)$ , the Modified Newton method converges in fewer iterations than the Truncated Newton method: this is expected since the starting point is far from the minimum, so the Newton system is solved approximately in the first steps: this lead to a very slow progress in the first iterations of the Truncated Newton method.

In both cases, despite the fact that the Truncated Newton method requires more iterations, its execution time is comparable or even lower than Modified Newton method. This is due to the fact that the Truncated Newton method requires less computational effort per iteration.

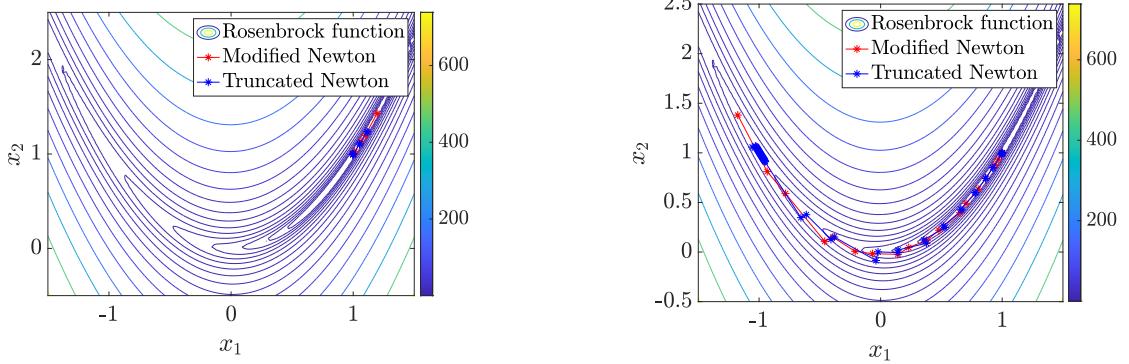


Figure 1: Convergence of the algorithms on the Rosenbrock function with different starting points and different optimization algorithms.

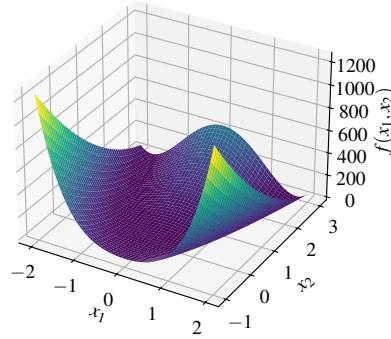


Figure 2: Surface plot of the 2-dimensional extended Rosenbrock function

### 3 Extended Rosenbrock function

The extended Rosenbrock function is a generalization of the Rosenbrock function to  $n$  dimensions, defined as follows. Figure 2 shows the surface plot of the 2-dimensional extended Rosenbrock function: notice that for  $n = 2$  it is identical to the standard Rosenbrock function, except for the  $\frac{1}{2}$  term.

$$F(x) = \frac{1}{2} \sum_{k=1}^n f_k^2(x), \quad f_k(x) = \begin{cases} 10(x_k^2 - x_{k+1}), & k \bmod 2 = 1 \\ x_{k-1} - 1, & k \bmod 2 = 0 \end{cases} \quad (5)$$

The minimum of the function is in a very flat valley which is easy to reach, but in practice it's harder to converge to a minimum, which makes the extended Rosenbrock function a challenging optimization problem. For convergence to happen for some points for  $n = 10^4$ , it has been necessary to adopt a higher constant factor  $c = 5$  for the modification of the Hessian in the Modified Newton method, as the default value  $c = 2$  was not enough to ensure convergence.

#### 3.1 Exact gradient and Hessian

The gradient of the extended Rosenbrock function is given by the following expression,

$$\frac{\partial F}{\partial x_k} = \begin{cases} 200(x_k^3 - x_k x_{k+1}) + (x_k - 1), & k \bmod 2 = 1 \\ -100(x_{k-1}^2 - x_k), & k \bmod 2 = 0 \end{cases} \quad (6)$$

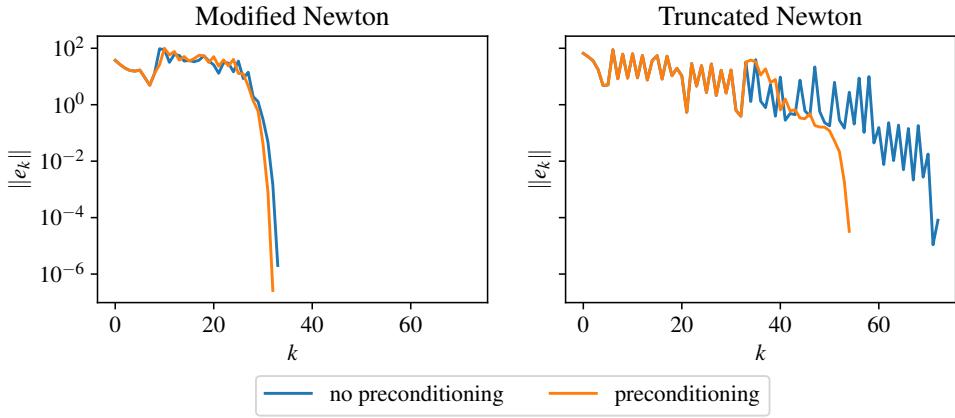


Figure 3: Estimate of the error for Modified Newton method and for the Truncated Newton method applied to the Extended Rosenbrock function with exact gradient and Hessian for  $n = 10^5$ , random point 1.

computation can be eased considering that component  $k$  depends only on  $f_k$  and  $f_{k+1}$  when  $k$  is odd, and only on  $f_{k-1}$  when  $k$  is even. The Hessian of the extended Rosenbrock function is given by the following expression.

$$\frac{\partial^2 F}{\partial x_k \partial x_j} = \begin{cases} 200(3x_k^2 - x_{k+1}) + 1, & j = k, k \bmod 2 = 1 \\ 100, & j = k, k \bmod 2 = 0 \\ -200x_k, & |k - j| = 1, k \bmod 2 = 1 \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

Notice that the Hessian is a sparse matrix, with only  $n$  non-zero elements on the diagonal and  $n/2$  non-zero elements on the first co-diagonal.

Table 2 shows the results for the *Modified Newton method* applied to the extended Rosenbrock function with exact gradient and Hessian. All attempts were successful, and the method converged in a small number of iterations, with a convergence rate that is close to 2 for all dimensions but  $10^5$ , where the convergence rate is smaller and close to 1. However, the time required to converge is significantly higher for the  $10^5$ -dimensional problem, which is expected due to the increased number of function evaluations required to compute the gradient and Hessian. We observe that being the number of iterations necessary to converge to the minimum of the function very low, the experimental convergence rate is not very reliable, as it is computed as the ratio between the number of iterations and the logarithm of the relative error.

Table 3 shows the results for the *Truncated Newton method* applied to the extended Rosenbrock function with exact gradient and Hessian. All attempts were successful, and the method converged in a small number of iterations. The computed experimental convergence rate is not reliable at all, presumably due to the small number of iterations and truncations of the iterative solver used to solve the Newton system. Moreover, when preconditioning is not adopted we get a negative convergence rate since  $\log\|e_k\|$  where  $e_k$  is the error at iteration  $k$  is not monotonic with respect to  $k$ , as shown in figure 3. However, when preconditioning is adopted and  $k$  is larger, i.e. for  $n = 10^5$ , the convergence rate is superlinear as expected from the truncated Newton method with a superlinear forcing term.

Figure 3 shows the estimate of the error for the Modified Newton method and for the Truncated Newton method applied to the Extended Rosenbrock function with exact gradient and Hessian for  $n = 10^5$  when starting from random point 1:

- for the *Modified Newton method*, in early iterates the estimated error is not monotonic, but it becomes monotonic after a few iterations regardless of preconditioning;
- for the *Truncated Newton method*, in early iterates the estimated error is not monotonic, but it becomes monotonic after a few iterations, but only if preconditioning is adopted.

Both in the case of the Modified Newton method and the Truncated Newton method, preconditioning improves performance of the optimization algorithms both in terms of number of iterations and time required to converge.

### 3.2 Finite differences gradient and Hessian

When applying 2, one can notice that the terms  $F(x + he_k)$  and  $F(x - he_k)$  only differ by terms  $f_k$  and  $f_{k+1}$  for  $k$  odd, by terms  $f_{k-1}$  for  $k$  even. Then to make function evaluations less expensive, we can define the following

Table 2: Results for Modified Newton method applied to Extended Rosenbrock with exact gradient and hessian, metrics are average metrics for successful attempts.

preconditioning dimension	iterations		convergence rate		time		success rate	
	False	True	False	True	False	True	False	True
3	31.91	28.91	1.99	1.51	0.05	0.02	1.00	1.00
4	32.36	29.18	2.10	2.04	0.11	0.08	1.00	1.00
5	26.50	26.00	1.10	1.00	0.68	0.53	1.00	1.00

Table 3: Results for Truncated Newton method applied to Extended Rosenbrock with exact gradient and hessian, metrics are average metrics for successful attempts.

preconditioning dimension	iterations		convergence rate		time		success rate	
	False	True	False	True	False	True	False	True
3	50.09	31.00	-2.66	11.55	0.01	0.01	1.00	1.00
4	57.27	34.73	-3.06	4.33	0.04	0.03	1.00	1.00
5	64.00	28.50	-5.58	1.62	0.45	0.23	1.00	1.00

function  $F_{fd,k}$ , which can be plugged in 2 in place of  $F$  yielding the same result.

$$F_{fd,k}(x) = \begin{cases} \frac{1}{2}f_k^2(x) + \frac{1}{2}f_{k+1}^2(x), & k \bmod 2 = 1 \\ \frac{1}{2}f_{k-1}^2(x), & k \bmod 2 = 0 \end{cases}$$

The same procedure can be applied for the Hessian, considering that:

- function evaluations to compute entry  $h_{k,k}$  differ only by  $f_k$  and  $f_{k+1}$  for  $k$  odd, and only on  $f_{k-1}$  for  $k$  even;
- function evaluations to compute entry  $h_{k,k+1}$  differ only by  $f_k$  and  $f_{k+1}$  for  $k$  odd.

Then to make function evaluations less expensive, we can define the functions  $F_{fd,k,k}$  and  $F_{fd,k,k+1}$ , which can be plugged in 3 in place of  $F$  yielding the same result to compute entries  $h_{k,k}$  and  $h_{k,k+1}$  respectively.

$$F_{fd,k,k}(x) = \begin{cases} \frac{1}{2}f_k^2(x) + \frac{1}{2}f_{k+1}^2(x), & k \bmod 2 = 1 \\ \frac{1}{2}f_{k-1}^2(x), & k \bmod 2 = 0 \end{cases}$$

$$F_{fd,k,k+1}(x) = \begin{cases} \frac{1}{2}f_k^2(x) + \frac{1}{2}f_{k+1}^2(x), & k \bmod 2 = 1 \\ 0, & k \bmod 2 = 0 \end{cases}$$

When plugging the functions  $F_{fd,k}$ ,  $F_{fd,k,k}$  and  $F_{fd,k,k+1}$  into 2 and 3 it's convenient to expand them so that the computation of the gradient and Hessian is not subject to numerical cancellation. After expanding the functions, the gradient and Hessian can be approximated as follows.

$$\frac{\partial F}{\partial x_k} \approx \begin{cases} 600h^2x_k - 100hx_{k+1} + \frac{1}{2}h + 350h^3 + 300hx_k^2, & k \bmod 2 = 1 \\ -100x_{k-1}^2 + 100x_k, & k \bmod 2 = 0 \end{cases}$$

$$\frac{\partial^2 F}{\partial x_k \partial x_j} \approx \begin{cases} 1200h_kx_k - 200x_{k+1} + 1 + 700h_k^2 + 600x_k^2, & j = k, k \bmod 2 = 1 \\ 100, & j = k, k \bmod 2 = 0 \\ -100h_kh_{k+1} - 200x_k, & |j - k| = 1, k \bmod 2 = 1 \\ 0, & \text{otherwise} \end{cases}$$

Tables 4 and 5 show the results for the *Modified Newton method* applied to the extended Rosenbrock function with absolute and specific finite differences respectively. All attempts were successful, and the method converged in a small number of iterations, with a convergence rate that is close to 2 for all dimensions, when a suitable choice of  $h$  is made. On the contrary, when a poor choice of  $h$  is made, the convergence rate is close to 1, which is expected since the convergence rate of the Newton method is 2: this happens when  $h = 10^{-2}$  or  $h = 10^{-4}$ , both for the constant increment and the specific increment methodologies. We notice that for a fixed dimension, performance improves as the stepsize  $h$  (be it constant or specific) decreases, as expected, up to a certain point

Table 4: Results for Modified Newton method applied to Extended Rosenbrock with absolute finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-02	149.91	150.82	1.00	1.00	0.07	0.07	1.00	1.00
	1e-04	33.18	32.09	1.01	1.00	0.02	0.02	1.00	1.00
	1e-06	32.09	30.18	1.96	2.09	0.03	0.02	1.00	1.00
	1e-08	32.00	30.18	2.04	2.23	0.02	0.02	1.00	1.00
	1e-10	32.00	30.18	2.00	2.23	0.02	0.02	1.00	1.00
	1e-12	32.00	30.18	2.00	2.23	0.02	0.02	1.00	1.00
4	1e-02	160.09	160.55	1.00	1.00	0.34	0.36	1.00	1.00
	1e-04	34.18	32.55	1.00	1.00	0.12	0.09	1.00	1.00
	1e-06	32.55	30.27	1.89	1.96	0.12	0.08	1.00	1.00
	1e-08	32.55	30.27	1.95	1.98	0.12	0.08	1.00	1.00
	1e-10	32.55	30.27	1.95	1.98	0.11	0.08	1.00	1.00
	1e-12	32.55	30.27	1.95	1.98	0.11	0.08	1.00	1.00
5	1e-02	169.36	169.82	1.00	1.00	3.26	3.36	1.00	1.00
	1e-04	34.82	33.73	1.00	1.00	1.08	0.79	1.00	1.00
	1e-06	34.00	31.73	2.13	2.00	1.17	0.75	1.00	1.00
	1e-08	33.27	31.73	2.60	2.10	1.17	0.73	1.00	1.00
	1e-10	33.45	31.73	1.93	2.10	1.13	0.73	1.00	1.00
	1e-12	33.45	31.73	1.90	2.10	1.15	0.73	1.00	1.00

Table 5: Results for Modified Newton method applied to Extended Rosenbrock with specific finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-02	139.82	141.45	1.00	1.00	0.06	0.06	1.00	1.00
	1e-04	33.73	32.27	1.00	1.00	0.02	0.02	1.00	1.00
	1e-06	32.09	30.18	1.89	2.06	0.02	0.02	1.00	1.00
	1e-08	32.00	30.18	2.04	2.23	0.02	0.02	1.00	1.00
	1e-10	32.00	30.18	2.00	2.23	0.02	0.02	1.00	1.00
	1e-12	32.00	30.18	2.00	2.23	0.02	0.02	1.00	1.00
4	1e-02	149.55	150.55	1.00	1.00	0.32	0.32	1.00	1.00
	1e-04	34.45	32.45	1.00	1.00	0.12	0.09	1.00	1.00
	1e-06	32.55	30.27	1.81	1.96	0.12	0.09	1.00	1.00
	1e-08	32.55	30.27	1.95	1.98	0.12	0.08	1.00	1.00
	1e-10	32.55	30.27	1.95	1.98	0.12	0.08	1.00	1.00
	1e-12	32.55	30.27	1.95	1.98	0.11	0.08	1.00	1.00
5	1e-02	159.18	160.27	1.00	1.00	3.05	3.17	1.00	1.00
	1e-04	34.64	34.45	1.00	1.00	1.09	0.80	1.00	1.00
	1e-06	33.45	31.73	2.56	2.00	1.16	0.74	1.00	1.00
	1e-08	33.09	31.73	1.84	2.10	1.13	0.73	1.00	1.00
	1e-10	33.45	31.73	1.90	2.10	1.12	0.74	1.00	1.00
	1e-12	33.45	31.73	1.90	2.10	1.14	0.77	1.00	1.00

Table 6: Results for Truncated Newton method applied to Extended Rosenbrock with absolute finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-02	182.00	179.27	1.00	1.00	0.02	0.04	1.00	1.00
	1e-04	54.36	37.91	1.00	1.00	0.01	0.01	1.00	1.00
	1e-06	52.82	36.36	-1.48	2.39	0.01	0.01	1.00	1.00
	1e-08	52.91	35.91	-2.12	2.86	0.01	0.01	1.00	1.00
	1e-10	53.36	35.45	-4.10	2.21	0.01	0.01	1.00	1.00
	1e-12	55.64	36.45	-2.49	3.01	0.01	0.01	1.00	1.00
4	1e-02	203.91	192.18	1.00	1.00	0.16	0.17	1.00	1.00
	1e-04	61.64	44.45	1.00	1.00	0.06	0.05	1.00	1.00
	1e-06	59.18	42.55	-1.35	4.36	0.06	0.05	1.00	1.00
	1e-08	61.64	42.09	-1.78	2.34	0.06	0.05	1.00	1.00
	1e-10	62.27	41.91	-2.46	3.14	0.06	0.05	1.00	1.00
	1e-12	61.73	43.00	-2.89	2.65	0.06	0.05	1.00	1.00
5	1e-02	224.55	206.82	1.00	1.00	1.84	1.83	1.00	1.00
	1e-04	74.64	42.91	1.00	1.01	0.73	0.41	1.00	1.00
	1e-06	79.09	49.27	-2.17	2.06	0.81	0.50	1.00	1.00
	1e-08	75.64	49.91	-2.83	7.31	0.75	0.47	1.00	1.00
	1e-10	78.00	51.64	-1.21	3.21	0.75	0.50	1.00	1.00
	1e-12	78.00	53.27	-2.11	2.50	0.77	0.52	1.00	1.00

where reducing the stepsize does not result in any improvement (i.e. time and iterations do not significantly decrease). This plateau is reached sooner when preconditioning is adopted and when the dimension is lower.

Tables 6 and 7 show the results for the *Truncated Newton method* applied to the extended Rosenbrock function with absolute and specific finite differences respectively. All attempts were successful, and the method converged in a small number of iterations, when a suitable choice of  $h$  is made. A larger number of iterations is required when a poor choice of  $h$  is made, namely when  $h = 10^{-2}$  or  $h = 10^{-4}$ , both for the constant increment and the specific increment methodologies, yielding a convergence rate of 1. For smaller values of  $h$ , the estimate of the convergence rate is not reliable at all and considerations made in subsection 3.1 for the Truncated Newton method applied to the extended Rosenbrock function with exact gradient and Hessian apply here as well. The plateau in performance as the stepsize  $h$  decreases is reached slower than in the case of the Modified Newton method, hinting that the Truncated Newton method highly benefits from a finer approximation of the derivatives.

## 4 Generalized Broyden tridiagonal function

The generalized Broyden tridiagonal function is defined as follows.

$$F(x) = \frac{1}{2} \sum_{i=1}^n f_k^2(x) \quad f_k(x) = (3 - 2x_k)x_k + 1 - x_{k-1} - x_{k+1} \quad (8)$$

Figure 4 shows the surface plot of the 2-dimensional generalized Broyden tridiagonal function. Notice that the area where the minimum lies is very flat, which makes it hard to converge to the minimum.

### 4.1 Exact gradient and Hessian

The gradient of the generalized Broyden tridiagonal function is given by the following expression,

$$\frac{\partial F}{\partial x_k} = \begin{cases} (3 - 4x_1)f_1(x) - f_2(x), & k = 1 \\ (3 - 4x_k)f_k(x) - f_{k+1}(x) - f_{k-1}(x), & 1 < k < n \\ (3 - 4x_n)f_n(x) - f_{n-1}, & k = n \end{cases} \quad (9)$$

Table 7: Results for Truncated Newton method applied to Extended Rosenbrock with specific finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-02	169.91	169.91	1.00	1.00	0.02	0.03	1.00	1.00
	1e-04	54.45	37.27	1.00	1.02	0.01	0.01	1.00	1.00
	1e-06	51.45	36.27	-2.69	2.35	0.01	0.01	1.00	1.00
	1e-08	53.64	36.27	-3.43	1.91	0.01	0.01	1.00	1.00
	1e-10	54.64	35.27	-2.94	2.01	0.01	0.01	1.00	1.00
	1e-12	53.82	35.91	-4.18	2.56	0.01	0.01	1.00	1.00
4	1e-02	191.09	182.64	1.00	1.00	0.15	0.17	1.00	1.00
	1e-04	64.27	42.55	1.00	1.01	0.06	0.05	1.00	1.00
	1e-06	62.18	43.00	-0.96	2.62	0.06	0.05	1.00	1.00
	1e-08	60.73	40.91	-2.95	2.87	0.06	0.05	1.00	1.00
	1e-10	60.82	42.73	-1.22	3.98	0.06	0.05	1.00	1.00
	1e-12	63.09	42.91	-2.01	2.51	0.06	0.05	1.00	1.00
5	1e-02	213.55	198.64	1.00	1.00	1.68	1.74	1.00	1.00
	1e-04	75.91	40.82	1.00	1.00	0.74	0.42	1.00	1.00
	1e-06	77.91	49.00	-2.96	6.91	0.77	0.49	1.00	1.00
	1e-08	77.64	50.18	-1.96	2.28	0.77	0.48	1.00	1.00
	1e-10	75.00	53.27	-4.07	3.38	0.75	0.52	1.00	1.00
	1e-12	75.09	52.45	-2.09	3.48	0.75	0.52	1.00	1.00

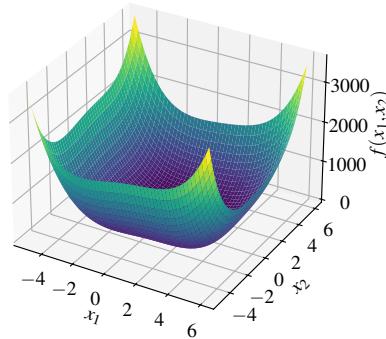


Figure 4: Surface plot of the 2-dimensional generalized Broyden tridiagonal function

Table 8: Results for Modified Newton method applied to Generalized Broyden with exact gradient and hessian, metrics are average metrics for succesful attempts.

preconditioning dimension	iterations		convergence rate		time		success rate	
	False	True	False	True	False	True	False	True
3	8.636	8.636	1.974	1.978	0.006	0.005	1.00	1.00
4	8.000	8.273	1.829	1.849	0.037	0.034	1.00	1.00
5	7.909	10.545	1.794	1.989	0.291	0.300	1.00	1.00

Table 9: Results for Truncated Newton method applied to Generalized Broyden with exact gradient and hessian, metrics are average metrics for succesful attempts.

preconditioning dimension	iterations		convergence rate		time		success rate	
	False	True	False	True	False	True	False	True
3	11.818	9.000	1.774	1.971	0.002	0.003	1.00	1.00
4	12.727	9.727	1.346	1.987	0.014	0.015	1.00	1.00
5	13.727	9.636	1.905	1.923	0.125	0.108	1.00	1.00

computation can be eased considering that component  $k$  depends only on  $f_k$ ,  $f_{k+1}$  and  $f_{k-1}$ . The Hessian of the generalized Broyden tridiagonal function is given by the following expression.

$$\frac{\partial^2 F}{\partial x_k \partial x_j} = \begin{cases} (3 - 4x_1)^2 - 4f_1(x) + 1, & k = j = 1 \\ (3 - 4x_k)^2 - 4f_k(x) + 2, & 1 < k = j < n \\ (3 - 4x_n)^2 - 4f_n(x) + 1, & k = j = n \\ 4x_k + 4x_{k+1} - 6, & |k - j| = 1 \\ 1, & |k - j| = 2 \\ 0, & \text{otherwise} \end{cases} \quad (10)$$

Notice that the Hessian is a banded matrix, with only  $n$  non-zero elements on the diagonal,  $n - 1$  non-zero elements on the first co-diagonal and  $n - 2$  non-zero elements on the second co-diagonal.

Table 8 shows the results for the *Modified Newton method* applied to the generalized Broyden tridiagonal function with exact gradient and Hessian. All attempts are succesful, and the method converges in a small number of iterations. Convergence rate is close to 2, which is the expected convergence rate for Newton's method. Preconditioning does not seem to have a significant impact on the performance of the method, compared to its impact on the performance on the extended Rosenbrock function due to the fact that the Hessian matrix for the Newton system is better scaled in the currently considered function.

Table 9 shows the results for the *Truncated Newton method* applied to the generalized Broyden tridiagonal function with exact gradient and Hessian. All attempts are succesful, and the method converges in a small number of iterations. When preconditioning is not adopted convergence rate is lower than 2, which is expected since we are solving the Newton system with a relative tolerance that is guaranteed to yield superlinear convergence. Differently from the Modified Newton method, the Truncated Newton method is more sensitive to preconditioning, which can be observed by the fact that the number of iterations is lower when preconditioning is adopted and convergence rate is almost 2.

Comparing performance of the two methods, the Truncated Newton method is 2 to 3 times faster than the Modified Newton method due to the fact that the Newton system is solved approximately, which is less computationally expensive than solving it exactly. However, the Truncated Newton methods requires more iterations to converge than the Modified Newton method, which is expected since the Modified Newton exactly solves the Newton system at each iteration.

## 4.2 Finite differences gradient and Hessian

When applying 2, one can notice that the terms  $F(x + he_k)$  and  $F(x - he_k)$  only differ by terms  $f_k$ ,  $f_{k+1}$  and  $f_{k-1}$ . Then to make function evaluations less expensive, we can define the following function  $F_{fd,k}$ , which can be plugged in 2 in place of  $F$  yielding the same result.

$$F_{fd,k}(x) = \frac{1}{2}f_k^2(x) + \frac{1}{2}f_{k+1}^2(x) + \frac{1}{2}f_{k-1}^2(x)$$

Table 10: Results for Modified Newton method applied to Generalized Broyden with absolute finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-04	8.545	8.545	2.016	2.010	0.006	0.006	1.00	1.00
	1e-06	8.545	8.636	1.964	1.979	0.006	0.006	1.00	1.00
	1e-08	8.636	8.636	1.974	1.978	0.006	0.005	1.00	1.00
	1e-10	8.636	8.636	1.974	1.978	0.006	0.006	1.00	1.00
	1e-12	8.636	8.636	1.973	1.978	0.006	0.006	1.00	1.00
4	1e-04	7.909	8.273	1.803	1.829	0.036	0.031	1.00	1.00
	1e-06	8.000	8.273	1.832	1.851	0.036	0.031	1.00	1.00
	1e-08	8.000	8.273	1.829	1.849	0.038	0.031	1.00	1.00
	1e-10	8.000	8.273	1.829	1.849	0.038	0.031	1.00	1.00
	1e-12	8.000	8.273	1.829	1.849	0.038	0.032	1.00	1.00
5	1e-04	7.818	10.455	1.751	2.045	0.321	0.344	1.00	1.00
	1e-06	7.909	10.545	1.796	1.992	0.333	0.332	1.00	1.00
	1e-08	7.909	10.545	1.794	1.989	0.324	0.337	1.00	1.00
	1e-10	7.909	10.545	1.794	1.989	0.314	0.325	1.00	1.00
	1e-12	7.909	10.545	1.794	1.989	0.320	0.328	1.00	1.00

The same procedure can be applied for the Hessian, considering that:

- function evaluations to compute entry  $h_{k,k}$  differ only by  $f_k$ ,  $f_{k+1}$  and  $f_{k-1}$ ;
- function evaluations to compute entry  $h_{k,k+1}$  differ only by  $f_k$  and  $f_{k+1}$ ;
- function evaluations to compute entry  $h_{k,k+2}$  differ only by  $f_{k-1}$ .

Then to make function evaluations less expensive, we can define the functions  $F_{fd,k,k}$ ,  $F_{fd,k,k+1}$ ,  $F_{fd,k,k+2}$ , which can be plugged in 3 in place of  $F$  yielding the same result to compute entries  $h_{k,k}$ ,  $h_{k,k+1}$  and  $h_{k,k+2}$  respectively.

$$\begin{aligned} F_{fd,k,k}(x) &= \frac{1}{2}f_k^2(x) + \frac{1}{2}f_{k-1}^2(x) + \frac{1}{2}f_{k+1}^2(x) \\ F_{fd,k,k+1}(x) &= \frac{1}{2}f_k^2(x) + \frac{1}{2}f_{k+1}^2(x) \\ F_{fd,k,k+2}(x) &= \frac{1}{2}f_{k-1}^2(x) \end{aligned}$$

When plugging the functions  $F_{fd,k}$ ,  $F_{fd,k,k}$ ,  $F_{fd,k,k+1}$  and  $F_{fd,k,k+2}$  into 2 and 3 it's convenient to expand them so that the computation of the gradient and Hessian is not subject to numerical cancellation as previously done for the extended Rosenbrock function in subsection 3.2.

Tables 10 and 11 show the results for the *Modified Newton method* applied to the generalized Broyden tridiagonal function with absolute and specific finite differences respectively. All attempts are successful, but attempts with  $h = 10^{-2}$  which don't converge within the fixed maximum number of iterations  $k_{max} = 1000$ , both for absolute and specific differences, for all dimensions. This is probably due to the fact that  $h = 10^{-2}$  is not a suitable increment for the finite differences method to approximate the gradient and Hessian of the generalized Broyden tridiagonal function. When the method converges, it does so in a small number of iterations with a convergence rate close to 2, which is the expected convergence rate for Newton's method.

Tables 12 and 13 show the results for the *Truncated Newton method* applied to the generalized Broyden tridiagonal function with absolute and specific finite differences respectively. All attempts are successful, but attempts with  $h = 10^{-2}$  which converge within the fixed maximum number of iterations  $k_{max} = 1000$  only once for  $n = 10^3$ , regardless of the type (absolute or specific) of finite differences adopted. When the method converges, it does so in a small number of iterations with a convergence rate that as expected is superlinear, almost quadratic for  $n = 10^5$ , even if preconditioning is not adopted. When preconditioning is adopted, the number of iterations is lower and the convergence rate is closer to 2, which is the expected convergence rate for Newton's method.

Also for this function, it's evident that the Truncated Newton method is faster than the Modified Newton method, but requires more iterations to converge. Moreover, the Truncated Newton method is more sensitive to preconditioning than the Modified Newton method.

Table 11: Results for Modified Newton method applied to Generalized Broyden with specific finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-04	8.545	8.636	1.977	2.000	0.006	0.006	1.00	1.00
	1e-06	8.545	8.636	1.964	1.979	0.006	0.005	1.00	1.00
	1e-08	8.636	8.636	1.974	1.978	0.007	0.007	1.00	1.00
	1e-10	8.636	8.636	1.973	1.978	0.006	0.006	1.00	1.00
	1e-12	8.636	8.636	1.973	1.978	0.007	0.006	1.00	1.00
4	1e-04	7.909	8.273	1.804	1.835	0.037	0.031	1.00	1.00
	1e-06	8.000	8.273	1.830	1.850	0.050	0.031	1.00	1.00
	1e-08	8.000	8.273	1.829	1.849	0.038	0.032	1.00	1.00
	1e-10	8.000	8.273	1.829	1.849	0.038	0.032	1.00	1.00
	1e-12	8.000	8.273	1.829	1.849	0.038	0.031	1.00	1.00
5	1e-04	7.818	10.545	1.752	2.020	0.315	0.325	1.00	1.00
	1e-06	7.909	10.545	1.795	1.990	0.314	0.326	1.00	1.00
	1e-08	7.909	10.545	1.794	1.989	0.329	0.329	1.00	1.00
	1e-10	7.909	10.545	1.794	1.989	0.325	0.327	1.00	1.00
	1e-12	7.909	10.545	1.794	1.989	0.313	0.325	1.00	1.00

Table 12: Results for Truncated Newton method applied to Generalized Broyden with absolute finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-02	13.000	NaN	1.051	NaN	0.002	NaN	0.09	NaN
	1e-04	11.818	9.091	1.773	2.009	0.003	0.003	1.00	1.00
	1e-06	11.818	9.000	1.774	1.972	0.003	0.003	1.00	1.00
	1e-08	11.818	9.000	1.774	1.971	0.003	0.003	1.00	1.00
	1e-10	11.818	9.000	1.774	1.971	0.003	0.003	1.00	1.00
	1e-12	11.818	9.000	1.774	1.971	0.003	0.004	1.00	1.00
4	1e-04	12.727	9.545	1.347	1.995	0.019	0.018	1.00	1.00
	1e-06	12.727	9.727	1.346	1.990	0.018	0.017	1.00	1.00
	1e-08	12.727	9.727	1.346	1.987	0.018	0.017	1.00	1.00
	1e-10	12.727	9.727	1.346	1.987	0.018	0.018	1.00	1.00
	1e-12	12.727	9.727	1.346	1.987	0.019	0.018	1.00	1.00
5	1e-04	13.727	9.545	1.905	1.903	0.160	0.132	1.00	1.00
	1e-06	13.727	9.636	1.905	1.926	0.162	0.133	1.00	1.00
	1e-08	13.727	9.636	1.905	1.923	0.159	0.137	1.00	1.00
	1e-10	13.727	9.636	1.905	1.923	0.161	0.132	1.00	1.00
	1e-12	13.727	9.636	1.905	1.923	0.156	0.132	1.00	1.00

Table 13: Results for Truncated Newton method applied to Generalized Broyden with specific finite differences, metrics are average metrics for successful attempts.

dimension	preconditioning h	iterations		convergence rate		time		success rate	
		False	True	False	True	False	True	False	True
3	1e-04	11.818	9.091	1.773	1.986	0.003	0.003	1.00	1.00
	1e-06	11.818	9.000	1.774	1.971	0.003	0.003	1.00	1.00
	1e-08	11.818	9.000	1.774	1.971	0.003	0.003	1.00	1.00
	1e-10	11.818	9.000	1.774	1.971	0.003	0.004	1.00	1.00
	1e-12	11.818	9.000	1.774	1.971	0.003	0.003	1.00	1.00
4	1e-04	12.727	9.545	1.347	1.982	0.018	0.017	1.00	1.00
	1e-06	12.727	9.727	1.346	1.989	0.017	0.018	1.00	1.00
	1e-08	12.727	9.727	1.346	1.987	0.018	0.018	1.00	1.00
	1e-10	12.727	9.727	1.346	1.987	0.017	0.017	1.00	1.00
	1e-12	12.727	9.727	1.346	1.987	0.018	0.018	1.00	1.00
5	1e-04	13.727	9.636	1.905	1.930	0.155	0.133	1.00	1.00
	1e-06	13.727	9.636	1.905	1.924	0.159	0.131	1.00	1.00
	1e-08	13.727	9.636	1.905	1.923	0.160	0.134	1.00	1.00
	1e-10	13.727	9.636	1.905	1.923	0.158	0.132	1.00	1.00
	1e-12	13.727	9.636	1.905	1.923	0.156	0.132	1.00	1.00

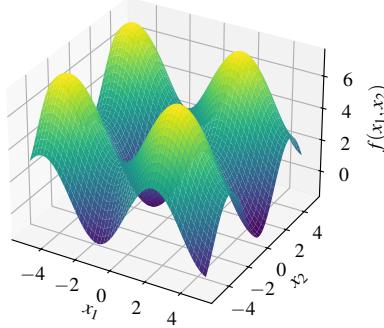


Figure 5: Surface plot of the 2-dimensional banded trigonometric function

## 5 Banded trigonometric function

The banded trigonometric function is defined as follows.

$$f(x) = \sum_{i=1}^n i[(1 - \cos x_i) + \sin x_{i-1} - \sin x_{i+1}] \quad (11)$$

Figure 5 shows the surface plot of the 2-dimensional banded trigonometric function. It has a highly oscillatory landscape with multiple peaks and valleys due to its sinusoidal terms. This results in a mix of locally smooth and rapidly changing regions, making optimization sensitive to initialization and prone to multiple local minima. The optimization runs for the banded trigonometric function were conducted using a high value for the Armijo condition parameter  $c_1 = 10^{-2}$  while the other parameters were kept the same as in the previous experiments. This was necessary since when running the modified Newton method, the too small required decrease led to stagnation.

Table 14: Results for Modified Newton method applied to Banded Trigonometric with exact gradient and hessian, metrics are average metrics for successful attempts. Results are given only with preconditioning adopted.

dimension	iterations	convergence rate	time	success rate
3	27.273	2.593	0.025	1.00
4	51.545	2.749	0.150	1.00
5	7.000	2.946	0.092	0.09

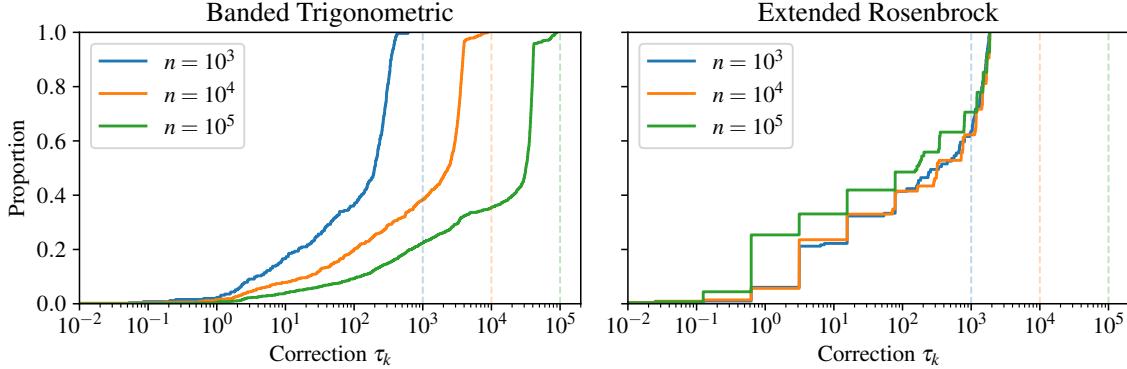


Figure 6: Empirical cumulative distribution function of  $\tau_k$  in banded trigonometric problem and extended Rosenbrock function with exact gradient and Hessian, considering all times a correction is needed across all runs. Dashed lines mark the values  $10^3$ ,  $10^4$  and  $10^5$ .

## 5.1 Exact gradient and Hessian

The gradient of the banded trigonometric function is given by the following expression.

$$\frac{\partial F}{\partial x_k} = \begin{cases} k \sin x_k + 2 \cos x_k, & 1 \leq k < n \\ n \sin x_n - (n-1) \cos x_n, & k = n \end{cases} \quad (12)$$

The Hessian of the banded trigonometric function is given by the following expression.

$$\frac{\partial^2 F}{\partial x_k \partial x_j} = \begin{cases} k \cos x_k - 2 \sin(x_k), & 1 \leq k = j < n \\ n \cos x_n + (n-1) \sin x_n, & k = j = n \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

Notice that the Hessian is a diagonal matrix, which makes the optimization problem easier to solve. However, the Hessian of the matrix has very distant eigenvalues due to the fact that the  $i$ -th diagonal entry is multiplied by  $i$ : the problem may become increasingly ill-conditioned as the dimension  $n$  increases. Due to this fact, we only perform optimization with preconditioning.

Table 14 shows the results for the *Modified Newton method* applied to the banded trigonometric function with exact gradient and Hessian. Modified newton method converges for all points when  $n = 10^3$  or  $n = 10^4$ , while attempts with  $n = 10^5$  yield only one success corresponding to the suggested starting point. When problem dimension is  $n = 10^5$ , attempts with randomly initialized points do not converge within  $k_{max} = 10^3$  iterations. This may be due to the fact that being the problem badly scaled, in case of a non positive-definite Hessian matrix, the modification of the Hessian matrix may happen with a  $\tau$  that is too large. Figure 6 compares the empirical cumulative distribution function of  $\tau_k$  for the banded trigonometric problem and the extended Rosenbrock function with exact gradient and Hessian. One can notice that while for the extended Rosenbrock function the maximum value of the correction does not depend on the dimension, for the banded trigonometric problem the maximum value of the correction increases with the dimension. Moreover, as suggested by the factor  $i$  that multiplies the  $i$ -th diagonal entry of the Hessian, the maximum correction needed for the dimension  $10^n$  is close to  $10^n$ . While a correction of this magnitude is acceptable for  $n = 10^3$  or  $n = 10^4$ , it is not for larger values of  $n$ , leading to a Hessian matrix that is too different from the original one and resulting in a poor descent direction. All failures happen because convergence has not been reached within  $k_{max} = 10^3$  iterations.

Table 15 shows the results for the *Truncated Newton method* applied to the banded trigonometric function with exact gradient and Hessian. All attempts are succesful, and the method converges in a small number of iterations with a large experimental rate of convergence. So, when preconditioning is adopted, the Truncated

Table 15: Results for Truncated Newton method applied to Banded Trigonometric with exact gradient and hessian, metrics are average metrics for successful attempts. Results are given only with preconditioning adopted.

dimension	iterations	convergence rate	time	success rate
3	14.091	2.811	0.009	1.00
4	20.455	2.743	0.028	1.00
5	25.000	2.776	0.301	1.00

Newton method is able to solve the problem efficiently unlike Modified Newton method. Probably adopting another matrix correction strategy that impacts less on the whole matrix, such as the one based on *Modified Cholesky factorization* proposed in [2], could lead to better results.

## 5.2 Finite differences gradient and Hessian

When applying 2, one can notice that the terms  $F(x + he_k)$  and  $F(x - he_k)$  only differ by summands with indices  $i = k$  and  $i = k \pm 1$ . Then to make function evaluations less expensive, we can define the following function  $F_{fd,k}$ , that can be plugged in 2 in place of  $F$  yielding the same result.

$$F_{fd,k} = \sum_{i=k-1}^{k+1} i[(1 - \cos x_i) + \sin x_{i-1} - \sin x_{i+1}]$$

Some of the terms in  $F_{fd,k}$  are constant and do not depend on  $x_k$ , so they will eventually subtract and we can further simplify the expression.

$$F_{fd,k} = \begin{cases} -k \cos x_k + 2 \sin x_k, & 1 \leq k < n \\ n \cos x_n - (n-1) \sin x_n, & k = n \end{cases}$$

The same procedure can be applied to the Hessian, considering that terms in 3 only differ by summands with indices  $i = k$  and  $i = k \pm 1$ . When plugging the function  $F_{fd,k}$  into 2 and 3 it's convenient to expand them so that the computation of the gradient and Hessian is not subject to numerical cancellation. After expanding the functions, the gradient and Hessian can be approximated as follows.

$$\frac{\partial F}{\partial x_k} \approx \begin{cases} \frac{2k \sin x_k \sin h + 4 \cos x_k \sin h}{2h}, & 1 \leq k < n \\ \frac{2n \sin x_n \sin h - 2(n-1) \cos x_n \sin h}{2h}, & k = n \end{cases}$$

$$\frac{\partial^2 F}{\partial x_k^2} \approx \begin{cases} (k \cos x_k - 2 \sin x_k) - h(k \sin x_k + 2 \cos x_k), & 1 \leq k < n \\ (n \cos x_n - (n-1) \sin x_n) - h(n \sin x_n - (n-1) \cos x_n), & k = n \end{cases}$$

The approximation for the gradient is obtained applying the trigonometric identities for the sum of angles. The approximation for the Hessian is obtained by applying the aforementioned trigonometric identities in order to isolate the following terms.

$$\cos(2h) - 2 \cos h \approx -1 - h^2 + \mathcal{O}(h^4), \quad \sin(2h) - 2 \sin h \approx -h^3 + \mathcal{O}(h^5)$$

To avoid numerical cancellation, it is necessary to substitute them with their Taylor expansions.

Tables 16 and 17 show the results for the *Modified Newton method* applied to the banded trigonometric function with absolute and specific finite differences, respectively. In terms of success rate, results are comparable to the ones obtained with exact gradient and Hessian. However, when specific finite differences are adopted, approximation with  $h = 10^{-2}$ ,  $h = 10^{-4}$  and  $h = 10^{-6}$  give poor results, with success rate for  $h = 10^{-2}$  being 0% for all dimension. On the other hand, when increments are constant, we get at least 9% success rate for all dimensions, all increments and 100% success rate when  $n = 10^4$  or  $n = 10^5$ . When absolute increment is adopted, all failures happen due to the fact that convergence has not been reached within  $k_{max} = 10^3$  iterations. When specific increment is adopted, most failures happen because the method does not converge within  $k_{max} = 10^3$  iterations, but in 1 case over 18 failures for  $n = 10^3$ , 4 cases over 24 failures for  $n = 10^4$  and in 3 cases over 61 failures for  $n = 10^5$ , failure happens because Armijo condition can't be satisfied.

Tables 18 and 19 show the results for the *Truncated Newton method* applied to the banded trigonometric function with absolute and specific finite differences, respectively. All attempts are successful, and the method converges in a small number of iterations with a large experimental rate of convergence. The results are similar to the ones obtained with exact gradient and Hessian, confirming that the Truncated Newton method is a suitable algorithm to tackle the banded trigonometric problem optimization.

Table 16: Results for Modified Newton method applied to Banded Trigonometric with finite differences and constant increment, metrics are average metrics for successful attempts. Results are given only with preconditioning adopted.

dimension	h	iterations	convergence rate	time	success rate
3	1e-02	27.364	2.676	0.018	1.00
	1e-04	27.545	2.743	0.014	1.00
	1e-06	25.909	2.807	0.013	1.00
	1e-08	27.818	2.467	0.015	1.00
	1e-10	27.182	2.770	0.014	1.00
	1e-12	27.182	2.817	0.015	1.00
4	1e-02	60.818	2.661	0.179	1.00
	1e-04	50.091	2.626	0.149	1.00
	1e-06	54.364	2.831	0.161	1.00
	1e-08	49.091	2.770	0.152	1.00
	1e-10	51.364	2.767	0.149	1.00
	1e-12	51.909	2.516	0.157	1.00
5	1e-02	9.000	2.705	0.129	0.09
	1e-04	8.000	2.828	0.113	0.09
	1e-06	6.000	2.856	0.085	0.09
	1e-08	6.000	2.749	0.088	0.09
	1e-10	7.000	2.948	0.098	0.09
	1e-12	7.000	2.946	0.099	0.09

Table 17: Results for Modified Newton method applied to Banded Trigonometric with finite differences and specific increment, metrics are average metrics for successful attempts. Results are given only with preconditioning adopted.

dimension	h	iterations	convergence rate	time	success rate
3	1e-04	162.000	0.984	0.050	0.36
	1e-06	26.818	2.389	0.014	1.00
	1e-08	27.000	2.610	0.015	1.00
	1e-10	27.091	2.597	0.014	1.00
	1e-12	27.273	2.820	0.015	1.00
4	1e-04	9.000	0.973	0.018	0.09
	1e-06	55.875	inf	0.161	0.73
	1e-08	50.636	2.791	0.152	1.00
	1e-10	51.091	2.617	0.152	1.00
	1e-12	51.273	2.684	0.152	1.00
5	1e-04	9.000	1.074	0.128	0.09
	1e-06	6.000	2.894	0.083	0.09
	1e-08	6.000	2.749	0.085	0.09
	1e-10	7.000	2.948	0.095	0.09
	1e-12	7.000	2.946	0.098	0.09

Table 18: Results for Truncated Newton method applied to Banded Trigonometric with finite differences and constant increment, metrics are average metrics for successful attempts. Results are given only with preconditioning adopted.

dimension	h	iterations	convergence rate	time	success rate
3	1e-02	14.909	2.723	0.007	1.00
	1e-04	14.091	2.615	0.006	1.00
	1e-06	14.091	2.531	0.005	1.00
	1e-08	14.091	2.770	0.006	1.00
	1e-10	14.000	2.794	0.005	1.00
	1e-12	14.000	2.779	0.006	1.00
4	1e-02	19.091	2.595	0.028	1.00
	1e-04	19.364	2.587	0.028	1.00
	1e-06	20.091	2.646	0.028	1.00
	1e-08	20.636	2.741	0.030	1.00
	1e-10	20.364	2.798	0.028	1.00
	1e-12	20.455	2.758	0.028	1.00
5	1e-02	23.000	2.471	0.301	0.91
	1e-04	25.909	2.648	0.315	1.00
	1e-06	25.727	2.612	0.327	1.00
	1e-08	25.636	2.762	0.313	1.00
	1e-10	25.364	2.754	0.312	1.00
	1e-12	25.273	2.699	0.311	1.00

Table 19: Results for Truncated Newton method applied to Banded Trigonometric with finite differences and specific increment, metrics are average metrics for successful attempts. Results are given only with preconditioning adopted.

dimension	h	iterations	convergence rate	time	success rate
3	1e-02	15.000	2.251	0.005	1.00
	1e-04	15.091	2.860	0.006	1.00
	1e-06	13.818	2.534	0.006	1.00
	1e-08	13.818	2.609	0.006	1.00
	1e-10	14.000	2.795	0.006	1.00
	1e-12	14.000	2.803	0.006	1.00
4	1e-02	19.909	2.118	0.027	1.00
	1e-04	19.182	2.666	0.026	1.00
	1e-06	20.455	2.851	0.028	1.00
	1e-08	20.182	2.650	0.028	1.00
	1e-10	19.818	2.811	0.027	1.00
	1e-12	19.909	2.699	0.027	1.00
5	1e-02	24.273	2.351	0.313	1.00
	1e-04	25.364	2.715	0.295	1.00
	1e-06	26.455	2.748	0.315	1.00
	1e-08	25.364	2.717	0.321	1.00
	1e-10	25.455	2.713	0.315	1.00
	1e-12	25.091	2.549	0.308	1.00

## 6 Conclusions

Our experiments provide insights into the practical performance of Newton-type methods for large-scale unconstrained optimization problems. The Modified Newton Method, while robust, can suffer from inefficiencies when dealing with ill-conditioned Hessians, particularly in the banded trigonometric function. On the other hand, the Truncated Newton Method, which relies on iterative solutions to the Newton system, generally achieves better computational efficiency, particularly for high-dimensional problems. However, it is more sensitive to preconditioning, which significantly influences its convergence rate and stability.

Finite difference approximations, though useful in the absence of exact derivatives, introduce additional computational costs and potential inaccuracies, particularly when the step size is not chosen carefully.

## A Code snippets

AP: [TODO]

### A.1 Code for method implementations

#### A.1.1 Modified Newton method

Listing 1: Implementation of the Modified Newton method

```
1 function [xk, fk, gradfk_norm, k, T, success, xseq] = ...
2     modified_newton(x0, f, gradf, Hessf, beta, ...
3         kmax, tolgrad, c1, rho, btmax, max_chol_iter, preconditioning, logging,
4         modification_coeff)
5 %
6 %
7 % INPUTS:
8 % x0 = n-dimensional column vector;
9 % f = function handle that describes a function R^n->R;
10 % gradf = function handle that describes the gradient of f;
11 % Hessf = function handle that describes the Hessian of f;
12 % kmax = maximum number of iterations permitted;
13 % tolgrad = value used as stopping criterion w.r.t. the norm of the gradient;
14 % c1 = the factor of the Armijo condition that must be a scalar in (0,1);
15 % rho = fixed factor, lesser than 1, used for reducing alpha0;
16 % btmax = maximum number of steps for updating alpha during the backtracking
17 % strategy.
18 %
19 % OUTPUTS:
20 % xk = the last x computed by the function;
21 % fk = the value f(xk);
22 % gradfk_norm = value of the norm of gradf(xk)
23 % k = index of the last iteration performed
24 % xseq = n-by-k matrix where the columns are the elements xk of the
25 % sequence
26 % btseq = 1-by-k vector where elements are the number of backtracking
27 % iterations at each optimization step.
28 %
29 if nargin < 14
30     modification_coeff = 2;
31 end
32
33 n = length(x0);
34 % Function handle for the armijo condition
35 farmijo = @(fk, alpha, c1_gradfk_pk) ...
36     fk + alpha * c1_gradfk_pk;
37
38 % Initializations
39 if logging
```

```

38     xseq = zeros(n, kmax);
39 else
40     xseq = [];
41 end
42 success = 1;
43
44 gradfkseq = zeros(1, kmax);
45 fkseq = zeros(1, kmax);
46 btseq = zeros(1, kmax);
47 pcgiterseq = zeros(1, kmax);
48 correctionseq = zeros(1, kmax);
49 errornormseq = zeros(1, kmax);
50
51 xk = x0;
52 fk = f(xk);
53 gradfk = gradf(xk);
54 k = 0;
55 gradfk_norm = norm(gradfk);
56
56 Hessfk = Hessf(xk);
57
58 while k < kmax && gradfk_norm >= tolgrad
    % Compute the descent direction as solution of
    % Hessf(xk) p = - gradf(xk)
    %%%%%% L.S. SOLVED WITH BACKSLASH (NOT USED) %%%%%%
    % pk = -Hessf(xk)\gradfk;
    %%%%%% L.S. SOLVED WITH pcg %%%%%%
    % For simplicity: default values for tol and maxit; no preconditioning
    % pk = pcg(Hessf(xk), -gradfk);
    % If you want to silence the messages about "solution quality", use
    % instead:
    [B, tau] = chol_with_addition(Hessfk, beta, modification_coeff,
        max_chol_iter);
    % B = eigenvalue_modification(Hessfk);
    % B = modchol_ldlt(Hessfk);
    % TODO Warning: Input tol may not be achievable by PCG - Try to use a
        bigger tolerance
    Hkm = Hessfk + B;
    if preconditioning
        try
            L = ichol(Hkm);
            [pk, ~, ~, iterk, ~] = pcg(Hkm, -gradfk, 1e-6, 1000, L, L');
        catch
            [pk, ~, ~, iterk, ~] = pcg(Hkm, -gradfk);
        end
    else
        [pk, ~, ~, iterk, ~] = pcg(Hkm, -gradfk);
    end
    %%%
    % Reset the value of alpha
    alpha = 1;
    %%%
    % Compute the candidate new xk
    xnew = xk + alpha * pk;
    % Compute the value of f in the candidate new xk
    fnew = f(xnew);
    %%%
    c1_gradfk_pk = c1 * gradfk' * pk;

```

```

96    bt = 0;
97    % Backtracking strategy:
98    % 2nd condition is the Armijo condition not satisfied
99    while bt < btmax && fnew > farmijo(fk, alpha, c1_gradfk_pk)
100        % Reduce the value of alpha
101        alpha = rho * alpha;
102        % Update xnew and fnew w.r.t. the reduced alpha
103        xnew = xk + alpha * pk;
104        fnew = f(xnew);

105        % Increase the counter by one
106        bt = bt + 1;
107    end
108    if bt == btmax && fnew > farmijo(fk, alpha, c1_gradfk_pk)
109        disp("Armijo condition could not be satisfied!")
110        success = 0;
111        break
112    end
113    errornormseq(k+1) = norm(xnew - xk);
114    % Update xk, fk, gradfk_norm
115    xk = xnew;
116    fk = fnew;
117    gradfk = gradf(xk);
118    gradfk_norm = norm(gradfk);
119    Hessfk = Hessf(xk);

120    % Increase the step by one
121    k = k + 1;

122    if logging
123        xseq(:, k) = xk;
124    end
125    gradfkseq(k) = gradfk_norm;
126    fkseq(k) = fk;
127    btseq(k) = bt;
128    pcgiterseq(k) = iterk;
129    correctionseq(k) = tau;
130
131 end

132 gradfkseq = gradfkseq(1:k);
133 fkseq = fkseq(1:k);
134 btseq = btseq(1:k);
135 pcgiterseq = pcgiterseq(1:k);
136 correctionseq = correctionseq(1:k);
137 errornormseq = errornormseq(1:k);
138 T = table(gradfkseq', fkseq', btseq', pcgiterseq', correctionseq',
139             errornormseq', ...
140             'VariableNames', {'gradient_norm', 'function_value', 'backtrack',
141                               'inner_iterations', 'correction', 'error_norm'});
142
143 if logging
144     xseq = xseq(:, 1:k);
145 end
146
147 if k >= kmax && gradfk_norm >= tolgrad
148     success = 0;
149 end
150
151 end

```

Listing 2: Implementation of *Cholesky with added multiple of the identity*

```

1 function [B, tau] = chol_with_addition(A, beta, coefficient, max_iter)
2
3 mindiag = min(diag(A));
4 if mindiag > 0
5     tau = 0;
6 else
7     tau = -mindiag + beta;
8 end
9
10 sizes = size(A);
11 n = sizes(1);
12
13 for i = 1:max_iter
14     try chol(A + tau*speye(n));
15         if i > 1
16             end
17             break
18         catch
19             tau = max(coefficient*tau, beta);
20         end
21     end
22 if i >= max_iter
23     tau = 0;
24     disp('Tau could not be found!')
25 end
26
27 B = speye(n) * tau;
28
29 end

```

### A.1.2 Truncated Newton method

Listing 3: Implementation of the Truncated Newton method

```

1 function [xk, fk, gradfk_norm, k, T, success, xseq] = ...
2     truncated_newton(x0, f, gradf, Hessf, ...
3         kmax, tolgrad, c1, rho, btmax, fterms, pcg_maxit, preconditioning,
4         logging)
5 %
6 % Function that performs the Inexact Newton optimization method, using
7 % backtracking strategy for the step-length selection.
8 %
9 % INPUTS:
10 % x0 = n-dimensional column vector;
11 % f = function handle that describes a function R^n->R;
12 % gradf = function handle that describes the gradient of f;
13 % Hessf = function handle that describes the Hessian of f;
14 % kmax = maximum number of iterations permitted;
15 % tolgrad = value used as stopping criterion w.r.t. the norm of the
16 % gradient;
17 % c1 = the factor of the Armijo condition that must be a scalar in (0,1);
18 % rho = fixed factor, lesser than 1, used for reducing alpha0;
19 % btmax = maximum number of steps for updating alpha during the
20 % backtracking strategy.
21 % fterms = function handle taking as input arguments k and gradfk, and
22 % returning the forcing term etak
22 % pcg_maxit = maximum number of iterations for the pcg solver
22 %

```

```

23 % OUTPUTS:
24 % xk = the last x computed by the function;
25 % fk = the value f(xk);
26 % gradfk_norm = value of the norm of gradf(xk)
27 % k = index of the last iteration performed
28 % xseq = n-by-k matrix where the columns are the elements xk of the
29 % sequence
30 % btseq = 1-by-k vector where elements are the number of backtracking
31 % iterations at each optimization step.
32 % pcgiterseq = 1-by-k vector where elements are the number of pcg
33 % iterations at each optimization step.
34 %
35 n = length(x0);
36 % Function handle for the armijo condition
37 farmijo = @(fk, alpha, c1_gradfk_pk) ...
38     fk + alpha * c1_gradfk_pk;
39
40 % Initializations
41 if logging
42     xseq = zeros(n, kmax);
43 else
44     xseq = [];
45 end
46 success = 1;
47
48 gradfkseq = zeros(1, kmax);
49 fkseq = zeros(1, kmax);
50 btseq = zeros(1, kmax);
51 pcgiterseq = zeros(1, kmax);
52 truncatedseq = zeros(1, kmax);
53 errornormseq = zeros(1, kmax);
54
55 xk = x0;
56 fk = f(xk);
57 gradfk = gradf(xk);
58 k = 0;
59 gradfk_norm = norm(gradfk);
60
61 while k < kmax && gradfk_norm >= tolgrad
62     % "INEXACTLY" compute the descent direction as approximated solution of
63     % Hessf(xk) p = - gradf(xk)
64
65     % TOLERANCE VARYING W.R.T. FORCING TERMS:
66     etak = fterms(k, gradfk);
67     % ATTENTION! We will use directly eta_k as tolerance in the pcg because
68     % this function looks at the RELATIVE RESIDUAL and not the RESIDUAL!
69
70     %%%%%% L.S. SOLVED WITH pcg %%%%%%
71     % For simplicity: default values for tol and maxit; no preconditioning
72     % pk = pcg(Hessf(xk), -gradfk, etak, pcg_maxit);
73     % If you want to silence the messages about solution "quality" use
74     % instead:
75     % [pk, flagk, relresk, iterk, resveck] = pcg(Hessf(xk), ...
76     % -gradfk, etak, pcg_maxit);
77     % [pk, ~, iterk] = cg(Hessf(xk), -gradfk, pcg_maxit, etak);
78     Hk = Hessf(xk);
79     if preconditioning
80         try
81             L = ichol(Hk);

```

```

82     [pk, ~, iterk, truncated] = cg_preconditioned(Hk, -gradfk,
83         pcg_maxit, etak, L);
84     catch
85         % If the preconditioner fails, we will use the default one
86         [pk, ~, iterk, truncated] = cg(Hk, -gradfk, pcg_maxit, etak);
87     end
88 else
89     [pk, ~, iterk, truncated] = cg(Hessf(xk), -gradfk, pcg_maxit, etak);
90 end
91
92 % Reset the value of alpha
93 alpha = 1;
94
95 % Compute the candidate new xk
96 xnew = xk + alpha * pk;
97 % Compute the value of f in the candidate new xk
98 fnew = f(xnew);
99
100 c1_gradfk_pk = c1 * gradfk' * pk;
101 bt = 0;
102 % Backtracking strategy:
103 % 2nd condition is the Armijo condition not satisfied
104 while bt < btmax && fnew > farmijo(fk, alpha, c1_gradfk_pk)
105     % Reduce the value of alpha
106     alpha = rho * alpha;
107     % Update xnew and fnew w.r.t. the reduced alpha
108     xnew = xk + alpha * pk;
109     fnew = f(xnew);
110
111     % Increase the counter by one
112     bt = bt + 1;
113 end
114 if bt == btmax && fnew > farmijo(fk, alpha, c1_gradfk_pk)
115     disp('Armijo condition could not be satisfied!')
116     success = 0;
117     break
118 end
119 errornormseq(k+1) = norm(xnew - xk);
120 % Update xk, fk, gradfk_norm
121 xk = xnew;
122 fk = fnew;
123 gradfk = gradf(xk);
124 gradfk_norm = norm(gradfk);
125
126 % Increase the step by one
127 k = k + 1;
128
129 if logging
130     xseq(:, k) = xk;
131 end
132 gradfkseq(k) = gradfk_norm;
133 fkseq(k) = fk;
134 btseq(k) = bt;
135 pcgiterseq(k) = iterk;
136 truncatedseq(k) = truncated;
137 end
138
139 gradfkseq = gradfkseq(1:k);
140 fkseq = fkseq(1:k);

```

```

141 btseq = btseq(1:k);
142 pcgiterseq = pcgiterseq(1:k);
143 truncatedseq = truncatedseq(1:k);
144 errornormseq = errornormseq(1:k);
145 T = table(gradfkseq', fkseq', btseq', pcgiterseq', truncatedseq',
146   errornormseq', ...
147   'VariableNames', {'gradient_norm', 'function_value', 'backtrack',
148   'inner_iterations', 'truncated', 'error_norm'});
149
150 if logging
151   xseq = xseq(:, 1:k);
152 end
153
154 if k >= kmax && gradfk_norm >= tolgrad
155   success = 0;
156 end
157
158 end

```

Listing 4: Implementation of conjugate gradient method with truncation

```

1 function [xk, k, relres, truncated] = ...
2   cg(A, b, kmax, tol)
3 % function [xk, k, relres] = ...
4 %   lab01_cg_linsys(A, b, x0, kmax, tol)
5 %
6 % Conjugate Gradient Method for linear systems.
7 %
8 % INPUTS:
9 % A : n-by-n symmetric, positive (semi-)definite matrix of the system
10 % b : column vector of n elements, known terms of the system
11 % x0 : columns vector of n elements, starting guess
12 % kmax : positive integer, maximum number of steps
13 % tol : positive real value, tolerance for relative residual
14
15 xk = zeros(size(b));
16 rk = b - A * xk;
17 pk = rk;
18
19 truncated = 0;
20 k = 0;
21
22 norm_b = norm(b);
23 relres = norm(rk) / norm_b;
24
25 while k < kmax && relres > tol
26   zk = A * pk;
27   if pk' * zk <= 0
28     if k == 0
29       xk = b;
30     end
31     truncated = 1;
32     % disp(['Stopped at iteration ', num2str(k)]);
33     return
34   end
35   alphak = (rk' * pk) / (pk' * zk);
36   xk = xk + alphak * pk;
37   rk = rk - alphak * zk;
38
39   betak = -(rk' * zk) / (pk' * zk);

```

```

40     pk = rk + betak * pk;
41
42     relres = norm(rk) / norm_b;
43     k = k + 1;
44 end
45
46 end

```

Listing 5: Implementation of preconditioned conjugate gradient method with truncation

```

1 function [xk, k, relres, truncated] = ...
2     cg_preconditioned(A, b, kmax, tol, L)
3 % function [xk, k, relres] = ...
4 %     lab01_cg_linsys(A, b, x0, kmax, tol)
5 %
6 % Conjugate Gradient Method for linear systems.
7 %
8 % INPUTS:
9 % A : n-by-n symmetric, positive (semi-)definite matrix of the system
10 % b : column vector of n elements, known terms of the system
11 % x0 : columns vector of n elements, starting guess
12 % kmax : positive integer, maximum number of steps
13 % tol : positive real value, tolerance for relative residual
14
15 xk = zeros(size(b));
16 rk = A * xk - b;
17 wk = L \ rk;
18 yk = L' \ wk;
19 pk = -yk;
20
21 k = 0;
22 truncated = 0;
23
24 norm_b = norm(b);
25 relres = norm(rk) / norm_b;
26
27 while k < kmax && relres > tol
28     zk = A * pk;
29     if pk' * zk <= 0
30         if k == 0
31             xk = b;
32         end
33         truncated = 1;
34         % disp(['Stopped at iteration ', num2str(k)]);
35         return
36     end
37     alphak = (rk' * yk) / (pk' * zk);
38     xk = xk + alphak * pk;
39     rk_new = rk + alphak * zk;
40
41     wk = L \ rk_new;
42     yk_new = L' \ wk;
43
44     betak = (rk_new' * yk_new) / (rk' * yk);
45     pk = -yk_new + betak * pk;
46
47     rk = rk_new;
48     yk = yk_new;
49
50     relres = norm(rk) / norm_b;

```

```

51     k = k + 1;
52 end
53
54 end

```

## A.2 Code for objective functions

### A.2.1 Rosenbrock function

Listing 6: Rosenbrock function

```

1 function f = rosenbrock(x)
2 f = 100*(x(2) - x(1)^2)^2 + (1 - x(1))^2;
3 end
4
5 function g = rosenbrock_gradient(x)
6 g = [-400*x(1)*(x(2) - x(1)^2) - 2*(1 - x(1));
7     200*(x(2) - x(1)^2)];
8 end
9
10 function H = rosenbrock_hessian(x)
11 H = [1200*x(1)^2 - 400*x(2) + 2, -400*x(1);
12     -400*x(1), 200];
13 end
14
15 x0_0 = [1.2; 1.2];
16 x0_1 = [-1.2; 1];
17 starting_points = {x0_0, x0_1};
18
19 kmax = 1e3;
20 tol = 1e-6;
21 % Backtracking parameters
22 rho = 0.5;
23 c1 = 1e-4;
24 btmax = 50;
25 % Modified Newton parameters
26 beta = 1e-3;
27 c = 2;
28 % Truncated Newton parameters
29 fterm_suplin = @(k, gradfk) min(0.5, sqrt(norm(gradfk)));
30
31 % Define the Rosenbrock function
32 rosenbrockFunction = @(x, y) (1 - x).^2 + 100 * (y - x.^2).^2;
33
34 % Create a grid of X and Y values
35 x = linspace(-1.5, 1.5, 100);
36 y = linspace(-0.5, 2.5, 100);
37 [X, Y] = meshgrid(x, y);
38 Z = rosenbrockFunction(X, Y);
39
40 % Run the optimization algorithms
41 for i = 1:length(starting_points)
42     % Set figure properties for readability in a paper
43     set(groot, 'defaultAxesFontSize', 30);
44     set(groot, 'defaultTextInterpreter', 'latex');
45     set(groot, 'defaultLegendInterpreter', 'latex');
46     set(groot, 'defaultColorbarTickLabelInterpreter', 'latex');
47     set(groot, 'defaultAxesTickLabelInterpreter', 'latex');
48     % Plot the Rosenbrock function contour
49     figure;

```

```

50 contourLevels = logspace(-0.5, 3.5, 20); % Logarithmically spaced
51     contour levels
52 contour(X, Y, Z, contourLevels);
53 xlabel('$x_1$', 'Interpreter', 'latex');
54 ylabel('$x_2$', 'Interpreter', 'latex');
55 colorbar;
56 % Starting point
57 x0 = starting_points{i};
58 disp(' ');
59 disp(['Starting point: ', mat2str(x0)]);
60 tic;
61 [xk, fk, gradfk_norm, k, T, success, xseq] = modified_newton( ...
62     x0, @rosenbrock, @rosenbrock_gradient, @rosenbrock_hessian, ...
63     beta, kmax, tol, c1, rho, btmax, 50, false, true, 2);
64 time = toc;
65 disp([9, 'Modified Newton: ', num2str(success), ', ', num2str(k), ', ', ...
66     num2str(gradfk_norm), ', ', num2str(fk)]);
67 disp([9, 'time: ', num2str(time), ' k = ', num2str(k)]);
68 disp([9, 'xk = ', mat2str(xk)]);
69 % Plot the descent path for Modified Newton
70 hold on;
71 plot(xseq(1, :), xseq(2, :), 'r*-', 'MarkerSize', 10);
72
73 tic;
74 [xk, fk, gradfk_norm, k, T, success, xseq] = truncated_newton( ...
75     x0, @rosenbrock, @rosenbrock_gradient, @rosenbrock_hessian, ...
76     kmax, tol, c1, rho, btmax, fterm_suplin, 100, false, true);
77 time = toc;
78 disp([9, 'Truncated Newton: ', num2str(success), ', ', num2str(k), ', ', ...
79     num2str(gradfk_norm), ', ', num2str(fk)]);
80 disp([9, 'time: ', num2str(time), ' k = ', num2str(k)]);
81 disp([9, 'xk = ', mat2str(xk)]);
82 % Plot the descent path for Truncated Newton
83 plot(xseq(1, :), xseq(2, :), 'b*-', 'MarkerSize', 10);
84 hold off;
85 legend('Rosenbrock function', 'Modified Newton', 'Truncated Newton');
86
87 % Save the figure
88 orient landscape;
89 set(gca, 'LooseInset', get(gca, 'TightInset'));
90 filename = ['../report/figures/rosenbrock_x0_', num2str(i)];
91 saveas(gcf, filename, 'pdf');
92 end

```

### A.2.2 Extended Rosenbrock function

Listing 7: Extended Rosenbrock function

```

1 function [f] = extended_rosenbrock(x)
2 n = length(x);
3 fk_odd = @(x) 10*(x(1:2:n-1).^2 - x(2:2:n));
4 fk_even = @(x) x(1:2:n-1) - 1;
5 f = 0.5*sum(fk_odd(x).^2 + fk_even(x).^2);
6 end

```

Listing 8: Gradient of the Extended Rosenbrock function

```

1 function [gradf] = extended_rosenbrock_gradf(x)
2 n = length(x);
3 gradf = zeros(n,1);

```

```

4 gradf(1:2:n-1) = 200*x(1:2:n-1).^3 - 200*x(2:2:n).*x(1:2:n-1) + x(1:2:n-1) -
5   1;
6 gradf(2:2:n) = -100*(x(1:2:n-1).^2 - x(2:2:n));
7 end

```

Listing 9: Hessian of the Extended Rosenbrock function

```

1 function [Hessf] = extended_rosenbrock_Hessf(x)
2 n = length(x);
3 Bin = zeros(n,3);
4 Bin(1:2:n-1, 2) = 200*(3*x(1:2:n-1).^2 - x(2:2:n)) + 1;
5 Bin(2:2:n, 2) = 100;
6 Bin(1:2:n-1, 1) = -200*x(1:2:n-1);
7 Bin(2:2:n-2, 1) = 0;
8 Bin(2:2:n, 3) = Bin(1:n-1, 1);
9 Hessf = spdiags(Bin, -1:1, n, n);
10 end

```

Listing 10: Approximated gradient of the Extended Rosenbrock function

```

1 function [gradf] = extended_rosenbrock_gradf_fd(x, h, relative)
2   n = length(x);
3   if relative
4     hs = h*abs(x);
5   else
6     hs = h*ones(n, 1);
7   end
8
9   gradf = zeros(n, 1);
10  gradf(1:2:n) = 2*x(1:2:n).*hs(1:2:n) - 2*hs(1:2:n)...
11    + 400*x(1:2:n).^3.*hs(1:2:n) + 400*x(1:2:n).*hs(1:2:n).^3 -
12    400*x(1:2:n).*x(2:2:n).*hs(1:2:n);
13  gradf(2:2:n) = -200.*hs(2:2:n).*x(1:2:n).^2 + 200*hs(2:2:n).*x(2:2:n);
14  gradf = gradf ./ (2*hs);
15 end

```

Listing 11: Approximated Hessian of the Extended Rosenbrock function

```

1 function [Hessf] = extended_rosenbrock_Hessf_fd(x, h, relative)
2   n = length(x);
3   Bin = zeros(n, 3);
4   if relative
5     hs = h*abs(x);
6   else
7     hs = h*ones(n, 1);
8   end
9
10  % Diagonal
11  Bin(1:2:n, 2) = 1200*hs(1:2:n).*x(1:2:n) - 200.*x(2:2:n)...
12    + 1 + 700*hs(1:2:n).^2 + 600*x(1:2:n).^2;
13  Bin(2:2:n, 2) = 100;
14  % Off-diagonal
15  Bin(1:2:n, 1) = -100*hs(1:2:n) - 200*x(1:2:n);
16  Bin(2:2:n, 3) = Bin(1:2:n, 1);
17  Hessf = spdiags(Bin, -1:1, n, n);
18 end

```

### A.2.3 Generalized Broyden tridiagonal function

Listing 12: Generalized Broyden tridiagonal function

```

1 function [f] = extended_rosenbrock(x)
2 n = length(x);
3 fk_odd = @(x) 10*(x(1:2:n-1).^2 - x(2:2:n));
4 fk_even = @(x) x(1:2:n-1) - 1;
5 f = 0.5*sum(fk_odd(x).^2 + fk_even(x).^2);
6 end

```

Listing 13: Gradient of the generalized Broyden tridiagonal function

```

1 % Problem 32
2 function [f] = generalized_broyden(x)
3     fk_function = @(x) (3-2*x).*x + 1 - [0; x(1:end-1)] - [x(2:end); 0];
4     fk = fk_function(x);
5     f = sum(fk.^2)/2;
6 end

```

Listing 14: Hessian of the generalized Broyden tridiagonal function

```

1 function [Hessf] = generalized_broyden_Hessf(x)
2     n = length(x);
3     fk_function = @(x) (3-2*x).*x + 1 - [0; x(1:end-1)] - [x(2:end); 0];
4     fk = fk_function(x);
5     Bin = zeros(n, 5);
6     % Main diagonal
7     Bin(2:n-1, 3) = -4*fk(2:n-1) + (3-4*x(2:n-1)).^2 + 2;
8     Bin(1, 3) = (3-4*x(1)).^2 - 4*fk(1) + 1;
9     Bin(n, 3) = (3-4*x(n)).^2 - 4*fk(n) + 1;
10    % First codiagonal
11    Bin(1:n-1, 2) = 4*x(1:n-1) + 4*x(2:n) - 6;
12    Bin(2:n, 4) = Bin(1:n-1, 2);
13    % Second codiagonal
14    Bin(1:n-2, 1) = 1;
15    Bin(3:n, 5) = 1;
16    % Create banded matrix
17    Hessf = spdiags(Bin, -2:2, n, n);
18 end

```

Listing 15: Approximated gradient of the generalized Broyden tridiagonal function

```

1 function [gradf] = generalized_broyden_gradf_fd(x, h, relative)
2     n = length(x);
3     if relative
4         hs = h*abs(x);
5     else
6         hs = h*ones(n, 1);
7     end
8     xm2 = [0; 0; x(1:end-2)];
9     xm1 = [0; x(1:end-1)];
10    xp1 = [x(2:end); 0];
11    xp2 = [x(3:end); 0; 0];
12    gradf = 2*xm1.^2 + 4*xm1.*x - 6*xm1 + 8*x.^3 - 18*x.^2 + 4*x.*xp1 +
13        8*x.*hs.^2 + 7.*x ...
14        + 2*xp1.^2 - 6*xp1 - 6*hs.^2 + xm2 + xp2 + 1;
15    gradf(1) = 6*x(1) - 6*x(2) + x(3) + 4*x(1)*x(2) + 8*x(1)*hs(1)^2 ...
16        - 18*x(1)^2 + 8*x(1)^3 + 2*x(2)^2 - 6*hs(1)^2 + 2;
17    gradf(n) = x(n-2) - 6*x(n-1) + 6*x(n) + 4*x(n-1)*x(n) + 8*x(n)*hs(n)^2
18        ...
19        + 2*x(n-1)^2 - 18*x(n)^2 + 8*x(n)^3 - 6*hs(n)^2 + 2;
end

```

Listing 16: Approximated Hessian of the generalized Broyden tridiagonal function

```

1 function [Hessf] = generalized_broyden_Hessf_fd(x, h, relative)
2 n = length(x);
3 if relative
4     hs = h*abs(x);
5 else
6     hs = h*ones(n, 1);
7 end
8 Bin = zeros(n, 5);
9
10 xm1 = [0; x(1:n-1)];
11 xp1 = [x(2:n); 0];
12
13 % Diagonal
14 Bin(1:n, 3) = 24*x.^2 + 48*x.*hs - 36*x + 28*hs.^2 - 36*hs + 4*xm1 + 4*xp1 +
    7;
15 Bin(1, 3) = 24*x(1)^2 + 48*x(1)*hs(1) - 36*x(1) + 28*hs(1)^2 - 36*hs(1) +
    4*x(2) + 6;
16 Bin(n, 3) = 24*x(n)^2 + 48*x(n)*hs(n) - 36*x(n) + 28*hs(n)^2 - 36*hs(n) +
    4*x(n-1) + 6;
17 % First codiagonal
18 Bin(1:n-1, 2) = 4*x(1:n-1) + 4*x(2:n) + 2*hs(1:n-1) + 2*hs(2:n) - 6;
19 Bin(2:n, 4) = Bin(1:n-1, 2);
20 % Second codiagonal
21 Bin(1:n-2, 1) = 1;
22 Bin(3:n, 5) = 1;
23
24 Hessf = spdiags(Bin, -2:2, n, n);
25 end

```

#### A.2.4 Banded trigonometric problem

Listing 17: Banded trigonometric problem

```

1 function [f] = banded_trigonometric(x)
2     n = length(x);
3     x_m1 = [0; x(1:n-1)];
4     x_p1 = [x(2:n); 0];
5     i = (1:n)';
6     f = sum(i.*(1-cos(x) + sin(x_m1) - sin(x_p1)));
7 end

```

Listing 18: Gradient of the banded trigonometric problem

```

1 function [gradf] = banded_trigonometric_gradf(x)
2     n = length(x);
3     i = (2:n-1)';
4     gradf = zeros(n, 1);
5     gradf(2:n-1) = i.*sin(x(2:n-1)) + 2*cos(x(2:n-1));
6     gradf(1) = sin(x(1)) + 2*cos(x(1));
7     gradf(n) = n*sin(x(n)) - (n-1)*cos(x(n));
8 end

```

Listing 19: Hessian of the banded trigonometric problem

```

1 function [Hessf] = banded_trigonometric_Hessf(x)
2     n = length(x);
3     diag_entries = zeros(n, 1);
4     i = (2:n-1)';

```

```

5 diag_entries(2:n-1) = i.*cos(x(2:n-1)) - 2*sin(x(2:n-1));
6 diag_entries(1) = cos(x(1)) - 2*sin(x(1));
7 diag_entries(n) = n*cos(x(n)) + (n-1)*sin(x(n));
8 Hessf = spdiags(diag_entries, 0, n, n);
9 end

```

Listing 20: Approximated gradient of the banded trigonometric problem

```

1 function [gradf] = banded_trigonometric_gradf_fd(x, h, relative)
2     n = length(x);
3     if relative
4         hs = h*abs(x);
5     else
6         hs = h*ones(n, 1);
7     end
8
9     i = (1:n-1)';
10
11    gradf = [
12        2*i.*sin(x(1:n-1)).*sin(hs(1:n-1)) + 4*cos(x(1:n-1)).*sin(hs(1:n-1));
13        2*n.*sin(x(n)).*sin(hs(n)) - 2*(n-1)*cos(x(n)).*sin(hs(n));
14    ];
15    gradf = gradf ./ (2*hs);
16 end

```

Listing 21: Approximated Hessian of the banded trigonometric problem

```

1 function [Hessf] = banded_trigonometric_Hessf_fd(x, h, relative)
2     n = length(x);
3     if relative
4         hs = h*abs(x);
5     else
6         hs = h*ones(n, 1);
7     end
8
9     i = (1:n)';
10    diagonal = (-i.*cos(x) + 2*sin(x)).*(-1) ...
11        + (i.*sin(x) + 2*cos(x)).*(-hs);
12    diagonal(n) = (-n.*cos(x(n)) - (n-1)*sin(x(n))).*(-1) ...
13        + (n.*sin(x(n)) - (n-1)*cos(x(n))).*(-hs(n));
14    Hessf = spdiags(diagonal, 0, n, n);
15 end

```

### A.3 Utility code for running experiments

Listing 22: Script to run all the experiments for all functions

```

1 clear all; close all; clc; %#ok

```

```

14 rng(seed);
15 disp('***** CHAINED ROSENBROCK *****')
16 codiags = 1;
17 test(@(x) extended_rosenbrock(x), @(x) extended_rosenbrock_gradf(x),...
18     @(x)extended_rosenbrock_Hessf(x), @(n)
19         extended_rosenbrock_initializer(n),...
20             @extended_rosenbrock_gradf_fd, @extended_rosenbrock_Hessf_fd, codiags, ...
21                 kmax, 1e-6, 1e-4, 0.5, 50, 100, 1e-3, fterm, 100,
22                     './results/extended_rosenbrock/', true, 5);
23
24 rng(seed);
25 disp('***** GENERALIZED BROYDEN *****')
26 codiags = 2;
27 test(@(x) generalized_broyden(x), @(x) generalized_broyden_gradf(x),...
28     @(x)generalized_broyden_Hessf(x), @(n)
29         generalized_broyden_initializer(n),...
30             @generalized_broyden_gradf_fd, @generalized_broyden_Hessf_fd, codiags, ...
31                 kmax, 1e-6, 1e-4, 0.5, 50, 100, 1e-3, fterm, 100,
32                     './results/generalized_broyden/', true, 2);
33
34 rng(seed);
35 disp('***** BANDED TRIGONOMETRIC *****')
36 codiags = 0;
37 test(@(x) banded_trigonometric(x), @(x) banded_trigonometric_gradf(x),...
38     @(x)banded_trigonometric_Hessf(x), @(n)
39         banded_trigonometric_initializer(n),...
40             @banded_trigonometric_gradf_fd, @banded_trigonometric_Hessf_fd,
41                 codiags, ...
42                     kmax, 1e-6, 1e-2, 0.5, 50, 100, 1e-3, fterm, 100,
43                         './results/banded_trigonometric/', true, 2);

```

Listing 23: Script to run all the experiments for a given function

```

1 function [] = test(f, gradf, Hessf, initializer, gradf_fd, Hessf_fd,
2     codiags, ...
3     kmax, tolgrad, c1, rho, btmax, chol_maxit, beta, fterms, pcg_maxit,
4     root_dir, findiff, tau_coeff)
5 %
6 % INPUTS
7 % n = dimension of the problem;
8 % f = function handle that describes a function R^n->R;
9 % gradf = function handle that describes the gradient of f;
10 % Hessf = function handle that describes the Hessian of f;
11 % initializer = function handle that generate starting point x0;
12 % kmax = maximum number of iterations permitted;
13 % tolgrad = value used as stopping criterion w.r.t. the norm of the gradient;
14 % c1 = the factor of the Armijo condition that must be a scalar in (0,1);
15 % rho = fixed factor, lesser than 1, used for reducing alpha0;
16 % btmax = maximum number of steps for updating alpha during the backtracking
17 % strategy.
18 % chol_maxit = maximum number of iterations for the Cholesky factorization
19 % beta = fixed factor, greater than 1, used for the Cholesky factorization;
20 % fterms = function handle taking as input arguments k and gradfk, and
21 % returning the forcing term etak
22 % pcg_maxit = maximum number of iterations for the pcg solver
23 % h = step size for the finite difference approximation of the gradient and
24 % Hessian.
25 %
26 % OUTPUTS

```

```

22 %
23 if ~exist(root_dir, 'dir')
24     mkdir(root_dir);
25 end
26 experiment = 1;
27
28 if nargin < 16
29     findiff = false;
30 end
31
32 for i=[3, 4, 5]
33     % Dimension of the problem
34     n=10^i;
35     % Initialization of the starting point
36     x0 = initializer(n);
37
38     % Run the optimization algorithm for 11 different starting points
39     for j=0:10
40         if j == 0
41             x0_j = x0;
42         else
43             x0_j = 2 * rand(n, 1) + x0 - 1;
44         end
45         disp([9, '* STARTING POINT: (dim:1e', num2str(i), ', test point #',
46               num2str(j), ')'])
47         disp([9, 9, '- EXACT GRADIENT AND HESSIAN ****'])
48         % Exact gradient and Hessian
49         for pre = [0, 1]
50             [fk_m, gradfk_norm_m, k_m, T_m, time_m, success_m, ...
51              fk_t, gradfk_norm_t, k_t, T_t, time_t, success_t] =
52                 run_optimization(x0_j, f, gradf, Hessf, beta, kmax(i),
53                               tolgrad, c1, rho, btmax, chol_maxit, fterms, pcg_maxit,
54                               pre, tau_coeff);
55             logger(root_dir, experiment, success_m, 0, pre, i, j, 1, 0, 0,
56                   fk_m, gradfk_norm_m, k_m, T_m, time_m);
57             experiment = experiment+1;
58             logger(root_dir, experiment, success_t, 1, pre, i, j, 1, 0, 0,
59                   fk_t, gradfk_norm_t, k_t, T_t, time_t);
60             experiment = experiment+1;
61         end
62         if findiff
63             % Finite difference gradient and Hessian for different values of
64             % h
65             for k=2:2:12
66                 h = 10^(-k);
67
68                 % Absolute
69                 disp([9, 9, '- ABSOLUTE FINITE DIFFERENCE GRADIENT AND
70                       HESSIAN - h=', num2str(h), ' ***'])
71                 % findiff_gradf = @(x) findiff_grad(f, x, h, 'c', false);
72                 findiff_gradf = @(x) gradf_fd(x, h, false);
73                 for pre = [0, 1]
74                     [fk_m, gradfk_norm_m, k_m, T_m, time_m, success_m, ...
75                      fk_t, gradfk_norm_t, k_t, T_t, time_t, success_t] =
76                         run_optimization(x0_j, f, findiff_gradf, @(x)
77                           Hessf_fd(x, h, false), beta, kmax(i), tolgrad,
78                           c1, rho, btmax, chol_maxit, fterms, pcg_maxit,
79                           pre, tau_coeff);
80                     logger(root_dir, experiment, success_m, 0, pre, i, j, 0,
81                           h, 1, fk_m, gradfk_norm_m, k_m, T_m, time_m);

```

```

69         experiment = experiment+1;
70         logger(root_dir, experiment, success_t, 1, pre, i, j, 0,
71                 h, 1, fk_t, gradfk_norm_t, k_t, T_t, time_t);
72         experiment = experiment+1;
73     end
74     % Relative
75     % findiff_gradf = @(x) findiff_grad(f, x, h, 'c', true);
76     findiff_gradf = @(x) gradf_fd(x, h, true);
77     disp([9, 9, '- RELATIVE FINITE DIFFERENCE GRADIENT AND
78           HESSIAN - h= ', num2str(h), ' ***'])
79     for pre = [0, 1]
80         [fk_m, gradfk_norm_m, k_m, T_m, time_m, success_m, ...
81          fk_t, gradfk_norm_t, k_t, T_t, time_t, success_t] =
82             run_optimization(x0_j, f, findiff_gradf, @(x)
83                               Hessf_fd(x, h, true), beta, kmax(i), tolgrad, c1,
84                               rho, btmax, chol_maxit, fterms, pcg_maxit, pre,
85                               tau_coeff);
86         logger(root_dir, experiment, success_m, 0, pre, i, j, 0,
87                 h, 0, fk_m, gradfk_norm_m, k_m, T_m, time_m);
88         experiment = experiment+1;
89         logger(root_dir, experiment, success_t, 1, pre, i, j, 0,
90                 h, 0, fk_t, gradfk_norm_t, k_t, T_t, time_t);
91         experiment = experiment+1;
92     end
93     end
94 end
95 end

```

Listing 24: Script to run experiments for a given configuration

```

1 function [fk_m, gradfk_norm_m, k_m, T_m, time_m, success_m, ...
2     fk_t, gradfk_norm_t, k_t, T_t, time_t, success_t] =
3     run_optimization(x0, f, gradf, Hessf, beta, ...
4     kmax, tolgrad, c1, rho, btmax, chol_maxit, fterms, pcg_maxit, pre,
5     tau_coeff)

6 disp([9, 9, 9, 'MODIFIED NEWTON *****'])
7 disp([9, 9, 9, 'PRECONDITIONING: ', num2str(pre)])
8 tic;
9 [~, fk_m, gradfk_norm_m, k_m, T_m, success_m, ~] = ...
10    modified_newton(x0, f, gradf, Hessf, beta, ...
11    kmax, tolgrad, c1, rho, btmax, chol_maxit, pre, false, tau_coeff);
12 time_m = toc;
13 disp([9, 9, 9, 'f(xk): ', num2str(fk_m)])
14 disp([9, 9, 9, 'gradfk_norm: ', num2str(gradfk_norm_m)])
15 disp([9, 9, 9, 'k: ', num2str(k_m)])
16 disp([9, 9, 9, '*****'])

17 disp([9, 9, 9, 'TRUNCATED NEWTON'])
18 disp([9, 9, 9, 'PRECONDITIONING: ', num2str(pre)])
19 tic;
20 [~, fk_t, gradfk_norm_t, k_t, T_t, success_t, ~] = ...
21    truncated_newton(x0, f, gradf, Hessf, ...
22    kmax, tolgrad, c1, rho, btmax, fterms, pcg_maxit, pre, false);
23 time_t = toc;
24 disp([9, 9, 9, 'f(xk): ', num2str(fk_t)])
25 disp([9, 9, 9, 'gradfk_norm: ', num2str(gradfk_norm_t)])
26 disp([9, 9, 9, 'k: ', num2str(k_t)])

```

```

27 disp([9, 9, 9, '*****'] )
28
29 end

```

Listing 25: Logger to CSV files

```

1 function [] = logger(root_dir, i, success, method, pre, dimension,
2   point_number, exact, h, absolute, fk, gradfk, k, T, time)
3   % Create the directory if it does not exist
4   if ~exist(root_dir, 'dir')
5     mkdir(root_dir);
6   end
7   % Save the results
8   writetable(T, [root_dir, 'experiment_', num2str(i), '.csv']);
9   % Append the experiment results to the log csv file
10  writematrix([i, success, method, pre, dimension, point_number, exact, h,
11    absolute, fk, gradfk, k, time], [root_dir, 'results.csv'],
12    'WriteMode', 'append');
13 end

```

## References

- [1] Ladislav Luksan and Jan Vlček. *Test Problems for Unconstrained Optimization*. Nov. 2003.
- [2] Jorge Nocedal and Stephen J. Wright. “Numerical optimization”. English (US). In: *Springer Series in Operations Research and Financial Engineering*. Springer Series in Operations Research and Financial Engineering. Springer Nature, 2006, pp. 1–664.