

Federated Multi-Task Learning

Federated Update of W

Algorithm 1 MOCHA: Federated Multi-Task Learning Framework

```
1: Input: Data  $\mathbf{X}_t$  from  $t = 1, \dots, m$  tasks, stored on one of  $m$  nodes, and initial matrix  $\mathbf{\Omega}_0$ 
2: Starting point  $\boldsymbol{\alpha}^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{v}^{(0)} := \mathbf{0} \in \mathbb{R}^b$ 
3: for iterations  $i = 0, 1, \dots$  do
4:   Set subproblem parameter  $\sigma'$  and number of federated iterations,  $H_i$ 
5:   for iterations  $h = 0, 1, \dots, H_i$  do
6:     for tasks  $t \in \{1, 2, \dots, m\}$  in parallel over  $m$  nodes do
7:       call local solver, returning  $\theta_t^h$ -approximate solution  $\Delta\boldsymbol{\alpha}_t$  of the local subproblem (4)
8:       update local variables  $\boldsymbol{\alpha}_t \leftarrow \boldsymbol{\alpha}_t + \Delta\boldsymbol{\alpha}_t$ 
9:       return updates  $\Delta\mathbf{v}_t := \mathbf{X}_t\Delta\boldsymbol{\alpha}_t$ 
10:    reduce:  $\mathbf{v}_t \leftarrow \mathbf{v}_t + \Delta\mathbf{v}_t$ 
11:   Update  $\mathbf{\Omega}$  centrally based on  $\mathbf{w}(\boldsymbol{\alpha})$  for latest  $\boldsymbol{\alpha}$ 
12: Central node computes  $\mathbf{w} = \mathbf{w}(\boldsymbol{\alpha})$  based on the latest  $\boldsymbol{\alpha}$ 
13: return:  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_m]$ 
```

Convergence Analysis

Definition (Per-Node-Per-Iteration-Approximation Parameter)

At each iteration h , we define the accuracy level of the solution calculated by node k to its subproblem as

$$\theta_k^h := \frac{G_k^{\sigma'}(\Delta\alpha_k^{(h)}; v^{(h)}, \alpha_k^{(h)}) - G_k^{\sigma'}(\Delta\alpha_k^\star; v^{(h)}, \alpha_k^{(h)})}{G_k^{\sigma'}(\mathbf{0}; v^{(h)}, \alpha_k^{(h)}) - G_k^{\sigma'}(\Delta\alpha_k^\star; v^{(h)}, \alpha_k^{(h)})}$$

$\theta_k^h \in [0, 1]$, $\theta_k^h = 1$ means that no updates to the subproblem are made at iteration h

Assumption

Let $\mathcal{H}_h := (\alpha^{(h)}, \dots, \alpha^{(1)})$ be the *dual vector history* until the beginning of iteration h , and define

$\Theta_k^h := \mathbb{E}[\theta_k^h | \mathcal{H}_h]$. For all tasks k and all iterations h , we assume $p_k^h := \mathbb{P}(\theta_k^h = 1) \leq p_{\max} < 1$ and

$\hat{\Theta}_k^h = \mathbb{E}[\theta_k^h | \mathcal{H}_h, \theta_k^h < 1] \leq \Theta_{\max} < 1$.