

On-Device Federated Learning via Second-Order Optimization with Over-the-Air Computation

Sheng Hua, Kai Yang, Yuanming Shi

School of Information Science and Technology
ShanghaiTech University



上海科技大学
ShanghaiTech University

Outline

Motivations

Problem Formulation

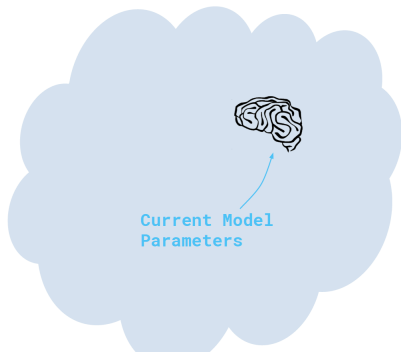
Proposed Algorithm

Simulation Results

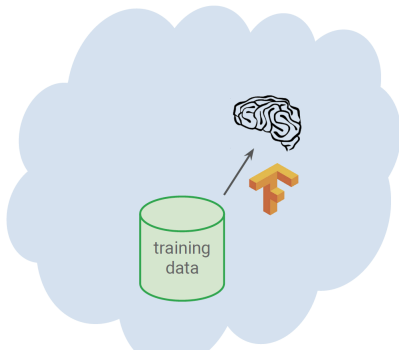
Summary

Cloud-Centric Machine Learning

The model lives in the cloud



We train models in the cloud



Mobile
Device

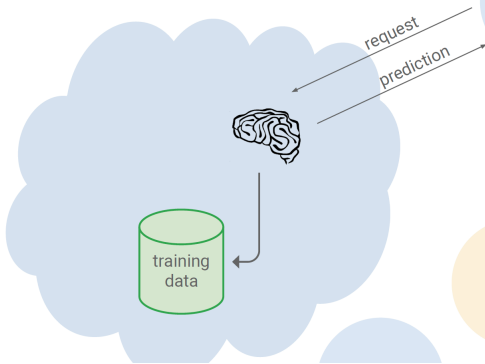
Current Model
Parameters



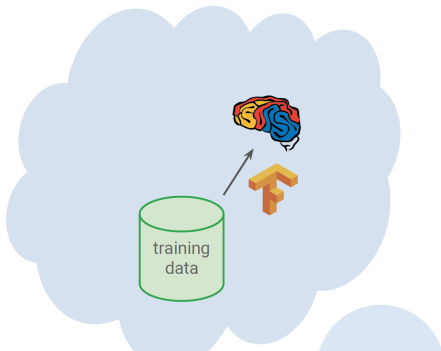
Make predictions in the cloud



Gather training data in the cloud



And make the models better



Why On-Device Learning?

- ▶ explosive growth in the volume of data on devices
- ▶ growing computation and storage capacity of devices
- ▶ privacy leakage, long delay
- ▶ ...

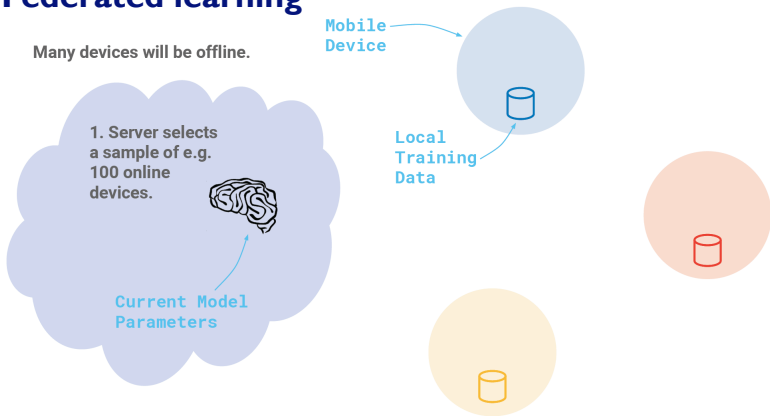
Why On-Device Learning?

- ▶ explosive growth in the volume of data on devices
- ▶ growing computation and storage capacity of devices
- ▶ privacy leakage, long delay
- ▶ ...

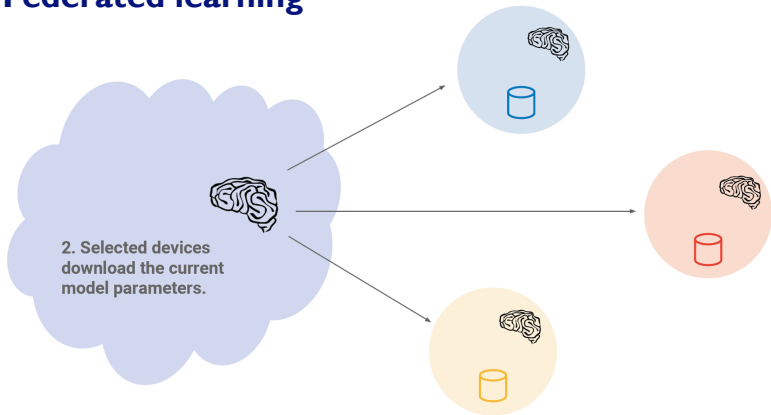
New Framework: Federated Learning

Federated learning

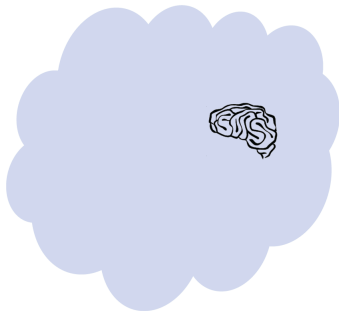
Many devices will be offline.



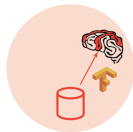
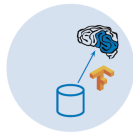
Federated learning



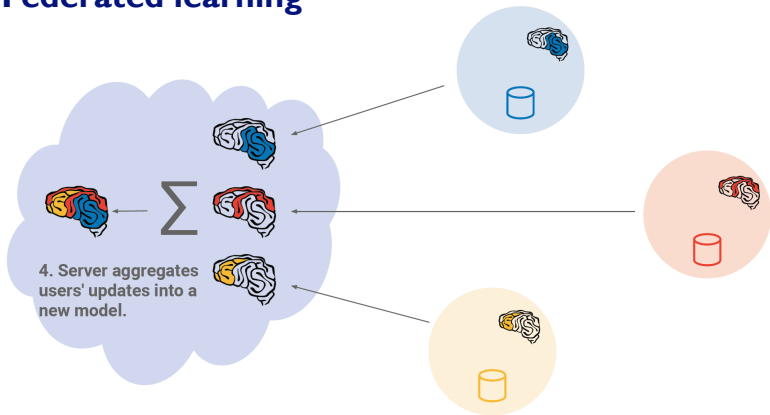
Federated learning



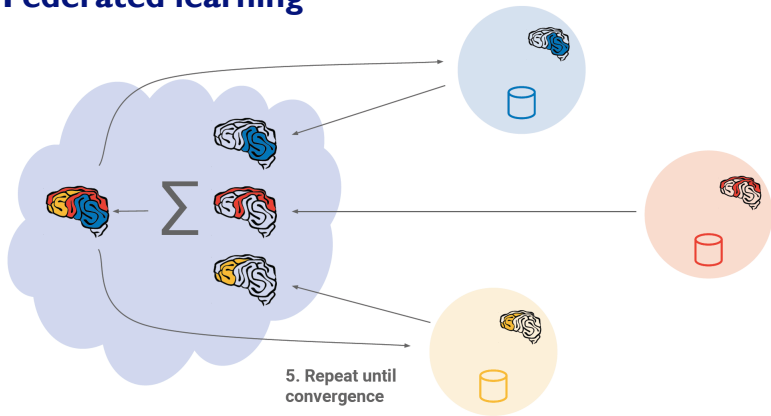
3. Devices compute an update using local training data



Federated learning

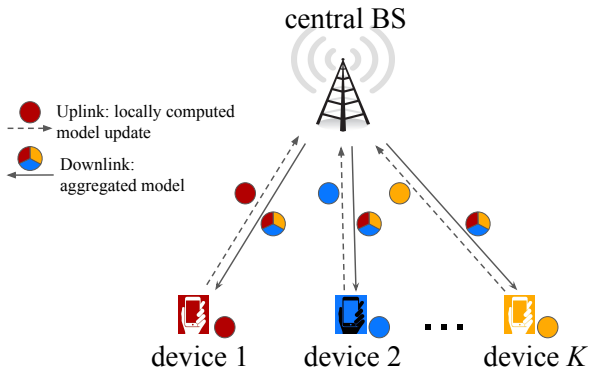


Federated learning



Federated Learning over Wireless Networks

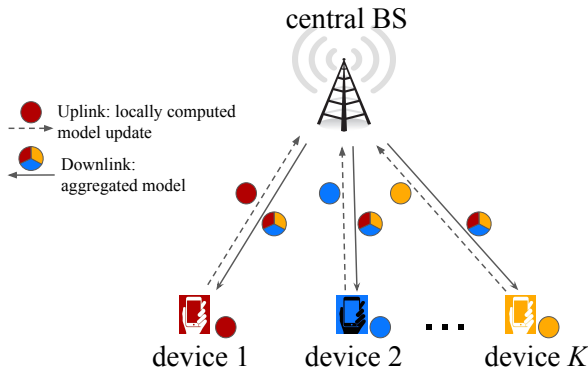
- **Goal:** train a shared global model via **wireless** federated computation



Q: How to efficiently aggregate models over wireless networks?

Federated Learning over Wireless Networks

- **Goal:** train a shared global model via **wireless** federated computation

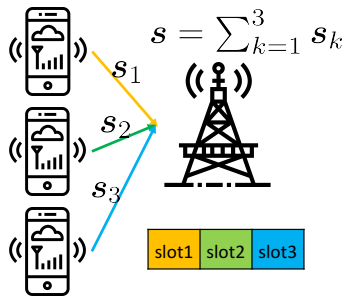


Q: How to efficiently aggregate models over wireless networks?

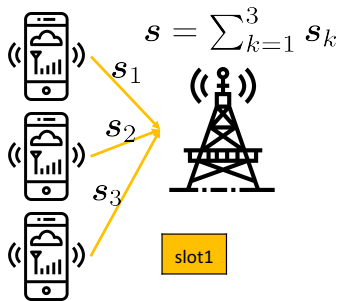
A: Via over-the-air computation.

Simple Illustration of Over-the-Air Computation

Communication Computation



Over-the-air Computation



The Training Procedures

► Challenges

- possible signal distortion in wireless communications
- slow convergence because a large number of iterations is required to train a satisfactory model

The Training Procedures

► Challenges

- possible signal distortion in wireless communications
- slow convergence because a large number of iterations is required to train a satisfactory model

► Our Work

- propose a difference-of-convex-functions (DC) algorithm to minimize signal distortion
- fasten model convergence by adopting the canonical Newton's method for local model updates

Outline

Motivations

Problem Formulation

Proposed Algorithm

Simulation Results

Summary

Notations

- ▶ One N -antenna base station (BS), K single-antenna mobile devices
- ▶ \mathcal{K} is the set of all devices
- ▶ the received signal at the BS after concurrent transmissions

$$\mathbf{y} = \sum_{k \in \mathcal{K}} \mathbf{h}_k b_k s_k + \mathbf{n}$$

- s_k : the representative signal; b_k : allocated power
- \mathbf{h}_k : the channel vector between the k -th device and the BS
- $\mathbf{n} \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I})$: the noise vector

Notations

- ▶ One N -antenna base station (BS), K single-antenna mobile devices
- ▶ \mathcal{K} is the set of all devices
- ▶ the received signal at the BS after concurrent transmissions

$$\mathbf{y} = \sum_{k \in \mathcal{K}} \mathbf{h}_k b_k s_k + \mathbf{n}$$

- s_k : the representative signal; b_k : allocated power
 - \mathbf{h}_k : the channel vector between the k -th device and the BS
 - $\mathbf{n} \sim \mathcal{CN}(\mathbf{0}, \sigma^2 \mathbf{I})$: the noise vector
- ▶ the signal after decoding at the BS

$$\hat{\mathbf{w}} = \frac{1}{\sqrt{\eta}} \mathbf{a}^H \mathbf{y} = \frac{1}{\sqrt{\eta}} \mathbf{a}^H \sum_{k \in \mathcal{K}} \mathbf{h}_k b_k s_k + \frac{1}{\sqrt{\eta}} \mathbf{a}^H \mathbf{n}$$

- η : a normalizing factor
- \mathbf{a} : the receive beamforming vector at the BS

Problem Formulation

- The distortion of decoded signal $\hat{\mathbf{w}}$ and ideal signal $\mathbf{w} = \sum_{k \in \mathcal{K}} s_k$ is measured by mean-square-error (MSE)

$$\begin{aligned}\text{MSE}(\hat{\mathbf{w}}, \mathbf{w}; \mathbf{a}) &= \mathbb{E}(|\hat{\mathbf{w}} - \mathbf{w}|^2) \\ &= \sum_k \left| \mathbf{a}^H \mathbf{h}_k b_k / \sqrt{\eta} - 1 \right|^2 + \sigma^2 \|\mathbf{a}\|^2 / \eta\end{aligned}$$

which can be further simplified as

$$\text{MSE}(\hat{\mathbf{w}}, \mathbf{w}; \mathbf{a}) = \frac{\|\mathbf{a}\|^2 \sigma^2}{\eta} = \frac{\|\mathbf{a}\|^2 \sigma^2}{P_0 \min_{k \in \mathcal{K}} \|\mathbf{a}^H \mathbf{h}_k\|^2}$$

by setting $b_k = \sqrt{\eta} \frac{(\mathbf{a}^H \mathbf{h}_k)^H}{\|\mathbf{a}^H \mathbf{h}_k\|^2}$ and $\eta = P_0 \min_{k \in \mathcal{K}} \|\mathbf{a}^H \mathbf{h}_k\|^2$.

Problem Formulation

- ▶ **Goal:** find the optimal receive beamforming vector \mathbf{a} to minimize MSE
- ▶ **Problem Formulation**

$$\underset{\mathbf{a} \in \mathbb{C}^N}{\text{minimize}} \left(\max_{k \in \mathcal{K}} \frac{\|\mathbf{a}\|^2}{\|\mathbf{a}^H \mathbf{h}_k\|^2} \right),$$

and it can be recast as

$$\begin{array}{ll} \underset{\mathbf{a} \in \mathbb{C}^N}{\text{minimize}} & \|\mathbf{a}\|^2 \\ \text{subject to} & \|\mathbf{a}^H \mathbf{h}_k\|^2 \geq 1, \forall k \in \mathcal{K}. \end{array}$$

Problem Formulation

- ▶ **Goal:** find the optimal receive beamforming vector \mathbf{a} to minimize MSE
- ▶ **Problem Formulation**

$$\underset{\mathbf{a} \in \mathbb{C}^N}{\text{minimize}} \left(\max_{k \in \mathcal{K}} \frac{\|\mathbf{a}\|^2}{\|\mathbf{a}^H \mathbf{h}_k\|^2} \right),$$

and it can be recast as

$$\begin{aligned} &\underset{\mathbf{a} \in \mathbb{C}^N}{\text{minimize}} && \|\mathbf{a}\|^2 \\ &\text{subject to} && \|\mathbf{a}^H \mathbf{h}_k\|^2 \geq 1, \forall k \in \mathcal{K}. \end{aligned}$$

- ▶ **Challenges:** a nonconvex quadratically constrained quadratic programming (QCQP)
- ▶ **Proposal:** low-rank optimization after matrix lifting

Outline

Motivations

Problem Formulation

Proposed Algorithm

Simulation Results

Summary

Low-Rank Optimization

► define $\mathbf{A} = \mathbf{a}\mathbf{a}^H$, $\mathbf{A} \succeq \mathbf{0}$ with $\text{rank}(\mathbf{A}) = 1$

► **Problem Rewrite**

$$\begin{array}{ll} \underset{\mathbf{A} \in \mathbb{C}^{N \times N}}{\text{minimize}} & \text{Tr}(\mathbf{A}) \\ \text{subject to} & \text{Tr}(\mathbf{A}\mathbf{H}_k) \geq 1, \forall k \in \mathcal{K} \\ & \mathbf{A} \succeq \mathbf{0}, \text{rank}(\mathbf{A}) = 1 \end{array}$$

Low-Rank Optimization

► define $\mathbf{A} = \mathbf{a}\mathbf{a}^H$, $\mathbf{A} \succeq \mathbf{0}$ with $\text{rank}(\mathbf{A}) = 1$

► **Problem Rewrite**

$$\begin{array}{ll} \underset{\mathbf{A} \in \mathbb{C}^{N \times N}}{\text{minimize}} & \text{Tr}(\mathbf{A}) \\ \text{subject to} & \text{Tr}(\mathbf{A}\mathbf{H}_k) \geq 1, \forall k \in \mathcal{K} \\ & \mathbf{A} \succeq \mathbf{0}, \text{rank}(\mathbf{A}) = 1 \end{array}$$

► **Challenges:** the nonconvex rank-one constraint

► **Proposal:** a DC representation for the rank-one constraint

DC Reformulation

► DC Representation

$$\text{rank}(\mathbf{A}) = 1 \iff \text{Tr}(\mathbf{A}) - \|\mathbf{A}\|_2 = 0, \text{Tr}(\mathbf{A}) > 0$$

► DC Reformulation

$$\begin{array}{ll} \underset{\mathbf{A} \in \mathbb{C}^{N \times N}}{\text{minimize}} & \text{Tr}(\mathbf{A}) + \beta (\text{Tr}(\mathbf{A}) - \|\mathbf{A}\|_2) \\ \text{subject to} & \text{Tr}(\mathbf{A}\mathbf{H}_k) \geq 1, \forall k \in \mathcal{K} \\ & \mathbf{A} \succeq \mathbf{0}, \text{Tr}(\mathbf{A}) > 0 \end{array}$$

DC Algorithm

- At iteration t , \mathbf{A}^t is obtained by solving subproblem

$$\begin{aligned} & \underset{\mathbf{A} \in \mathbb{C}^{N \times N}}{\text{minimize}} && (1 + \beta) \text{Tr}(\mathbf{A}) - \beta \langle \partial \|\mathbf{A}^t\|_2, \mathbf{A} \rangle \\ & \text{subject to} && \text{Tr}(\mathbf{A}\mathbf{H}_k) \geq 1, \forall k \in \mathcal{K} \\ & && \mathbf{A} \succeq \mathbf{0}, \text{Tr}(\mathbf{A}) > 0 \end{aligned} \quad ,$$

where $\partial \|\mathbf{A}^t\|_2$ is one of subgradients of the spectral norm at point \mathbf{A}^t , and $\langle \cdot, \cdot \rangle$ is the inner product of two matrices defined as $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{Real}(\text{Tr}(\mathbf{X}^H \mathbf{Y}))$

DC Algorithm

- At iteration t , \mathbf{A}^t is obtained by solving subproblem

$$\begin{aligned} & \underset{\mathbf{A} \in \mathbb{C}^{N \times N}}{\text{minimize}} && (1 + \beta) \text{Tr}(\mathbf{A}) - \beta \langle \partial \|\mathbf{A}^t\|_2, \mathbf{A} \rangle \\ & \text{subject to} && \text{Tr}(\mathbf{A}\mathbf{H}_k) \geq 1, \forall k \in \mathcal{K} \quad , \\ & && \mathbf{A} \succeq \mathbf{0}, \text{Tr}(\mathbf{A}) > 0 \end{aligned}$$

where $\partial \|\mathbf{A}^t\|_2$ is one of subgradients of the spectral norm at point \mathbf{A}^t , and $\langle \cdot, \cdot \rangle$ is the inner product of two matrices defined as $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{Real}(\text{Tr}(\mathbf{X}^H \mathbf{Y}))$

- Repeat the above DC algorithm until convergence for a feasible \mathbf{A} with exact rank-one; then \mathbf{a} is obtained via singular value decomposition (SVD).

Outline

Motivations

Problem Formulation

Proposed Algorithm

Simulation Results

Summary

A Quick Review

► Challenges

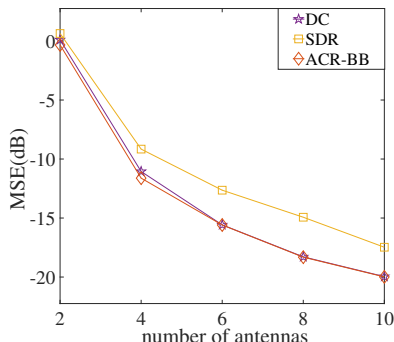
- possible signal distortion in wireless communications
- slow convergence because a large number of iterations is required to train a satisfactory model

► Our Work

- propose a difference-of-convex-functions (DC) algorithm to minimize signal distortion
- fasten model convergence by adopting the canonical Newton's method for local model updates

Simulation Results

- $K = 5$, results averaged over 100 independently generated channel realizations



SDR[Sidiropoulos et al.'06]: convexify the nonconvex QCQP by simply dropping the rank-one constraint

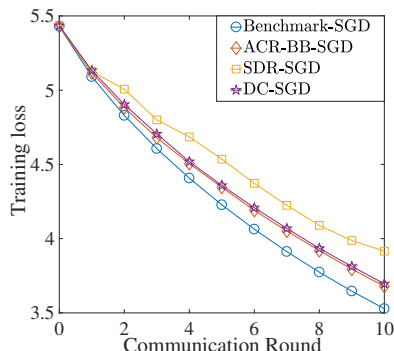
ACR-BB[Lu et al.'17]: a global approach for QCQP based on the branch-and-bound algorithm

Remark

- The proposed DC algorithm achieves nearly-optimal performance on minimizing MSE.

Simulation Results

- Classification experiment over CIFAR10 dataset
 - train a softmax classifier via the distributed stochastic gradient descent (SGD)



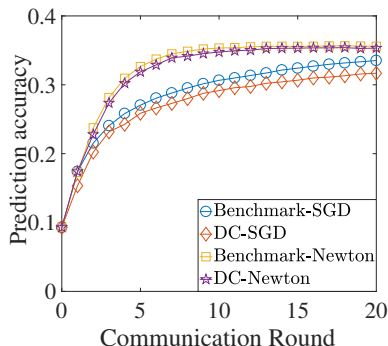
Benchmark: ideal transmission with no aggregation errors, i.e., $\text{MSE} = 0\text{dB}$

Remark

- The aggregation errors significantly degrade performance.

Simulation Results

- Classification experiment over CIFAR10 dataset
 - train a softmax classifier via the canonical Newton's method



Remark

- The Newton's method significantly fasten the convergence of trained model and is much more robust to the aggregation errors.

Outline

Motivations

Problem Formulation

Proposed Algorithm

Simulation Results

Summary

Concluding Remarks

- ▶ We propose a **second-order based** model update method for on-device federated learning.
- ▶ We develop a low-rank approach to support over-the-air computation, followed by a **novel DC algorithm**.
- ▶ We demonstrate the connection between aggregation errors and model convergence behaviors through experiments, i.e., large aggregation error \implies slow convergence
- ▶ Second order methods benefit from two aspects:
 1. reduce the communication burden because much less communication rounds are required for convergence
 2. be **much more robust to aggregation errors** therefore errors result in very limited performance loss

Thanks !

huasheng@shanghaitech.edu.cn