

---

# Sparsified Distributed Adaptive Learning with Error Feedback

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 To be completed...

## 2 1 Introduction

3 Some related work:

4 [18] develops variant of signSGD (as a biased compression schemes) for distributed optimization.  
5 Contributions are mainly on this error feedback variant. In [26], the authors provide theoretical  
6 results on the convergence of sparse Gradient SGD for distributed optimization (we want that for  
7 AMS here). [27] develops a variant of distributed SGD with sparse gradients too. Contributions  
8 include a memory term used while compressing the gradient (using top k for instance). Speeding up  
9 the convergence in  $\frac{1}{T^3}$ .

## 10 2 Preliminaries

### 11 Sparse Optimization Methods.

12 **Distributed Learning.** When a large number of compute engines is available, being able to  
13 train global machine learning models while mutualizing the available and *decentralized* source of  
14 computation has been a growing focus for the community.

15 Decentralized optimization methods include methods such as ADMM [6], Distributed Subgradient  
16 Descent [24], Dual Averaging [11], Prox-PDA [14], GNSD [21], and Choco-SGD [20].

17 A recent work [7], which focuses on adaptive gradient methods, namely the Adam [19] and the  
18 AMSGrad [25] optimization methods, develops a decentralized variant of gradient based and adap-  
19 tive methods in the context of gossip protocols. To date, very few contributions provided attempt  
20 to efficiently run adaptive gradient method in such a distributed setting. Apart from [7], (author?)  
21 [23] proposes a decentralized version of AMSGrad [25] which provably satisfies some non-standard  
22 regret. Though, no sparsified variants of them have been proposed for practical purposes nor been  
23 studied in the literature.

24 **Compression-Based Distributed Optimization.** While the capabilities of the compute powers  
25 is exploding, the communication complexity between either the central server and the decentralized  
26 workers or among workers is becoming ineffectively large [9, 22]. Gradient sparsification con-  
27 stitutes one popular method to induce sparsity through the optimization procedure and reduce the  
28 number of bits transmitted at each iteration. Extensive works have studied this technique to improve  
29 the communication efficiency of SGD-based methods such as distributed SGD. This large class of  
30 sparsification techniques include gradient quantization leveraging quantized vector of gradients in  
31 the communication phase [2, 29, 16, 28, 13, 8, 15], gradient sparsification generally selection top

k components of the vector to be communicated, see [27, 1], or variants of the particular SGD algorithm such as low-precision SGD [4, 18] proposing a trade-off between communication cost and precision, and signSGD [10, 30] where only the signs of the gradient vectors are communicated. Most of these works apply to the SGD method [5] as a prototype where a novel method and some convergence results are presented with a rate of  $\mathcal{O}(\frac{1}{\sqrt{T}})$  where  $T$  denotes the total number of iterations, see [3], thus achieving the same rate as plain SGD, see [12, 17].

Yet these communication reduction techniques, still presents a negative dependence on the number of workers, typically a linear dependence. Hence the need for even more efficient techniques which constitutes the object of our paper.

### 3 Method

Consider standard synchronous distributed optimization setting. AMSGrad is used as the prototype, and the local workers is only in charge of gradient computation.

#### 3.1 TopK AMSGrad with Error Feedback

The key difference (and interesting part) of our TopK AMSGrad compared with the following arxiv paper “Quantized Adam”<https://arxiv.org/pdf/2004.14180.pdf> is that, in our model only gradients are transmitted. In “QAdam”, each local worker keeps a local copy of moment estimator  $m$  and  $v$ , and compresses and transmits  $m/v$  as a whole. Thus, that method is very much like the sparsified distributed SGD, except that  $g$  is changed into  $m/v$ . In our model, the moment estimates  $m$  and  $v$  are computed only at the central server, with the compressed gradients instead of the full gradient. This would be the key (and difficulty) in convergence analysis.

---

#### Algorithm 1 SPARS-AMS for Federated Learning

---

```

1: Input: parameter  $\beta_1, \beta_2$ , learning rate  $\eta_t$ .
2: Initialize: central server parameter  $\theta_0 \in \Theta \subseteq \mathbb{R}^d$ ;  $e_{t,i} = 0$  the error accumulator for each
   worker; sparsity parameter  $k$ ;  $N$  local workers;  $m_0 = 0, v_0 = 0, \hat{v}_0 = 0$ 
3: for  $t = 1$  to  $T$  do
4:   parallel for worker  $i \in [n]$  do:
5:     Receive model parameter  $\theta_{t-1}$  from central server
6:     Compute stochastic gradient  $g_{t,i}$  at  $\theta_t$ 
7:     Compute  $\tilde{g}_{t,i} = \text{TopK}(g_{t,i} + e_{t,i}, k)$ 
8:     Update the error  $e_{t+1,i} = e_{t,i} + g_{t,i} - \tilde{g}_{t,i}$ 
9:     Send  $\tilde{g}_{t,i}$  back to central server
10:  end parallel
11:  Central server do:
12:     $\bar{g}_t = \frac{1}{N} \sum_{i=1}^N \tilde{g}_{t,i}$ 
13:     $m_t = \beta_1 m_{t-1} + (1 - \beta_1) \bar{g}_t$ 
14:     $v_t = \beta_2 v_{t-1} + (1 - \beta_2) \bar{g}_t^2$ 
15:     $\hat{v}_t = \max(v_t, \hat{v}_{t-1})$ 
16:    Update global model  $\theta_t = \theta_{t-1} - \eta_t \frac{m_t}{\sqrt{\hat{v}_t}}$ 
17: end for

```

---

#### 3.2 Convergence Analysis

Several mild assumptions to make: Nonconvex and smooth loss function, unbiased stochastic gradient, bounded variance of the gradient, bounded norm of the gradient, control of the distance between the true gradient and its sparse variant.

Check [7] starting with single machine and extending to distributed settings (several machines).

##### 3.2.1 Single machine

Under the centralized setting, the goal is to derive an upper bound to the second order moment of the gradient of the objective function at some iteration  $T_f \in [1, T]$ .

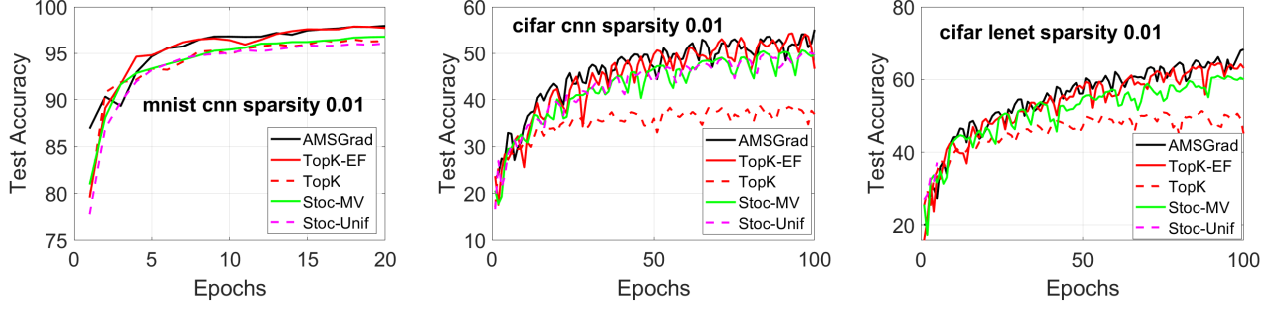


Figure 1: Test accuracy.

60 We begin by making the following assumptions.

61 We first define multiple auxiliary sequences. For the first moment, define

$$\begin{aligned}\bar{m}_t &= m_t + \mathcal{E}_t, \\ \mathcal{E}_t &= \beta_1 \mathcal{E}_{t-1} + (1 - \beta_1)(e_{t+1} - e_t),\end{aligned}$$

62 such that

$$\begin{aligned}\bar{m}_t &= \bar{m}_t + \mathcal{E}_t \\ &= \beta_1(m_t + \mathcal{E}_t) + (1 - \beta_1)(\bar{g}_t + e_{t+1} - e_1) \\ &= \beta_1 \bar{m}_{t-1} + (1 - \beta_1)g_t.\end{aligned}$$

### 63 3.2.2 Multiple machine

## 64 4 Experiments

65 Our proposed TopK-EF with AMSGrad matches that of full AMSGrad, in distributed learning.  
66 Number of local workers is 20. Error feedback fixes the convergence issue of using solely the  
67 TopK gradient.

## 68 5 Conclusion

## References

- [1] Alham Fikri Aji and Kenneth Heafield. Sparse communication for distributed gradient descent. *arXiv preprint arXiv:1704.05021*, 2017.
- [2] Dan Alistarh, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. Qsgd: Communication-efficient sgd via gradient quantization and encoding. In *Advances in Neural Information Processing Systems*, pages 1709–1720, 2017.
- [3] Dan Alistarh, Torsten Hoefer, Mikael Johansson, Sarit Khirirat, Nikola Konstantinov, and Cédric Renggli. The convergence of sparsified gradient methods. *arXiv preprint arXiv:1809.10505*, 2018.
- [4] 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.
- [5] Léon Bottou and Olivier Bousquet. The tradeoffs of large scale learning. In J. C. Platt, D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 161–168. Curran Associates, Inc., 2008.
- [6] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, Jonathan Eckstein, et al. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends® in Machine learning*, 3(1):1–122, 2011.
- [7] Congliang Chen, Li Shen, Haozhi Huang, Qi Wu, and Wei Liu. Quantized adam with error feedback. *arXiv preprint arXiv:2004.14180*, 2020.
- [8] Yongjian Chen, Tao Guan, and Cheng Wang. Approximate nearest neighbor search by residual vector quantization. *Sensors*, 10(12):11259–11273, 2010.
- [9] Trishul Chilimbi, Yutaka Suzue, Johnson Apacible, and Karthik Kalyanaraman. Project adam: Building an efficient and scalable deep learning training system. In *Symposium on Operating Systems Design and Implementation*, pages 571–582, 2014.
- [10] Christopher De Sa, Matthew Feldman, Christopher Ré, and Kunle Olukotun. Understanding and optimizing asynchronous low-precision stochastic gradient descent. In *Proceedings of the 44th Annual International Symposium on Computer Architecture*, pages 561–574, 2017.
- [11] John C Duchi, Alekh Agarwal, and Martin J Wainwright. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control*, 57(3):592–606, 2011.
- [12] Saeed Ghadimi and Guanghui Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Journal on Optimization*, 23(4):2341–2368, 2013.
- [13] Farzin Haddadpour, Mohammad Mahdi Kamani, Mehrdad Mahdavi, and Viveck Cadambe. Trading redundancy for communication: Speeding up distributed sgd for non-convex optimization. In *International Conference on Machine Learning*, pages 2545–2554. PMLR, 2019.
- [14] Mingyi Hong, Davood Hajinezhad, and Ming-Min Zhao. Prox-pda: The proximal primal-dual algorithm for fast distributed nonconvex optimization and learning over networks. In *International Conference on Machine Learning*, pages 1529–1538, 2017.
- [15] Herve Jegou, Matthijs Douze, and Cordelia Schmid. Product quantization for nearest neighbor search. *IEEE transactions on pattern analysis and machine intelligence*, 33(1):117–128, 2010.
- [16] Peng Jiang and Gagan Agrawal. A linear speedup analysis of distributed deep learning with sparse and quantized communication. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 2530–2541, 2018.

- 113 [17] Belhal Karimi, Blazej Miasojedow, Eric Moulines, and Hoi-To Wai. Non-asymptotic analysis  
114 of biased stochastic approximation scheme. In *Conference on Learning Theory*, pages 1944–  
115 1974. PMLR, 2019.
- 116 [18] Sai Praneeth Karimireddy, Quentin Rebjock, Sebastian U Stich, and Martin Jaggi. Error feed-  
117 back fixes signsgd and other gradient compression schemes. *arXiv preprint arXiv:1901.09847*,  
118 2019.
- 119 [19] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*  
120 *preprint arXiv:1412.6980*, 2014.
- 121 [20] Anastasia Koloskova, Sebastian U Stich, and Martin Jaggi. Decentralized stochastic optimiza-  
122 tion and gossip algorithms with compressed communication. In *International Conference on*  
123 *Machine Learning*, pages 3478–3487, 2019.
- 124 [21] Songtao Lu, Xinwei Zhang, Haoran Sun, and Mingyi Hong. Gnsd: A gradient-tracking based  
125 nonconvex stochastic algorithm for decentralized optimization. In *2019 IEEE Data Science*  
126 *Workshop (DSW)*, pages 315–321, 2019.
- 127 [22] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Ar-  
128 cas. Communication-efficient learning of deep networks from decentralized data. In *Artificial*  
129 *Intelligence and Statistics*, pages 1273–1282. PMLR, 2017.
- 130 [23] Parvin Nazari, Davoud Ataee Tarzanagh, and George Michailidis. Dadam: A consensus-  
131 based distributed adaptive gradient method for online optimization. *arXiv preprint*  
132 *arXiv:1901.09109*, 2019.
- 133 [24] Angelia Nedic and Asuman Ozdaglar. Distributed subgradient methods for multi-agent opti-  
134 mization. *IEEE Transactions on Automatic Control*, 54(1):48, 2009.
- 135 [25] Sashank J Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond.  
136 In *International Conference on Learning Representations*, 2018.
- 137 [26] Shaohuai Shi, Kaiyong Zhao, Qiang Wang, Zhenheng Tang, and Xiaowen Chu. A convergence  
138 analysis of distributed sgd with communication-efficient gradient sparsification. In *IJCAI*,  
139 pages 3411–3417, 2019.
- 140 [27] Sebastian U Stich, Jean-Baptiste Cordonnier, and Martin Jaggi. Sparsified sgd with memory.  
141 In *Advances in Neural Information Processing Systems*, pages 4447–4458, 2018.
- 142 [28] Jianqiao Wangni, Jialei Wang, Ji Liu, and Tong Zhang. Gradient sparsification for  
143 communication-efficient distributed optimization. *arXiv preprint arXiv:1710.09854*, 2017.
- 144 [29] Wei Wen, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Tern-  
145 grad: Ternary gradients to reduce communication in distributed deep learning. *arXiv preprint*  
146 *arXiv:1705.07878*, 2017.
- 147 [30] Guandao Yang, Tianyi Zhang, Polina Kirichenko, Junwen Bai, Andrew Gordon Wilson, and  
148 Chris De Sa. Swalp: Stochastic weight averaging in low precision training. In *International*  
149 *Conference on Machine Learning*, pages 7015–7024. PMLR, 2019.

