

Vladimir I. NORKIN¹

STOCHASTIC GENERALIZED GRADIENT METHODS FOR TRAINING NONCONVEX NONSMOOTH NEURAL NETWORKS ²

September 30, 2019

Abstract. The paper observes a similarity between the stochastic optimal control of discrete dynamical systems and the learning multilayer neural networks. It focuses on contemporary deep networks with nonconvex nonsmooth loss and activation functions. The machine learning problems are treated as nonconvex nonsmooth stochastic optimization problems. As a model of nonsmooth nonconvex dependences, the so-called generalized differentiable functions are used. The backpropagation method for calculating stochastic generalized gradients of the learning quality functional for such systems is substantiated basing on Hamilton-Pontryagin formalism. Stochastic generalized gradient learning algorithms are extended for training nonconvex nonsmooth neural networks. The performance of a stochastic generalized gradient algorithm is illustrated in the linear multiclass classification problem.

Keywords: Machine learning, deep learning, multilayer neural networks, nonsmooth nonconvex optimization, stochastic optimization, stochastic generalized gradient.

Introduction

The machine learning problem consists of the identification of parameters of a neural network model, e.g., neural weights, using a set of input-output observations. The training task is formulated as the task of minimizing some smooth loss functional (empirical risk), which measures the average forecast error of the neural network model.

Methods of training (identification) of large neural network models are discussed in many articles and monographs [1 – 11]. To train deep (i.e., multilayer) neural networks, the stochastic gradient method, and its modifications are mainly used [9 – 11], adopted from the theory of stochastic approximation [12] and stochastic programming [13 – 15], since only they are practically applicable for training such networks. The stochastic gradient of the risk functional is a random vector whose mathematical expectation approximates the gradient of the target functional, and the stochastic gradient descent method is an iterative method for changing the desired model parameters in the direction of the stochastic (anti-) gradient.

¹ V.M. Glushkov Institute of Cybernetics of the National Academy of Sciences of Ukraine, Kiev & Faculty of Applied Mathematics of the National Technical University of Ukraine “Igor Sikorsky Kyiv Polytechnic Institute”
Email: vladimir.norkin@gmail.com

² The work was partially supported by grant CPEA-LT-2016/10003 funded by the Norwegian Agency for International Cooperation and Quality Enhancement in Higher Education (Diku).

To solve smooth neural network training problems, the backpropagation of error technique, BackProp method, is widely used [1 – 8], i.e. a special method for calculating gradients of the target functional over various parameters. The history of discovery, development, and application of the BackProp method was studied in [8]. Nonsmooth machine learning tasks arise when using nonsmooth (module type) indicators of the quality of training, when applying nonsmooth regularizations, and also when using nonsmooth (for example, piecewise linear, ReLU, etc.) activation functions in multilayer neural [5, Section 6.3.1], [6, Section 3.3], [8]. Such functions give rise to essentially nonconvex nonsmooth functionals of the quality of learning, and the question arises of the convergence of the stochastic generalized gradient descent method in solving such problems. This problem has been relatively recently recognized and is already being considered in the literature [16 – 21]. However, it is usually assumed using the Clarke stochastic subgradients [22] of the optimized functional but the problem of their calculation for deep networks is not profoundly discussed.

In this paper, we extend the BackProp method to calculating stochastic gradients of nonconvex nonsmooth problems for training multilayer neural networks and formulate the method in terms of stochastic generalized gradients of the nonsmooth Hamilton-Pontryagin function. As a model of nonsmooth nonconvex dependencies, we use the so-called generalized differentiable functions [23, 24]. We also consider an important version of the BackProp method for training the so-called recurrent neural networks, i.e. networks with feedbacks and memory [5, Section 10].

In this paper, we show that the convergence of the stochastic generalized gradient method follows from the earlier results of the theory of nonconvex nonsmooth stochastic optimization [24 - 28]. In [16], as descent directions, Clarke’s stochastic generalized gradients of the optimized risk functional are used. However, the question remains what kind of objects the backpropagation method calculates in the case of a nonsmooth nonconvex functional and whether these objects can be used for optimization purposes. It may be supposed that the BackProp method calculates the (stochastic) Clarke’s subgradients of the optimized function but this takes place only in the case of the so-called subdifferentially regular functions [22, Section 2.3], that can be not the case. In this connection, it was proposed in [24], [27 – 30] to randomize the method of generalized gradient descent, namely, to calculate gradients not at the current iteration point but at a random close point, where the Lipschitz function is almost always differentiable.

Thus, although the problems of learning deep smooth neural networks have been studied for a long time, there are several new aspects related to the nonsmoothness of networks that still require discussion:

Nonconvexity and nonsmoothness of the optimized risk functional;

Methods for calculating stochastic (generalized) gradients for nonsmooth nonconvex networks;

A convergence of the stochastic gradient method in the nonconvex nonsmooth case;

The method parameters control and the method modifications for solving nonconvex nonsmooth problems;

Multixtremal nature of learning tasks;

The possibility of retraining a neural network model.

The purpose of this article is to apply results of the theory of nonconvex nonsmooth stochastic programming to machine learning problems and to discuss the peculiarities of the application of the stochastic (generalized) gradient method for these problems. In particular, we illustrate the application of the stochastic generalized gradients method to the problem of multiclass linear classifications.

1. Nonconvex nonsmooth learning problems and calculation of stochastic generalized gradients

Let us consider a standard neural network model. Let the network consists of m layers of neurons, each layer $i \in \{1, \dots, m\}$ has n_i neurons with numbers $j = 1, \dots, n_i$ and each of them has n_{i-1} inputs and one output. In the initial layer, there are n_1 neurons, each neuron of this layer has n_0 common inputs and one output. The outputs of the neurons of each layer go to the inputs of the neurons of the next layer. The output layer of the network may consist of one or more neurons.

In the theory of neural networks, the standard mathematical model of neuron (i, j) is some smooth activation function $g_i^j(x_i, w_{ij}, v_{ij})$ (e.g., the logistic sigmoid, the hyperbolic tangent, and etc. [5, Section 6.3.2; 6]). which expresses the dependence of the output signal $x_{(i+1)j}$ of neuron

$$(i, j) \text{ on the input signal } x_i, \text{ for example, } x_{(i+1)j} = g_i^j(x_i, w_{ij}, v_{ij}) = \left(1 + \exp\left\{-\langle x_i, w_{ij} \rangle - v_{ij}\right\}\right)^{-1},$$

where $x_i \in \mathbf{R}^{n_{i-1}}$ is a common input of all neurons in layer i ; $w_{ij} \in \mathbf{R}^{n_{i-1}}$ and $v_{ij} \in (-\infty, +\infty)$ are the individual weight vector and the activation level of neuron $j \in \{1, \dots, n_i\}$ in layer i ; the expression $\langle x_i, w_{ij} \rangle$ denotes the scalar product of vectors x_i and w_{ij} . The weights w_{ij} and thresholds v_{ij} may satisfy constraints $w_{ij} \in W_{ij}$, $v_{ij} \in V_{ij}$. Here notation like \mathbf{R}^n is used for n -dimensional Euclidian vector space.

Nonsmooth machine learning tasks arise when using nonsmooth indicators of the quality of learning, when applying nonsmooth regularization functions, and when using nonsmooth (for

example, piecewise linear) activation functions in multilayer neural networks, for example, $g_i^j(x_i, w_{ij}, v_{ij}) = \max\{-1, \min\{1, \langle x_i, w_{ij} \rangle + v_{ij}\}\}$ [5, Section 6.3.3].

For example, piecewise linear activation functions are essentially used in a dynamic brain model with positive BSB feedbacks (brain-state-in-box) [31, section 14.10, p. 884].

In [5], the problem of non-differentiability is informally discussed, caused by a nonsmooth activation function, e.g. the function of linear rectification (the positive part of the argument) $g(z) = \max\{0, z\}$ and its generalizations $g(z) = \max\{\alpha z, \beta z\}$, $g(z) = \max_{i \in I} z_i$, and others [5, Section 6.3, p. 169; Subsection 6.3.1, p. 170; Section 6.6, p. 197]. The use of piecewise linear activation functions instead of the sigmoidal functions significantly improved the quality of direct neural networks [3, 32].

Note that activation functions themselves can be random, for example, neurons can accidentally fall into the so-called sleep (drop out [5, 6]) state, i.e. produce a zero output signal: $g_i^j(x_i, w_{ij}, v_{ij}, \omega_{ij}) = \omega_{ij} \cdot g_i^j(x_i, w_{ij}, v_{ij})$, where ω_{ij} is an additional random parameter taking values 1 or 0 with probabilities p_{ij} and $1 - p_{ij}$. We assume that the random parameters $\{\omega_{ij}\}$ are independent and combined into a vector $\omega = \{\omega_{ij}\}$ that takes values from a finite set Ω .

In what follows, we assume that the activation functions $g_i^j(x_i, w_{ij}, v_{ij}, \omega_{ij})$ of neurons $j = 1, \dots, n_i$ in each layer i for any fixed value of ω_{ij} are generalized differentiable over their variables (x_i, w_{ij}, v_{ij}) in the sense of the following definition, which covers all practical examples.

Definition 1 [23, 24, 33]. *A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^1$ is called generalized differentiable at point $\bar{z} \in \mathbb{R}^n$, if in some ε -neighborhood $\{z \in \mathbb{R}^n : \|z - \bar{z}\| < \varepsilon\}$ of the point \bar{z} it is defined an upper semicontinuous at \bar{z} multivalued mapping $\partial f(\cdot)$ with convex compact values $\partial f(z)$ and such that the following expansion holds true:*

$$f(z) = f(\bar{z}) + \langle d, z - \bar{z} \rangle + o(\bar{z}, z, d), \quad (1)$$

where $d \in \partial f(z)$, $\langle \cdot, \cdot \rangle$ denotes a scalar product of two vectors, and the remainder term $o(\bar{z}, z, d)$ satisfies the condition: $\lim_{k \rightarrow \infty} o(\bar{z}, z^k, d^k) / \|z^k - \bar{z}\| = 0$ for all sequences $d^k \in \partial f(z^k)$, $z^k \rightarrow \bar{z}$ as $k \rightarrow \infty$. A function f is called generalized differentiable if it is generalized differentiable at each point $\bar{z} \in \mathbb{R}^n$; the mapping $\partial f(\cdot)$ is called the generalized gradient mapping of the function f ; the set $\partial f(z)$ is called a generalized gradient set of the function $f(\cdot)$ at point z ; vectors $d \in \partial f(z)$ are called generalized gradients of the function $f(\cdot)$ at point z .

Properties of generalized differentiable functions were studied in details in [23, 24, 33]. Any generalized differentiable function $f(z)$ is locally Lipschitzian and its Clark subdifferential $\partial_C f(z)$ [22] is the minimal (with respect to inclusion) generalized gradient mapping for $f(z)$, i.e. for all $z \in \mathbf{R}^n$ it takes place $\partial f(z) \supseteq \partial_C f(z)$ and for almost all $z \in \mathbf{R}^n$ it holds $\partial f(z) = \partial_C f(z)$ [24, Theorem 1.10]. The class of generalized differentiable functions contains continuously differentiable, convex, concave, weakly convex and weakly concave [26], semismooth [34], and some other piece-wise smooth functions [35], and is closed with respect to the operations of maximum, minimum, superposition and mathematical expectation (see [23, 24, 33, 36]).

Suppose there is a (training) set $\{(x_1^s \in \mathbf{R}^{n_0}, y_{m+1}^s \in \mathbf{R}^{n_m}), s = 1, \dots, S\}$ of observations of a network inputs-outputs. The standard training (identification) task for the network with the training quality criterion $\varphi(x_{m+1}^s, y_{m+1}^s)$ (for example, $\varphi(x_{m+1}^s, y_{m+1}^s) = \|x_{m+1}^s - y_{m+1}^s\|^2$) and regularization is as follows:

$$J(\{w_{ij}\}, \{v_{ij}\}) = \mathbf{E}_\omega \frac{1}{S} \sum_{s=1}^S \varphi(x_{m+1}^s, y_{m+1}^s) \rightarrow \min_{\{u_{ij} \in W_{ij}\}, \{v_{ij} \in V_{ij}\}}, \quad (2)$$

where $x_{m+1}^s \in \mathbf{R}^{n_k}$ is the vector of outputs of the last network layer for a training example s ; $y_{m+1}^s \in \mathbf{R}^{n_k}$ is a known, generally speaking, multidimensional vector of observations of the network outputs; $\|w_{ij}\|$ denotes a norm of the vector w_{ij} ; $\lambda > 0$, $\alpha \geq 1$; \mathbf{E}_ω is the mathematical expectation operator over ω ; the sequence of layers' outputs $\{x_i^s = (x_{i1}^s, \dots, x_{in_i}^s)^T, i = 2, \dots, m+1\}$ for a given first layer input $x_1^s \in \mathbf{R}^{n_0}$ is given by the relations

$$x_{(i+1)j}^s = g_i^j(x_i^s, w_{ij}, v_{ij}, \omega_{ij}), \quad j = 1, \dots, n_i; \quad i = 1, \dots, m. \quad (3)$$

The empirical criterion $J(\{w_{ij}\}, \{v_{ij}\})$ in (2) can be interpreted as the mathematical expectation of the random quantity $\varphi_{m+1}(x_{m+1}^s, y_{m+1}^s)$ over a discrete random variable $\theta = (s, \omega)$ that takes values in the set $\Theta = \{1, \dots, S\} \times \Omega$.

In machine learning, together with (2), regularized problems are considered [5, Ch. 7], [6, Section 4.1]:

$$J(\{w_{ij}\}, \{v_{ij}\}) = \mathbf{E}_\omega \frac{1}{S} \sum_{s=1}^S \varphi(x_{m+1}^s, y_{m+1}^s) + \sum_{i=1}^m \lambda_i \sum_{j=1}^{n_i} \left(\|w_{ij}\|^\alpha + \|v_{ij}\|^\alpha \right) \rightarrow \min_{\{u_{ij} \in W_{ij}\}, \{v_{ij} \in V_{ij}\}} \quad (4)$$

with smooth ($\alpha = 2$) and nonsmooth ($\alpha = 1$) regularizing terms $\|w_{ij}\|^\alpha, \|v_{ij}\|^\alpha$, and (penalty) parameters $\lambda_i \geq 0$ for layers $i = 1, \dots, m$. Regularization, on the one hand, improves the

conditionality of the problem, and on the other hand, suppresses the influence of excess neurons in the network.

Moreover, possibly, the training examples may contain not only the input and output of a network (for example, features and labels of objects) $\{(x_1^s, y_{m+1}^s), s = 1, \dots, S\}$, but also may include additional intermediate features $y_i^s \in \mathbb{R}^{n_i}$, $i \in I \subset \{2, \dots, m\}$, which can be used to improve the learning of the intermediate layers of the network, i.e. training examples may take the form of sequences $\{(x_1^s, \{y_i^s, i \in I\}, y_{m+1}^s), s = 1, \dots, S\}$. Then the criterion of the quality of training takes the following form:

$$J(\{w_{ij}\}, \{v_{ij}\}) = \mathbf{E}_\omega \frac{1}{S} \sum_{s=1}^S \sum_{i \in I} \varphi(x_i^s, y_i^s) + \sum_{i=1}^m \lambda_i \sum_{j=1}^{n_i} (\|w_{ij}\|^\alpha + \|v_{ij}\|^\alpha) + \mathbf{E}_\omega \frac{1}{S} \sum_{s=1}^S \varphi(x_{m+1}^s, y_{m+1}^s) \rightarrow \min_{\{u_{ij} \in W_{ij}\}, \{v_{ij} \in V_{ij}\}}. \quad (5)$$

That's why, next, we consider the following general network training task:

$$J(u) = \mathbf{E}_\theta \sum_{i=1}^m \varphi_i(x_i(\theta), u_i) + \mathbf{E}_\theta \varphi_{m+1}(x_{m+1}(\theta)) \rightarrow \min_{u \in U} \quad (6)$$

subject to constraints (satisfied for all values of the random parameter $\theta \in \Theta$):

$$x_{i+1}(\theta) = g_i(x_i(\theta), u_i, \theta) = \left\{ g_i^j(x_i(\theta), u_{ij}, \theta) \right\}_{j=1}^{n_i}, \quad i = 1, \dots, m; \quad x_1(\theta) \in \mathbb{R}^{n_0}. \quad (7)$$

Here $u = (u_1, \dots, u_m) \in \mathbb{R}^l$ ($l = \sum_{i=1}^m n_i$) is the vector of all adjusted parameters; $x_i = (x_{i1}, \dots, x_{in_i})^T$ is the input vector for neurons in layer i ; u_{ij} is the vector of the adjusted parameters of neuron (i, j) ; $u_i = \{u_{ij}\}_{j=1}^{n_i-1}$ is the vector of the adjusted parameters of all neurons in layer i ; g_i^j is the activation function of neuron j in layer i ; $g_i = \{g_i^j\}_{j=1}^{n_i}$ is the vector activation function of the neurons in layer i ; $x_1(\theta) \in \mathbb{R}^{n_0}$ is a random vector of input signals to the network; θ is a random vector parameter that defines the distribution of input signals and influences on the propagation of signals through the network; \mathbf{E}_θ denotes the sign of the mathematical expectation over θ .

In problems (2) - (5) $u_{ij} = (w_{ij}, v_{ij})$, the role of the random parameter θ is played by the random pair $\theta = (s, \omega)$; here $x_1(\theta) = x_1^s$, $\varphi_{m+1}(x_{m+1}(\theta)) = \varphi(x_{m+1}(s, \omega), y_{m+1}^s)$, and

$$\varphi_i(x_i(\theta), u_i, \theta) = \begin{cases} \varphi(x_i(s, \omega), y_i^s) + \lambda_i \sum_{j=1}^{n_i} (\|w_{ij}\|^\alpha + \|v_{ij}\|^\alpha) & i \in I \\ \lambda_i \sum_{j=1}^{n_i} (\|w_{ij}\|^\alpha + \|v_{ij}\|^\alpha), & i \notin I. \end{cases}$$

We make the following assumptions.

Assumptions. Suppose that in problem (6), (7) the functions $\varphi_i(x_i, u_i)$, $g_i^j(x_i, u_{ij}, \theta)$, and $\varphi_{m+1}(x_{m+1})$ are generalized-differentiable over the totality of their arguments, respectively, over

(x_i, u_i) , (x_i, u_{ij}) , and x_{m+1} (under fixed θ). Here, the activation function $g_i^j(x_i, u_{ij}, \theta)$ can be of a general form, i.e. optionally, the function g_i^j may depend not on all elements of the vector x_i and the dimension of the vector of the adjustable parameters u_{ij} may not coincide with the dimension of the vector of inputs x_i . The random parameter $\theta \in \Theta$ is a random variable defined on some probability space.

Note that in the literature (see, for example, [16 – 21]), for the purpose of training neural networks, it is proposed to use (stochastic) Clarke subgradients of the risk functional $J(u)$ but these subgradients are relatively simple to calculate only for subdifferentially regular Lipschitz functions [7, §2.3, §2.7], and for general nonconvex nonsmooth functions, their calculation may be a problem.

The next theorem exploits a similarity between optimal control problems for discrete dynamical systems and multilayer neural networks, and formalizes a method for calculating stochastic generalized gradients in the problem of training a nonconvex nonsmooth neural network. It extends the well-known method of “backpropagation of the error” (BackProp) [1 - 5] to nonconvex nonsmooth learning problems.

First, we introduce the following notation. For arbitrary generalized differentiable (by the totality of variables) vector functions $g_i(x, u) \in \mathbb{R}^{n_i}$ with arguments $x = (x^1, \dots, x^n)^T \in \mathbb{R}^n$, $u = (u^1, \dots, u^l)^T \in \mathbb{R}^l$, we denote the matrices:

$$g_{ix} = \begin{pmatrix} g_{ix^1}^1 & \dots & g_{ix^n}^1 \\ \dots & \dots & \dots \\ g_{ix^1}^{n_i} & \dots & g_{ix^n}^{n_i} \end{pmatrix} = \begin{pmatrix} g_{ix}^1 \\ \dots \\ g_{ix}^{n_i} \end{pmatrix}, \quad g_{iu} = \begin{pmatrix} g_{iu^1}^1 & \dots & g_{iu^l}^1 \\ \dots & \dots & \dots \\ g_{iu^1}^{n_i} & \dots & g_{iu^l}^{n_i} \end{pmatrix} = \begin{pmatrix} g_{iu}^1 \\ \dots \\ g_{iu}^{n_i} \end{pmatrix};$$

and for arbitrary generalized differentiable (over the totality of arguments) scalar functions $f_i(x, u)$, $x \in \mathbb{R}^n$, $u \in \mathbb{R}^l$, and $\varphi_{m+1}(x)$, $x \in \mathbb{R}^n$, let us introduce vectors:

$$(f_{ix})^T = (f_{ix^1}, \dots, f_{ix^n}), \quad (f_{iu})^T = (f_{iu^1}, \dots, f_{iu^l}), \quad (\varphi_{kx})^T = (\varphi_{kx^1}, \dots, \varphi_{kx^n}),$$

where $(f_{ix}, f_{iu})^T$, $(g_{ix}^j, g_{iu}^j)^T$ are some generalized gradients of the functions $f_i(\cdot, \cdot)$, $g_i^j(\cdot, \cdot, \theta)$; $\varphi_{(m+1)x}(\cdot)$ is some generalized gradient of the function φ_{m+1} ; the expression $(\cdot)^T$ denotes the transposition of the matrix (\cdot) .

Theorem 1. *Under the assumptions made, the objective function $J(u)$ of problem (6), (7) is generalized differentiable with respect to variables $u = (u_1 \in \mathbb{R}^{n_1}, \dots, u_m \in \mathbb{R}^{n_m})$, and the vector*

$$h_u(u, \theta) = \left(h_{1u_1}(x_1(\theta), \psi_1(\theta), u_1), h_{2u_2}(x_2(\theta), \psi_2(\theta), u_2), \dots \right. \\ \left. \dots, h_{mu_m}(x_{mm}(\theta), \psi_m(\theta), u_m) \right)^T, \quad (8)$$

is a stochastic generalized gradient of the function $J(u)$ at a given point u , i.e. $\mathbf{E}_\theta h_u(u, \theta) \in \partial J(u)$, where $h_i(x_i, \psi_i, u_i) = f_i(x_i, u_i) + g_i^T(x_i, u_i) \cdot \psi_i$, $i=1, \dots, m$, is a discrete (over i) Hamilton-Pontryagin function; the vector

$$h_{iu_i}(x_i(\theta), \psi_i(\theta), u_i) = \left(h_{iu_i^1}(x_i(\theta), \psi_i(\theta), u_i), \dots, h_{iu_i^{n_i}}(x_i(\theta), \psi_i(\theta), u_i) \right)^T = \\ = \left(f_{iu_i^1}(x_i(\theta), u_i) + g_{iu_i^1}^T(x_i(\theta), u_i) \cdot \psi_i(\theta), \dots, f_{iu_i^{n_i}}(x_i(\theta), u_i) + g_{iu_i^{n_i}}^T(x_i(\theta), u_i) \cdot \psi_i(\theta) \right)^T \in \mathbf{R}^{n_i} \quad (9)$$

is the u_i -component of a generalized gradient of the function $h_i(\cdot, \psi_i, \cdot)$, $i=1, \dots, m$; $x(\theta) = (x_1(\theta), \dots, x_{m+1}(\theta))$ is a discrete random trajectory of process (7), corresponding to the given sequence of parameters $(u_1, \dots, u_m) = u$ and the random initial data $x_1(\theta) \in \mathbf{R}^{n_0}$. Here, the random sequence of auxiliary (conjugate) vector functions $(\psi_1(\theta), \dots, \psi_m(\theta)) = \psi(\theta)$ is determined through the backpropagation equations (adopted from the Pontryagin maximum principle): $\psi_m(\theta) = \varphi_{(m+1)x_{m+1}}(x_{m+1}(\theta))$,

$$\psi_{i-1}(\theta) = h_{ix_i}(x_i(\theta), \psi_i(\theta), u_i) = f_{ix_i}(x_i(\theta), u_i) + g_{ix_i}^T(x_i(\theta), u_i) \cdot \psi_i(\theta), \quad i = m, m-1, \dots, 2; \quad (10)$$

$(f_{ix}(x_i(\theta), u_i), f_{iu}(x_i(\theta), u_i))^T$, $(g_{ix}^j(x_i(\theta), u_i), g_{iu}^j(x_i(\theta), u_i))^T$ are some generalized gradients of the functions $f_i(\cdot, \cdot)$, $g_i^j(\cdot, \cdot, \theta)$ at the point $(x_i(\theta), u_i)$ and $\varphi_{(m+1)x}(x_{m+1}(\theta))$ is some generalized gradient of the function φ_{m+1} at the point $x_{m+1}(\theta)$, which are used in (9), (10).

Proof. Note that process (7) can be formally treated as a stochastic dynamic system in discrete time $i=1, \dots, m+1$ with states x_i , control parameters u_i , a given initial state $x_1(\theta)$, and the optimality criterion (6). The stochasticity of system (7) is generated by the random input $x_1(\theta)$, a random mechanism of dropping out of neurons, and, possibly, by other factors. So the present theorem is a particular case of a similar Theorem 7 from [37], established for discrete stochastic dynamic system. Similar to the proofs of Theorems 6, 7 from [37] and using relations (7), the vectors x_i , $i=2, \dots, m+1$, can be sequentially excluded from the formulation of the optimization problem (6). Then under the sign of summation in (6) there remains some complex composite function

$$f(u, \theta) = \sum_{i=1}^m \varphi_i(\tilde{x}_i(u_1, \dots, u_{i-1}, \theta), u_i) + \varphi_{m+1}(\tilde{x}_{m+1}(u_1, \dots, u_m, \theta)),$$

which depends on optimization variables u and where $\tilde{x}_i(u_1, \dots, u_{i-1}, \theta)$ are complex compound functions of their arguments. And since the class of generalized differentiable functions is closed with respect to compositions, then under the made assumptions this function $f(u, \theta)$ becomes generalized differentiable with respect to u for each $\theta \in \Theta$. The mathematical expectation \mathbf{E}_θ (in this case, summation) does not move out from the class of generalized differentiable functions, therefore the function $J(u) = \mathbf{E}_\theta f(u, \theta)$ is also generalized differentiable with respect to u . Now, similar to the proofs of Theorems 6 - 8 from [37], applying the rules of differentiation of the sum, the chain rule of differentiation of complex generalized differentiable functions (which are analogues to the rules of differentiation of smooth and convex functions) [23, 24, 33], and introducing auxiliary variables (10), after some algebraic manipulations (see [37, Theorem 6]), we obtain formula (8) for stochastic generalized gradients $h_u(u, \theta)$ of the function $J(u)$. The proof is complete.

Formulas (8), (9) use procedure (7) of direct calculation of the trajectory of motion x_1, x_2, \dots, x_{m+1} and the reverse calculation (10) of auxiliary conjugate variables ψ_m, \dots, ψ_1 , essentially adopted from the Pontryagin maximum principle [38, 39]. Thus, the vector $h_u(u, \theta)$ is a stochastic generalized gradient of the function $J(u)$ such that $\mathbf{E}_\theta h_u(u, \theta) \in \partial J(u)$, and it can be used in stochastic generalized gradient methods for minimizing the nonsmooth functional $J(u)$. Note that the set of generalized gradients $\partial J(u)$ may be wider than the Clarke's subdifferential $\partial_c J(u)$, and it may turn out that $\mathbf{E}_\theta h_u(u, \theta) \notin \partial_c J(u)$. To take this possibility into account, in the next section we modify the standard stochastic gradient descent method by introducing an artificial randomization into it.

Neural networks may have a complex and heterogeneous structure. However, introducing additional dummy neurons, such networks can be reduced to a canonical multilayer form and for them, one can apply the formulas of Theorem 1.

Let us consider an important special case of networks, in which the adjusted parameters are the same for each layer. For example, a network can consist of identical neurons, or some identical neurons are added to each layer in the already-trained network, or the network consists of duplicates of the same layer. Then, similarly to (6), (7), the training task consists in solving the following problem:

$$J(u) = \mathbf{E}_\theta \sum_{i=1}^m \varphi_i(x_i(\theta), u) + \mathbf{E}_\theta \varphi_{m+1}(x_{m+1}(\theta)) \rightarrow \min_{u=(u_1, \dots, u_l) \in U}, \quad (11)$$

$$x_{i+1}(\theta) = g_i(x_i(\theta), u, \theta) = \left\{ g_i^j(x_i(\theta), u, \theta) \right\}_{j=1}^{n_i}, \quad i = 1, \dots, m; \quad x_1(\theta) \in \mathbf{R}^{n_0}. \quad (12)$$

Theorem 2. *Under the assumptions made, the objective function $J(u)$ of problem (11), (12) is generalized differentiable with respect to variables $u = (u_1, \dots, u_l)$, and the vector*

$$\left(\sum_{i=1}^m h_i(x_i(\theta), \psi_i(\theta), u) \right)_u = \left(\sum_{i=1}^m h_{iu_1}(x_i(\theta), \psi_i(\theta), u), \sum_{i=1}^m h_{iu_2}(x_i(\theta), \psi_i(\theta), u), \dots, \sum_{i=1}^m h_{iu_l}(x_i(\theta), \psi_i(\theta), u) \right)^T, \quad (13)$$

is a stochastic generalized gradient of the function $J(u)$ at point u , i.e. $\mathbf{E}_\theta \left(\sum_{i=1}^m h_i(x_i(\theta), \psi_i(\theta), u) \right)_u \in \partial J(u)$, where $h_i(x_i, \psi_i, u) = f_i(x_i, u) + g_i^T(x_i, u) \cdot \psi_i$, $i = 1, \dots, m$, is a discrete (over i) Hamilton-Pontryagin function; $x(\theta) = (x_1(\theta), \dots, x_{m+1}(\theta))$ is a discrete random trajectory of process (12) that corresponds to the vector parameter u and the random initial data $x_1(\theta) \in \mathbf{R}^{n_0}$. Here, the random sequence of auxiliary (conjugate) vector functions $(\psi_1, \dots, \psi_m) = \psi$ is determined through the backpropagation equations:

$$\psi_{i-1}(\theta) = h_{ix_i}(x_i(\theta), \psi_i(\theta), u) = f_{ix_i}(x_i(\theta), u) + g_{ix_i}^T(x_i(\theta), u) \cdot \psi_i(\theta),$$

$$i = m, m-1, \dots, 2, \quad \psi_m(\theta) = \varphi_{(m+1)x_{m+1}}(x_{m+1}(\theta)).$$

This theorem is an analogue of Theorem 6 from [37] and is proved similarly.

2. The method of stochastic generalized gradient descent and its variants

The stochastic gradient descent method is the main method for training deep neural networks, firstly, because of enormous dimensions of such networks and, secondly, due to the regularizing properties of the method. The properties of the stochastic gradient method and its modifications were studied in details in cases of smooth and convex optimized functions [9 - 11]. In this article, we analyze this method as applied to nonsmooth nonconvex problems of machine learning.

In the notation and assumptions of the previous section, the learning task (6), (7) consists in optimizing the complex function of the mathematical expectation $J(u) = \mathbf{E}_\theta f(u, \theta)$, where

$$f(u, \theta) = \sum_{i=1}^m \varphi_i(\tilde{x}_i(u_1, \dots, u_{i-1}, \theta), u_i) + \varphi_{m+1}(\tilde{x}_{m+1}(u_1, \dots, u_m, \theta)).$$

As indicated in Theorem 1, we can assume that the integrand $f(u, \theta)$ is generalized differentiable with respect to variables $u = (u_1, \dots, u_m)$ for each fixed θ ; then the function $J(u) = \mathbf{E}_\theta f(u, \theta)$ is also generalized differentiable with respect to u , and its stochastic generalized gradients can be calculated by formulas (8), (10).

We now consider the (randomized) method of stochastic generalized gradient descent to minimize on a convex set U a generalized differentiable mathematical expectation function,

$$J(u) = \mathbf{E}_\theta f(u, \theta) \rightarrow \min_{u \in U} \quad (14)$$

The method has the form:

$$u^{k+1} \in \Pi_U \left(u^k - \rho_k d^k \right) = \arg \min_{u \in U} \left(\langle d^k, u - u^k \rangle + \frac{1}{2\rho_k} \|u - u^k\|^2 \right), \quad (15)$$

$$d^k = d(\tilde{u}^k, \theta^k) \in \partial_u f(\tilde{u}^k, \theta^k), \quad \|\tilde{u}^k - u^k\| \leq \delta_k, \quad k = 0, 1, \dots, \quad (16)$$

where k denotes the iteration number; $\Pi_U(\cdot)$ is the projection operator on a convex feasible set U ; $d(u, \theta)$ is a (u, θ) -measurable section (see [36] for details) of a generalized gradient map $\partial_u f(u, \theta)$ of the generalized differentiable random function $f(\cdot, \theta)$; $\{\theta^k, k = 0, 1, \dots\}$ are independent identically distributed observations of a random variable θ ; points $\{\tilde{u}^k\}$ are randomly sampled from the sets $\{u : \|u - u^k\| \leq \delta_k\}$; non-negative quantities ρ_k, δ_k can depend on $\{u^0, \dots, u^{k-1}\}$ but must be measurable with respect to the σ -algebra $\sigma\{u^0, \dots, u^{k-1}\}$ and with probability one must satisfy the conditions:

$$\lim_{k \rightarrow \infty} \delta_k = \lim_{k \rightarrow \infty} \rho_k = 0, \quad \sum_{k=0}^{\infty} \rho_k = +\infty, \quad \sum_{k=0}^{\infty} \rho_k^2 < +\infty. \quad (17)$$

Here, the randomization consists in calculating the current generalized gradient $d(\tilde{u}^k, \theta^k) \in \partial_u f(\tilde{u}^k, \theta^k)$ not at the current iteration u^k but at a random close point \tilde{u}^k , where $\|\tilde{u}^k - u^k\| \leq \delta_k$. A similar idea of randomization of the generalized gradient descent algorithm is used in [24, 29, 30]. If $\delta_k \equiv 0$, then method (15) - (17) turns into the usual method of stochastic generalized gradient descent. If $U = \mathbb{R}^n$, then the projection operation $\Pi_U(\cdot)$ in method (15) is absent. We denote $N_U(u)$ the cone of normals to the set U at a given point u .

Theorem 3 (convergence with probability 1 of the non-convex randomized method of stochastic generalized gradients [28, section 2], [29, Theorem 5]. *Under conditions (17), for almost all trajectories of the process (15), (16), the minimum (in terms of the function J) limit points of the sequence $\{u^k\}$ belong to the set $U_C^* = \{u \in U : 0 \in \partial_C J(u) + N_U(u)\}$ of points satisfying the Clarke's necessary optimality conditions [22], and all limit points of the numerical sequence $\{J(u^k)\}$ constitute an interval in the set $J_C^* = J(U_C^*)$. If the set F_C^* does not contain intervals (for*

example, F_C^* is finite or countable), then all limit points of the sequence $\{u^k\}$ belong to the connected subset of U_C^* , and there is a limit $\lim_{k \rightarrow \infty} J(u^k) \in J_C^*$.

If in algorithm (15) - (17) all $\delta_k \equiv 0$, then the statement of Theorem 3 holds for the set $U^* = \{u \in U : 0 \in \partial J(u) + N_U(u)\}$ [27, Theorem 5.1]. The convergence of the randomized method of generalized gradient descent for deterministic problems is shown in [27, Remark 4.2], [28], [40, Remark 4.2]. Similar results on the convergence of the (nonrandomized) method of stochastic generalized gradient descent using generalized Clarke gradients for piecewise smooth Lipschitz (Whitney stratifiable) functions were recently obtained (by another, differential inclusion method) in [16].

Many other stochastic methods of convex optimization were considered in [6, 7, 9, 11] (methods with averaging the trajectory, averaging generalized gradients, ravine step methods, heavy ball, and others). In [24], these methods were extended to problems of non-convex nonsmooth stochastic optimization, in particular, the stochastic ravine step method for solving problem (14) has the form:

$$\begin{aligned} y^0 &= u^0, \\ y^{k+1} &= \Pi_U(u^k - \rho_k d^k), \\ u^{k+1} &= y^{k+1} + \lambda_k (y^{k+1} - y^k), \quad k = 0, 1, \dots, \end{aligned}$$

where $d^k = d(u^k, \theta^k) \in \partial_u f(u^k, \theta^k)$; $\{\theta^k, k = 0, 1, \dots\}$ are independent identically distributed observations of a random variable θ ; numbers ρ_k, λ_k satisfy conditions:

$$0 \leq \rho_{k+1} \leq \rho_k \leq \rho, \quad \sum_{k=0}^{+\infty} \rho_k = +\infty, \quad \sum_{k=0}^{+\infty} \rho_k^2 < +\infty; \quad 0 \leq \lambda_k \leq \lambda < 1.$$

This is a stochastic analog of the deterministic ravine method [41], which, when optimizing smooth convex functions, has a high convergence rate of the order $O(1/k^2)$ [42]. A geometric interpretation of the method shows that it moves (descends) along the ‘‘gullies’’ of the minimized function or the boundary of the feasible region. In [6, 7, 11], [43 – 45], adaptive step adjustments in stochastic gradient optimization methods are considered.

The randomized method also admits the following interpretation [24]. We introduce the so-called smoothed functions:

$$J_k(u) = \frac{1}{V_{\delta_k}} \int_{\|\tilde{u}-u\| \leq \delta_k} J(\tilde{u}) d\tilde{u}, \quad (18)$$

где V_{δ_k} - объем δ_k -окрестности нуля. Если ввести случайный вектор \tilde{u}^k , равномерно распределенный в δ_k -окрестности точки u , то сглаженную функцию $J_k(u)$ и ее градиент $\nabla J_k(u)$ можно представить соответственно в виде $J_k(u) = \tilde{E}J(\tilde{u}^k)$ и $\nabla J_k(u) = \tilde{E}\partial J(\tilde{u}^k)$, где \tilde{E} обозначает математическое ожидание по \tilde{u}^k , $\tilde{E}\partial J(\tilde{u}^k)$ обозначает математическое ожидание случайного многозначного отображения $(u, \tilde{u}^k) \rightarrow \partial J(\tilde{u}^k)$.

Thus, the randomization in method (15) - (17) plays a threefold role: on the one hand, it allows us to narrow the convergence set of the generalized gradient method to the set $U_c^* \subseteq U^*$, and on the other hand, it provides to method (15) - (17) some global properties due to the fact that it minimizes the sequence of smoothed functions $J_k(u)$. Besides, the randomization prevents the method to sticking at critical points that are not local minima. In case $\delta_k = \delta$ the randomized generalized gradient method (15) - (17) becomes a stochastic gradient method for minimizing the same non-changing smoothed function $J_k(u)$. To strengthen global properties of method (15) - (17) it is possible to use the estimate of the gradient of the smoothed function (18) using several independent realizations of a random point \tilde{u} from δ_k -vicinity of the current point u^k .

Note that the use of smoothed functions (18) with the norm $\|u\| = \max_i |u_i|$ allows us to construct stochastic finite-difference minimization methods, i.e. methods without calculating subgradients [24, 46]. A review of randomized gradient algorithms for optimizing smooth functions is also available in [47].

3. The method of stochastic generalized gradients in the linear classification problem

As an illustration of the application of the method of stochastic generalized gradient descent for training nonsmooth neural networks, we consider the classical problem of multiple classifications of objects based on a family of precedents [4, 7], [48 – 50]. For classification purposes, the methods of nonsmooth optimization were used in [51 – 53], and the method of stochastic generalized gradients was used in [7]. The linear classifier can be used as an output layer of a deep neural network in transfer learning models [5, section 15.2; 6]. In this section, we use the stochastic generalized gradient method to solve the multiclass classification problem with a new learning/loss function.

The abstract setting of the problem has the following form. Let a graph be given and fixed, which is used for the representation of classification objects. Such objects will be the subsets of edges (or vertices, or both) of this graph. Each object can have several forms, i.e. it can be represented by different collections of edges. In this model, edges can be interpreted as features of

the objects, and each instance of an object can be interpreted as a set of specific features. Obviously, instead of edges, graph vertices can be used to represent objects. Let the edges of the graph be numbered. We will encode objects with (0,1)-strings with the number of elements equal to the number of edges. Moreover, 0 or 1 at a certain position in the string means that this edge is not used or is used to represent the object. Thus, a set of training examples is encoded by a set of (0,1)-strings with labels of belonging to one or another class. A graphical interpretation is used to visualize objects. Representation of objects in the form of (0,1)-lines can be noisy, i.e. instead of 0 and 1 in the string at their positions can be random numbers with values from zero to one. By noising, we can expand the original training set arbitrarily largely. Examples may include unclassified objects as a separate class, therefore, without loss of generality, we assume that all examples belong to one or another class. In numerical experiments, noisy codes of stylized numbers and letters were used.

Let us denote $x = (x_0 = 1, x_1, \dots, x_n)$ the vector of features of a presented object, where for convenience the feature $x_0 = 1$ is introduced, and $x_i \in \{0, 1\}$ or $x_i \in [0, 1]$. We introduce a numerical (m, n) -matrix $W = \{w_{ij}\}$, where m is the number of classes, with rows $w_i = (w_{i0}, w_{i1}, \dots, w_{in})$, which must be determined based on training examples. The linear classification method $k(x)$ determines the class of the object x under the found weights W according to the formula:

$$k(x) \in \arg \max_{1 \leq i \leq m} \langle x, w_i \rangle.$$

The matrix of weights W can be found in many ways. For example, it can be found by minimizing the smooth convex cross entropy function:

$$F(W) = - \sum_{s=1}^S \ln \frac{\exp \langle w_{i_s}, x^s \rangle}{\sum_{s'=1}^S \exp \langle w_i, x^{s'} \rangle} \rightarrow \min_W, \quad (19)$$

where S is the number of training examples; $x^s = (x_0^s, x_1^s, \dots, x_n^s)$ denotes the feature vector of the example s ; i_s is the known class number (label) for the example s . For comparison, we also use the classic learning function:

$$\Omega(W) = \sum_{s=1}^S \left(\max_{1 \leq i \leq m, i \neq i_s} \langle w_i, x^s \rangle - \langle w_{i_s}, x^s \rangle \right) \rightarrow \min_W. \quad (20)$$

This article also proposes to use the following nonsmooth convex learning function:

$$\Phi(W) = \sum_{s=1}^S \left(\max_{1 \leq i \leq m} \langle w_i, x^s \rangle - \langle w_{i_s}, x^s \rangle \right) \rightarrow \min_W. \quad (21)$$

Compared to $\Omega(W)$, the function $\Phi(W)$ is advantageous in that its optimal (zero) value is known (for linearly separable data). The advantage of the function $\Phi(W)$ compared to $F(W)$ is in its simplicity. It's obvious that $\Phi(W) \geq 0$. If the classes are linearly separable, i.e. there exists a matrix W^* such that $i_s \in \arg \max_{1 \leq i \leq m} \langle w_i^*, x^s \rangle$ for all $s = 1, \dots, S$, then $\min_W \Phi(W) = 0$. Note that the set of minima of the function $\Phi(W)$ is not empty and includes non-trivial subsets of the space of all $(m \times n)$ -matrices, in particular, all matrices such that their rows coincide belong to the set of minima of the function $\Phi(W)$. Therefore, not all minima of the problem $\min_W \Phi(W)$ are suitable for solving the classification problem. The tasks of minimizing the functions $F(W)$, $\Phi(W)$, and $\Omega(W)$ can be solved by the method of stochastic subgradients. Denote $\Phi^s(W) = \max_{1 \leq i \leq m} \langle w_i, x^s \rangle - \langle w_{i_s}, x^s \rangle = \langle w_{i_s^*}, x^s \rangle - \langle w_{i_s}, x^s \rangle$, where the index i_s^* is such that $\max_{1 \leq i \leq m} \langle w_i, x^s \rangle = \langle w_{i_s^*}, x^s \rangle$. Subgradients of the function $\Phi^s(W)$ have the following components:

$$\Phi_{w_{ij}}^s(W) = \begin{cases} x_j^s, & \text{если } i = i_s^* \text{ и } k \neq i_s, \\ -x_j^s, & \text{если } i = i_s \text{ и } k \neq i_s^*, \\ 0, & \text{если } i = i_s^* = i_s \text{ или } (i \neq i_s^* \text{ и } i \neq i_s); \end{cases} \quad i = 1, \dots, m; \quad j = 0, 1, \dots, n.$$

If an index s (the number of a training example) is chosen randomly and equiprobably, then the vector with the components

$$\left(\Phi_{w_{10}}^s(W), \dots, \Phi_{w_{1n}}^s(W), \Phi_{w_{20}}^s(W), \dots, \Phi_{w_{2n}}^s(W), \dots, \Phi_{w_{m0}}^s(W), \dots, \Phi_{w_{mn}}^s(W) \right)^T = \Phi_w^s(W)$$

is called the stochastic subgradient of the function $\Phi(\cdot)$ at a point W .

The randomized method with averaging stochastic subgradients for solving problem (21) has the form:

$$w_{ij}^{k+1} = w_{ij}^k - \rho_k \cdot \frac{1}{L} \sum_{l=1}^L \Phi_{w_{ij}}^{s_l}(\tilde{W}^{kl}), \quad i = 1, \dots, m, \quad j = 0, 1, \dots, n, \quad (22)$$

where $k = 0, 1, \dots$ denotes the iteration number of the algorithm; $\{s^k, k = 0, 1, \dots\}$ are independently and equally likely taken numbers of training examples; $W^k = \{w_{ij}^k\}$ is $(m \times n)$ -matrix with components w_{ij}^k ; $\tilde{W}^{kl} = \{\tilde{w}_{ij}^{kl}\}$ is $(m \times n)$ -matrix with components \tilde{w}_{ij}^{kl} ; $\{\Phi_{w_{ij}}^{s_k}(\tilde{W}^{kl})\}$ are components of the the stochastic generalized gradient $\Phi_w^{s_k}(\tilde{W}^{kl})$ at a random point \tilde{W}^{kl} such that $\|W^k - \tilde{W}^{kl}\| \leq \delta_k$; W^0 is an initial set of weights; L is the number of subgradients averaged at each

iteration. In numerical experiments, it was set $\rho_k = \rho_0/k^{1/2}$, $\delta_k = \rho_0/k^{1/3}$. In the classical stochastic subgradient method, it is assumed that $\delta_k \equiv 0$, $L = 1$. When solving problems (19), (20) in method (22), instead of $\Phi_w^{s_k}$ the subgradients of the functions F and Ω are used.

Figures 1-6 show the results of training linear classifiers by the stochastic subgradient method using learning functions $\Phi(W)$, $\Omega(W)$, and $F(W)$. These figures show typical examples of performance of the constructed linear classifiers, i.e. they show the fractions of correctly recognized examples in the training and test samples, as a function of the number of iterations of the stochastic subgradient method. The results of the numerical experiments indicate that the linear classification with the function Φ is more effective than using classical functions F and Ω .

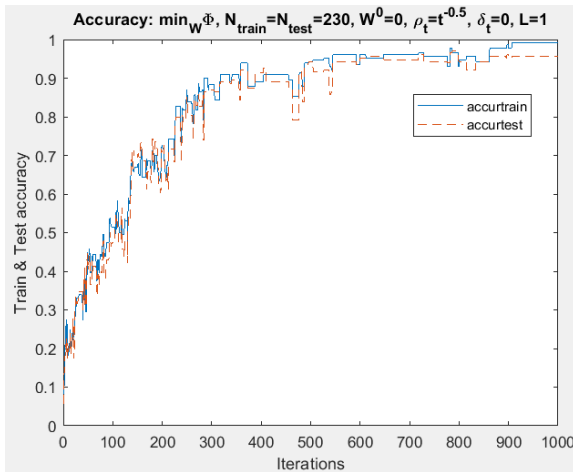


Fig. 1. The performance of the stochastic subgradient method on the function Φ under zero initial weights.

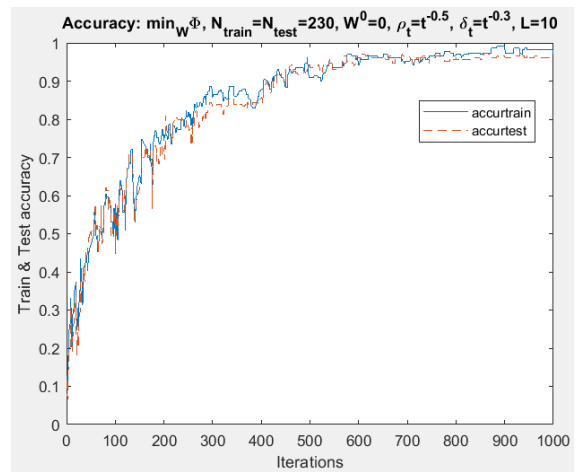


Fig.2. The performance of the randomized stochastic subgradient Method on the function Φ under zero initial weights.

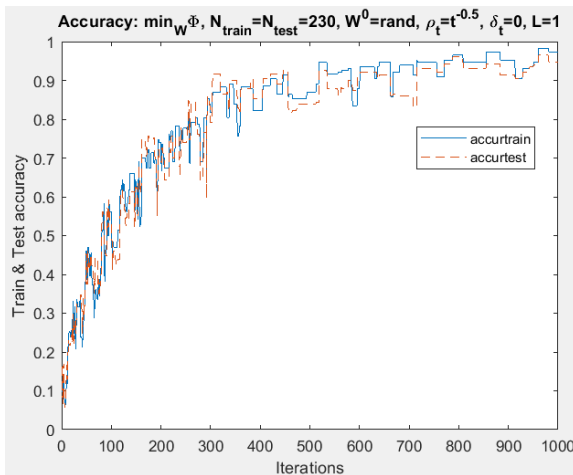


Fig. 3. The performance of the stochastic subgradient method on the function Φ under some random initial weights.

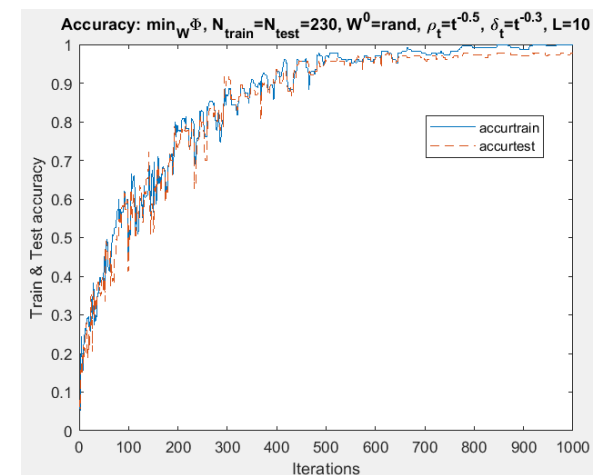


Fig. 4. The performance of the randomized stochastic subgradient method on the function Φ under random initial weights.

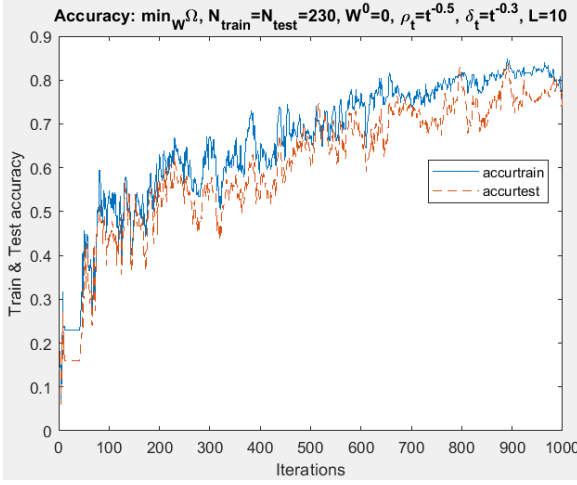


Fig. 5. The performance of the stochastic subgradient method on the function Ω under zero initial weights.

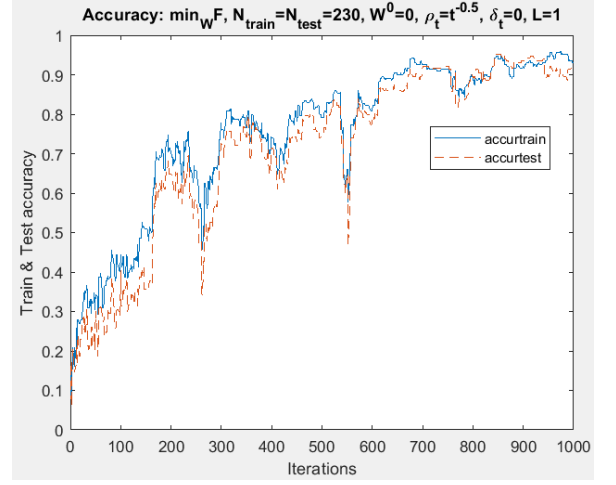


Fig. 6. The performance of the randomized stochastic subgradient method on the function F under zero initial weights.

Conclusions

In the present work, the following results were obtained:

The well-known method of backpropagation of errors is extended to nonconvex nonsmooth machine learning problems:

The randomized stochastic generalized gradient method is substantiated for the training nonsmooth nonconvex deep neural networks:

For the linear classification problem, a variant of the error function is proposed and the advantages of the method of stochastic generalized gradients are demonstrated.

It is of interest to extend the methods of block-coordinate [11] and asynchronous [17, 20] stochastic gradient descent to general nonconvex nonsmooth machine learning problems, which are smoothed out through artificial randomization.

References

1. Rumelhart, D.E., Hinton, G.E., Williams, R.J. (1986). Learning Representations by Back-Propagating Errors. *Nature*. Vol. 323, pp. 533-536. DOI: <https://doi.org/10.1038/323533a0>
2. LeCun, Y.A, Bottou, L., Orr G.B., and Muller K.-R. (2012). Efficient BackProp. G. Montavon et al. (Eds.): NN: Tricks of the Trade, 2nd edn., LNCS 7700, pp. 9 - 48. Berlin, Heidelberg: Springer-Verlag, 2012.
3. LeCun, Y., Bengio, Y., Hinton, G. (2015). Deep learning. 436. *Nature*. Vol. 521. DOI:10.1038/nature14539

4. Flach, P. (2012). Machine Learning. The Art and Science of Algorithms that Make Sense of Data. Cambridge University Press.
5. Goodfellow, I., Bengio, Y., Courville, A. (2016). Deep Learning, MIT Press.
6. Nikolenko, S., Kadurin, A., Arhangel'skaia, E. (2018). Deep Learning. Diving in the world of neural networks. St. Petersburg: Piter. (In Russian).
7. Vorontsov, K.V. Machine learning (a year course). Access (11.07.2019): <http://www.machinelearning.ru/wiki/> (In Russian).
8. Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks*, 61, pp. 85–117. DOI: <http://dx.doi.org/10.1016/j.neunet.2014.09.003>
9. Bottou, L., Curtis, F.E., Nocedal, J. (2018). Optimization Methods for Large-Scale Machine Learning. *SIAM Rev.*, 60, No. 2, pp. 223–311. DOI: <https://doi.org/10.1137/16M1080173>
10. Zhang, C., Liao, Q., Rakhlin, A., Miranda, B., Golowich, N., Poggio, T. (2018). Theory of Deep Learning IIb: Optimization Properties of SGD. *CBMM Memo No. 072*. Center for Brains, Minds, and Machines, McGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA, 2018. 9 p. arXiv:1801.02254v1 [cs.LG] 7 Jan 2018
11. Newton, D., Yousefian, F., Pasupathy, R. (2018). Stochastic Gradient Descent: Recent Trends. *INFORMS TutORials in Operations Research*, Published online: 19 Oct 2018, pp. 193-220. Doi: <https://doi.org/10.1287/educ.2018.0191>
12. Robbins, H. and Monro, S. (1951). A Stochastic Approximation Method. *The Annals of Mathematical Statistics*. Vol. 22(3). P. 400-407.
13. Ermoliev, Y.M. (1976). Methods of stochastic programming. Moscow: Nauka. (In Russian).
14. Nemirovski, A., Juditsky, A., Lan, G., and Shapiro, A. (2009). Robust Stochastic Approximation Approach to Stochastic Programming. *SIAM J. on Optimization*. Vol. 19(4). P. 1574-1609.
15. Shapiro, A., Dentcheva, D., and Ruszczyński, A. (2009). Lectures on Stochastic Programming: Modeling and Theory. Philadelphia: SIAM.
16. Davis, D., Drusvyatskiy, D., Kakade, S., and Lee, J.D. (2019). Stochastic subgradient method converges on tame functions. *Found. Comput. Math.* P. 1-36. DOI: <https://doi.org/10.1007/s10208-018-09409-5>
17. Zhu, R., Niu, D., and Li, Z. (2018). Asynchronous stochastic proximal methods for nonconvex nonsmooth optimization. *arXiv preprint arXiv:1802.08880* .

18. Li, Z., and Li, J. (2018). A simple proximal stochastic gradient method for nonsmooth nonconvex optimization. In: *Advances in Neural Information Processing Systems*. 2018. P. 5564-5574.
19. Majewski, S., Miasojedow, B., and Moulines, E. (2018). Analysis of nonsmooth stochastic approximation: the differential inclusion approach. *arXiv preprint arXiv:1805.01916* .
20. Kungurtsev, V., Egan, M., Chatterjee, B., Alistarh, D. (2019). Asynchronous Stochastic Subgradient Methods for General Nonsmooth Nonconvex Optimization. *arXiv preprint arXiv:1905.11845 [math.OC]*
21. Davis, D., Drusvyatskiy, D. (2019). Stochastic Model-Based Minimization of Weakly Convex Functions. *SIAM J. Optim.* 2019. Vol. 29(1), pp. 207-239. DOI: <https://doi.org/10.1137/18M1178244>
22. Clarke, F.H. (1990). Optimization and nonsmooth analysis, Vol. 5 of Classics in Applied Mathematics. Philadelphia, PA: Society for Industrial and Applied Mathematics (SIAM), 2nd ed.
23. Norkin, V.I. (1980). Generalized differentiable functions. *Cybernetics*, Vol. 16, No. 1, pp. 10-12. DOI: <https://doi.org/10.1007/BF01099354>
24. Mikhalevich, V.S., Gupal, A.M., Norkin, V.I. (1987). Methods of Nonconvex Optimization, Moscow: Nauka. (In Russian).
25. Nurminskii, E.A. (1974). Minimization of nondifferentiable functions in the presence of noise. *Cybernetics*. Vol. 10(4), pp. 619-621. DOI: <https://doi.org/10.1007/BF01071541>
26. Nurminski, E.A. (1979). Numerical methods for solving stochastic minimax problems. Kyiv: Naukova Dumka. (In Russian).
27. Ermoliev, Y.M., Norkin, V.I. (1998). Stochastic generalized gradient method for solving nonconvex nonsmooth stochastic optimization problems. *Cybern. Syst. Anal.* Vol. 34, No. 2, pp. 196-215. DOI: <https://doi.org/10.1007/BF02742069>
28. Norkin, V.I. (1998). Stochastic methods for solving nonconvex stochastic optimization problems and their applications. *Extended abstract of the Doctor Thesis*. Kyiv: V.M.Glushkov Institute of Cybernetics of the National Academy of Sciences of Ukraine. 32 p. Access (22.07.2019): <http://library.nuft.edu.ua/ebook/file/01.05.01%20Norkin%20VI.pdf> (In Ukrainian).
29. Ermoliev, Y.M., Norkin, V.I. (2003). Solution of nonconvex nonsmooth stochastic optimization problems. *Cybern. Syst. Anal.*, 39, No. 5, pp. 701-715. DOI: <https://doi.org/10.1023/B:CASA.0000012091.84864.65>

30. Burke, J., Lewis, A., and Overton, M. (2005). A robust gradient sampling algorithm for nonsmooth nonconvex optimization. *SIAM J. Opt.* Vol. 25, pp. 751-779.
31. Haykin, S. (2006). *Neural Networks. A Comprehensive Foundation. Second Edition.* Prentice Hall, New Jersey.
32. Jarrett, K., Kavukcuoglu, K., Ranzato, M., and LeCun, Y. (2009). What is the best multi-stage architecture for object recognition? 2009 IEEE 12th International Conference on Computer Vision (ICCV). IEEE. P. 2146-2153. DOI: 10.1109/iccv.2009.5459469
33. Norkin, V.I. (1978). Nonlocal minimization algorithms of nondifferentiable functions. *Cybernetics*. Vol. 14, No. 5, pp. 704-707. DOI: <https://link.springer.com/article/10.1007/BF01069307>
34. Mifflin, R. (1977). An algorithm for constrained optimization with semi-smooth functions. *Math. Oper. Res.* 1977. Vol. 2, No. 2, pp. 191-207.
35. Bolte, J., Daniilidis, A., Lewis, A. (2009). Tame functions are semismooth. *Math. Program., Ser. B.* Vol. 117, pp. 5–19. DOI: 10.1007/s10107-007-0166-9
36. Norkin, V.I. (1986). Stochastic generalized-differentiable functions in the problem of nonconvex nonsmooth stochastic optimization. *Cybernetics*. Vol. 22, No.6, pp. 804-809. DOI: <https://doi.org/10.1007/BF01068698>
37. Norkin, V.I. (2019). Generalized gradients in problems of dynamic optimization, optimal control, and machine learning. *Preprint*. Kyiv: V.M.Glushkov Institute of Cybernetics of the National Academy of Sciences of Ukraine. Access (2019.09.19): http://www.optimization-online.org/DB_HTML/2019/09/7374.html
38. Pontryagin, L.S., Boltianski, V.G., Gamkelidze, R.V., Mischenko, E.F. (1962). *The Mathematical Theory of Optimal Processes.* John Wiley.
39. Boltianski, V.G. (1973). *Optimal Control of Discrete Systems.* Moscow: Nauka. (In Russian).
40. Ermoliev, Yu.M. and Norkin, V.I. (1997). Stochastic generalized gradient method with application to insurance risk management. *Interim Report IR-97-021.* Laxenburg (Austria): Int. Inst. for Appl. Syst. Anal. Laxenburg, Austria. 19 p. Access: <http://pure.iiasa.ac.at/id/eprint/5270/1/IR-97-021.pdf>
41. Gel'fand, I.M., Tzeitlin, M.L. (1961). The principle of the nonlocal search in automatic optimization systems. *Dokl. Akad. Nauk SSSR*. Vol. 137(2), pp. 295–298.
42. Nesterov, Y. (1983). A method of solving a convex programming problem with convergence rate $O(1/k^2)$. *Soviet Math. Dokl.*, Vol. 27(2), pp. 372–376.
43. Urjas'ev, S.P. (1980). Step control for direct stochastic-programming methods. *Cybernetics*. Vol. 16(6), pp. 886-890. DOI: <https://doi.org/10.1007/BF01069063>

44. Urjas'ev, S.P. (1986). Stochastic quasi-gradient algorithms with adaptively controlled parameters. *Working paper* WP-86-32. Laxenburg (Austria): Int. Inst. for Appl. Syst. Anal. 27 p. Access: <http://pure.iiasa.ac.at/id/eprint/2827/1/WP-86-032.pdf>
45. Urjas'ev, S.P. (1990). Adaptive algorithms of stochastic optimization and game theory. Nauka, Moscow. (In Russian).
46. Gupal, A.M. (1979). Stochastic Methods for Solution of Nonsmooth Extremal Problems, Naukova Dumka, Kiev. (In Russian).
47. Granichin, O.N., Polyak, B.T. (2003). Randomized algorithms of optimization and estimation under almost arbitrary noise. Nauka, Moscow. (In Russian).
48. Vapnik, V.N. (1998). Statistical learning theory. New York: Wiley & Sons.
49. Shlezinger, M., Glavach, V. (2004). Ten lectures on the statistical and structural recognition. Nakova dumka, Kyiv. (In Russian).
50. Yuan, G.-X., Ho, C.-H., Lin, C.-J. (2012). Recent Advances of Large-Scale Linear Classification. *Proceedings of the IEEE*. Vol. 100, Issue 9, pp. 2584-2603. DOI: 10.1109/JPROC.2012.2188013
51. Laptin, Yu.P., Zhuravlev, Yu.I., Vinogradov, A.P. (2011). Empirical risk minimization and problems of constructing linear classifiers. *Cybern. Syst. Anal.* Vol. 47(4), pp. 640-648. DOI: <https://doi.org/10.1007/s10559-011-9344-0>
52. Zhuravlev, Yu.I., Laptin, Yu.P., Vinogradov, A.P., Zhurbenko, N.G., Lykhovyd, O.P., and Berezovskyi, O.A. (2017). Linear classifiers and selection of informative features. *Pattern Recognition and Image Analysis*. Vol. 27, No. 3, pp. 426-432. DOI: <https://doi.org/10.1134/S1054661817030336>
53. Zhuravlev, Yu.I., Laptin, Yu.P., Vinogradov, A.P. (2014). A comparison of some approaches to classification problems, and possibilities to construct optimal solutions efficiently. *Pattern Recognition and Image Analysis*. Vol. 24, No. 2, pp. 189-195. DOI: <https://doi.org/10.1134/S1054661814020175>