

Collaborative Machine Learning at the Wireless Edge with Blind Transmitters

Mohammad Mohammadi Amiri*, Tolga M. Duman[†], Deniz Gündüz*

*Electrical and Electronic Engineering Department, Imperial College London, London SW7 2BT, U.K.

[†]Department of Electrical and Electronics Engineering, Bilkent University, Ankara 06800, Turkey

Email: *{m.mohammadi-amiri15, d.gunduz}@imperial.ac.uk, [†]duman@ee.bilkent.edu.tr

Abstract—We study wireless collaborative machine learning (ML), where mobile edge devices, each with its own dataset, carry out distributed stochastic gradient descent (DSGD) over-the-air with the help of a wireless access point acting as the parameter server (PS). At each iteration of the DSGD algorithm wireless devices compute gradient estimates with their local datasets, and send them to the PS over a wireless fading multiple access channel (MAC). Motivated by the additive nature of the wireless MAC, we propose an analog DSGD scheme, in which the devices transmit scaled versions of their gradient estimates in an uncoded fashion. We assume that the channel state information (CSI) is available only at the PS. We instead allow the PS to employ multiple antennas to alleviate the destructive fading effect, which cannot be cancelled by the transmitters due to the lack of CSI. Theoretical analysis indicates that, with the proposed DSGD scheme, increasing the number of PS antennas mitigates the fading effect, and, in the limit, the effects of fading and noise disappear, and the PS receives aligned signals used to update the model parameter. The theoretical results are then corroborated with the experimental ones.

I. INTRODUCTION

With the growing prevalence of Internet of things (IoT) devices, constantly collecting information about various physical phenomena, and the growth in the number and processing capability of mobile edge devices (phones, tablets, smart watches and activity monitors), there is a growing interest in enabling distributed machine learning (ML) to learn from data distributed across mobile devices. Centralized ML techniques are often developed, assuming that the datasets are offloaded to a central processor. In the case of wireless edge devices, centralized ML techniques are not desirable, since offloading such massive amounts of data to a central cloud may be too costly in terms of both energy and privacy.

In many ML problems, the goal is to minimize a loss function, $F(\theta)$, where $\theta \in \mathbb{R}^d$ captures the model parameters to be optimized. The loss function $F(\theta)$ represents the average of empirical loss functions computed at different data samples with respect to model parameter θ , $F(\theta) = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{u} \in \mathcal{B}} f(\theta, \mathbf{u})$, where \mathcal{B} is the set of available data points, and \mathbf{u} represents a data sample and its label.

We assume that an iterative stochastic gradient descent (SGD) algorithm is used to minimize the loss function $F(\theta)$, in which the model parameter vector at iteration t , θ_t , is updated according to the stochastic gradient $\mathbf{g}(\theta_t)$. SGD allows parallelization across multiple mobile devices. In distributed SGD (DSGD), devices process data locally with

respect to a globally consistent parameter vector, and send their gradient estimates to the parameter server (PS). To be more precise, at iteration t , device m computes the gradient estimate $\mathbf{g}_m(\theta_t) \triangleq \frac{1}{|\mathcal{B}_m|} \sum_{\mathbf{u} \in \mathcal{B}_m} \nabla f(\theta_t, \mathbf{u})$ with respect to its local dataset \mathcal{B}_m and model parameter θ_t , and sends the result to the PS. Having M devices in the system, the PS updates the model parameter vector according to

$$\theta_{t+1} = \theta_t - \eta_t \frac{1}{M} \sum_{m=1}^M \mathbf{g}_m(\theta_t), \quad (1)$$

where η_t denotes the learning rate at iteration t , and shares the result with the devices for the computations at the following iterations. Although parallelism reduces the computation load at each device, communication from the devices to the PS becomes the main performance bottleneck [1]–[5], particularly for wireless edge learning due to limited bandwidth and power.

Several architectures have been proposed in recent years to employ computational capabilities of edge devices, and train an ML model collaboratively with the help of a remote PS. However, these works ignore the physical characteristics of the communication channel from the devices to the PS, and consider interference-and-error-free links with a fixed capacity, which is hard to guarantee in most wireless environments.

Collaborative ML taking into account the physical layer channel characteristics has recently been studied in [6]–[9]. These works consider a wireless multiple access channel (MAC) from the edge devices to the PS, and propose over-the-air computation to average gradient vectors or estimated model parameters at the PS. In [6] the authors focus on bandwidth efficient learning, and employ gradient sparsification followed by linear projection to design a communication efficient DSGD algorithm. This scheme has been extended to the fading MAC model in [9]. Distributed ML over a wireless fading MAC is studied in [7], where the wireless devices employ power allocation with perfect channel state information (CSI) to align the received signals at the PS. A single-input multiple-output (SIMO) wireless fading MAC is studied in [8], where a beamforming technique is designed to maximize the number of devices participating in each iteration, while keeping the quality of the received signal at the PS above the specified threshold level.

Our goal in this paper is to enable distributed learning over a wireless fading MAC, while removing the requirement of CSI at the transmitters (CSIT). This will be achieved by

employing multiple antennas at the PS. Similarly to [6]–[9] we considering uncoded transmission of gradient estimates and over-the-air computation. We design a receive beamformer at the PS in order to mitigate the fading effect and align the desired signals. We analytically show that the proposed scheme alleviates the destructive effects of interference and noise terms at the PS thanks to the utilization of multiple antennas, and, in the limit, due to channel hardening, it boils down to a deterministic channel with identical gains from all the devices. This result is validated by numerical experiments, where we investigate the impact of the number of antennas on the performance of the proposed scheme with no CSIT. It is worth noting that the CSI requirements of over-the-air computation with a multi-antenna receiver was also studied in [10]. The authors proposed a scheme that encodes the information on the energy of the transmitter signals, and hence, limited only to positive values, but requires CSI neither at the transmitters nor at the PS. Performance of this no-CSI scheme for DSGD will be studied in the extended version of this paper.

Notations: \mathbb{R} and \mathbb{C} represent the sets of real and complex values, respectively. We denote entry-wise complex conjugate of vector \mathbf{x} by $(\mathbf{x})^*$, and $\text{Re}\{\mathbf{x}\}$ and $\text{Im}\{\mathbf{x}\}$ return entry-wise real and imaginary components of \mathbf{x} , respectively. For \mathbf{x} and \mathbf{y} with the same dimension, $\mathbf{x} \cdot \mathbf{y}$ returns their inner product. We denote a zero-mean normal distribution with variance σ^2 by $\mathcal{N}(0, \sigma^2)$, and $\mathcal{CN}(0, \sigma^2)$ represents a circularly symmetric complex normal distribution with real and imaginary terms each distributed according to $\mathcal{N}(0, \sigma^2/2)$. We let $[i] \triangleq \{1, \dots, i\}$. We denote the cardinality of set \mathcal{X} by $|\mathcal{X}|$, and l_2 norm of vector \mathbf{x} by $\|\mathbf{x}\|_2$.

II. SYSTEM MODEL

We consider M devices, where device m has access to a local dataset \mathcal{B}_m , and employs SGD to compute the gradient estimate $\mathbf{g}_m(\boldsymbol{\theta}_t) \in \mathbb{R}^d$ at iteration t , $m \in [M]$. These local gradient estimates are transmitted to the PS, equipped with K antennas, through a wireless shared medium. The PS updates the model parameter based on its received signal, and shares it with all the devices over an error-free shared link, so that all the devices have a globally consistent model parameter.

We model the shared wireless channel from the edge devices to the PS as a wireless fading MAC, where OFDM is used to divide the available bandwidth into s subchannels, $s \leq d$ (in practice, we typically have $s \ll d$). We assume that N OFDM symbols can be transmitted over each subchannel at each iteration of DSGD algorithm. The received vector corresponding to the n -th OFDM symbol in iteration t at the k -th antenna of the PS is given by

$$\mathbf{y}_k^n(t) = \sum_{m=1}^M \mathbf{h}_{m,k}^n(t) \cdot \mathbf{x}_m^n(t) + \mathbf{z}_k^n(t), \quad k \in [K], \quad (2)$$

where $\mathbf{x}_m^n(t)$ is the n -th symbol of dimension s transmitted by the m -th device, $\mathbf{h}_{m,k}^n(t) \in \mathbb{C}^s$ denotes the vector of channel gains from device m to the k -th PS antenna, $m \in [M]$, and $\mathbf{z}_k^n(t) \in \mathbb{C}^s$ represents the circularly symmetric complex white Gaussian noise at the k -th antenna of the PS, $n \in [N]$. The

i -th entry of channel vector $\mathbf{h}_{m,k}^n(t)$, denoted by $h_{m,k,i}^n(t)$, is distributed according to $\mathcal{CN}(0, \sigma_h^2)$, $i \in [s]$, and different entries of $\mathbf{h}_{m,k}^n(t)$ can be correlated, while the channel gains are assumed to be independent and identically distributed (i.i.d.) across PS antennas, OFDM symbols, and wireless devices, $k \in [K]$, $n \in [N]$, $m \in [M]$. Similarly, different entries of noise vector $\mathbf{z}_k^n(t)$ can be correlated, and its i -th entry, denoted by $z_{k,i}^n(t)$, distributed according to $\mathcal{CN}(0, \sigma_z^2)$, $i \in [s]$, $k \in [K]$, $n \in [N]$. Noise vectors are also assumed to be i.i.d. across PS antennas and OFDM symbols. We consider the following average power constraint imposed at each wireless device assuming a total of T iterations of the DSGD algorithm:

$$\frac{1}{NT} \sum_{t=1}^T \sum_{n=1}^N \mathbb{E} [\|\mathbf{x}_m^n(t)\|_2^2] \leq \bar{P}, \quad \forall m \in [M], \quad (3)$$

where the expectation is taken with respect to the randomness of the communication channel.

We assume that the PS has perfect CSI, while there is no CSI at the wireless devices. At each iteration, the goal at the PS is to estimate the average of the gradient estimates, $\frac{1}{M} \sum_{m=1}^M \mathbf{g}_m(\boldsymbol{\theta}_t)$, denoted by $\hat{\mathbf{g}}(\boldsymbol{\theta}_t)$, and update the model parameter as in (1) at the end of each iteration based on the received symbols $\mathbf{y}_k^1(t), \dots, \mathbf{y}_k^N(t)$, $\forall k$, and its knowledge of the CSI $\mathbf{h}_{m,k}^n(t)$, $\forall k, n, m$.

We note that the PS is interested in the average of the gradient estimates computed by the devices rather than each individual estimate. Motivated by the additive nature of the wireless MAC, we consider an analog approach similarly to [6]–[9], where the devices transmit their gradient estimates simultaneously without employing any channel coding.

III. ANALOG DSGD WITHOUT CSIT

At iteration t of DSGD, device m transmits its gradient estimate $\mathbf{g}_m(\boldsymbol{\theta}_t) \in \mathbb{R}^d$ over $N = \lceil d/2s \rceil$ OFDM symbols across s subchannels in an uncoded manner, $m \in [M]$. We denote the i -th entry of $\mathbf{g}_m(\boldsymbol{\theta}_t)$ by $g_{m,i}(\boldsymbol{\theta}_t)$, $i \in [d]$, and define, for $n \in [N]$, $m \in [M]$,

$$\mathbf{g}_{m,\text{re}}^n(\boldsymbol{\theta}_t) \triangleq [g_{m,2(n-1)s+1}(\boldsymbol{\theta}_t), \dots, g_{m,(2n-1)s}(\boldsymbol{\theta}_t)]^T, \quad (4a)$$

$$\mathbf{g}_{m,\text{im}}^n(\boldsymbol{\theta}_t) \triangleq [g_{m,(2n-1)s+1}(\boldsymbol{\theta}_t), \dots, g_{m,2ns}(\boldsymbol{\theta}_t)]^T, \quad (4b)$$

$$\mathbf{g}_m^n(\boldsymbol{\theta}_t) \triangleq \mathbf{g}_{m,\text{re}}^n(\boldsymbol{\theta}_t) + j\mathbf{g}_{m,\text{im}}^n(\boldsymbol{\theta}_t), \quad (4c)$$

where $j \triangleq \sqrt{-1}$, and we zero-pad $\mathbf{g}_m(\boldsymbol{\theta}_t)$ to have length $2sN$. The i -th entry of $\mathbf{g}_m^n(\boldsymbol{\theta}_t)$ is then given by

$$g_{m,i}^n(\boldsymbol{\theta}_t) = g_{m,2(n-1)s+i}(\boldsymbol{\theta}_t) + jg_{m,(2n-1)s+i}(\boldsymbol{\theta}_t), \quad \text{for } i \in [s], n \in [N], m \in [M]. \quad (5)$$

According to (4), we have

$$\mathbf{g}_m(\boldsymbol{\theta}_t) = [\mathbf{g}_{m,\text{re}}^1(\boldsymbol{\theta}_t), \mathbf{g}_{m,\text{im}}^1(\boldsymbol{\theta}_t), \dots, \mathbf{g}_{m,\text{re}}^N(\boldsymbol{\theta}_t), \mathbf{g}_{m,\text{im}}^N(\boldsymbol{\theta}_t)]^T, \quad (6)$$

with $N = \lceil d/2s \rceil$. At the n -th OFDM symbol of iteration t , device m sends

$$\mathbf{x}_m^n(t) = \alpha_t \mathbf{g}_m^n(t), \quad n \in [N], m \in [M]. \quad (7)$$

Accordingly, the average transmit power depends on α_t , and is evaluated as follows:

$$\frac{1}{NT} \sum_{t=1}^T \alpha_t^2 \sum_{n=1}^N \|\mathbf{g}_m^n(t)\|_2^2 \leq \bar{P}. \quad (8)$$

The PS observes the following signal at its k -th antenna, for $k \in [K], n \in [N]$:

$$\mathbf{y}_k^n(t) = \alpha_t \sum_{m=1}^M \mathbf{h}_{m,k}^n(t) \cdot \mathbf{g}_m^n(t) + \mathbf{z}_k^n(t). \quad (9)$$

Having known the CSI, the PS combines the signals at different antennas in the following form:

$$\mathbf{y}^n(t) \triangleq \frac{1}{K} \sum_{k=1}^K \left(\sum_{m=1}^M \mathbf{h}_{m,k}^n(t) \right)^* \cdot \mathbf{y}_k^n(t), \quad (10)$$

whose i -th entry is given by

$$y_i^n(t) = \frac{1}{K} \sum_{k=1}^K \sum_{m=1}^M (h_{m,k,i}^n(t))^* y_{k,i}^n(t), \quad (11)$$

where $y_{k,i}^n(t)$ denotes the i -th entry of $\mathbf{y}_k^n(t)$, $i \in [s], n \in [N]$.

By substituting $y_{k,i}^n(t)$, given in (9), it follows that

$$\begin{aligned} y_i^n(t) &= \alpha_t \underbrace{\sum_{m=1}^M \left(\frac{1}{K} \sum_{k=1}^K |h_{m,k,i}^n(t)|^2 \right)}_{\text{signal term}} g_{m,i}^n(\boldsymbol{\theta}_t) \\ &+ \underbrace{\frac{\alpha_t}{K} \sum_{k=1}^K \sum_{m=1}^M \sum_{m'=1, m' \neq m}^M (h_{m,k,i}^n(t))^* h_{m',k,i}^n(t) g_{m',i}^n(\boldsymbol{\theta}_t)}_{\text{interference term}} \\ &+ \underbrace{\sum_{m=1}^M \left(\frac{1}{K} \sum_{k=1}^K (h_{m,k,i}^n(t))^* \right)}_{\text{noise term}} z_{k,i}^n(t). \end{aligned} \quad (12)$$

There are three terms with $y_i^n(t)$ specified by signal, interference, and noise terms, respectively, in (12). With the law of large numbers, as the number of antennas at the PS $K \rightarrow \infty$, the signal term approaches

$$y_{i,\text{sig}}^n(t) \triangleq \alpha_t \sigma_h^2 \sum_{m=1}^M g_{m,i}^n(\boldsymbol{\theta}_t), \quad i \in [s], n \in [N], \quad (13)$$

from which the PS can recover

$$\frac{1}{M} \sum_{m=1}^M g_{m,2(n-1)s+i}(\boldsymbol{\theta}_t) = \frac{\text{Re}\{y_{i,\text{sig}}^n(t)\}}{\alpha_t M \sigma_h^2}, \quad (14a)$$

$$\frac{1}{M} \sum_{m=1}^M g_{m,(2n-1)s+i}(\boldsymbol{\theta}_t) = \frac{\text{Im}\{y_{i,\text{sig}}^n(t)\}}{\alpha_t M \sigma_h^2}. \quad (14b)$$

However, the interference term in (12) does not allow the exact recoveries of $\frac{1}{M} \sum_{m=1}^M g_{m,2(n-1)s+i}(\boldsymbol{\theta}_t)$ and $\frac{1}{M} \sum_{m=1}^M g_{m,(2n-1)s+i}(\boldsymbol{\theta}_t)$ from $y_i^n(t)$, which is observed at the PS. To analyze the interference term, we first define, for $i \in [s], n \in [N]$,

$$\mathbf{h}_i^n(t) \triangleq \frac{1}{K} \sum_{k=1}^K \sum_{m=1}^M \sum_{m'=1, m' \neq m}^M (h_{m,k,i}^n(t))^* h_{m',k,i}^n(t). \quad (15)$$

It is easy to verify that the mean and the variance of $\mathbf{h}_i^n(t)$ are

given by

$$\mathbb{E}[\mathbf{h}_i^n(t)] = 0, \quad (16a)$$

$$\mathbb{E}[|\mathbf{h}_i^n(t)|^2] = \frac{M(M-1)\sigma_h^4}{K}, \quad (16b)$$

respectively. We note that the gradient values computed at each iteration are independent of the channel realizations experienced during the same iteration. Accordingly, by fixing the gradient values, from the analysis in (16), we conclude that the interference term in (12) has zero mean and a variance that scales with M^2/K . Thus, for a fixed number of wireless devices M , the variance of the interference term in (12) approaches zero as $K \rightarrow \infty$. In practice, it is feasible to employ sufficiently large number of antennas at the PS exploiting massive multiple-input multiple-output (MIMO) systems [11].

According to the above analysis, the PS estimates $\frac{1}{M} \sum_{m=1}^M g_{m,2(n-1)s+i}(\boldsymbol{\theta}_t)$ and $\frac{1}{M} \sum_{m=1}^M g_{m,(2n-1)s+i}(\boldsymbol{\theta}_t)$, for $i \in [s], n \in [N]$, through

$$\hat{g}_{2(n-1)s+i}(\boldsymbol{\theta}_t) = \frac{\text{Re}\{y_i^n(t)\}}{\alpha_t M \sigma_h^2}, \quad (17a)$$

$$\hat{g}_{(2n-1)s+i}(\boldsymbol{\theta}_t) = \frac{\text{Im}\{y_i^n(t)\}}{\alpha_t M \sigma_h^2}, \quad (17b)$$

respectively. It then utilizes the estimated vector $\hat{\mathbf{g}}(\boldsymbol{\theta}_t) \triangleq [\hat{g}_1(\boldsymbol{\theta}_t), \dots, \hat{g}_d(\boldsymbol{\theta}_t)]^T$, which can provide a good estimate of the actual average of gradients if a sufficiently large number of PS antennas are employed, to update the model parameters.

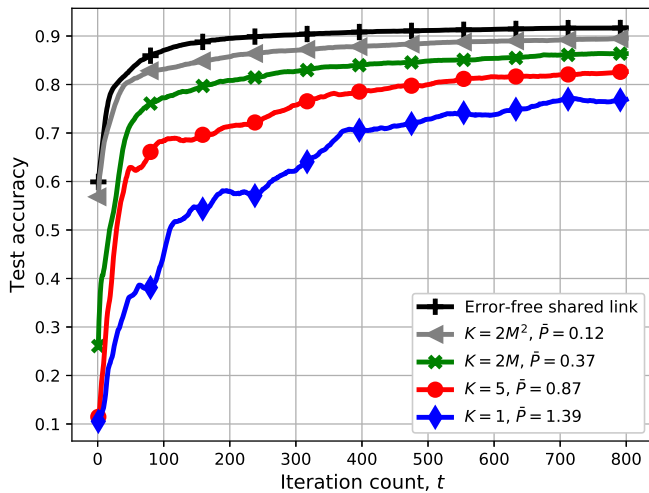
Remark 1. We note that with SGD the empirical variances of the gradient estimates decay over time and approach zero asymptotically [3], [4], [6], [12], [13]. Thus, for robust communication of the gradient estimates against noise at each iteration of the DSGD algorithm, it is reasonable to increase the power allocation factor α_t over time.

Remark 2. We remark that the main focus in this paper is to develop techniques to perform a DSGD algorithm at the wireless edge with no CSIT. We propose to employ multiple antennas at the PS, which can help to mitigate the effect of fading, and, in the limit, align the received signals at the PS. We can further employ some of the existing schemes in the literature providing more efficient communication over the limited bandwidth wireless MAC, such as the idea of linear projection proposed in [6]. We leave the analysis of such combined techniques to future work.

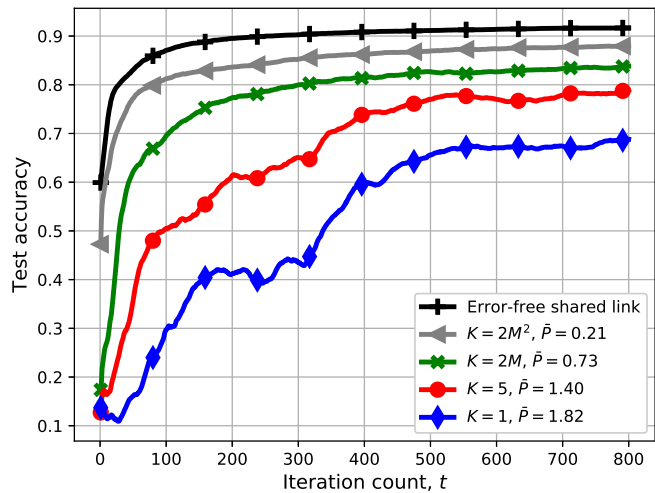
IV. NUMERICAL EXPERIMENTS

Here we evaluate the performance of the proposed analog DSGD algorithm with no CSI available at the wireless devices. We are particularly interested in investigating the impact of the number of PS antennas on the performance of the proposed scheme. We run experiments on MNIST dataset [14] with 60000 training and 10000 test samples, and train a single layer neural network with $d = 7850$ parameters utilizing ADAM optimizer [15]. We train the network for $T = 800$ iterations.

We consider $M = 20$ wireless devices in the system. To have a realistic model of data distribution across the devices



(a) Noise variance, $\sigma_z^2 = 20$



(b) Noise variance, $\sigma_z^2 = 100$

Fig. 1: Test accuracy of the proposed multi-antenna analog DSGD algorithm without CSIT for different number of antennas values ($K \in \{1, 5, 2M, 2M^2\}$) and noise variances σ_z^2 .

for the wireless edge learning model, we assume that each device has access to 1000 training data samples selected at random from the training dataset. Thus, some of the training data samples are not assigned to any device, and the data samples across different devices may not be independent. For simplicity, we assume that the s channel gains associated with each OFDM symbol from each device to each PS antenna are i.i.d., and $\sigma_h^2 = 1$. The performance is measured as the accuracy with respect to the test samples based on the updated model parameters at each DSGD iteration.

For numerical comparison, we also consider the benchmark scenario, in which the PS receives the actual average of the gradient estimates $\frac{1}{M} \sum_{m=1}^M \mathbf{g}_m(\boldsymbol{\theta}_t)$, and updates the parameter vector according to this noiseless observation at each DSGD iteration. We refer to this as the error-free shared link scenario, and its accuracy can serve as an upper bound on the performance of the proposed analog DSGD scheme.

In Fig. 1 we illustrate the performance of the proposed analog DSGD scheme with no CSIT for different K values and different noise levels. We consider $K \in \{1, 5, 2M, 2M^2\}$, and investigate the performance of the proposed scheme for $\sigma_z^2 = 20$ and $\sigma_z^2 = 100$ in Figures 1a and 1b, respectively. We also include the performance of the error-free shared link scenario. We set the power allocation factor $\alpha_t = 1 + t/1000$, $t \in [T]$, and for simplicity, we assume that $s = d/2$ resulting in $N = 1$. We note that, for a fixed power allocation α_t , $\forall t$, the value of s does not have any impact on the accuracy of the considered schemes; instead, any change in s scales the average transmit power, whose value is proportional to N . As it can be seen, employing more antennas at the PS results in a higher accuracy with the improvement more highlighted when the noise level is higher. This is due to the fact that increasing K mitigates the effects of both the interference and noise terms, inferred from (12). Thus, the advantage of having more PS antennas is more pronounced when the channel is noisier.

For example, even when $\sigma_z^2 = 100$, the proposed scheme with $K = 2M^2$ PS antennas and average power $\bar{P} = 0.21$ provides a slightly smaller accuracy than that of the error-free shared link scenario; this result indicates the success of the proposed scheme in mitigating the noise term even when the ratio \bar{P}/σ_z^2 is relatively small. We further observe that, compared to having a single-antenna PS, the accuracy improves by exploiting even a few antennas at the PS, e.g., $K = 5$, where the improvement is much higher when the channel is noisier, i.e., $\sigma_z^2 = 100$ case. We note that, with all the other parameters fixed, the required average transmit power reduces with K , which verifies a faster convergence rate with higher K resulting in a faster reduction in the empirical gradients' variances over time. The same observation is made by reducing σ_z^2 from 100 to 20 while all the other parameters are fixed.

V. CONCLUSIONS

We have studied DSGD at the wireless edge, where wireless devices compute the gradient estimates based on their available limited datasets, and transmit their estimates to the PS over a wireless fading MAC. To make the model more realistic, we have assumed that the devices do not have CSI for the underlying fast fading channel. With the goal of recovering the average gradient estimates at the PS, we have developed an analog DSGD technique, where the effect of fading, which cannot be cancelled at the transmitters due to the lack of CSIT, is alleviated by employing multiple antennas at the PS. Theoretical analysis, corroborated with numerical results, indicates that, with the proposed approach, increasing the number of PS antennas provides a better estimate of the average gradients through a better alignment of the desired signals, as well as elimination of the interference and noise terms. Asymptotically, the proposed DSGD scheme guarantees, despite the lack of CSIT, that the wireless MAC becomes deterministic, and both the fading and noise effects disappear.

REFERENCES

- [1] D. Alistarh, D. Grubic, J. Z. Li, R. Tomioka, and M. Vojnovic, “QSGD: Communication-efficient SGD via randomized quantization and encoding,” in *NIPS*, Long Beach, CA, Dec. 2017, pp. 1709–1720.
- [2] F. Seide, H. Fu, J. Droppo, G. Li, and D. Yu, “1-bit stochastic gradient descent and its application to data-parallel distributed training of speech DNNs,” in *INTERSPEECH*, Singapore, Sep. 2014, pp. 1058–1062.
- [3] S. Gupta, A. Agrawal, K. Gopalakrishnan, and P. Narayanan, “Deep learning with limited numerical precision,” in *ICML*, Jul. 2015.
- [4] N. Strom, “Scalable distributed DNN training using commodity gpu cloud computing,” in *INTERSPEECH*, 2015.
- [5] M. Mohammadi Amiri and D. Gündüz, “Computation scheduling for distributed machine learning with stragglers,” *arXiv:1810.09992 [cs.DC]*, May 2019.
- [6] M. Mohammadi Amiri and D. Gündüz, “Machine learning at the wireless edge: Distributed stochastic gradient descent over-the-air,” *arXiv:1901.00844 [cs.DC]*, Jan. 2019.
- [7] G. Zhu, Y. Wang, and K. Huang, “Low-latency broadband analog aggregation for federated edge learning,” *arXiv:1812.11494 [cs.IT]*, Jan. 2019.
- [8] K. Yang, T. Jiang, Y. Shi, and Z. Ding, “Federated learning via over-the-air computation,” *arXiv:1812.11750 [cs.LG]*, Jan. 2019.
- [9] M. Mohammadi Amiri and D. Gündüz, “Over-the-air machine learning at the wireless edge,” in *Proc. IEEE Int. Workshop on Signal Process. Advances in Wireless Commun. (SPAWC)*, Cannes, France, Jul. 2019.
- [10] M. Goldenbaum and S. Stanczak, “On the channel estimation effort for analog computation over wireless multiple-access channels,” *IEEE Wireless Commun. Lett.*, vol. 3, no. 3, pp. 261–264, Jun. 2014.
- [11] F. Rusek et al., “Scaling up MIMO: Opportunities and challenges with very large arrays,” *IEEE Signal Process. Mag.*, vol. 30, no. 1, pp. 40–60, Jan. 2013.
- [12] L. Bottou, “Large-scale machine learning with stochastic gradient descent,” in *Proc. COMPSTAT*, 2010, pp. 177–187.
- [13] T. Lin, S. U. Stich, and M. Jaggi, “Don’t use large mini-batches, use local SGD,” *arXiv:1808.07217v3 [cs.LG]*, Oct. 2018.
- [14] Y. LeCun, C. Cortes, and C. Burges, “The MNIST database of handwritten digits,” <http://yann.lecun.com/exdb/mnist/>, 1998.
- [15] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv:1412.6980v9 [cs.LG]*, Jan. 2017.