Imperfect CSI: A Key Factor of Uncertainty to Over-the-Air Federated Learning

Jiacheng Yao, Graduate Student Member, IEEE, Zhaohui Yang, Member, IEEE,

Wei Xu, Senior Member, IEEE, Dusit Niyato, Fellow, IEEE, and

Xiaohu You, Fellow, IEEE

Abstract

Over-the-air computation (AirComp) has recently been identified as a prominent technique to enhance communication efficiency of wireless federated learning (FL). This letter investigates the impact of channel state information (CSI) uncertainty at the transmitter on an AirComp enabled FL (AirFL) system with the truncated channel inversion strategy. To characterize the performance of the AirFL system, the weight divergence with respect to the ideal aggregation is analytically derived to evaluate learning performance loss. We explicitly reveal that the weight divergence deteriorates as $O(1/\rho^2)$ as the level of channel estimation accuracy ρ vanishes, and also has a decay rate of $O(1/K^2)$ with the increasing number of participating devices, K. Building upon our analytical results, we formulate the channel truncation threshold optimization problem to adapt to different ρ , which can be solved optimally. Numerical results verify the analytical results and show that a lower truncation threshold is preferred with more accurate CSI.

Index Terms

Federated learning (FL), over-the-air computation (AirComp), imperfect channel state information (CSI)

J. Yao, W. Xu, and X. You are with the National Mobile Communications Research Laboratory (NCRL), Southeast University, Nanjing 210096, China ({jcyao, wxu, xhyu}@seu.edu.cn).

Zhaohui Yang is with the Zhejiang Lab, Hangzhou 311121, China, and also with the College of Information Science and Electronic Engineering, Zhejiang University, Hangzhou, Zhejiang 310027, China (yang_zhaohui@zju.edu.cn).

Dusit Niyato is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 308232 (dniyato@ntu.edu.sg).

I. INTRODUCTION

Federated learning (FL), a distributed machine learning paradigm, has been regarded as a promising technique to support ubiquitous intelligence in the beyond fifth-generation (B5G) wireless networks [1], [2]. In a wireless FL system, the distributed devices, orchestrated by a parameter server (PS), iteratively train a shared learning model through the exchange of model parameters rather than the raw data, thereby protecting data privacy [3], [4]. However, due to the frequent uplink transmissions of model parameters from a large number of devices, the communication overhead and latency of FL become excessively high, which hinders its deployment in resource-constrained wireless networks.

To facilitate communication-efficient FL design, over-the-air computation (AirComp) has been greatly adopted for effective uplink model transmission [5]–[7]. By exploiting the waveform superposition nature of multiple access (MAC) channels, simultaneous model transmission and over-the-air model aggregation can be achieved, which can reduce the communication latency and save the uplink communication bandwidth substantially. In [5], a truncated channel inversion scheme was proposed to combat deep fadings in an AirComp-aided FL (AirFL), and the fundamental trade-offs between communication and learning was discussed. Then in [6], the power control strategy was further optimized to alleviate the impacts brought by AirComp errors. Considering the constraint of limited wireless communication resources, device selection and power control were jointly optimized to minimize the accuracy loss for AirFL in [7].

However, most of the existing AirFL scheme design and resource allocation optimization optimistically assumed the availability of perfect channel state information (CSI) at the transmitter, which is hardly to acquire in practice especially in a wireless network. More importantly, unlike traditional communication systems, the CSI imperfection in the AirFL system brings a severe impact. To be concrete, considering transmit power constraints, the users in deep fading should be truncated and therefore not participate in AirComp. Moreover, to achieve the uniform model aggregation, channel inversion should be performed at the transmitter. Considering the CSI uncertainties, the model aggregation is perturbed due to the imperfect truncation decision and channel inversion, resulting in the deterioration in learning performance. In [8], the authors considered a bounded CSI error and analyzed the impact of imperfect CSI on the convergence rate of FL. However, few effort has been endeavored to explicitly analyze in theory the aggregation distortion and accuracy loss brought by CSI uncertainty. To the best of our knowledge, there is

no theoretical guidance on channel truncation strategy of imperfect CSI.

Against the above backgrounds, we focus on an AirFL system adopting the truncated channel inversion scheme, where only partial CSI is available at the PS. We theoretically characterize the aggregation distortion due to imperfect CSI and the corresponding channel truncation, and derive an upper bound of the weight divergence of the aggregated gradient to evaluate the degradation of learning performance. Our results unrevil that as the level of channel estimation accuracy ρ decreases, the weight divergence enlarges at the order of $1/\rho^2$. The analytical result further suggests that increasing the number of participating devices, K, help decrease the weight divergence as $O(1/K^2)$ and can completely eliminate the impact of CSI imperfection. Moreover, based on the derived analytical results, we derive the optimal truncation threshold as a function of channel estimation uncertainty and system SNR. Numerical tests are conducted to verify the effectiveness of performance analysis and truncation threshold optimization.

II. SYSTEM MODEL OF AIRFL

A. Federated Learning Model

We consider a typical FL algorithm, where a shared machine learning model is trained via the collaboration between a central PS and K distributed devices. Let \mathcal{D}_k denote the local dataset owned by the kth device. The local loss function of model parameters, w, at the device k is defined as

$$F_k(\boldsymbol{w}, \mathcal{D}_k) = \frac{1}{|\mathcal{D}_k|} \sum_{\boldsymbol{u} \in \mathcal{D}_k} \mathcal{L}(\boldsymbol{w}, \boldsymbol{u}),$$
(1)

where u is a data sample and $\mathcal{L}(w, u)$ represents the sample-wise loss function. Without loss of generality, we assume that the size of all local datasets is the same, i.e., $|\mathcal{D}_k| = D$, $\forall k$. Then, the global loss function over all the datasets is given by

$$F(\boldsymbol{w}) = \frac{1}{K} \sum_{k=1}^{K} F_k(\boldsymbol{w}, \mathcal{D}_k).$$
(2)

The goal of the FL is to find the optimal model parameters, denoted by w^* , to minimize the global loss function in (2).

To effectively handle this problem, we apply the widely used FL algorithm in [8]. Specifically, in the *m*th round of the FL algorithm, the PS firstly broadcasts the up-to-date global parameter w_m to all devices. With the received global model w_m and their local datasets, each device



Fig. 1. An architecture of AirFL with one PS and K devices.

runs a stochastic gradient descent (SGD) algorithm on a local mini-batch to compute the local gradient, which follows

$$\boldsymbol{g}_{m}^{k} \triangleq \nabla F_{k}\left(\boldsymbol{w}_{m}, \mathcal{D}_{k,m}\right) = \frac{1}{|\mathcal{D}_{k,m}|} \sum_{\boldsymbol{u} \in \mathcal{D}_{k,m}} \mathcal{L}(\boldsymbol{w}_{m}, \boldsymbol{u}),$$
(3)

where $\mathcal{D}_{k,m}$ is the mini-batch selected from \mathcal{D}_k . Next, all devices report the local gradients in (3) to the PS. Upon receiving all the local gradients, PS performs the update as

$$\boldsymbol{w}_{m+1} = \boldsymbol{w}_m - \eta \boldsymbol{g}_m, \tag{4}$$

where η denotes the learning rate and

$$\boldsymbol{g}_m \triangleq \frac{1}{K} \sum_{k=1}^{K} \boldsymbol{g}_m^k.$$
 (5)

The FL algorithm iterates (3) and (4) until convergence.

B. Over-the-Air Computation for FL

In practice, we adopt the AirComp method for model aggregation in wireless networks, shown in Fig. 1. We express the channel between the kth device and the PS as $d_k^{-\frac{\alpha}{2}}h_k$, where d_k denotes the distance between the PS and device k, α is the large scale path loss exponent, and h_k represents the small-scale fading of the channel. Assume that the channels are independent Rayleigh fading channels, i.e., $h_k \sim C\mathcal{N}(0, 1)$. In general, the small-scale fading of the channel cannot be perfectly estimated at devices. Denote the channel estimate at device k by \hat{h}_k and a relationship between h_k and \hat{h}_k can be modelled as $h_k = \rho \hat{h}_k + \sqrt{1 - \rho^2} v_k$, where $\rho \in (0, 1]$ represents the correlation coefficient between h_k and \hat{h}_k , and $v_k \sim C\mathcal{N}(0, 1)$ is the error independent of \hat{h}_k . Note that ρ directly corresponds to the level of channel estimation accuracy and $\rho = 1$ implies the availability of perfect CSI. To overcome the negative impact of deep fading, a must truncated channel inversion scheme is headed for the uplink transmission [5]. To be concrete, only when $|\hat{h}_k|^2$ exceeds a predetermined threshold, $\gamma_{\rm th}$, the device is activated to transmit its gradient to PS. Accordingly, the received signal at the PS follows

$$\boldsymbol{y} = \sum_{k \in \mathcal{S}_m} d_k^{-\frac{\alpha}{2}} h_k \beta_k \boldsymbol{g}_m^k + \boldsymbol{z}_m,$$
(6)

where S_m represents the set of activated devices, β_k is the pre-processing factor for device k, and $z_m \sim C\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ is the additive Gaussian noise with power σ^2 . Given the imperfect CSI, the pre-processing factor for device k is chosen as $\beta_k = \frac{\zeta \lambda d_k^{\alpha/2} \hat{h}_k^*}{K|\hat{h}_k|^2}$ [7], where ζ is a scaling factor for ensuring the transmit power constraint and λ is a compensation constant for ensuring unbiasedness of the gradient estimation. For simplicity, we consider the uniform transmit power budget P_{max} at each device and choose the factor ζ to guarantee

$$\mathbb{E}\left[\left\|\beta_k \boldsymbol{g}_m^k\right\|^2\right] \le P_{\max}.$$
(7)

At the receiver, by scaling y with $\frac{1}{\zeta}$ and taking the real part, an estimate of the actual gradient in (5) is given by

$$\hat{\boldsymbol{g}}_m = \frac{1}{K} \sum_{k=1}^{K} \xi_k \boldsymbol{g}_m^k + \bar{\boldsymbol{z}}_m, \tag{8}$$

where $\bar{z}_m \triangleq \frac{\Re\{z_m\}}{\zeta}$ is the equivalent noise, and ξ_k is given by

$$\xi_{k} = \begin{cases} \lambda \frac{\Re\{h_{k}^{*}\hat{h}_{k}\}}{|\hat{h}_{k}|^{2}} & |\hat{h}_{k}|^{2} \ge \gamma_{\text{th}}, \\ 0 & |\hat{h}_{k}|^{2} < \gamma_{\text{th}}. \end{cases}$$
(9)

By comparing (8) and (5), the distortion in the gradient estimation comes from two aspects, i.e., the coefficient distortion ξ_k caused by the imperfect CSI, and the scaled additive Gaussian noise \bar{z}_m . Also, we notice that the expectation of ξ_k determines whether the gradient estimation is unbiased, and the variance of ξ_k and \bar{z}_m measure the gradient estimation distortion, which brings notable deterioration in convergence performance [4].

III. PERFORMANCE ANALYSIS AND OPTIMIZATION

In this section, we theoretically capture the impact of the imperfect CSI on the performance of AirFL. Based on the analytical results, we further optimize the truncation threshold $\gamma_{\rm th}$ with respect to ρ and system SNR.

A. Performance Analysis for AirFL

Firstly, we need to determine the value of the compensation constant λ for achieving an unbiased gradient estimation.

Lemma 1: In order to ensure the unbiasedness of gradient transmission, the compensation constant truncation and imperfect CSI is chosen by $\lambda = \frac{e^{\gamma_{th}}}{\rho}$.

Proof: Please refer to Appendix A.

Then, to facilitate the performance analysis for AirFL, the following *Lemma 2* derives the variance of AirComp parameters ξ_k , which directly reflects the mean squared error (MSE) of AirComp [9].

Lemma 2: The variance of ξ_k with unit mean is given by

$$\mathbb{E}\left[(\xi_k - 1)^2\right] = e^{\gamma_{\rm th}} - \frac{1 - \rho^2}{2\rho^2} \mathrm{Ei}(-\gamma_{\rm th}) e^{2\gamma_{\rm th}} - 1,$$
(10)

where $Ei(\cdot)$ denotes the exponential integral function.

Proof: Please refer to Appendix B.

Remark 1: Note that $\lim_{\gamma_{th}\to 0} Ei(\gamma_{th}) = -\infty$. Hence, from (10), regardless of the level of channel estimation accuracy, the truncation is necessary to avoid unbounded variance.

Remark 2: It is worth noting that for other power control schemes of AirComp, e.g., that in [6], the method of theoretical analysis still applies through treating them as special truncated channel inversion schemes. Without loss of generality, we therefore consider the most commonly adopted truncated channel inversion scheme as a general analysis.

Next, to evaluate the accuracy loss caused by the imperfect model aggregation, we choose a popular performance metric as the expected weight divergence with respect to \hat{g}_m and g_m [10], defined by $\Delta^2 = \mathbb{E} \left[\|\hat{g}_m - g_m\|^2 \right]$. It is worth noting that Δ^2 corresponds to the MSE of the

gradient estimation at the PS and directly reflects the accuracy of gradient estimation via the AirComp, which determines the convergence performance. To pave the way for performance analysis, we make the following widely used assumption [6].

Assumption: The stochastic gradients on random batches are uniformly bounded, i.e., $\mathbb{E}\left[\left\|\boldsymbol{g}_{m}^{k}\right\|^{2}\right] \leq G^{2}$. And the obtained global gradient, \boldsymbol{g}_{m} , is unbiased and variance bounded, i.e.,

$$\mathbb{E}[\boldsymbol{g}_m] = \nabla F(\boldsymbol{w}_m), \quad \mathbb{E}[\|\boldsymbol{g}_m - \nabla F(\boldsymbol{w}_m)\|] \le \delta^2.$$
(11)

Then, the weight divergence can be accurately characterized under this general assumption.

Theorem 1: The weight divergence, Δ^2 , is bounded by

$$\Delta^{2} \leq \frac{G^{2}}{K^{2}} \left(e^{\gamma_{\rm th}} - \frac{1 - \rho^{2}}{2\rho^{2}} \operatorname{Ei}(-\gamma_{\rm th}) e^{2\gamma_{\rm th}} - 1 + \frac{\sigma^{2} \max_{k} \{d_{k}^{\alpha}\} e^{2\gamma_{\rm th}}}{2P_{\max}\rho^{2}\gamma_{\rm th}} \right).$$
(12)

Proof: Please refer to Appendix C.

Remark 3: According to (12), the imperfect channel estimation deteriorates the learning performance with the order of $\frac{1}{\rho^2}$. It validates our statement that the accurate channel estimation is a key to the AirFL system.

Remark 4: Especially for high SNR regime, i.e., $\frac{P_{\text{max}}}{\sigma^2} \to \infty$, the weight divergence Δ^2 is dominated by the impact of imperfect CSI rather than noise. It implies that, for a given level of channel estimation accuracy, the accuracy loss caused by imperfect CSI can no longer be compensated by increasing the transmit power while only weakens the impact of the noise. This is the key observation that is different from the impact of CSI error in pure communication systems for data recovery.

Remark 5: By direct inspection of (12), as the increase of the number of devices, K, the weight divergence decreases as $O(1/K^2)$ and eventually tends towards zero, i.e., the impact of imperfect CSI is completely eliminated. This phenomenon can be qualitatively explained by the law of Large Numbers, that is, the randomness of aggregation distortion is eliminated when the participating devices tend to be infinite many.

Then, starting from (12), the convergence performance of FL is characterized in the following theorem.

Theorem 2: Suppose the loss function F is L-Lipschitz with respect to w and the learning rate satisfies $\eta < \frac{2}{L}$. The convergence of FL algorithm is bounded by

$$\frac{1}{M}\sum_{m=0}^{M-1}\mathbb{E}\left[\left\|\nabla F(\boldsymbol{w}_{m})\right\|^{2}\right] \leq \frac{1}{M}\left(\frac{F(\boldsymbol{w}_{0}) - \mathbb{E}\left[F(\boldsymbol{w}_{M})\right]}{\eta - \frac{L\eta^{2}}{2}} + \frac{ML\eta(\Delta^{2} + \delta^{2})}{2 - L\eta}\right).$$
(13)

Proof: Please refer to Appendix D.

This theorem implies that the convergence is guaranteed with sufficiently large M and the gap to the optimality converges to $\frac{L\eta(\Delta^2+\delta^2)}{2-L\eta}$, which linearly increasing with respect to Δ^2 .

B. Optimization of the Truncation Threshold

According to the result in *Theorem 1*, we find that he impacts of learning algorithms and the wireless transmission are decoupled. Hence, the truncation threshold optimization can be isolated from the specific learning algorithms and parameters, thus being defined as follows:

$$\max_{\gamma_{\rm th}>0} \quad h(\gamma_{\rm th}) \triangleq e^{\gamma_{\rm th}} - k_1 {\rm Ei}(-\gamma_{\rm th}) e^{2\gamma_{\rm th}} + k_2 \frac{e^{2\gamma_{\rm th}}}{\gamma_{\rm th}},\tag{14}$$

where $k_1 \triangleq \frac{1-\rho^2}{2\rho^2}$, and $k_2 \triangleq \frac{\sigma^2 \max_k \{d_k^{\alpha}\}}{2P_{\max}\rho^2}$ are positive constants.

Theorem 3: The objective function in (14) is convex.

Proof: Please refer to Appendix E.

Based on convexity of $h(\cdot)$, the optimal value of γ_{th} can be easily obtained from a bisection method with low complexity. Specifically, the derivative of h(x) is

$$h'(x) = e^{x} - k_{1} \frac{e^{x}}{x} - 2k_{1} \text{Ei}(-x)e^{2x} + k_{2}e^{2x} \frac{2x-1}{x^{2}}.$$
(15)

Since $\lim_{x\to 0} h'(x) < 0$ and $\lim_{x\to\infty} h'(x) > 0$, the unique zero point of h'(x), i.e., the optimal solution of γ_{th} , can be found through the bisection search.

IV. SIMULATION RESULTS

In this section, we provide simulation results to verify the performance analysis and truncation threshold optimization. We train a multi-layer perceptron (MLP) on the popular MNIST dataset via the AirFL algorithm. The distance d_k is uniformly distributed over (0, 500) m. Unless



Fig. 2. Variance versus ρ and $\gamma_{\rm th}$.



Fig. 3. Test accuracy versus different truncation thresholds.

otherwise specified, the other parameters are set as: K = 10, $\alpha = 2.2$, $P_{\text{max}} = 0.1$ W, $\sigma^2 = -40$ dBm, and $\eta = 0.005$.

Fig. 2 depicts the numerical variance of ξ_k obtained from Monte-Carlo simulations, compared with the theoretical result in (10). It shows that the numerical results matches well with the theoretical results, which verifies our analysis. Moreover, the channel estimation accuracy level, ρ , imposes more significant impacts on the variance than the truncation threshold.

In Fig. 3, we evaluate the impact of truncation threshold optimization on learning performance and compare it with other schemes. The four benchmark schemes are described as: "Communication oriented" and "Computation oriented" schemes represent γ_{th} is optimized to minimize the noise related and computation related term in (12), respectively; "Fixed γ_{th} " represents the truncation threshold is set as a constant [5]; The "Full power" scheme represents the transmitter does not perform power control and only compensates channel phase offset [11]. For all the tested setups, the proposed optimization method outperforms all the benchmarks due to the joint consideration of communication and computation. It is observed that the test accuracy first improves and then deteriorates with γ_{th} . This is because with the increase of γ_{th} , the performance is first limited by noise and then limited by less participating devices. Also, for larger ρ , a lower truncation threshold is preferred. Moreover, compared with the full power scheme, the proposed power control strategy successfully alleviates the impact of data heterogeneity, leading to a much prominent performance gain.

V. CONCLUSION

In this paper, we theoretically analyzed the performance of AirFL with imperfect CSI and optimized a channel truncation strategy. The analytical results revealed the importance of accurate channel estimation for AirFL. Our results can also be extended to performance analysis and optimization for other power control schemes of AirFL.

APPENDIX A

PROOF OF LEMMA 1

To determine λ , we start with the expectation of ξ_k , expressed as

$$\mathbb{E}\left[\xi_k\right] = \lambda \mathbb{E}\left[\frac{\Re\{h_k^* \hat{h}_k\}}{|\hat{h}_k|^2} \middle| |\hat{h}_k|^2 \ge \gamma_{\rm th}\right] \Pr\left\{|\hat{h}_k|^2 \ge \gamma_{\rm th}\right\},\tag{16}$$

which should be equal to 1 to guarantee an unbiased gradient estimation in (8). Considering that h_k and its estimate \hat{h}_k are correlated, we first introduce a new random variable to tackle with this difficulty, which follows

$$x \triangleq \frac{1}{\sqrt{1-\rho^2}} \left(\frac{\Re\{h_k^* \hat{h}_k\}}{|\hat{h}_k|^2} - \rho \right) = \frac{\Re\{v_k^* \hat{h}_k\}}{|\hat{h}_k|^2},$$
(17)

where v_k and \hat{h}_k are uncorrelated Gaussian variables with zero mean and unit variance. Then, we have

$$\mathbb{E}\left[x\left||\hat{h}_{k}|^{2} \geq \gamma_{\mathrm{th}}\right] = \mathbb{E}\left[\frac{v_{k}^{*}\hat{h}_{k} + v_{k}\hat{h}_{k}^{*}}{2|\hat{h}_{k}|^{2}}\right||\hat{h}_{h}|^{2} \geq \gamma_{\mathrm{th}}\right] = 0.$$

$$(18)$$

Then, by comparing (16) and (17), through some linear transformations, we arrive at $\mathbb{E}[\xi_k] = \lambda e^{-\gamma_{th}}\rho = 1$, which implies that $\lambda = e^{\gamma_{th}}/\rho$ and the proof completes.

APPENDIX B

PROOF OF LEMMA 2

We derive the variance of ξ_k by using the form of conditional expectation as

$$\mathbb{E}\left[\left(\xi_{k}-1\right)\right]^{2} = \mathbb{E}\left[\left.\left(\frac{\Re\{h_{k}^{*}\hat{h}_{k}\}e^{\gamma_{th}}}{|\hat{h}_{k}|^{2}\rho}-1\right)^{2}\right||\hat{h}_{k}|^{2} \ge \gamma_{th}\right] \Pr\left\{|\hat{h}_{k}|^{2} \ge \gamma_{th}\right\} + \Pr\left\{|\hat{h}_{k}|^{2} < \gamma_{th}\right\} \\
= \frac{e^{\gamma_{th}}(1-\rho^{2})}{\rho^{2}}\mathbb{E}\left[\left(x-c\right)^{2}\right|y \le -\gamma_{th}\right] + 1 - e^{-\gamma_{th}},$$
(19)

where $y \triangleq -|\hat{h}_k|^2$ and $c \triangleq \frac{\rho(1-e^{\gamma_{\text{th}}})}{e^{\gamma_{\text{th}}}\sqrt{1-\rho^2}}$. To calculate the conditional expectation, we first need to find the joint distribution of x and y. The joint cumulative distribution function (CDF) of x and y equals

$$F_{xy}(t,\gamma) = \Pr\left\{\frac{\Re\{v_k^*\hat{h}_k\}}{|\hat{h}_k|^2} < t, -|\hat{h}_k|^2 < \gamma\right\}$$

= $\Pr\left\{v_k^*\hat{h}_k + v_k\hat{h}_k^* - 2t|\hat{h}_k|^2 < 0, -|\hat{h}_k|^2 < \gamma\right\}$
= $\Pr\left\{\boldsymbol{z}^H \boldsymbol{A}_1 \boldsymbol{z} < 0, \, \boldsymbol{z}^H \boldsymbol{A}_2 \boldsymbol{z} < \gamma\right\},$ (20)

where $\boldsymbol{z} \triangleq [\hat{h}_k, v_k]^H$, and

$$\boldsymbol{A}_{1} \triangleq \begin{bmatrix} -2t & 1 \\ 1 & 0 \end{bmatrix}, \quad \boldsymbol{A}_{2} \triangleq \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix}.$$
(21)

According to [12, Eq. (3.2c.5)], the joint moment generating function (MGF) of the two quadratic forms, $z_1 \triangleq \boldsymbol{z}^H \boldsymbol{A}_1 \boldsymbol{z}$ and $z_2 \triangleq \boldsymbol{z}^H \boldsymbol{A}_2 \boldsymbol{z}$, follows

$$\mathcal{M}_{z_1, z_2}(s_1, s_2) = \det \left(\mathbf{I} - s_1 \mathbf{A}_1 - s_2 \mathbf{A}_2 \right)^{-1}.$$
 (22)

Applying the inverse Laplace transformation, we express the probability in (20) as

$$\Pr\left\{\boldsymbol{z}^{H}\boldsymbol{A}_{1}\boldsymbol{z} < 0, \, \boldsymbol{z}^{H}\boldsymbol{A}_{2}\boldsymbol{z} < \gamma\right\}$$
$$= \frac{1}{(2\pi i)^{2}} \int_{\epsilon_{1}-i\infty}^{\epsilon_{1}+i\infty} \int_{\epsilon_{2}-i\infty}^{\epsilon_{2}+i\infty} \frac{\mathrm{e}^{\gamma s_{2}}}{s_{1}s_{2}} \mathcal{M}_{z_{1},z_{2}}(s_{1},s_{2}) \mathrm{d}s_{2} \mathrm{d}s_{1}$$

$$\stackrel{\text{(a)}}{=} \frac{1}{2\pi i} \int_{\epsilon_{1}-i\infty}^{\epsilon_{1}+i\infty} \frac{1}{s_{1}(1+2ts_{1}-s_{1}^{2})} \mathrm{d}s_{1} - \frac{1}{2\pi i} \int_{\epsilon_{1}-i\infty}^{\epsilon_{1}+i\infty} \frac{1}{s_{1}(1+2ts_{1}-s_{1}^{2})} \mathrm{e}^{-(1+2ts_{1}-s_{1}^{2})\gamma} \mathrm{d}s_{1}$$

$$\stackrel{\text{(b)}}{=} \frac{t+\sqrt{1+t^{2}}}{2\sqrt{1+t^{2}}} + \frac{1}{2\pi i} \int_{\epsilon-i\infty}^{\epsilon+i\infty} \frac{(t+\sqrt{t^{2}+s+1})}{2s(s+1)\sqrt{t^{2}+s+1}} \mathrm{d}s$$

$$\stackrel{\text{(c)}}{=} \frac{t}{2\sqrt{1+t^{2}}} \left(1 - \operatorname{Erf}\left(\sqrt{-\gamma(1+t^{2})}\right)U(-\gamma)\right) + \frac{\mathrm{e}^{\gamma}}{2} \operatorname{Erfc}\left(-\sqrt{-\gamma}t\right)U(-\gamma),$$

$$(23)$$

where (a) comes from Eq. (5.2.4) in supplements of [13], (b) exploits the Cauchy's residue theorem and $s \triangleq s_1^2 - 2ts_1 - 1$, (c) is due to Eq. (5.3.7) in supplements of [13], $\text{Erf}(\cdot)$, $\text{Erfc}(\cdot)$, and $U(\cdot)$ represent the error function, the complementary error function and the Heaviside function, respectively. From (20) and (23) and by taking the derivative of joint CDF, the joint probability density function (PDF) of x and y equals

$$f_{xy}(t,\gamma) = \sqrt{-\frac{\gamma}{\pi}} e^{\gamma(1+t^2)} U(-\gamma).$$
(24)

From (24), we calculate the conditional expectation in (19) as

$$\mathbb{E}\left[\left(x-c\right)^{2}\middle| y \leq -\gamma_{\rm th}\right] = \int_{-\infty}^{\infty} \frac{\int_{-\infty}^{-\gamma_{\rm th}} f_{xy}(t,\gamma) \mathrm{d}\gamma}{\Pr\left\{y \leq -\gamma_{\rm th}\right\}} \mathrm{d}t$$

$$= e^{\gamma_{\rm th}} \int_{-\infty}^{-\gamma_{\rm th}} \int_{-\infty}^{\infty} (t-c)^{2} \sqrt{-\frac{\gamma}{\pi}} e^{\gamma(1+t^{2})} \mathrm{d}t \mathrm{d}\gamma$$

$$\stackrel{(a)}{=} e^{\gamma_{\rm th}} \int_{-\infty}^{-\gamma_{\rm th}} \frac{2c^{2}\gamma - 1}{2\gamma} e^{\gamma} \mathrm{d}\gamma$$

$$\stackrel{(b)}{=} c^{2} - \frac{1}{2} \mathrm{Ei}(-\gamma_{\rm th}) e^{\gamma_{\rm th}},$$
(25)

where (a) exploits [14, Eq. (3.462.8)] and [14, Eq. (3.321.1)] and the fact that $\int_{-\infty}^{\infty} t e^{\gamma t^2} dt = 0$. The equality in (b) comes from the definition of the exponential integral function, Ei(·). Applying (25) into (19), we complete the proof.

APPENDIX C

PROOF OF THEOREM 1

The weight divergence is reformulated as

$$\mathbb{E}\left[\left\|\hat{\boldsymbol{g}}_{m}-\boldsymbol{g}_{m}\right\|^{2}\right] = \mathbb{E}\left[\left\|\frac{1}{K}\sum_{k=1}^{K}(\xi_{k}-1)\boldsymbol{g}_{m}^{k}+\bar{\boldsymbol{z}}_{m}\right\|^{2}\right]$$
$$\stackrel{\text{(a)}}{=}\frac{1}{K^{2}}\sum_{k=1}^{K}\mathbb{E}\left[(\xi_{k}-1)^{2}\right]\mathbb{E}\left[\left\|\boldsymbol{g}_{m}^{k}\right\|^{2}\right]+\mathbb{E}\left[\left\|\bar{\boldsymbol{z}}_{m}\right\|^{2}\right],$$
(26)

where (a) is due to the zero mean and independence between $\xi_k - 1$, $\forall k$. As for the noise term, recall that $\bar{z}_m = \frac{\Re\{z_m\}}{\zeta}$ and we have $\mathbb{E}\left[\|\bar{z}_m\|^2\right] = \frac{\sigma^2}{2\zeta^2}$. According to the transmit power constraint in (7), the scaling factor ζ must satisfy

$$\max_{k \in \mathcal{S}_m} \left\{ \frac{\zeta^2 \lambda^2 d_k^{\alpha}}{K^2 |\hat{h}_k|^2} \mathbb{E} \left[\left\| \boldsymbol{g}_m^k \right\|^2 \right] \right\} \le P_{\max}.$$
(27)

Note that for all $k \in S_m$, we have $|\hat{h}_k|^2 \ge \gamma_{\text{th}}$. Combining (27) with the value of λ and the bound assumption of $\mathbb{E}\left[\left\|\boldsymbol{g}_m^k\right\|^2\right] \le G^2$, ζ is set as $\zeta = \frac{K\rho\sqrt{P_{\max}\gamma_{\text{th}}}}{G_{\max}\left\{d_k^{\alpha/2}\right\}e^{\gamma_{\text{th}}}}$. Then, combining all the derived results, we obtain (12) and complete the proof.

APPENDIX D Proof of Theorem 2

Under the general assumption, we have

$$\mathbb{E}\left[F(\boldsymbol{w}_{m+1}) - F(\boldsymbol{w}_{m})\right] \\
\stackrel{(a)}{\leq} \mathbb{E}\left[-\eta(\nabla F(\boldsymbol{w}_{m}))^{T}\hat{\boldsymbol{g}}_{m} + \frac{L\eta^{2}}{2}\|\hat{\boldsymbol{g}}_{m}\|^{2}\right] \\
\stackrel{(b)}{=} -\eta\mathbb{E}\left[\|\nabla F(\boldsymbol{w}_{m})\|^{2}\right] + \frac{L\eta^{2}}{2}\mathbb{E}\left[\|\hat{\boldsymbol{g}}_{m} - \boldsymbol{g}_{m} + \boldsymbol{g}_{m} - \nabla F(\boldsymbol{w}_{m}) + \nabla F(\boldsymbol{w}_{m})\|^{2}\right] \\
\stackrel{(c)}{=} -\left(\eta - \frac{L\eta^{2}}{2}\right)\mathbb{E}\left[\|\nabla F(\boldsymbol{w}_{m})\|^{2}\right] + \frac{L\eta^{2}}{2}\mathbb{E}\left[\|\hat{\boldsymbol{g}}_{m} - \boldsymbol{g}_{m}\|^{2}\right] + \frac{L\eta^{2}}{2}\mathbb{E}\left[\|\boldsymbol{g}_{m} - \nabla F(\boldsymbol{w}_{m})\|^{2}\right] \\
\stackrel{(d)}{\leq} -\left(\eta - \frac{L\eta^{2}}{2}\right)\mathbb{E}\left[\|\nabla F(\boldsymbol{w}_{m})\|^{2}\right] + \frac{L\eta^{2}(\Delta^{2} + \delta^{2})}{2},$$
(28)

where (a) is due to the fact that $F(\cdot)$ is *L*-Lipschitz and the definition of w_{m+1} , (b) comes from *Lemma 1*, (c) uses the assumption, and (d) exploits *Theorem 1*. By summing (28) from m = 0 to m = M - 1, we have

$$\frac{1}{M}\sum_{m=0}^{M-1}\mathbb{E}\left[\left\|\nabla F(\boldsymbol{w}_{m})\right\|^{2}\right] \leq \frac{1}{M}\left(\frac{F(\boldsymbol{w}_{0}) - \mathbb{E}\left[F(\boldsymbol{w}_{M})\right]}{\eta - \frac{L\eta^{2}}{2}} + \frac{ML\eta(\Delta^{2} + \delta^{2})}{2 - L\eta}\right), \quad (29)$$

which holds for any small $\eta < \frac{2}{L}$.

APPENDIX E

PROOF OF THEOREM 3

We check the second derivative of h(x) as

$$h''(x) = e^{x} + \frac{k_{1}e^{x}}{x^{2}} \left(-4x^{2}e^{x} \operatorname{Ei}(-x) - 3x + 1 \right) + \frac{2k_{2}e^{x}(2x^{2} - 2x + 1)}{x^{3}}.$$
 (30)

It is obvious that the first and the third terms in (30) are positive for positive x. To prove the nonnegativity of the second term, we use the inequalities [15, Eq. (5.1.20)]

$$-\mathrm{Ei}(-x) > \frac{1}{2}\mathrm{e}^{-x}\ln\left(1+\frac{2}{x}\right) > \mathrm{e}^{-x}\frac{1}{x+1},\tag{31}$$

which yields to $-4x^2 e^x Ei(-x) - 3x + 1 > \frac{(x-1)^2}{x+1} \ge 0$. Then, we conclude that h''(x) > 0 for x > 0 and hence h(x) is convex.

REFERENCES

- [1] W. Xu *et al.*, "Edge learning for B5G networks with distributed signal processing: Semantic communication, edge computing, and wireless sensing," *IEEE J. Sel. Topics Signal Process.*, vol. 17, no. 1, pp. 9–39, Jan. 2023.
- [2] G. Zhu *et al.*, "Pushing AI to wireless network edge: An overview on integrated sensing, communication, and computation towards 6G," *Sci. China Inf. Sci.*, vol. 66, no. pp. 130301:1–19, Feb. 2023.
- [3] Z. Yang *et al.*, "Energy efficient federated learning over wireless communication networks," *IEEE Trans. Wireless Commun.*, vol. 20, no. 3, pp. 1935–1949, Mar. 2021.
- [4] K. Yang et al., "Federated learning via over-the-air computation," IEEE Trans. Wireless Commun., vol. 19, no. 3, pp. 2022–2035, Mar. 2020.
- [5] G. Zhu, Y. Wang, and K. Huang, "Broadband analog aggregation for low-latency federated edge learning," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 491–506, Jan. 2020.
- [6] X. Cao et al., "Transmission power control for over-the-air federated averaging at network edge," IEEE J. Sel. Areas Commun., vol. 40, no. 5, pp. 1571–1586, May 2022.
- [7] W. Guo *et al.*, "Joint device selection and power control for wireless federated learning," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 8, pp. 2395–2410, Aug. 2022.
- [8] G. Zhu *et al.*, "One-bit over-the-air aggregation for communication-efficient federated edge learning: Design and convergence analysis," *IEEE Trans. Wireless Commun.*, vol. 20, no. 3, pp. 2120–2135, Mar. 2021.
- [9] W. Zhang *et al.*, "Worst-case design for RIS-aided over-the-air computation with imperfect CSI," *IEEE Commun. Lett.*, vol. 26, no. 9, pp. 2136–2140, Sept. 2022.
- [10] J. Yao et al., "GoMORE: Global model reuse for rescource-constrained wireless federated learning," IEEE Wireless Lett., early access, 2023. Doi: 10.1109/LWC.2023.3281881.
- [11] X. Fan *et al.*, "BEV-SGD: Best effort voting SGD against byzantine attacks for analog-aggregation-based federated learning over the air," *IEEE Internet Things J.*, vol. 9, no. 19, pp. 18946–18959, Oct. 2022.
- [12] A. Mathai and S. B. Provost, Quadratic Forms in Random Variables: Theory and Applications. New York, NY, USA: Marcel Dekker, 1992.

- [13] A. Polyanin and A. Manzhirov, *Handbook of Integral Equations: Second Edition*, Handbooks of Mathematical Equations. Taylor & Francis, 2008.
- [14] I. S. Gradshteyn and I. M. Ryzhik, Table of Integrals, Series and Products. New York: Academic Press, 6th ed, 2000.
- [15] M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables.* New York, NY, USA: Academic, 1972.