Distributed Online Convex Optimization with Compressed Communication

Anonymous Author(s) Affiliation Address email

Abstract

We consider a distributed online convex optimization problem when streaming 1 data are distributed among computing agents over a connected communication 2 network. Since the data are high-dimensional or the network is large-scale, com-3 munication load can be a bottleneck for the efficiency of distributed algorithms. 4 To tackle this bottleneck, we apply the state-of-art data compression scheme to 5 the fundamental GD-based distributed online algorithms. Three algorithms with 6 difference-compressed communication are proposed for full information feedback 7 (DC-DOGD), one-point bandit feedback (DC-DOBD), and two-point bandit feed-8 back (DC-DO2BD), respectively. We obtain regret bounds explicitly in terms 9 of the time horizon, compression ratio, decision dimension, agent number, and 10 network parameters. Our algorithms are proved to be no-regret and match the same 11 12 regret bounds, w.r.t. the time horizon, with their uncompressed versions for both convex and strongly convex losses. Numerical experiments are given to validate 13 the theoretical findings and illustrate that the proposed algorithms can effectively 14 reduce the total transmitted bits for distributed online training compared with the 15 uncompressed baseline. 16

17 **1 Introduction**

Online optimization has attracted considerable attention in recent decades, for its remarkable ap-18 plications in machine learning tasks such as spam filtering, dictionary learning, ad. selection, and 19 so on [1, 2, 3]. In such online tasks, data are revealed incrementally, and decisions must be made 20 before all data are available. When the streaming data are collected at multiple agents, the distributed 21 online optimization over a multi-agent network is considered, where data storage and processing are 22 performed in the agents [4, 5]. It is often impractical to communicate data among different agents 23 from multiple concerns such as privacy and bandwidth utilization. Also, there is no center agent for 24 global coordination. In such settings, each agent relies on its own data to run an algorithm while 25 26 communicating decisions with its immediate neighbors.

To be specific, this paper considers the distributed online convex optimization (DOCO) prob-27 lem over an N-agent network. The objective is to minimize the accumulated system-wide loss 28 $\min_{x \in \mathcal{K}} \sum_{t=1}^{T} \sum_{i=1}^{N} f_i^t(x)$, where the local convex loss function f_i^t is formed by the data arriving 29 at time t in agent i, and $\mathcal{K} \subset \mathbb{R}^d$ is a convex feasible set. Note that the loss information is revealed to 30 agent i after its decision x_i^t is made. Generally, there are two basic types of information feedback 31 that agents can possess. One is the full information feedback, where agents have access to the loss 32 functions. The other is the bandit feedback, where agents can only possess the values of the loss 33 function at points around the decision. At each time step, agents choose the decisions based on their 34 local feedback and neighbors' information. To measure the performance of an algorithm, the (static) 35 regret is frequently used, which compares the cumulative loss of online decisions and the loss of the 36

best decision chosen in hindsight through all the time horizons. The regret of node $j \in \mathcal{V}$ is defined 37 as $R(j,T) = \sum_{t=1}^{T} \sum_{i=1}^{N} f_i^t(x_j^t) - \min_{x \in \mathcal{K}} \sum_{t=1}^{T} \sum_{i=1}^{N} f_i^t(x)$. An algorithm is called *no-regret* [6] if the average regret over T goes to zero as T is large, which means the online decision updated 38 39 by the streaming data is not far from the best decision chosen in hindsight. Distributed no-regret 40 online algorithms have been widely studied in recent years [7, 8, 9]. 41

Although distributed algorithms are theoretically feasible, most of them are not practical as the model 42 size gets large, since communication cost can be a bottleneck for efficiency. In distributed training 43 tasks, agents can be powful microcomputers, while their communication network may be with low 44 bandwidth. The information exchange over the network is pretty slow compared with the computation 45 taking place in agents [10]. Thus, communication compression techniques are of significance for 46 practical implementations. 47

There have been many attempts to combine compressors with distributed optimization algorithms. 48 A straightforward idea is the direct compression scheme, while algorithms with this simple scheme 49 fail to converge even for the distributed average consensus problem [11, 12]. As an improvement, 50 extrapolation compression scheme and difference compression scheme are proposed [13]. Along 51 this line, quite a number of studies successfully extend the distributed optimization algorithms with 52 compressors and meanwhile maintain the convergence rate [14, 15, 16]. However, distributed online 53 optimization with compression is still an area that has not been fully exploited. [17] proposed ECD-54 AMSGrad algorithm that extended the AMSGrad to the distributed online setting with extrapolation 55 compression, while only empirical results were given without theoretical analysis. The key open 56 problem in this area is 57

whether it is possible to design provably no-regret distributed online algorithms 58 that work with compressors. 59

Contributions In this work, we answer the above question in the affirmative. We apply the 60 difference compression scheme to the fundamental GD-based distributed online algorithms. Although 61 the idea of such combination is simple, the underlying algorithm design and theoretical principle are 62 challenging since the compression error, projection error, and consensus error will be coupled. Our 63 contributions are summarized as follows: 64

• We propose communication-efficient distributed online algorithms, which consist of difference 65 compression, γ -gossip consensus, gradient descent, and projection, for the cases of full infor-66 mation feedback (DC-DOGD), one-point bandit feedback (DC-DOBD), and two-point bandit 67 feedback (DC-DO2BD), respectively. We make the technical advance to combine the difference 68 compression scheme with the projection scheme. Through proper design, the errors can be 69 estimated and controlled with the consensus stepsize γ and the gradient descent stepsizes. 70

• We analyze the regret bounds of the proposed algorithms for convex and strongly convex losses, 71 respectively, which are established explicitly in terms of the time horizon T, compression ratio 72 ω , decision dimension d, agent number N, and the parameters of the communication graph \mathcal{G} , as 73 simplified and summarized in Table 1. The obtained regret bounds are in accordance with that of 74 [18] w.r.t T, N, d. To the best of our knowledge, the proposed algorithms are the first distributed 75 online algorithms with theoretical no-regret guarantees for ω -contracted compressors. 76

• We give exhaustive experiments to illustrate the performance of the proposed algorithms. 77 Compared with the uncompressed algorithm DAOL [7], the proposed algorithms can reduce the 78 total transmitted bits for distributed online training. Moreover, DC-DOGD and DC-DO2BD 79 significantly outperform the algorithm ECD-AMSGrad [17]. 80

Table 1: Regret bounds in different settings		
Settings	convex losses	strongly convex losses
Full information	$\mathcal{O}\left(\left(\omega^{-2}N^{1/2}+\omega^{-4} ight)N\sqrt{T} ight)$	$\mathcal{O}\left(\left(\omega^{-2}N^{1/2}+\omega^{-4}\right)N\ln(T)\right)$
One-point bandit	$\mathcal{O}\left(\left(\omega^{-2}N^{1/2}+\omega^{-4}\right)^{1/2}Nd^{1/2}T^{3/4}\right)$	$\mathcal{O}\left(\left(\omega^{-2}N^{1/2}+\omega^{-4}\right)^{1/3}Nd^{2/3}T^{2/3}\ln^{1/3}(T)\right)$
Two-point bandit	$\mathcal{O}\left(\left(\omega^{-2}N^{1/2}+\omega^{-4} ight)Nd\sqrt{T} ight)^{-1}$	$\mathcal{O}\left(\left(\omega^{-2}N^{1/2}+\omega^{-4}\right)Nd^2\ln(T)\right)$

T 1 1 **D** 1 . 1.00

Related Work Distributed online convex optimization has received numerous attention in recent 81 years. Many basic algorithms have been extended to distributed settings. For example, [7] proposed 82 a distributed online subgradient algorithm over a static connected directed network and achieved 83 the regrets $\mathcal{O}(\sqrt{T})$ and $\mathcal{O}(\ln(T))$ for convex and strongly convex losses, which is in line with the 84 regrets of classical centralized online algorithms [19, 20]. Then [8] studied DOCO with long-term 85 constraints over a time-varying network and achieved the regrets $\mathcal{O}(T^{3/4})$ and $\mathcal{O}(T^{2/3} \ln^{1/3}(T))$ for 86 convex and strongly convex losses in the one-point bandit feedback. As for two-point bandit feedback, 87 the regrets $\mathcal{O}(\sqrt{T})$ and $\mathcal{O}(\ln(T))$ were established in [9] for convex and strongly convex losses, 88 89 which are the same as the centralized two-bandit algorithms [21]. [18] comprehensively studied DOCO over Erdős-Rényi random networks in full gradient feedback, one-point bandit feedback, and 90 two-point bandit feedback, and gave regret bounds. Along the line of [18], this paper aims to further 91 introduce compressed communication strategies, while preserving the regret bounds. 92

Recently, combining distributed optimization algorithms with compressors has seen a dramatic 93 rise in interest. Traditional compressors include the quantization and sparsification. Quantization 94 is to reduce the precision of each element, such as 1-bit SGD [22], SignSGD [23], GSGD [24], 95 and so on. Sparsification is to transmit only a few elements of the vectors, for instance, Top_k 96 [25], Rand_k [26] and Threshold_v [27]. Also there are hybrid compressors combining quantization 97 with sparsification, to name a few, SketchML [28], 3LC [29], Qsparse-local-SGD [30], etc. The 98 way to apply compressors is called a compression scheme. The most widely used compression 99 schemes in distributed optimization are extrapolation compression and difference compression [13]. 100 Extrapolation compression allows agents to compress the extrapolation between the last two local 101 states. Decentralized PSGD with extrapolation compression (ECD-PSGD) [13] was proved to 102 converge sublinearly and match the rate of its uncompressed case (D-PSGD). Difference compression, 103 which is also called CHOCO [14] or innovation compression [31], allows agents to add replicas of 104 neighboring states and compress the state-difference. There have been extensive successful designs 105 combining distribute optimization algorithms with difference compression, to name a few, DCD-106 PSGD [13] (based on PSGD), CHOCO-SGD [14] (base on gossip SGD), SPARQ-SGD [15] (based 107 on event-trigger), C-GT [16] (base on gradient tracking), and COLD[31] (based on NIDS), etc. Also, 108 the idea of difference compression have been widely adopted in federated learning [32, 33, 34]. In 109 this work, we use difference compression scheme to design communication-efficient distributed 110 online algorithms. 111

The results about the distributed online optimization with compression are quite limited. [17] poposed the ECD-AMSGrad algorithm which combined AMSGrad with extrapolation compression. Actually, the AMSGrad algorithm may not be a good choice for DOCO algorithm design since although AMSGrad itself is proved no-regret [35], a considerable performance gap still exists between AMSGrad and SGD [36]. Besides, the introduction of compression errors will further worsen the algorithm such that ECD-AMSGrad will lose the no-regret performance (seen Section 5). In this paper, we focus on the fundamental GD-based algorithms and give no-regret guarantees.

119 2 Full Information Feedback

In this section, we first introduce the multi-agent network and the compressor we use, and then propose a communication-efficient distributed online algorithm for the DOCO with full information feedback. Expected regret bounds will be given for both convex and strongly convex losses.

Graph The multi-agent network is described by an undirected graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{1, \ldots, N\}$ is the set of nodes, representing the set of agents, and $\mathcal{E} \subset \mathcal{V} \times \mathcal{V}$ is the set of edges. Let $A = [a_{ij}] \in \mathbb{R}^{N \times N}$ be the connectivity matrix of \mathcal{G} such that $a_{ij} = a_{ji}$. If $(v_i, v_j) \in \mathcal{E}$, then v_i and v_j can exchange information, and $a_{ij} = 0$ otherwise. It is worth noting that the communication is node-to-node in our distributed setting, and there is no central node. The graph in this paper satisfies the following assumption.

Assumption 1. The communication graph G is undirected and connected. Its connectivity matrix 130 $A \in [0,1]^{N \times N}$ is a symmetric doubly stochastic matrix.

131 **Compressor** A compressor $Q(\cdot) : \mathbb{R}^d \to \mathbb{R}^d$ is a mapping whose output can be usually encoded 132 with fewer bits than its input. In this paper, we consider a broad class of compressors with the

following general property, which has been widely considered in distributed optimization with 133 compression [14, 15, 31]. 134

Assumption 2. For some $\omega \in (0, 1]$, Q satisfies $\mathbb{E}_Q ||Q(x) - x||^2 \le (1 - \omega) ||x||^2, \forall x \in \mathbb{R}^d$, where 135 \mathbb{E}_Q denotes the expectation over the internal randomness of Q. 136

Compressors satisfy the above assumption are called ω -contracted, which include many important 137 compressors, such as Sparsification $Rand_k$ and Top_k [26], Random quantization $QSGD_s$ [24], 138

Randomized gossip [14], etc. 139

2.1 Algorithm design 140

In the full information feedback, the loss function f_i^t is revealed to node *i* at time *t* after the decision x_i^t is made. Then node *i* has access to the gradient value $\nabla f_i^t(x_i^t)$ and can use $g_i^t = \nabla f_i^t(x_i^t)$ to 141 142 make the next decision x_i^{t+1} . We propose the DC-DOGD algorithm as shown in Algorithm 1, which 143 is based on DAOL [7] and Memory-efficient CHOCO-GOSSIP [14]. The DC-DOGD algorithm 144 consists of two main parts: difference compressed communication (steps 2 and 3) and local decision 145 update (steps 4 and 5). 146

Algorithm 1 Distributed Online Gradient Descent with Difference Compression (DC-DOGD)

Input: consensus stepsize γ , GD stepsizes $\{\eta_t\}_{t=1}^T$, time T **Initialize:** set $x_i^1 = \mathbf{0}, \hat{x}_i^1 = \mathbf{0}, s_i^1 = \mathbf{0}$, for each node $i \in \mathcal{V}$.

- 1: for t = 1 to T 1 do in parallel for each node $i \in \mathcal{V}$
- 2:
- Compress the difference vector $q_i^t = Q(x_i^t \hat{x}_i^t)$ and update the replica $\hat{x}_i^{t+1} = \hat{x}_i^t + q_i^t$. Send q_i^t to its neighbors and receive q_j^t from all its neighbors $j \in \mathcal{N}_i$. Update the estimate of 3: the consensus decision by $s_i^{t+1} = s_i^t + \sum_{j \in \mathcal{N}_i} a_{ij} q_j^t$. Receive the full information feedback and calculate $g_i^t = \nabla f_i^t(x_i^t)$.
- 4:
- 5: Update its decision variable as follows

$$x_i^{t+1} = P_{\mathcal{K}} \left(x_i^t + \gamma (s_i^{t+1} - \hat{x}_i^{t+1}) - \eta_t g_i^t \right), \tag{1}$$

where $P_{\mathcal{K}}$ denotes the Euclidean projection, i.e., $P_{\mathcal{K}}(x) = \operatorname{argmin}_{y \in \mathcal{K}} ||x - y||$.

Output: $\{x_i^t\}_{t=1}^T$

The insights of introducing the variables \hat{x}_i^t and using the difference compression are as below. 147 Assume that $x^* \neq 0$ without loss of generality. Then, if node *i* transmits the directly compressed 148 information $Q(x_i^t)$ to its neighbors, the compression error $Q(x_i^t) - x_i^t$ will not vanish for $t \to \infty$. 149 The accumulation of compression errors makes the algorithm fail to converge. Instead, we compress 150 something that goes to zero. Let node i and all its neighbors keep an auxiliary variable \hat{x}_i^t locally, 151 which acts as a replica of x_i^t . Whenever node *i* updates its decision variable x_i^t , node *i* calculates the difference $x_i^t - \hat{x}_i^t$, compresses the difference $q_i^t = Q(x_i^t - \hat{x}_i^t)$, and sends the compressed information 152 153 q_i^t to its neighbors. After that, node i and all its neighbors update the local replica $\hat{x}_i^{t+1} = \hat{x}_i^t + q_i^t$. 154 When all nodes are reaching a consensus optimal decision, the updates of local decisions are small, 155 and the differences between the replica variables and the true decision variables are also small. Then 156 the compression errors are expected to vanish. 157

Local decision variables update through the gradient descent, γ -gossip, and projection, in order to 158 minimize the local loss function, keep consensus with neighbors, and remain in the feasible set \mathcal{K} , 159 respectively. For each node i, s_i^1 is initialized to **0**, and thus, $s_i^t = \sum_{j=1}^N a_{ij} \hat{x}_j^t$. Recall that \hat{x}_i^t tracks x_i^t , then s_i^t acts as node i's estimate of the consensus decision at time t. The γ -gossip protocol is 160 161 adopted to renovate the decision variable towards the consensus decision. The consensus stepsize 162 $\gamma \in (0, 1]$ is tunable to control the consensus speed, which also plays a crucial role in controlling the 163 compression error. 164

If there is no compression, i.e., using exact communication, then \hat{x}_i^{t+1} turns out to be x_i^t , and s_i^{t+1} becomes $\sum_{j=1}^N a_{ij} x_j^t$. Besides, take $\gamma = 1$, and then (1) reduces to 165 166

$$x_i^{t+1} = P_{\mathcal{K}}\left(x_i^t + \gamma \sum_{j \in \mathcal{N}_i} a_{ij}(x_j^t - x_i^t) - \eta_t \nabla f_i^t(x_i^t)\right) = P_{\mathcal{K}}\left(\sum_{j \in \mathcal{N}_i} a_{ij}x_j^t - \eta_t \nabla f_i^t(x_i^t)\right),$$

- which is the DAOL algorithm in [7].
- **Remark 1.** s_i^t is introduced for the memory-efficiency. Eq. (1) is equivalent to the update rule

$$x_{i}^{t+1} = P_{\mathcal{K}}\left(x_{i}^{t} + \gamma \sum_{j \in \mathcal{N}_{i}} a_{ij}(\hat{x}_{j}^{t+1} - \hat{x}_{i}^{t+1}) - \eta_{t}g_{i}^{t}\right).$$
(2)

If we adopt the update rule (2) together with $\hat{x}_j^{t+1} = \hat{x}_j^t + q_j^t$ for $j \in \mathcal{N}_i \cup \{i\}$ instead of steps 3 and 5, then each node have to store deg(i) + 2 vectors, namely, x_i, \hat{x}_i and $\hat{x}_j, j \in \mathcal{N}_i$, which is memory-consuming.

172 2.2 Regret bounds

- We consider the following assumptions, which are widely used in the studies of distributed online optimization [1, 18, 37].
- Assumption 3. The convex set \mathcal{K} is bounded with diameter D, i.e., $||x y|| \le D$, $\forall x, y \in \mathcal{K}$.

Assumption 4. For each $i \in \mathcal{V}$ and t = 1, 2, ..., T, the loss function f_i^t is convex and differentiable with bounded gradient over \mathcal{K} , i.e., $\max_{i,t,x} \|\nabla f_i^t(x)\| \le G$.

Assumption 5. For each $i \in \mathcal{V}$ and t = 1, 2, ..., T, the loss function f_i^t is μ -strongly convex over \mathcal{K} with the parameter $\mu > 0$, i.e., $f_i^t(x) - f_i^t(y) \ge \langle x - y, \nabla f_i^t(y) \rangle + \frac{\mu}{2} ||x - y||^2$, $\forall x, y \in \mathcal{K}$.

Suppose that the eigenvalues of the symmetric doubly stochastic connectivity matrix A are $1 = |\lambda_1(A)| > |\lambda_2(A)| \ge \cdots \ge |\lambda_N(A)|$. Define the spectral gap $\delta := 1 - |\lambda_2(A)| \in (0, 1]$ and the spectral radius of the Laplacian matrix $\beta := ||I_N - A||_2 \in [0, 2]$. Then we give the expected regret bounds of Algorithm 1 for convex and strongly convex losses, respectively.

Theorem 1. Let common Assumptions 1 and 2 hold. Consider Algorithm 1 with the consensus stepsize

$$\gamma = \frac{3\delta^3\omega^2(\omega+1)}{48(\delta^2 + 18\delta\beta^2 + 36\beta^2)\beta^2(\omega+2)(1-\omega) + 4\delta^2(\beta^2 + \beta)(\omega+2)(1-\omega)\omega + 6\delta^3\omega},$$
 (3)

(*i*) (Convex case) Under Assumptions 3 and 4, take the gradient descent stepsize $\eta_t = \frac{D}{G\sqrt{t+c}}$ for a constant $c \ge \frac{8}{3\gamma\delta}$. Then for each $j \in \mathcal{V}$ and $T \ge 1$,

$$\mathbb{E}_{Q}\left[\mathbf{R}(j,T)\right] \leq \left(\frac{1}{2} + 8\sqrt{3}\left(\sqrt{N} + 2\sqrt{3}\gamma^{-1}\delta^{-1}\right)\left(1 + \gamma^{-1}\delta^{-1} + \omega^{-1}\right)\right)NGD\sqrt{T+c}.$$
 (4)

(*ii*) (Strongly convex case) Under Assumptions 4 and 5, take the gradient descent stepsize $\eta_t = \frac{1}{\mu(t+c)}$ for a constant $c \ge \frac{16}{3\gamma\delta}$. Then for each $j \in \mathcal{V}$ and $T \ge 1$,

$$\mathbb{E}_{Q}\left[\mathbf{R}(j,T)\right] \le 4\sqrt{3}\left(\sqrt{N} + 2\sqrt{3}\gamma^{-1}\delta^{-1}\right)\left(1 + \gamma^{-1}\delta^{-1} + \omega^{-1}\right)NG^{2}\mu^{-1}\ln(T+c).$$
 (5)

The proof ideas are as follows. Firstly, we estimate the general regret bounds for each node, which depend on the consensus error, the projection error, the compression error, and the gradient descent stepsize. Then comes the key points that we analyze the coupled relationship between the errors, and bound them with the consensus stepsize γ and the GD stepsize η_t . Finally, we choose proper γ and η_t to obtain Theorem 1. Complete proofs are attached to Appendix B.

The consensus stepsize γ chosen in (3) depends on the compression ratio ω and the communication graph paremeters δ and β . Notice that γ is an increasing function with respect to ω , and $\gamma|_{\omega=0} =$ $0, \gamma|_{\omega=1} = 1$. Thus, $\gamma \in (0, 1]$ for $\omega \in (0, 1]$. If there is no compression ($\omega = 1$), then $\gamma = 1$, and Algorithm 1 exactly reduces to DAOL [7], as mentioned in the algorithm design.

Theorem 1 shows that Algorithm 1 achieves the regret bounds $\mathcal{O}((\omega^{-2}N^{1/2} + \omega^{-4})N\sqrt{T})$ and $\mathcal{O}((\omega^{-2}N^{1/2} + \omega^{-4})N\ln(T))$ for convex losses and strongly convex losses, respectively. The results suggest that

- Algorithm 1 is no-regret in both convex case and strongly convex case, since the time averaged regret $\mathbb{E}_{Q}[\mathrm{R}(j,T)]/T \to 0$ for $T \to \infty$. The obtained regret bounds $\mathcal{O}(\sqrt{T})$ and $\mathcal{O}(\ln(T))$ are in
- accordance with that of the centralized online algorithms in the respective cases [19, 20].

• The node averaged regret $\mathbb{E}_{Q}[R(j,T)]/N$ increases with N, which in line with the results in [18]. 205

• As the compression ratio ω decreases, fewer bits are needed for node-to-node communication in 206

each iteration, while more iteration rounds are needed to reach the desired regret. ω can be used 207 to balance the iteration rounds and the transmitted bits in each iteration from multiple concerns

208 such as the bandwidth and agent computation capability. In practice, we can choose a proper ω 209

to minimize the total transmitted bits or minimize the overall training time. 210

3 **One-point Bandit Feedback** 211

In this section, we apply difference compression to DOCO with one-point bandit feedback. We 212 propose DC-DOBD algorithm, which basically follows DC-DOGD, except for the gradient estimation. 213

Algorithm 2 Distributed Online One-point Bandit Gradient Descent with Difference Compression (DC-DOBD)

Input: consensus stepsize γ , GD stepsizes $\{\eta_t\}_{t=1}^T$, time T, exploration parameter ϵ , shrinkage parameter ζ

Initialize: set $x_i^1 = \mathbf{0}, \hat{x}_i^1 = \mathbf{0}, s_i^1 = \mathbf{0}$, for each node $i \in \mathcal{V}$.

1: for t = 1 to T - 1 do in parallel for each node $i \in \mathcal{V}$

- Compress the difference vector $q_i^t = Q(x_i^t \hat{x}_i^t)$ and update $\hat{x}_i^{t+1} = \hat{x}_i^t + q_i^t$. Spread q_i^t and receive $q_j^t, j \in \mathcal{N}_i$. Update $s_i^{t+1} = s_i^t + \sum_{j \in \mathcal{N}_i} a_{ij} q_j^t$. 2:
- 3:
- 4:
- Receive the one-point bandit feedback and construct $g_i^t = \frac{d}{dt} f_i^t (x_i^t + \epsilon u_i^t) u_i^t$. Update the decision variable $x_i^{t+1} = P_{(1-\zeta)\mathcal{K}} \left(x_i^t + \gamma(s_i^{t+1} \hat{x}_i^{t+1}) \eta_t g_i^t \right)$. 5:

Output: $\{x_i^t\}_{t=1}^T$

In the one-point bandit feedback, after making the decision x_i^t at time t, agent i can query the loss 214 function value at one point around x_i^t and use the feedback to construct the gradient estimator g_i^t . 215 Like the procedure in [38], let agent i choose a unit-norm vector $u_i^t \in \mathbb{R}^d$ uniformly at random, 216

query the value of f_i^t at the point $y_i^t = x_i^t + \epsilon u_i^t$, and calculate $g_i^t = \frac{d}{\epsilon} f_i^t(y_i^t) u_i^t$. Since the loss 217 function f_i^t is defined in the set \mathcal{K} , we slightly modify the projection in (1) as $P_{(1-\zeta)\mathcal{K}}$ to ensure 218 the query point $y_i^t \in \mathcal{K}$. Algorithm 2 actually performs the gradient descent on the function 219 $\hat{f}_i^t(x) = \mathbb{E}_{u \in \mathcal{B}} [f_i^t(x + \epsilon u)]$ restricted to the convex set $(1 - \zeta)\mathcal{K}$. It has been shown by [38] that 220 $\mathbb{E}[g_i^t] = \nabla \hat{f}_i^t(x_i^t)$. In the bandit setting, Assumptions 3 and 4 are modified as follows, which are 221 commonly used in online bandit optimization [38, 21, 18]. 222

Assumption 6. The convex set \mathcal{K} contains the ball of radius r centered at the origin, and is contained 223 in the ball of radius R, i.e., $r\mathcal{B} \subseteq \mathcal{K} \subseteq R\mathcal{B}, \mathcal{B} = \{u \in \mathbb{R}^d : ||u|| \leq 1\}.$ 224

Assumption 7. For each $i \in \mathcal{V}$ and t = 1, 2, ..., T, the loss function f_i^t is convex and *l*-lipschitz continuous in \mathcal{K} , i.e., $|f_i^t(x) - f_i^t(y)| \le l ||x - y||, \forall x, y \in \mathcal{K}$. 225 226

Assumptions 6 and 7 lead to an uniform upper bound on the function value, i.e., there exists a constant 227

B > 0 such that $\max_{x,i,t} |f_i^t(x)| \le B$. Then we establish the expected regret bounds of Algorithm 2 228 for convex and strongly convex losses, respectively. 229

Theorem 2. Let common Assumptions 1, 2, 6 and 7 hold. Consider Algorithm 2 with the consensus 230 stepsize γ chosen in (3). Denote 231

$$H = 4\sqrt{3} \left(\sqrt{N} + 2\sqrt{3}\gamma^{-1}\delta^{-1}\right) \left(1 + \gamma^{-1}\delta^{-1} + \omega^{-1}\right).$$
 (6)

(i) (Convex case) Take the gradient descent stepsize $\eta_t = \frac{2R\epsilon}{dB\sqrt{t+c}}$ for a constant $c \geq \frac{8}{3\gamma\delta}$, $\epsilon =$ 232 $\sqrt{1}$

233
$$\left(\frac{(1+4H)dBR}{2(l+\frac{B}{r})}\right)^2 \frac{(T+c)^{\frac{1}{4}}}{T^{\frac{1}{2}}} and \zeta = \frac{\epsilon}{r}.$$
 Then for each $j \in \mathcal{V}$ and $T \ge 1$,

$$\mathbb{E}\left[\mathrm{R}(j,T)\right] \le 2NT^{\frac{1}{2}}(T+c)^{\frac{1}{4}}\sqrt{2(1+4H)(l+B/r)dBR}.$$
(7)

(ii) (Strongly convex case) With additional Assumptions 5, take the gradient descent stepsize
$$\eta_t =$$

235
$$\frac{1}{\mu(t+c)}$$
 for a constant $c \ge \frac{16}{3\gamma\delta}$, $\epsilon = \left(\frac{Hd^2B^2\ln(T+c)}{(l+\frac{B}{r})\mu T}\right)^{\frac{1}{3}}$ and $\zeta = \frac{\epsilon}{r}$. Then for each $j \in \mathcal{V}$ and $T \ge 1$,
 $\mathbb{E}\left[\mathrm{R}(j,T)\right] \le 3N \left(Hd^2B^2\mu^{-1}\right)^{\frac{1}{3}} (l+B/r)^{\frac{2}{3}} T^{\frac{2}{3}} \ln^{\frac{1}{3}}(T+c).$
(8)

Theorem 2 shows that Algorithm 2 is also no-regret, and it achieves the regret bounds 236 $\mathcal{O}(d^{1/2}N^{5/4}T^{3/4})$ and $\mathcal{O}(d^{2/3}N^{7/6}T^{2/3}\ln^{1/3}(T))$ for convex losses and strongly convex losses, respectively, which match the bounds obtained by [18]. The regrets are scaled with 237 238 $(\omega^{-2}N^{1/2} + \omega^{-4})^{1/2}$ and $(\omega^{-2}N^{1/2} + \omega^{-4})^{1/3}$ for convex and strongly convex losses, which in-239 dicates that the influence of the compression ratio ω on the regret bounds in the one-point bandit 240 setting is less than that in the full information setting. 241

Two-point Bandit Feedback 4 242

In the two-point bandit feedback, agent i can query the loss function values at two points around 243

 x_i^t . Like the procedure in [21], let agent *i* bick a unit-norm vector $u_i^t \in \mathbb{R}^d$ uniformly at random, query the values of f_i^t at $y_{i,1}^t = x_i^t + \epsilon u_i^t$ and $y_{i,2}^t = x_i^t - \epsilon u_i^t$, and estimate the gradient as $g_i^t = \frac{d}{2\epsilon} \left(f_i^t(y_{i,1}^t) - f_i^t(y_{i,2}^t) \right) u_i^t$. Then we obtain DC-DO2BD as a variant of DC-DOBD. 244

245

246

Algorithm 3 Distributed Online Two-point Bandit Gradient Descent with Difference Compression (DC-DO2BD)

Input: consensus stepsize γ , GD stepsizes $\{\eta_t\}_{t=1}^T$, time T, exploration parameter ϵ , shrinkage parameter ζ

Initialize: set $x_i^1 = \mathbf{0}, \hat{x}_i^1 = \mathbf{0}, s_i^1 = \mathbf{0}$, for each node $i \in \mathcal{V}$.

- 1: for t = 1 to T 1 do in parallel for each node $i \in \mathcal{V}$
- Compress the difference vector $q_i^t = Q(x_i^t \hat{x}_i^t)$ and update $\hat{x}_i^{t+1} = \hat{x}_i^t + q_i^t$. Spread q_i^t and receive q_j^t , $j \in \mathcal{N}_i$. Update $s_i^{t+1} = s_i^t + \sum_{j \in \mathcal{N}_i} a_{ij} q_j^t$. 2:
- 3:
- Receive the two-point feedback and construct $g_i^t = \frac{d}{2\epsilon} \left(f_i^t (x_i^t + \epsilon u_i^t) f_i^t (x_i^t \epsilon u_i^t) \right) u_i^t$. Update the decision variable $x_i^{t+1} = P_{(1-\zeta)\mathcal{K}} \left(x_i^t + \gamma(s_i^{t+1} \hat{x}_i^{t+1}) \eta_t g_i^t \right)$. 4:
- 5:

Output: $\{x_i^t\}_{t=1}^T$

- In the two-point bandit setting, the regret of node $j \in \mathcal{V}$ is modified as $R_2(j,T) =$ 247 $\sum_{t=1}^{T} \sum_{i=1}^{N} \frac{f_i^t(y_{i,1}^t) + f_i^t(y_{i,2}^t)}{2} - \sum_{t=1}^{T} \sum_{i=1}^{N} f_i^t(x^*).$ 248
- Theorem 3. Let common Assumptions 1, 2, 6 and 7 hold. Consider Algorithm 3 with the consensus 249 stepsize γ chosen in (3) and H defined in (6). 250

(i) (Convex case) Take the gradient descent stepsize $\eta_t = \frac{2R}{dl\sqrt{t+c}}$ for a constant $c \geq \frac{8}{3\gamma\delta}$, $\epsilon = \frac{1}{\sqrt{T}}$ 251 and $\zeta = \frac{\epsilon}{r}$. Then for each $j \in \mathcal{V}$ and $T \geq 1$, 252

$$\mathbb{E}\left[\mathrm{R}_{2}(j,T)\right] \le (1+4H)RNdl\sqrt{T+c} + (3+2R/r)Ndl\sqrt{T}.$$
(9)

(ii) (Strongly convex case) With additional Assumptions 5, take the gradient descent stepsize $\eta_t =$ 253 $\frac{1}{\mu(t+c)}$ for a constant $c \geq \frac{16}{3\gamma\delta}$, $\epsilon = \frac{\ln(T)}{T}$ and $\zeta = \frac{\epsilon}{r}$. Then for each $j \in \mathcal{V}$ and $T \geq 1$, 254

$$\mathbb{E}\left[\mathrm{R}_{2}(j,T)\right] \leq \mu^{-1} N d^{2} l^{2} H \ln(T+c) + (3+2R/r) N dl \ln(T).$$
(10)

Theorem 3 shows that Algorithm 3 achieves $\mathcal{O}((\omega^{-2}N^{1/2} + \omega^{-4})Nd\sqrt{T})$ and $\mathcal{O}((\omega^{-2}N^{1/2} + \omega^{-4})Nd\sqrt{T})$ 255 ω^{-4})Nd² ln(T)) regret bounds for convex and strongly convex losses, respectively, which recovers 256 the regret bounds $\mathcal{O}(\sqrt{T})$ (convex) and $\mathcal{O}(\ln(T))$ (strongly convex) in the full information case, 257 while the constants are larger than those of Theorem 1. 258

Numerical Experiments 5 259

In this section, we evaluate the three proposed algorithms on a real-world online problem. A 260 prominent example is the diabetes prediction task, which aims to build a model to diagnose diabetes 261 through several risk factors. Consider the distributed online regularized logistic regression with the 262 local loss function 263

$$f_{i}^{t}(x) = \sum_{j=1}^{S} \log\left(1 + \exp\left(-b_{i,j}^{t}\left\langle a_{i,j}^{t}, x\right\rangle\right)\right) + \frac{\mu}{2} \left\|x\right\|^{2},$$
(11)

where μ is the regularization parameter, and a batch data samples $\{(a_{i,j}, b_{i,j})\}_{j=1}^{S}$ are revealed to 264 agent i at time t. We adopt diabetes-binary-BRFSS2015 dataset with 70692 instances, 21 features, 265 and 2 labels from Kaggle¹. Here, $a_{i,j} \in \mathbb{R}^d$ with d = 21, and $b_{i,j} \in \{-1,1\}$. We standardize the 266 data samples and distribute them evenly among N agents under the sorted setting, i.e., each agent 267 only gets data samples from one class. The connected communication network $\mathcal{G}(N, M)$ with N 268 nodes and M edges is generated randomly by tool NetworkX [39], and then we use the Metropolis 269 rule [40] to construct the connectivity matrix A to satisfy Assumption 1. We repeat each experiment 270 ten times and depict the mean curve 2 . The choice of parameters is given in Appendix E. 271

Comparison experiment We run our algorithms DC-DOGD, DC-DOBD, DC-DO2BD, and make comparisons with ECD-AMSGrad [17], for the convex case ($\mu = 0$) and strongly convex case ($\mu = 1$). The compressor type, the compression ratio, and the communication network are kept the same. Take the setting of QSGD₂ with $\omega = 0.3$ over $\mathcal{G}(9, 18)$ as an example. We plot the time averaged maximum regret $SR(T) := \max_j R(j,T)/T$ versus the time horizon T and versus the total number of transmitted bits in Fig. 1, where the best solution in the hindsight x^* is obtained by *LogisticRegression* optimizer from scikit-learn [41].

Fig. 1 shows that the time averaged regrets of DC-DOGD, DC-DOBD, and DC-DO2BD go to 279 zero as T goes to infinity, which is in agreement with the theoretical results that our algorithms are 280 no-regret. Among the three proposed algorithms, the one-bandit feedback has the worst performance, 281 while using two-bandit information can improve the performance and even reach that of the full 282 283 information feedback. ECD-AMDGrad gets deteriorated in the first few steps because the inverse of the second raw moment estimation is large and this algorithm does not have a projection to restrict 284 variables. Then, ECD-AMDGrad declines fast, while its time-average regret can not reach zero. 285 Clearly, DC-DOGD and DC-DO2BD significantly outperform ECD-AMDGrad. 286



Figure 1: Comparison of algorithms DC-DOGD, DC-DOBD, DC-DO2BD, and ECD-AMSGrad with $QSGD_2$, $\omega = 0.3$, $\mathcal{G}(9, 18)$.

Impact of compression ratio and compressor type Fixing the compressor type (Top_k) and the 287 graph $\mathcal{G}(9, 18)$, we run DC-DOGD with different compression ratios ($\omega = 0.05, 0.1, 0.5$) for strongly 288 convex losses³. As the baseline we consider DAOL [7], which is with exact communication and is 289 the special case of DC-DOGD with $\omega = 1$. The greater the compression degree (less ω), the more 290 iteration rounds are needed to reach the ε average regret as Fig. 2a shows, while the total transmitted 291 292 bits are actually the fewer as Fig. 2b shows. DC-DOGD with $\omega = 0.5$ performs almost as good as DAOL while using $2 \times$ less total bits to reach the ε average regret. DC-DOGD with $\omega = 0.05$ have 293 approximately $10 \times$ reduction on bits to reach the ε average regret compared with DAOL. 294

Then, we fix the compression ratio ($\omega = 0.3$) and the graph $\mathcal{G}(9, 18)$, and run DC-DOGD with different compressors (Top_k, Rand_k, RGossip_p, GSGD_s) for strongly convex losses. Figs. 2c and 2d show that Rand_k and RGossip_p have almost the same performance. It is expected, since in Rand_k each element of the vector has the probability $\omega = k/d$ to be chosen to be transmitted, which is equivalent to randomly transmitting the whole vector with the probability $p = \omega$. Besides, Figs. 2c

¹The data set is from https://www.kaggle.com/code/encode0/diabetes-prediction-and-risk-factors-evaluation.

²All experiments are performed on a 64-bit Windows platform with the Intel(R) Core(TM) i7-6850K 3.6Ghz CPU. The codes are provided in the supplementary materials.

³Since the trajectories in the convex case and strongly convex case share similar trends, we only present the experiment results in the strongly convex case, for space limitation.

and 2d show that Top_k has better performance than $Rand_k$, which is in line with the intuition that

the largest k coordinates contain more useful information than arbitrary k coordinates. In addition, quantization $GSGD_s$ performs better in reducing the total transmitted bits than sparsification under the same compression ratio.



Figure 2: The impact of compression ratio and compressor type for DC-DOGD over $\mathcal{G}(9, 18)$ in the strongly convex case.

303

Impact of topology and node number The network topology concerns the parameters δ and β . 304 which influence the choice of the consensus stepsize γ as well as the algorithm performance. A 305 simpler topology with fewer edges needs fewer bits to transmit information in each iteration, while 306 more iteration rounds are needed for decision consensus. Thus, there will be a tradeoff. We take 307 DC-DOGD in the strongly convex case as an example. Fixing the compressor (Top_1) with the 308 compression ratio ($\omega = 0.05$), we assess DC-DOGD over three basic topologies (ring, $\mathcal{G}(N, 2N)$), 309 full connected). Fig. 3a shows that the full connected graph uses the smallest iteration round to reach 310 the ε average regret, while Fig. 3b illustrates that the total numbers of transmitted bits to reach the ε 311 average regret are close. 312

Finally, we let the node number N vary from to 10 to 50, and run DC-DOGD, DC-DOBD, DC-

DO2BD with the same compressor Top_2 , the same compression ratio $\omega = 0.1$, over the full connected graph, for convex losses and strongly convex losses. We plot the node averaged regret $AR(T) := \max_{j \in [N,T]/N} N$ versus the node number N in Figs. 3c and 3d.



Figure 3: The impact of topology and node number.

316

317 6 Conclusions

In this paper, we considered DOCO with the full information feedback, one-point and two-points 318 bandit feedback. We designed provably no-regret distributed online algorithms that work with ω -319 contracted compressors. The obtained regret bounds for both convex and strongly convex losses 320 matched those of uncompressed algorithms in the literature. We further assessed the influence of the 321 compressor type and the compression ratio ω on the regrets, and showed that ω can be used to balance 322 the iteration rounds and the transmitted bits according to the bandwidth. The limitation of this work is 323 that the obtained regret bounds show high order inverse dependence on the compression ratio, which 324 are pretty conservative and may be further improved. We believe this paper is an important step in 325 this direction. Future research includes designing provably no-regret distributed online algorithms 326 with other compression schemes such as extrapolation compression. 327

328 **References**

- [1] Elad Hazan et al. Introduction to online convex optimization. *Foundations and Trends* (R) *in Optimization*, 2(3-4):157–325, 2016.
- [2] David Sculley and Gabriel M Wachman. Relaxed online svms for spam filtering. In *Proceedings of the 30th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*,
 pages 415–422, 2007.
- [3] Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online dictionary learning for sparse
 coding. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pages 689–696,
 2009.
- [4] Deming Yuan, Alexandre Proutiere, and Guodong Shi. Distributed online linear regressions. *IEEE Transactions on Information Theory*, 67(1):616–639, 2020.
- [5] Wenpeng Zhang, Peilin Zhao, Wenwu Zhu, Steven CH Hoi, and Tong Zhang. Projection-free distributed
 online learning in networks. In *International Conference on Machine Learning*, pages 4054–4062. PMLR,
 2017.
- [6] Niranjan Srinivas, Andreas Krause, Sham Kakade, and Matthias Seeger. Gaussian process optimization in
 the bandit setting: No regret and experimental design. In *Proceedings of the 27th International Conference on Machine Learning*, number CONF. Omnipress, 2010.
- [7] Feng Yan, Shreyas Sundaram, SVN Vishwanathan, and Yuan Qi. Distributed autonomous online learn ing: Regrets and intrinsic privacy-preserving properties. *IEEE Transactions on Knowledge and Data Engineering*, 25(11):2483–2493, 2012.
- [8] Deming Yuan, Alexandre Proutiere, and Guodong Shi. Distributed online optimization with long-term
 constraints. *IEEE Transactions on Automatic Control*, 67(3):1089–1104, 2022.
- [9] Jueyou Li, Chaojie Li, Wenwu Yu, Xiaomei Zhu, and Xinghuo Yu. Distributed online bandit learning in
 dynamic environments over unbalanced digraphs. *IEEE Transactions on Network Science and Engineering*,
 8(4):3034–3047, 2021.
- [10] Dan Alistarh. A brief tutorial on distributed and concurrent machine learning. In *Proceedings of the 2018 ACM Symposium on Principles of Distributed Computing*, pages 487–488, 2018.
- [11] Tuncer Can Aysal, Mark J Coates, and Michael G Rabbat. Distributed average consensus with dithered
 quantization. *IEEE Transactions on Signal Processing*, 56(10):4905–4918, 2008.
- Ruggero Carli, Fabio Fagnani, Paolo Frasca, Tom Taylor, and Sandro Zampieri. Average consensus on networks with transmission noise or quantization. In 2007 European Control Conference (ECC), pages 1852–1857. IEEE, 2007.
- [13] Hanlin Tang, Shaoduo Gan, Ce Zhang, Tong Zhang, and Ji Liu. Communication compression for
 decentralized training. *Advances in Neural Information Processing Systems*, 31, 2018.
- [14] Anastasia Koloskova, Sebastian Stich, and Martin Jaggi. Decentralized stochastic optimization and gossip
 algorithms with compressed communication. In *International Conference on Machine Learning*, pages
 3478–3487. PMLR, 2019.
- [15] Navjot Singh, Deepesh Data, Jemin George, and Suhas Diggavi. Sparq-sgd: Event-triggered and compressed communication in decentralized optimization. *IEEE Transactions on Automatic Control*, 2022.
- [16] Yiwei Liao, Zhuorui Li, Kun Huang, and Shi Pu. Compressed gradient tracking methods for decentralized
 optimization with linear convergence. *arXiv preprint arXiv:2103.13748*, 2021.
- [17] Guangxia Li, Jia Liu, Xiao Lu, Peilin Zhao, Yulong Shen, and Dusit Niyato. Decentralized online
 learning with compressed communication for near-sensor data analytics. *IEEE Communications Letters*,
 25(9):2958–2962, 2021.
- [18] Jinlong Lei, Peng Yi, Yiguang Hong, Jie Chen, and Guodong Shi. Online convex optimization over
 erdős-rényi random networks. *Advances in neural information processing systems*, 33:15591–15601, 2020.
- [19] Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceed- ings of the 20th International Conference on Machine Learning*, pages 928–936, 2003.
- [20] Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization.
 Machine Learning, 69(2):169–192, 2007.

- [21] Alekh Agarwal, Ofer Dekel, and Lin Xiao. Optimal algorithms for online convex optimization with
 multi-point bandit feedback. In *Proceedings of the 23rd Annual Conference on Learning Theory*, pages
 28–40. Citeseer, 2010.
- [22] Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. 1-bit stochastic gradient descent and its
 application to data-parallel distributed training of speech dnns. In *Fifteenth Annual Conference of the International Speech Communication Association*. Citeseer, 2014.
- [23] Jeremy Bernstein, Yu-Xiang Wang, Kamyar Azizzadenesheli, and Animashree Anandkumar. signsgd:
 Compressed optimisation for non-convex problems. In *International Conference on Machine Learning*,
 pages 560–569. PMLR, 2018.
- [24] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication efficient sgd via gradient quantization and encoding. *Advances in Neural Information Processing Systems*,
 30, 2017.
- [25] Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. In
 Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages
 440–445, 2017.
- [26] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory. *Advances in Neural Information Processing Systems*, 31, 2018.
- [27] Aritra Dutta, El Houcine Bergou, Ahmed M Abdelmoniem, Chen-Yu Ho, Atal Narayan Sahu, Marco Canini,
 and Panos Kalnis. On the discrepancy between the theoretical analysis and practical implementations
 of compressed communication for distributed deep learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 3817–3824, 2020.
- Jiawei Jiang, Fangcheng Fu, Tong Yang, and Bin Cui. Sketchml: Accelerating distributed machine learning
 with data sketches. In *Proceedings of the 2018 International Conference on Management of Data*, pages
 1269–1284, 2018.
- [29] Hyeontaek Lim, David G Andersen, and Michael Kaminsky. 3lc: Lightweight and effective traffic
 compression for distributed machine learning. *Proceedings of Machine Learning and Systems*, 1:53–64,
 2019.
- [30] Debraj Basu, Deepesh Data, Can Karakus, and Suhas Diggavi. Qsparse-local-sgd: Distributed sgd with
 quantization, sparsification and local computations. *Advances in Neural Information Processing Systems*,
 32, 2019.
- [31] Jiaqi Zhang, Keyou You, and Lihua Xie. Innovation compression for communication-efficient distributed
 optimization with linear convergence. *arXiv preprint arXiv:2105.06697*, 2021.
- 410 [32] Peter Richtárik, Igor Sokolov, and Ilyas Fatkhullin. Ef21: A new, simpler, theoretically better, and 411 practically faster error feedback. *Advances in Neural Information Processing Systems*, 34, 2021.
- [33] Constantin Philippenko and Aymeric Dieuleveut. Preserved central model for faster bidirectional compression in distributed settings. *Advances in Neural Information Processing Systems*, 34, 2021.
- [34] Amirhossein Reisizadeh, Aryan Mokhtari, Hamed Hassani, Ali Jadbabaie, and Ramtin Pedarsani. Fedpaq:
 A communication-efficient federated learning method with periodic averaging and quantization. In International Conference on Artificial Intelligence and Statistics, pages 2021–2031. PMLR, 2020.
- [35] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- [36] Nitish Shirish Keskar and Richard Socher. Improving generalization performance by switching from adam
 to sgd. *arXiv preprint arXiv:1712.07628*, 2017.
- [37] Elad Hazan, Alexander Rakhlin, and Peter Bartlett. Adaptive online gradient descent. *Advances in Neural Information Processing Systems*, 20, 2007.
- [38] Abraham D Flaxman, Adam Tauman Kalai, and H Brendan McMahan. Online convex optimization in
 the bandit setting: gradient descent without a gradient. In *Proceedings of the 16th Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 385–394, 2005.
- [39] Aric Hagberg, Pieter Swart, and Daniel S Chult. Exploring network structure, dynamics, and function
 using networkx. Technical report, Los Alamos National Lab.(LANL), Los Alamos, NM (United States),
 2008.

- [40] Lin Xiao and Stephen Boyd. Fast linear iterations for distributed averaging. Systems & Control Letters,
 53(1):65–78, 2004.
- [41] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel,
 Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in
 python. *the Journal of Machine Learning Research*, 12:2825–2830, 2011.

434 Checklist

435	1. For all authors
436	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's
437	contributions and scope? [Yes]
438	(b) Did you describe the limitations of your work? [Yes] See Section 6 Conclusions.
439	(c) Did you discuss any potential negative societal impacts of your work? [No] This
440	work is a theoretical finding to solve the communication bottleneck of the distributed
441	algorithms. Therefore, this work does not present foreseeable societal impacts.
442	(d) Have you read the ethics review guidelines and ensured that your paper conforms to
443	them? [Yes]
444	2. If you are including theoretical results
445 446	(a) Did you state the full set of assumptions of all theoretical results? [Yes] See Assumptions 1, 2, 3, 4, 5, 6, and 7.
447	(b) Did you include complete proofs of all theoretical results? [Yes] We present the proof
448	ideas in the paper and put the complete proofs in supplemental materials due to space
449	limitation.
450	3. If you ran experiments
451	(a) Did you include the code, data, and instructions needed to reproduce the main experi-
452	mental results (either in the supplemental material or as a URL)? [Yes] Codes are in
453	the supplemental material. Data are from Kaggle with URL. Instructions are included
454	in Section 5.
455	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
456	were chosen)? [Yes] See Section 5 and Appendix E
457	(c) Did you report error bars (e.g., with respect to the random seed after running experi-
458	ments multiple times)? [No] We repeat each experiment ten times and only depict the
459	mean curves to make the figures easier to read. Actually, our experiment results with
460	different random seeds are similar.
461 462	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See Section 5 footnote 2
463	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
464	(a) If your work uses existing assets, did you cite the creators? [Yes] We use the tools
465	NetworkX and scikit-learn with citation. We use the open dateset from Kaggle with
466	URL.
467	(b) Did you mention the license of the assets? [Yes] See Appendix F
468	(c) Did you include any new assets either in the supplemental material or as a URL? [No]
469	(d) Did you discuss whether and how consent was obtained from people whose data you're
470	using/curating? [No]
471	(e) Did you discuss whether the data you are using/curating contains personally identifiable
472	information or offensive content? [No]
473	5. If you used crowdsourcing or conducted research with human subjects
474	(a) Did you include the full text of instructions given to participants and screenshots, if
475	applicable? [N/A]
476	(b) Did you describe any potential participant risks, with links to Institutional Review
477	Board (IRB) approvals, if applicable? [N/A]
478	(c) Did you include the estimated hourly wage paid to participants and the total amount
479	spent on participant compensation? [N/A]