

Stochastic Gradient Push for Distributed Deep Learning

+ HHHFL: Hierarchical Heterogeneous
Horizontal Federated Learning for
Electroencephalography

Tao Shen

Zhejiang University

November 16, 2019

Stochastic Gradient Push for Distributed Deep Learning

Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, Mike Rabbat

Distributed Optimization

Problem formulation

$$\begin{aligned} \min_{\mathbf{x}_i \in \mathbb{R}^d, i=1, \dots, n} \quad & \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\xi_i \sim D_i} F_i(\mathbf{x}_i; \xi_i) \\ \text{subject to} \quad & \mathbf{x}_i = \mathbf{x}_j, \forall i, j = 1, \dots, n \end{aligned}$$

Two Principle

1. Fit Local Model and Data (Training)
2. Fit Local Model and Other Model (Consensus)

Consensus

Approximate distributed averaging

$$\frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^{(0)} \quad \mathbf{y}_i^{(0)} \in \mathbb{R}^d \quad \mathbf{Y}^{(0)} \in \mathbb{R}^{n \times d}$$

$$\mathbf{y}_i^{(k+1)} = \sum_{j=1}^n p_{i,j}^{(k)} \mathbf{y}_j^{(k)}$$

$$\mathbf{Y}^{(k+1)} = \mathbf{P}^{(k)} \mathbf{Y}^{(k)} \quad \mathbf{P}^{(k)} \in \mathbb{R}^{n \times n}$$

Doubly-stochastic Matrix

$$\mathbf{P} = \left(\begin{array}{cc} 0.2 & 0.8 \\ 0.8 & 0.2 \end{array} \right), \left(\begin{array}{cc} 0.3 & 0.7 \\ 0.7 & 0.3 \end{array} \right), \left(\begin{array}{cc} 0.5 & 0.5 \\ 0.5 & 0.5 \end{array} \right)$$

$$\lim_{K \rightarrow \infty} \prod_{k=0}^K \mathbf{P}^{(k)} = \{p(i,j) = \frac{1}{n}\} \quad \mathbf{y}_i^{(\infty)} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^{(0)}$$

Consensus

Column-stochastic Matrix

$$\mathbf{Y}^{(k+1)} = \mathbf{P}^{(k)} \mathbf{Y}^{(k)} \quad \mathbf{P}^{(k)} \in \mathbb{R}^{n \times n}$$

$$\mathbf{P} = \begin{pmatrix} 0.2 & 0.3 \\ 0.8 & 0.7 \end{pmatrix}, \begin{pmatrix} 0.3 & 0.4 \\ 0.7 & 0.6 \end{pmatrix}, \begin{pmatrix} 0.23 & 0.23 \\ 0.76 & 0.76 \end{pmatrix}$$

$$\lim_{K \rightarrow \infty} \prod_{k=0}^K \mathbf{P}^{(k)} = \boldsymbol{\pi} \mathbf{1}^\top$$

$$\mathbf{Y}^{(\infty)} = \boldsymbol{\pi} \left(\mathbf{1}^\top \mathbf{Y}^{(0)} \right)$$

$$y_i^{(\infty)} = \pi_i \sum_{j=1}^n y_j^{(0)}$$

Add a scalar parameter $w_i^{(k)}$

Consensus

Column-stochastic Matrix

$$\lim_{K \rightarrow \infty} \prod_{k=0}^K \mathbf{P}^{(k)} = \boldsymbol{\pi} \mathbf{1}^\top$$

$$\mathbf{w}^{(k+1)} = \mathbf{P}^{(k)} \mathbf{w}^{(k)} \quad w_i^{(0)} = 1$$

$$\mathbf{w}^{(\infty)} = \boldsymbol{\pi} \left(\mathbf{1}^\top \mathbf{w}^{(0)} \right)$$

$$w_i^{(\infty)} = \pi_i n$$

De-biased ratio

$$\frac{y_i^{(\infty)}}{w_i^{(\infty)}} = \frac{1}{n} \sum_{i=1}^n y_i^{(0)}$$

Stochastic Gradient Push

Algorithm 1 Stochastic Gradient Push (SGP)

Require: Initialize $\gamma > 0$, $\mathbf{x}_i^{(0)} = \mathbf{z}_i^{(0)} \in \mathbb{R}^d$ and $w_i^{(0)} = 1$ for all nodes $i \in \{1, 2, \dots, n\}$

- 1: **for** $k = 0, 1, 2, \dots, K$, at node i , **do**
 - 2: Sample new mini-batch $\xi_i^{(k)} \sim \mathcal{D}_i$ from local distribution
 - 3: Compute mini-batch gradient at $\mathbf{z}_i^{(k)}$: $\nabla \mathbf{F}_i(\mathbf{z}_i^{(k)}; \xi_i^{(k)})$
 - 4: $\mathbf{x}_i^{(k+\frac{1}{2})} = \mathbf{x}_i^{(k)} - \gamma \nabla \mathbf{F}_i(\mathbf{z}_i^{(k)}; \xi_i^{(k)})$
 - 5: Send $(p_{j,i}^{(k)} \mathbf{x}_i^{(k+\frac{1}{2})}, p_{j,i}^{(k)} w_i^{(k)})$ to out-neighbors;
 receive $(p_{i,j}^{(k)} \mathbf{x}_j^{(k+\frac{1}{2})}, p_{i,j}^{(k)} w_j^{(k)})$ from in-neighbors
 - 6: $\mathbf{x}_i^{(k+1)} = \sum_j p_{i,j}^{(k)} \mathbf{x}_j^{(k+\frac{1}{2})}$
 - 7: $w_i^{(k+1)} = \sum_j p_{i,j}^{(k)} w_j^{(k)}$
 - 8: $\mathbf{z}_i^{(k+1)} = \mathbf{x}_i^{(k+1)} / w_i^{(k+1)}$
 - 9: **end for**
-

HHHFL: Hierarchical Heterogeneous Horizontal Federated Learning for Electroencephalography

DashanGao, CeJu, Xiguang Wei, Yang Liu, Tianjian Chen, Qiang Yang

Heterogeneous Data

Electroencephalography (EEG)

Heterogeneous & Privacy

Data Heterogeneity

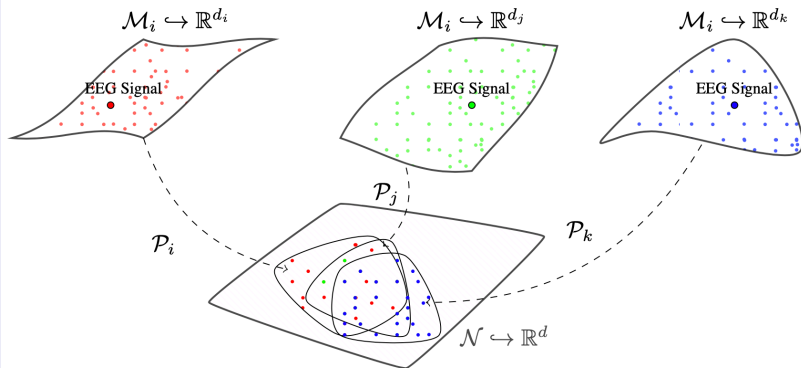
Heterogeneous Domain Adaptation

Privacy-Preserving

Federated Learning

Heterogeneous Domain Adaptation

Heterogeneous Domain Adaptation



$\mathcal{P}_i?$

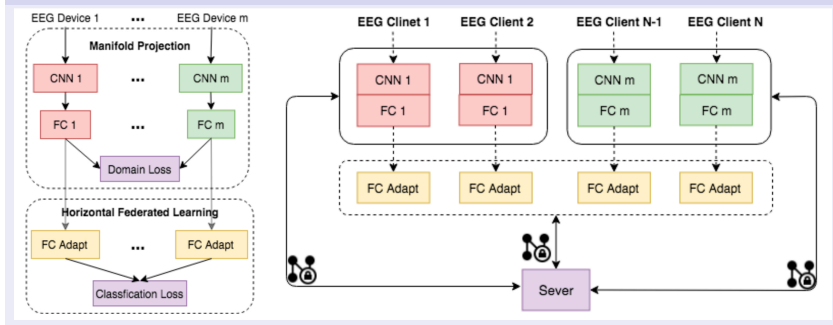
Neural Network Approach

Heterogeneous Domain Adaptation

Loss Function

Classification Loss & Domain Loss

Heterogeneous Domain Adaptation & Federated Learning



Maximum Mean Discrepancy (MMD)

Maximum Mean Miscrepancy (MMD)

$$\mathcal{L} := \mathcal{L}_C(X_{EEG}, Y) + \sum_{1 \leq i < j \leq m} \lambda_{i,j} \cdot \text{MMD}^2(Q_i, Q_j)$$

$$\text{MMD}(Q_i, Q_j) := \left\| \mathbb{E}_{\mathcal{P}_i(X_i) \sim Q_i} \psi(\mathcal{P}_i(X_i)) - \mathbb{E}_{\mathcal{P}_j(X_j) \sim Q_j} \psi(\mathcal{P}_j(X_j)) \right\|_{\mathcal{H}}$$
$$\mathcal{P}_i(\{x_1, \dots, x_{N_i}\}) \sim Q_i, \quad \psi : \mathcal{N} \longrightarrow \mathcal{H}$$

Reproducing Kernel Hilbert Space (RKHS)

Maximum Mean Discrepancy (MMD)

MMD in Euclidean Space

$$MMD[F, p, q] := \sup_{f \in F} (E_{x \sim p}[f(x)] - E_{y \sim q}[f(y)])$$

$$MMD[F, X, Y] := \sup_{f \in F} \left(\frac{1}{m} \sum_{i=1}^m f(x_i) - \frac{1}{n} \sum_{j=1}^n f(x_j) \right)$$

F is rich enough but not too much

Maximum Mean Discrepancy (MMD)

Reproducing Kernel Hilbert Space (RKHS)

$$f(x) = \langle f, \phi(x) \rangle_{\mathcal{H}} \quad f \rightarrow f(x)$$

$$E_p[f(x)] = \langle f, E_p[\phi(x)] \rangle_{\mathcal{H}}$$

MMD in RKHS

$$\begin{aligned} MMD[F, p, q] &= \sup_{\|f\|_{\mathcal{H}} \leq 1} E_p[f(x)] - E_q[f(y)] \\ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} E_p[\langle \phi(x), f \rangle_{\mathcal{H}}] - E_q[\langle \phi(y), f \rangle_{\mathcal{H}}] \\ &= \sup_{\|f\|_{\mathcal{H}} \leq 1} \langle \mu_p - \mu_q, f \rangle_{\mathcal{H}} \\ &= \|\mu_p - \mu_q\|_{\mathcal{H}} \end{aligned}$$

Maximum Mean Discrepancy (MMD)

MMD in RKHS

$$\begin{aligned} \text{MMD}^2[F, p, q] &:= \langle \mu_p - \mu_q, \mu_p - \mu_q \rangle_{\mathcal{H}} \\ &= \langle \mu_p, \mu_q \rangle_{\mathcal{H}} + \langle \mu_q, \mu_p \rangle_{\mathcal{H}} - 2 \langle \mu_p, \mu_q \rangle_{\mathcal{H}} \\ &= E_p \langle \phi(x), \phi(x') \rangle_{\mathcal{H}} + E_q \langle \phi(y), \phi(y') \rangle_{\mathcal{H}} \\ &\quad - 2 E_{p,q} \langle \phi(x), \phi(y) \rangle_{\mathcal{H}} \end{aligned}$$

Universal Kernel

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$$

MMD in Euclidean Space

$$\text{MMD}^2[F, p, q] = \frac{1}{m(m-1)} \sum_{i \neq j}^m k(x_i, x_j) + \frac{1}{n(n-1)} \sum_{i \neq j}^n k(y_i, y_j) - \frac{2}{mn} \sum_{i,j=1}^{m,n} k(x_i, y_j)$$

The End