

# Learning Supervised PageRank with Gradient-Based and Gradient-Free Optimization Methods

**Authored by:**

Gleb Gusev  
Lev Bogolubsky  
Pavel Dvurechensky  
Alexander Gasnikov  
Yurii Nesterov  
Andrei M. Raigorodskii  
Aleksey Tikhonov  
Maksim Zhukovskii

## **Abstract**

In this paper, we consider a non-convex loss-minimization problem of learning Supervised PageRank models, which can account for features of nodes and edges. We propose gradient-based and random gradient-free methods to solve this problem. Our algorithms are based on the concept of an inexact oracle and unlike the state-of-the-art gradient-based method we manage to provide theoretically the convergence rate guarantees for both of them. Finally, we compare the performance of the proposed optimization methods with the state of the art applied to a ranking task.

## **1 Paper Body**

In this paper, we consider a non-convex loss-minimization problem of learning Supervised PageRank models, which can account for features of nodes and edges. We propose gradient-based and random gradient-free methods to solve this problem. Our algorithms are based on the concept of an inexact oracle and unlike the state-of-the-art gradient-based method we manage to provide theoretically the convergence rate guarantees for both of them. Finally, we compare the performance of the proposed optimization methods with the state of the art applied to a ranking task.

1

INTRODUCTION

The most acknowledged methods of measuring importance of nodes in graphs are based on random walk models. Particularly, PageRank [18], HITS [11], and their variants [8, 9, 19] are originally based on a discrete-time Markov random walk on a link graph. Despite undeniable advantages of PageRank and its mentioned modifications, these algorithms miss important aspects of the graph that are not described by its structure. In contrast, a number of approaches allows to account for different properties of nodes and edges between them by encoding them in restart and transition probabilities (see [3, 4, 6, 10, 12, 20, 21]). These properties may include, e.g., the statistics about users' interactions with the nodes (in web graphs [12] or graphs of social networks [2]), types of edges (such as URL redirecting in web graphs [20]) or histories of nodes' and edges' changes [22]. In the general ranking framework called Supervised PageRank [21], weights of nodes and edges in a graph are linear combinations of their features with coefficients as the model parameters. The existing optimization method [21] of learning these parameters and the optimizations methods proposed in the presented paper have two levels. On the lower level, the following problem is solved: to estimate the value of the loss function (in the case of zero-order oracle) and its derivatives (in the case of first-order oracle) for a given parameter vector. On the upper level, the estimations obtained on the lower level of the optimization methods (which we also call inexact oracle information) are used for tuning the parameters by an iterative algorithm. Following [6], the authors of Supervised PageRank consider a non-convex loss-minimization problem for learning the parameters and solve 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

it by a two-level gradient-based method. On the lower level of this algorithm, an estimation of the stationary distribution of the considered Markov random walk is obtained by classical power method and estimations of derivatives w.r.t. the parameters of the random walk are obtained by power method introduced in [23, 24]. On the upper level, the obtained gradient of the stationary distribution is exploited by the gradient descent algorithm. As both power methods give imprecise values of the stationary distribution and its derivatives, there was no proof of the convergence of the state-of-the-art gradient-based method to a stationary point. The considered non-convex loss-minimization problem [21] can not be solved by existing optimization methods such as [16] and [7] due to presence of constraints for parameter vector and the impossibility to calculate the exact value of the loss function. Moreover, standard global optimization methods can not be applied, because they require unbiased estimations of the loss function. In our paper, we propose two two-level methods to solve the problem [21]. On the lower level of these methods, we use the linearly convergent method [17] to calculate an approximation to the stationary distribution of Markov random walk. We show that this method allows to approximate the value of the loss function at any given accuracy and has the lowest proved complexity bound among methods proposed in [5]. We develop a gradient method for general constrained non-convex optimization problems with inexact oracle, estimate its convergence rate to the stationary point of the problem. We exploit this gradient method on the upper level of the two-level algorithm for

learning Supervised PageRank. Our contribution to the gradient-free methods framework consists in adapting the approach of [16] to the case of constrained optimization problems when the value of the function is calculated with some known accuracy. We prove a convergence theorem for this method and exploit it on the upper level of the second two-level algorithm. Another contribution consists in investigating both for the gradient and gradient-free methods the trade-off between the accuracy of the lower-level algorithm, which is controlled by the number of iterations of method in [17] and its generalization (for derivatives estimation), and the computational complexity of the two-level algorithm as a whole. Finally, we estimate the complexity of the whole two-level algorithms for solving the loss-minimization problem with a given accuracy. In the experiments, we apply our algorithms to learning Supervised PageRank on a real ranking task. Summing up, both two-level methods, unlike the state-of-the-art [21], have theoretical guarantees on convergence rate, and outperform it in the ranking quality in experiments. The main advantages of the first gradient-based algorithm: the guarantees of a convergence do not require the convexity, this algorithm has less input parameters than gradient-free one. The main advantage of the second gradient-free algorithm is that it avoids calculating the derivative for each element of a large matrix.

2

## MODEL DESCRIPTION

We consider the following random walk on a directed graph  $\mathcal{G} = (V, E)$  introduced in [21]. Assume 1 that each node  $i \in V$  and each edge  $i \rightarrow j \in E$  is represented by a vector of features  $V_i \in \mathbb{R}^{m_1}$  and  $m_2$  a vector of features  $E_{ij} \in \mathbb{R}^{m_2}$  respectively. A surfer starts from a random page  $v_0$  of a seed set  $U \subset V$ . The restart probability that  $v_0 = i$  equals  $h_1$ ,  $V_i \in U$ ,  $h_1$ ,  $V_i$

$$[P_0]_i = P$$

$i \in U$

(2.1)

and  $[P_0]_i = P_0$  for  $i \in V \setminus U$ , where  $h_1 \in \mathbb{R}^{m_1}$  is a parameter, which conducts the random walk. We assume that  $h_1$ ,  $V_i$  should be non-zero. At each step, the surfer makes a restart with probability  $\alpha \in (0, 1)$  (originally [18],  $\alpha = 0.15$ ) or traverses an outgoing edge (makes a transition) with probability  $1 - \alpha$ . In the former case, the surfer chooses a vertex according to the distribution  $P_0$ . In the latter one, the transition probability of traversing an edge  $i \rightarrow j \in E$  is  $h_2$ ,  $E_{ij} \in \mathbb{R}^{m_2}$ ,  $h_2$ ,  $E_{ij}$

$$[P]_{i,j} = P$$

(2.2)

where  $h_2 \in \mathbb{R}^{m_2}$  is a parameter and the current position  $i$  has non-zero outdegree, and  $[P(\alpha)]_{i,j} = [P_0(\alpha)]_j$  for all  $j \in V$  if the outdegree of  $i$  is zero (thus the surfer always makes a restart in this case). We assume that  $h_2$ ,  $E_{ij}$  is non-zero for all  $i$  with non-zero outdegree. 2

By Equations 2.1 and 2.2, the total probability of choosing vertex  $j \in V$  conditioned by the surfer being at vertex  $i$  equals  $[P_0(\alpha)]_j + (1 - \alpha)[P(\alpha)]_{i,j}$ , where  $\alpha = (h_1, h_2)^T$  and we use  $P_0(\alpha)$  and  $P(\alpha)$  to express the dependence of  $P_0$ ,  $P$  on the parameters. The stationary distribution  $\pi(\alpha) \in \mathbb{R}^p$  of the

described Markov process is a solution of the system  $\pi = \pi P + (1 - \alpha)P$   $T(\pi)$ .

(2.3)

In this paper, we learn an algorithm, which ranks nodes  $i$  according to scores  $[f(\pi)]_i$ . Let  $Q$  be a set of queries and a set of nodes  $V_q \subseteq V$  is associated to each query  $q$ . For example, vertices in  $V_q$  may represent web pages visited by users after submitting query  $q$ . For each  $q \in Q$ , some nodes of  $V_q$  are manually judged by relevance labels  $1, \dots, \ell$ . Our goal is to learn the parameter vector  $\pi$  of a ranking algorithm  $f_q = f(\pi)$  which minimizes the discrepancy of its ranking scores  $[f_q]_i, i \in V_q$ , from the assigned labels. We consider the square loss function [12, 21, 22] —

$$\frac{1}{2} \sum_{q \in Q} \sum_{i \in V_q} (A_{q,i} - f_q(i))^2. \quad (2.4)$$

Each row of matrix  $A_q \in \mathbb{R}^{V_q \times \ell}$  corresponds to some pair of pages  $i_1, i_2 \in V_q$  such that the label of  $i_1$  is strictly greater than the label of  $i_2$  (we denote by  $r_q$  the number of all such pairs from  $V_q$  and  $p_q := |V_q|$ ). The  $i_1$ -th element of this row is equal to 1,  $i_2$ -th element is equal to 1, and all other elements are equal to 0. Vector  $x_+$  has components  $[x_+]_i = \max\{x_i, 0\}$ . To make ranking scores (2.3) query-dependent, we assume that  $f$  is defined on a query-dependent graph  $G_q = (V_q, E_q)$  with query-dependent feature vectors  $V_{ij}, i \in V_q, E_{ij}, i \in j \in E_q$ . For example, these features may reflect different aspects of query-page relevance. For a given  $q \in Q$ , we consider all the objects related to the graph  $G_q$  introduced above:  $U_q := U, \pi_q := \pi, P_q := P, f_q := f$ . In this way, the ranking scores  $f_q$  depend on query via the query-dependent features, but the parameters of the model  $\pi$  and  $P$  are not query-dependent. In what follows, we use the following notations throughout the paper:  $n_q := |U_q|, m = m_1 + m_2, r = \max_{q \in Q} r_q, p = \max_{q \in Q} p_q, n = \max_{q \in Q} n_q, s = \max_{q \in Q} s_q$ , where  $s_q = \max_{i \in V_q} |\{j : i \in j \in E_q\}|$ . In order to guarantee that the probabilities in (2.1) and (2.2) are correctly defined, we need to appropriately choose a set  $\pi$  of possible values of parameters  $\pi$ . We choose some  $\pi_0$  and  $R \geq 0$  such that  $1 - \pi_0 = \{\pi \in \mathbb{R}^m : \pi \geq \pi_0, \pi \leq R\}$  lies in the set of vectors with positive components  $\pi \geq 0$ . In this paper, we solve the following loss-minimization problem:  $\min_{\pi \in \pi_0 + R \cdot \mathbb{B}} f(\pi), \pi_0 = \{\pi \in \mathbb{R}^m : \pi \geq \pi_0, \pi \leq R\}$ .

(2.5)

???

NUMERICAL CALCULATION OF  $f(\pi)$  AND  $\nabla f(\pi)$

3

Our goal is to provide methods for solving Problem 2.5 with guarantees on rate of convergence and complexity bounds. The calculation of the values of  $f(\pi)$  and its gradient  $\nabla f(\pi)$  is problematic, since it requires to calculate those for  $|Q|$  vectors  $f_q(\pi)$  defined by Equation 2.3. While the exact values are impossible to derive in general, existing methods provide estimations of  $f_q(\pi)$  and its  $d f_q(\pi)$  derivatives  $d f$  in an iterative way with a trade-off between time and accuracy. To be able to guarantee convergence of our optimization algorithm in this inexact oracle setting, we consider numerical

methods that calculate approximation for  $q(\cdot)$  and its derivatives with any required accuracy. We have analysed state-of-the-art methods summarized in the review [5] and power method used in [18, 2, 21] and have found that the method [17] is the most suitable. It constructs a sequence  $k$  and outputs  $q(\cdot, N)$  by the following rule (integer  $N \geq 0$  is a parameter):  $q_0 = q_0(\cdot)$ ,

$$\begin{aligned} k+1 &= P_q T(\cdot) k, \\ q(\cdot, N) &= \\ N \times \sum_{k=0}^{N-1} (1 - \frac{1}{N})^k q_k(\cdot) & \quad (3.1) \end{aligned}$$

As probabilities  $[q_0(\cdot)]_i, i = 1, \dots, n, [P_q(\cdot)]_{i,i}, i = 1, \dots, n$  are scale-invariant ( $q_0(\cdot) = q_0(\cdot), P_q(\cdot) = P_q(\cdot)$ ), in our experiments, we consider the set  $\mathcal{Q} = \{q : \sum_{k=1}^n q_k \leq 0.99\}$ , where  $q$  is the vector of all ones, that has large intersection with the simplex  $\{q : \sum_{k=1}^n q_k = 1\}$ .

**Lemma 1.** Assume that, for some  $\epsilon > 0$ , Method 3.1 with  $N = \lceil \ln 8r / \epsilon \rceil$  is used to calculate  $P - Q - \frac{1}{N} q(\cdot, N)$  and  $k_2$  satisfies the vector  $q(\cdot, N)$  for every  $q \in \mathcal{Q}$ . Then  $f(q, \epsilon) = \sum_{k=1}^n k(Aq - f(q, \epsilon)) f(q) - \epsilon$ . Moreover, the calculation of  $f(q, \epsilon)$  requires not more than  $O(3\text{mps} + 3\text{ps}N + 6r)$  a.o. The proof of Lemma 1 is in Supplementary Materials. Let  $p_i(\cdot)$  be the  $i$ -th column of the matrix  $P_q T(\cdot)$ . Our generalization of the method [17] for  $d q(\cdot)$  is

for any  $q \in \mathcal{Q}$  is the following. Choose some non-negative integer  $N_1$  and  $q(\cdot, N_2)$  calculate  $q(\cdot, N_1)$  using (3.1). Choose some  $N_2 \geq 0$ , calculate  $k, k = 0, \dots, N_2$  and  $q$  calculation of

$$\begin{aligned} q_0 &= q \\ p_q X d q_0(\cdot) d p_i(\cdot) + (1 - \frac{1}{N_2}) [q(\cdot, N_1)]_i, T d q d T i=1 \\ q(\cdot, N_2) &= q \\ k+1 &= P_q T(\cdot) k, \\ N_2 \times \sum_{k=0}^{N_2-1} (1 - \frac{1}{N_2})^k q_k(\cdot) & \quad (3.2) \end{aligned}$$

$n_1, n_2$  In what follows,  $\|A\|_1 = \sum_{i=1}^n \sum_{j=1}^n |a_{ij}|$ , computable constant (see Supplementary Materials). Assume Lemma 2. Let  $\epsilon$  be some explicitly  $\epsilon$  that Method 3.1 with  $N_1 = \lceil \ln 24r / \epsilon \rceil$  is used for every  $q \in \mathcal{Q}$  to calculate the vector  $q(\cdot, N_1)$  and Method 3.2, 3.3 with  $N_2 = \lceil \ln 8r / \epsilon \rceil$  is used for every  $q \in \mathcal{Q}$  to calculate the  $q_2$

$T q(\cdot, N_2)$  (3.3). Then the vector  $g(q, \epsilon) = 2 P - Q - \frac{1}{N_2} q(\cdot, N_2)$   $AT_q(Aq - q(\cdot, N_1)) +$  matrix  $q = 1 - Q$

satisfies  $\|g(q, \epsilon) - f(q, \epsilon)\|_1 \leq \epsilon$ . Moreover, the calculation of  $g(q, \epsilon)$  requires not more than  $O(10\text{mps} + 3\text{ps}N_1 + 3\text{mps}N_2 + 7r)$  a.o. The proof of Lemma 2 can be found in Supplementary Materials.

## RANDOM GRADIENT-FREE OPTIMIZATION METHODS

In this section, we first describe general framework of random gradient-free methods with inexact oracle and then apply it for Problem 2.5. Lemma 1 allows to control the accuracy of the inexact zero-order oracle and hence apply random gradient-free methods with inexact oracle. 4.1

### GENERAL FRAMEWORK

Below we extend the framework of random gradient-free methods [1, 16, 7] for the situation of presence of uniformly bounded error of unknown nature in the value of an objective function in general optimization problem. Unlike [16], we consider a constrained optimization problem and a randomization on a Euclidean sphere which seems to give better large deviations bounds and doesn't need the assumption that the objective function can be calculated at any point of  $\mathbb{R}^m$ . Let  $E$  be a  $m$ -dimensional vector space and  $E^*$  be its dual. In this subsection, we consider a general function  $f : E^* \rightarrow \mathbb{R}$  and denote its argument by  $x$  or  $y$  to avoid confusion with other sections. We denote the value of linear function  $g : E^* \rightarrow \mathbb{R}$  at  $x \in E$  by  $hg, x$ . We choose some norm  $\|\cdot\|$  in  $E$  and say that  $f \in \text{CL}_{1,1}(k, \ell)$  iff  $|f(x) - f(y)| \leq \ell \|x - y\|$  and  $|f(x)| \leq k$  for all  $x \in E$ . The problem of our interest is to find  $\min_{x \in X} f(x)$ , where  $f \in \text{CL}_{1,1}(k, \ell)$ ,  $X$  is a closed convex set and there exists a number  $D > 0$ ,  $\ell > 0$  such that  $\text{diam} X := \max_{x, y \in X} \|x - y\| \leq D$ . Also we assume that  $f$  is the inexact zero-order oracle for  $f(x)$  returns a value  $f^?(x, \xi) = f(x) + \xi(x)$ , where  $\xi(x)$  is the error satisfying for some  $\epsilon > 0$  (which is known)  $|\xi(x)| \leq \epsilon$  for all  $x \in X$ . Let  $x^* \in \arg \min_{x \in X} f(x)$ . Denote  $f^* = \min_{x \in X} f(x)$ . Unlike [16], we define the biased gradient-free oracle  $g^?(x, \xi) = m^{-1} \sum_{i=1}^m (f(x + \xi_i), \xi_i)$  where  $m$  is a random vector uniformly distributed over the unit sphere  $S = \{\xi \in \mathbb{R}^m : \|\xi\| = 1\}$ ,  $\epsilon$  is a smoothing parameter. 4

**Algorithm 1** Gradient-type method Input: Point  $x_0 \in X$ , stepsize  $h > 0$ , number of steps  $M$ . Set  $k = 0$ . repeat Generate  $\xi_k$  and calculate corresponding  $g^?(x_k, \xi_k)$ . Calculate  $x_{k+1} = \Pi_X(x_k + hg^?(x_k, \xi_k))$  ( $\Pi_X(\cdot)$  is Euclidean projection onto the set  $X$ ). Set  $k = k + 1$ . until  $k \geq M$  Output: The point  $y_M = \arg \min_{x \in \{x_0, \dots, x_M\}} f(x)$ . Theorem 1. Let  $f \in \text{CL}_{1,1}(k, \ell)$  and convex. Assume that  $x^* \in \text{int} X$ , and the sequence  $x_k$  is generated by Algorithm 1 with  $h = 8mL$ . Then for any  $M \geq 0$ , we have  $E f(y_M) \leq f^* + 8mLD^2 M^{-1} +$

$$\begin{aligned} &+ \text{vector } \epsilon. \\ &\leq 2L(m+8) \epsilon \\ &+ \\ &\leq mD^2 \epsilon \\ &+ \\ &\leq 2mL^2 \epsilon. \end{aligned}$$

Here  $\xi_k = (\xi_k^1, \dots, \xi_k^m)$  is the history of realizations of the

The full proof of the theorem is in Supplementary Materials. 4.2

### SOLVING THE LEARNING PROBLEM

Now, we apply the results of Subsection 4.1 to solve Problem 2.5. Note that presence of constraints and oracle inexactness do not allow to directly apply the results of [16]. We assume that there is a local minimum  $x^*$ , and  $f$  is

a small vicinity of  $x^*$ , in which  $f(x)$  (2.4) is convex (generally speaking, it is nonconvex). We choose the desired accuracy  $\epsilon$  for  $f^*$  (the optimal value) approximation in the sense that  $E_M \leq f(y_M) - f^* \leq \epsilon$ . In accordance with Theorem 1,  $M$  gives the number of steps  $M$  of Algorithm 1, the value of  $\epsilon$ , the value of the required accuracy  $\epsilon$  of the inexact zero-order oracle. The value  $\epsilon$ , by Lemma 1, gives the number of steps  $N$  of Method 3.1 required to calculate a  $\epsilon$ -approximation  $f^*(\epsilon, \epsilon)$  for  $f^*$ . Then the inexact zero-order oracle  $f^*(\epsilon, \epsilon)$  is used to make Algorithm 1 step. Theorem 1 and the choice of the feasible set  $X$  to be a Euclidean ball make it natural to choose  $k_2$ -norm in the space  $\mathbb{R}^m$  of parameter  $x$ . It is easy to see that in this norm  $\text{diam}(X) \leq 2R$ . Algorithm 2 in Supplementary Materials is a formal record of these ideas. The most computationally hard on each iteration of the main cycle of this method are calculations of  $f^*(x_k + \epsilon_k, \epsilon)$ ,  $f^*(x_k, \epsilon)$ . Using Lemma 1, we obtain the complexity of each iteration and the following result, which gives the complexity of Algorithm 2. Theorem 2. Assume that the set  $X$  in (2.5) is chosen in a way such that  $f(x)$  is convex on  $X$  and some  $x^* \in \arg \min_{x \in X} f(x)$  belongs also to  $\text{int}(X)$ . Then the mean total number of arithmetic operations of the Algorithm 2 for the accuracy  $\epsilon$  (i.e. for the inequality  $E_M \leq f(y_M) - f^* \leq \epsilon$  to hold) is not more than  $\frac{1}{\epsilon} \left( \frac{1}{L} \ln \frac{1}{\epsilon} + 8 \right) \frac{R^2}{\epsilon} + m + \ln \frac{1}{\epsilon} + 6$ .

5

## GRADIENT-BASED OPTIMIZATION METHODS

In this section, we first develop a general framework of gradient methods with inexact oracle for non-convex problems from rather general class and then apply it for the particular Problem 2.5. Lemma 1 and Lemma 2 allow to control the accuracy of the inexact first-order oracle and hence apply proposed framework.

### 5.1 GENERAL FRAMEWORK

In this subsection, we generalize the approach in [7] for constrained non-convex optimization problems. Our main contribution consists in developing this framework for an inexact first-order oracle and unknown "Lipschitz constant" of this oracle. We consider a composite optimization problem of the form  $\min_{x \in X} \{f(x) := f(x) + h(x)\}$ , where  $X \subseteq \mathbb{R}^n$  is a closed convex set,  $h(x)$  is a simple convex function, e.g.  $\|x\|_1$ . We assume that  $f(x)$  is

a general function endowed with an inexact first-order oracle in the following sense. There exists a number  $L \in (0, +\infty)$  such that for any  $\epsilon \geq 0$  and any  $x \in X$  one can calculate  $f^*(x, \epsilon) \in \mathbb{R}$  and  $g^*(x, \epsilon) \in \mathbb{R}^n$  satisfying  $L - f(y) \leq (f^*(x, \epsilon) - h(y) - g^*(x, \epsilon), y - x) \leq \|x - y\|^2 + \epsilon$ . (5.1) for all  $y \in X$ . The constant  $L$  can be considered as "Lipschitz constant" because for the exact firstorder oracle for a function  $f \in \text{CL}_{1,1}(k, k)$  Inequality 5.1 holds with  $\epsilon = 0$ . This is a generalization of the concept of  $(\epsilon, L)$ -oracle considered in [25] for convex problems. We choose a prox-function  $d(x)$  which is continuously differentiable and 1-strongly convex on  $X$  with respect to  $k \geq k$ . This means that for any  $x, y \in X$   $d(y) \geq d(x) + h^*d(x), y - x \geq \frac{1}{k} \|y - x\|^2$ . We define also the corresponding Bregman distance  $V(x, z) = d(x) - d(z) - h^*d(z), x \in X, z \in X$ . Algorithm 2 Adaptive projected gradient algorithm Input: Point  $x_0 \in X$ ,

number  $L_0 \leq 0$ . Set  $k = 0$ ,  $z = +\infty$ . repeat Set  $M_k = L_k$ ,  $\text{flag} = 0$ . repeat ? Set  $\epsilon = 16M$ . Calculate  $f^*(x_k, \epsilon)$  and  $g^*(x_k, \epsilon)$ .  $k$  Find  $w_k = \arg \min_{x \in Q} \{h^*(g(x_k, \epsilon), x) + M_k V(x, x_k) + h(x)\}$  and calculate  $f^*(w_k, \epsilon)$ . ? ? If the inequality  $f(w_k, \epsilon) \leq f^*(x_k, \epsilon) + h^*$  holds,  $g(x_k, \epsilon)$ ,  $w_k \leftarrow x_k + M_2 k$   $k w_k \leq x_k^2 + 8M_k$  set  $\text{flag} = 1$ . Otherwise set  $M_k = 2M_k$ . until  $\text{flag} = 1$  Set  $x_{k+1} = w_k$ ,  $L_{k+1} = M_2 k$ . If  $k M_k (x_k - x_{k+1}) \leq z$ , set  $z = k M_k (x_k - x_{k+1})$ ,  $K = k$ . Set  $k = k + 1$ . until  $z \leq \epsilon$  Output: The point  $x_{K+1}$ . Theorem 3. Assume that  $f(x)$  is endowed with the inexact first-order oracle in a sense (5.1) and that there exists a number  $\epsilon \leq \epsilon^*$  such that  $f^*(x) \leq \epsilon$  for all  $x \in X$ . Then after  $M$  iterations of Algorithm 2 it holds that  $k M K (x_K - x_{K+1}) \leq 4L(\epsilon(x + 2\epsilon))$ . Moreover, the total number of  $M + 1$  inexact oracle calls is not more than  $2M + 2 \log_2 2L/L_0$ . The full proof of the theorem is in Supplementary Materials. 5.2

#### SOLVING THE LEARNING PROBLEM

In this subsection, we return to Problem 2.5 and apply the results of the previous subsection. Note that we can not directly apply the results of [7] due to the inexactness of the oracle. For this problem,  $h(\epsilon) \leq 0$ . It is easy to show that in 1-norm  $\text{diam}(\epsilon) \leq 2Rm$ . For any  $\epsilon \leq 0$ , Lemma 1 with  $\epsilon_1 = 2\epsilon$  allows us to obtain  $f^*(\epsilon, \epsilon_1)$  such that inequality  $-f^*(\epsilon, \epsilon_1) \leq f(\epsilon) - \epsilon_1$  holds and Lemma 2 with  $\epsilon_2 = 4R\epsilon m$  allows us to obtain  $g^*(\epsilon, \epsilon_2)$  such that inequality  $k \leq g(\epsilon, \epsilon_2) \leq f(\epsilon)k \leq \epsilon_2$  holds. Similar to [25], since  $f \in \text{CL}_{1,1}(k \leq k^2)$ , these two inequalities lead to Inequality 5.1 for  $f^*(\epsilon, \epsilon_1)$  in the role of  $f^*(x, \epsilon)$ ,  $g^*(\epsilon, \epsilon_2)$  in the role of  $g^*(x, \epsilon)$  and  $k \leq k^2$  in the role of  $k \leq k$ . We choose the desired accuracy  $\epsilon$  for approximating the stationary point of Problem 2.5. This accuracy gives the required accuracy  $\epsilon$  of the inexact first-order oracle for  $f(\epsilon)$  on each step of the inner cycle of the Algorithm 2. Knowing the value  $\epsilon_1 = 2\epsilon$  and using Lemma 1, we choose the number of steps  $N$  of Method 3.1 and thus approximate  $f(\epsilon)$  with the required accuracy  $\epsilon_1$  by  $f^*(\epsilon, \epsilon_1)$ . Knowing the value  $\epsilon_2 = 4R\epsilon m$  and using Lemma 2, we choose the number of steps  $N_1$  of Method 3.1 and the number of steps  $N_2$  of Method 3.2, 3.3 and obtain the approximation  $g^*(\epsilon, \epsilon_2)$  of  $f(\epsilon)$  with the required accuracy  $\epsilon_2$ . Then we use the inexact first-order oracle ( $f^*(\epsilon, \epsilon_1)$ ,  $g^*(\epsilon, \epsilon_2)$ ) to perform a step of Algorithm 2. Since  $\epsilon$  is the Euclidean ball, it is natural to set  $E = Rm$  and  $k \leq k \leq k^2$ , 6

choose the prox-function  $d(\epsilon) = 12k \leq k^2$ . Then the Bregman distance is  $V(\epsilon, \epsilon) = \text{Algorithm 4 in Supplementary Materials is a formal record of the above ideas.}$

$1 \leq k \leq$

$\epsilon \leq k^2$ .

The most computationally consuming operations of the inner cycle of Algorithm 4 are calculations of  $f^*(\epsilon_k, \epsilon_1)$ ,  $f^*(\epsilon_k, \epsilon_1)$  and  $g^*(\epsilon_k, \epsilon_2)$ . Using Lemma 1 and Lemma 2, we obtain the complexity of each iteration. From Theorem 3 we obtain the following result, which gives the complexity of Algorithm 4. Theorem 4. The total number of arithmetic operations in Algorithm 4 for the accuracy  $\epsilon$  (i.e. for the 2 inequality  $k M K (\epsilon_K - \epsilon_{K+1}) \leq \epsilon$  to hold) is not more than





next step and the current step are less than  $10^{-5}$  on the test sets. 2

yandex.com

7

In Table 1, we present the performances of the optimization algorithms in terms of the loss function  $f$  (2.4). We also compare the algorithms with the untuned Supervised PageRank ( $\alpha = 0 = \text{em}$ ). On Figure 1, we give the outputs of the optimization algorithms on each iteration of the upper levels of the learning processes on the test set Q32, similar results were obtained for the sets Q12, Q22. Q12

Q22

Q32

Meth. PR GBN GFN GBP 50s. GBP 100s. GBP 200s. GBP 500s.

	loss	steps	loss	steps	loss	steps	loss	steps	loss	steps	loss	steps
12	.00295	12	.00274	106	.00297	106	.00292	106	.00282	16	.00307	31
	.00282	8	.00307	16	.00295	20	.00283	4	.00308	7	.00295	9
	.00283	2	.00308	2	.00295							

3 Table 1: Comparison of the algorithms on the test sets.

Figure 1: Values of the loss function on each iteration of the optimization algorithms on the test set Q32. GFN significantly outperforms the state-of-the-art algorithms on all test sets. GBN significantly outperforms the state-of-the-art algorithm on Q12 (we obtain the p-values of the paired t-tests for all the above differences on the test sets of queries, all these values are less than 0.005). However, GBN requires less iterations of the upper level (until it stops) than GBP for step sizes 50 and 100 on Q22, Q32. Finally, we show that Nesterov-Nemirovski method converges to the stationary distribution faster than the power method (in supplementary materials, on Figure 2, we demonstrate the dependencies of the value of the loss function on Q11 for both methods of computing the untuned Supervised PageRank  $\alpha = 0 = \text{em}$ ).

7

## CONCLUSION

We propose a gradient-free optimization method for general convex problems with inexact zero-order oracle and an adaptive gradient method for possibly non-convex general composite optimization problems with inexact first-order oracle. For both methods, we provide convergence rate analysis. We also apply our new methods for known problem of learning a web-page ranking algorithm. Our new algorithms not only outperform existing algorithms, but also are guaranteed to solve this learning problem. In practice, this means that these algorithms can increase the reliability and speed of a search engine. Also, to the best of our knowledge, this is the first time when the ideas of random gradient-free and gradient optimization methods are combined with some efficient method for huge-scale optimization using the concept of an inexact oracle. Acknowledgments The research by P. Dvurechensky and A. Gasnikov presented in Section 4 of this paper was conducted in IITP RAS and supported by the Russian Science Foundation grant (project 14-50-00150), the research presented in Section 5 was supported by RFBR.

8

## 2 References

- [1] A. Agarwal, O. Dekel and L. Xiao, Optimal algorithms for online convex optimization with multi-point bandit feedback, 2010, 23rd Annual Conference on Learning Theory (COLT). [2] L. Backstrom and J. Leskovec, Supervised random walks: predicting and recommending links in social networks, 2011, WSDM. [3] Na Dai and Brian D. Davison, Freshness Matters: In Flowers, Food, and Web Authority, 2010, SIGIR. [4] N. Eiron, K. S. McCurley and J. A. Tomlin, Ranking the web frontier, 2004, WWW. [5] A. Gasnikov and D. Dmitriev, Efficient randomized algorithms for PageRank problem, *Comp. Math. & Math. Phys.*, 2015, 55(3): 1?18. [6] B. Gao, T.-Y. Liu, W. W. Huazhong, T. Wang and H. Li, Semi-supervised ranking on very large graphs with rich metadata, 2011, KDD. [7] S. Ghadimi, G. Lan, Stochastic first- and zeroth-order methods for non-convex stochastic programming, *SIAM Journal on Optimization*, 2014, 23(4), 2341?2368. [8] T. H. Haveliwala, Efficient computation of PageRank, Stanford University, 1999. [9] T. H. Haveliwala, Topic-Sensitive PageRank, 2002, WWW. [10] G. Jeh and J. Widom, Scaling Personalized Web Search, 2003, WWW. [11] J. M. Kleinberg, Authoritative sources in a hyperlinked environment, 1998, SODA. [12] Y. Liu, B. Gao, T.-Y. Liu, Y. Zhang, Z. Ma, S. He, H. Li, BrowseRank: Letting Web Users Vote for Page Importance, 2008, SIGIR. [13] J. Matyas, Random optimization, *Automation and Remote Control*, 1965, 26: 246?253. [14] Yu. Nesterov, *Introductory Lectures on Convex Optimization*, Springer, 2004, New York. [15] Yu. Nesterov, Efficiency of coordinate descent methods on huge-scale optimization problems, *SIAM Journal on Optimization*, 2012, 22(2): 341?362. [16] Yu. Nesterov and V. Spokoiny, *Random Gradient-Free Minimization of Convex Functions*, *Foundations of Computational Mathematics*, 2015, 1?40. [17] Yu. Nesterov and A. Nemirovski, Finding the stationary states of Markov chains by iterative methods, *Applied Mathematics and Computation*, 2015, 255: 58?65. [18] L. Page, S. Brin, R. Motwani, T. Winograd, The PageRank citation ranking: Bringing order to the web, Stanford InfoLab, 1999. [19] M. Richardson and P. Domingos, The intelligent surfer: Probabilistic combination of link and content information in PageRank, 2002, NIPS. [20] M. Zhukovskii, G. Gusev, P. Serdyukov, URL Redirection Accounting for Improving Link-Based Ranking Methods, 2013, ECIR. [21] M. Zhukovskii, G. Gusev, P. Serdyukov, Supervised Nested PageRank, 2014, CIKM. [22] M. Zhukovskii, A. Khropov, G. Gusev, P. Serdyukov, Fresh BrowseRank, 2013, SIGIR. [23] A. L. Andrew, Convergence of an iterative method for derivatives of eigensystems, *Journal of Computational Physics*, 1978, 26: 107?112. [24] A. Andrew, Iterative computation of derivatives of eigenvalues and eigenvectors, *IMA Journal of Applied Mathematics*, 1979, 24(2): 209?218. [25] O. Devolder, F. Glineur, Yu. Nesterov, First-order methods of smooth convex optimization with inexact oracle, *Mathematical Programming*, 2013, 146(1): 37?75. [26] Yu. Nesterov, B.T. Polyak, Cubic regularization of Newton method and its global performance, *Mathematical Programming*, 2006, 108(1) 177?205. [27] Yu. Nesterov, Gradient methods for minimizing composite functions, *Mathematical Programming*, 2012, 140(1) 125?161.

