FEDERAL STATE AUTONOMOUS EDUCATIONAL INSTITUTION
OF HIGHER EDUCATION
ITMO UNIVERSITY

# Report

on the practical task No. 3

«Algorithms for unconstrained nonlinear optimization.

First- and second-order methods»

Performed by:

Putnikov Semyon

J4132c

Accepted by:

Dr Petr Chunaev

St. Petersburg

2020

## 1. Goal

The use of first- and second-order methods (Gradient Descent, Conjugate Gradient Descent, Newton's method and Levenberg-Marquardt algorithm) in the tasks of unconstrained nonlinear optimization.

## 2. Formulation of the problem

Generate random numbers $\alpha \in (0, 1)$ and $\beta \in (0, 1)$. Furthermore, generate the noisy data $\{x_k, y_k\}$, where $k = 0, \ldots, 100$, according to the following rule:

$$y_k = \alpha x_k + \beta + \delta_k, \quad x_k = \frac{k}{100}$$

where $\delta_k\ N(0, 1)$ are values of a random variable with standard normal distribution. Approximate the data by the following linear and rational functions:

(1) $F(x, a, b) = ax + b$    (linear approximant),

(2) $F(x, a, b) = \frac{a}{1+bx}$    (rational approximant),

by means of least squares through the numerical minimization (with precision $\varepsilon = 0.001$) of the following function:

$$D(a, b) = \sum_{k=0}^{100} (F(x_k, a, b) - y_k)^2$$

To solve the minimization problem, use the methods of Gradient Descent, Conjugate Gradient Descent, Newton's method and Levenberg-Marquardt algorithm. If necessary, set the initial approximations and other parameters of the methods. Visualize the data and the approximants obtained separately for each type of approximant. Analyze the results obtained (in terms of number of iterations, precision, number of function evaluations, etc.) and compare them with those from Task 2 for the same dataset.

# 3. Brief theoretical part

## 3.1 Gradient Descent method

Gradient descent is based on the observation that if $f(x)$ is defined and differentiable in a neighbourhood of a point $a$, then $f(x)$ decreases fastest in a neighbourhood of $a$ in the direction of $-\nabla_a f$. One obtains the following formula:

$$a_{n+1} = a_n - \gamma \nabla_{a_n} f$$

for $\gamma \in \mathbb{R}_+$ small enough, then $f(a_n) \geq f(a_{n+1})$. With this observation in mind, one starts with a guess $a_0$ for a local minimum of $f$, and considers the sequence $\{a_n\}$ such that

$$a_{n+1} = a_n - \gamma \nabla_{a_n} f, n \geq 0$$

Here the value of the step size $\gamma_n$ may be non-fixed and changed at every iteration (many possible ways to choose).

## 3.2 Conjugate Gradient Descent method

Given a function $f(x), x \in \mathbb{R}^n$ and an initial approximation $a_0$, one starts in the steepest descent direction:

$$\Delta a_0 = -\nabla_{a_0} f$$

Find the step length $\alpha_0 := argmin_\beta f(a_0 + \alpha \Delta a_0)$ and the next point $a_1 = a_0 + \alpha \Delta a_0$. After this iteration, the following steps constitute one iteration of moving along a subsequent conjugate direction $s_n$, where $s_0 = \Delta a_0$:

- Calculate the steepest direction $\Delta a_n = -\nabla_{a_n} f$.

- Compute $\beta_n$ according to certain formulas.

- Update the conjugate direction $s_n = \Delta a_n + \beta_n s_{n-1}$.

- Find $\alpha_n = argmin_\alpha f(a_n + \alpha s_n)$.

- Update the position: $a_{n+1} = a_n + \alpha_n s_n$.

The choice of $\beta_n$ due to Fletcher-Reeves:

$$\beta_n^{FR} = \frac{\Delta a_n^T \Delta a_n}{\Delta a_{n-1}^T \Delta a_{n-1}}$$

The choice of $\beta_n$ due to Polak-Ribiere:

$$\beta_n^{PR} = \frac{\Delta a_n^T (\Delta a_n - \Delta a_{n-1})}{\Delta a_{n-1}^T \Delta a_{n-1}}$$

## 3.3 Newton's method

Let $f : \mathbb{R} \to \mathbb{R}$ be convex and twice-differentiable. Find the roots of $f'$ by constructing a sequence $a_n$ from $a_n$ initial guess $a_0$ s.t. $a_n \to x^*$ as $n \to \infty$, where $f'(x^*) = 0$, i.e. $x^*$ is a stationary point of $f$.

From the Taylor expansion of $f$ near $a_n$ (think that $x^* \approx a_n + \Delta a$),

$$f(a_n + \Delta a) \approx T_f(\Delta a) := f(a_n) + f'(a_n)\Delta a + \frac{1}{2}f''(a_n)(\Delta a)^2$$

Use this quadratic functions as approximants to $f$ in a neighbourhood of $a_n$ and find their minimum points (take into account that $f''(x) > 0$):

$$0 = \frac{dT_f(\Delta a)}{d\Delta a} = f'(a_n) + f''(a_n)\Delta a \Rightarrow \Delta a = -\frac{f'(a_n)}{f''(a_n)}$$

Incrementing $a_n$ by this $\Delta a$ yields a point closer to $x^*$:

$$a_{n+1} = a_n + \Delta a = a_n - \frac{f'(a_n)}{f''(a_n)}$$

It is proved that for the chosen class of $f$, $a_n \to x^*$ as $n \to \infty$.

## 3.4 Levenberg-Marquardt algorithm

The application of LMA is the **least-squares curve fitting problem**: given a set $(x_i, y_i)_{i=1}^m$, find the parameters $\beta$ (column vector) of the model curve

$f(x, \beta)$ so that the sum of the squares of the deviations $S(\beta)$ is minimized:

$$arg \min_{\beta} S(\beta) \equiv arg \min_{\beta} \sum_{i=1}^{m} [y_i - f(x_i, \beta)]^2$$

Start with an initial guess for $\beta$. In each iteration step, the parameter vector $\beta$ is replaced by a new estimate $\beta + \Delta\beta$. To determine $\Delta\beta$, the function $f(x_i, \beta + \Delta\beta)$ is approximated by its linearization:

$$f(x_i, \beta + \Delta\beta) \approx f(x_i, \beta) + J_i \Delta\beta, J_i = (\nabla_{x_i} f(\beta))^T$$

The sum $S(\beta)$ has its minimum at a zero gradient with respect to $\beta$. The above first-order approximation of $f(x_i, \beta + \Delta\beta)$ gives

$$S(\beta + \Delta\beta) \approx \sum_{i=1}^{m} [y_i - f(x_i, \beta) - J_i \Delta\beta]^2$$

or in vector notation,

$$S(\beta + \Delta\beta) \approx [y - f(\beta)]^T [y - f(\beta)] - 2[y - f(\beta)]^T J \Delta\beta + \Delta\beta^T J^T J \Delta\beta,$$

where $J$ is the Jacobian matrix, whose i-th row equals $J_i$ , and where $f(\beta)$ and $y$ are vectors with i-th component $f(x_i, \beta)$ and $y_i$ , respectively.

Taking the derivative of $S(\beta + \Delta\beta)$ with respect to $\Delta\beta$ and setting to zero gives

$$(J^T J)\Delta\beta = J^T [y - f(\beta)],$$

that is in fact a system of linear equations with respect to $\Delta\beta$.

The system may be replaced by

$$(J^T J + \lambda I)\Delta\beta = J^T [y - f(\beta)],$$

where **I** is the identity matrix, giving the increment $\Delta\beta$ to the estimated parameter vector $\beta$.

# 4. Results

All results of calcutions we can see at Figures 1.

| function name | method name | arg | nit | fval |
|---|---|---|---|---|
| **D_linear** | **Conjugate Gradient Descent** | [0.02827500644475886, 0.8196527910802767] | 3 | 2.419561 |
| | **Gradient Descent** | [0.028803538677273345, 0.8213947577182771] | 62 | 2.423613 |
| | **Levenberg-Marquardt algorithm** | [0.024987446073156807, 0.8223670879803155] | 30 | 2.421611 |
| | **Newton's method** | [0.028270324897951105, 0.8196474478765564] | 4 | 2.419561 |
| **D_rational** | **Conjugate Gradient Descent** | [0.8198442384048378, -0.0332645308552266] | 9 | 2.419716 |
| | **Gradient Descent** | [0.7897523756700222, -0.10005675731880896] | 13 | 2.719580 |
| | **Levenberg-Marquardt algorithm** | [0.8542103260654681, -0.01258092802532218] | 10 | 3.113180 |
| | **Newton's method** | [0.8226017164621675, -0.02672775149589726] | 3 | 2.422361 |

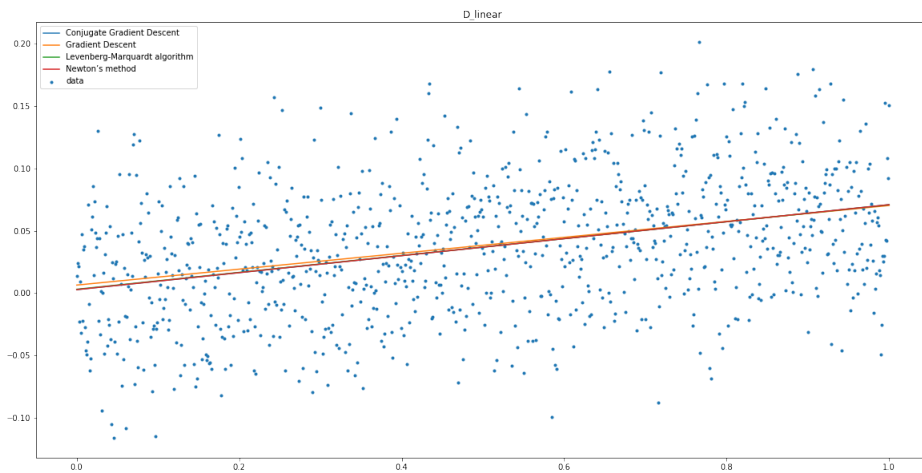**Figure 1:** Results of calcutions.



**Figure 2:** Linear approximant. Methods of Gradient Descent, Conjugate Gradient Descent, Newton's method and Levenberg-Marquardt algorithm.

All methods in case of liniar approximant (Figure 2)function converged to global minimum. Only one optimization method coverege to good approximation of rational approximant function (Figure 3).

Since the results of methods with linear function more representative, we can evaluate number of iterations. Most simplest method (Gradient Descent) worked for the greatest amount iterations. Newton and conjugate gradient
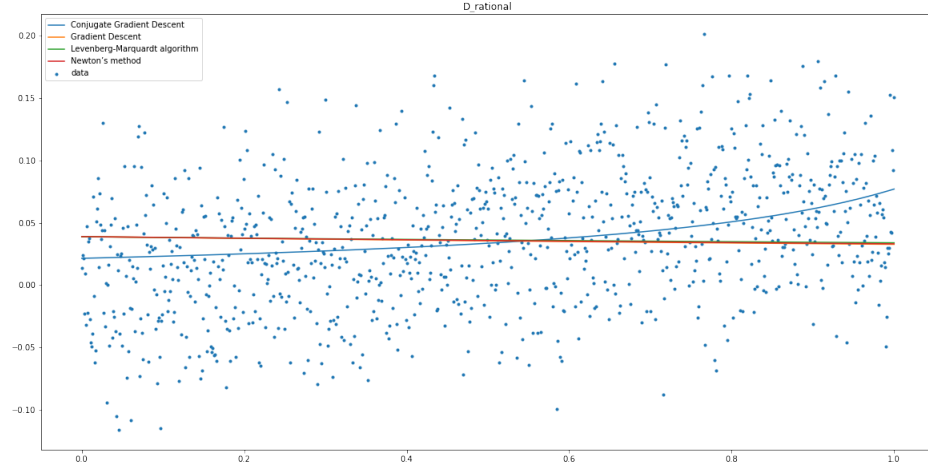
**Figure 3:** Rational approximant. Methods of Gradient Descent, Conjugate Gradient Descent, Newton's method and Levenberg-Marquardt algorithm.

methods are coverage with small number of iteration. Levenberg-Marquardt evaluate function more times than Newton or conjuate gradient method.

If we will compare this results with results of second lab, we see that first- and second-order methods give better result with less amount of iterations.