

CS-E4740 - Federated Learning

FL Design Principle

Assoc. Prof. Alexander Jung

Spring 2025

Playlist



Glossary



Course Site



Table of Contents

Formulating FL as Optimization

Computational Aspects

Statistical Aspects

Interpretations

Table of Contents

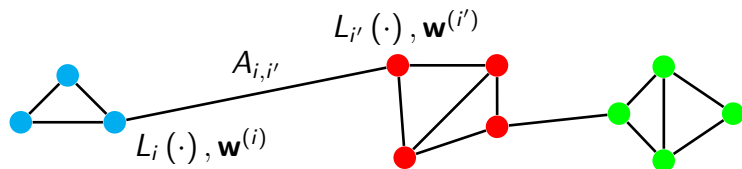
Formulating FL as Optimization

Computational Aspects

Statistical Aspects

Interpretations

An FL Network



- ▶ FL network consisting of devices $i=1, \dots, n$.
- ▶ Some i, i' connected by an edge with weight $A_{i,i'} > 0$.
- ▶ Device i learns model params. $\mathbf{w}^{(i)} \in \mathbb{R}^d$.
- ▶ Usefulness of $\mathbf{w}^{(i)}$ measured by some local loss, e.g.,

$$L_i(\mathbf{w}^{(i)}) := \frac{1}{m_i} \sum_{r=1}^{m_i} \left(y^{(i,r)} - (\mathbf{w}^{(i)})^T \mathbf{x}^{(i,r)} \right)^2.$$

FL via Regularization

- ▶ Each node carries a linear model $h^{(\mathbf{w}^{(i)})}(\mathbf{x}) := \mathbf{x}^T \mathbf{w}^{(i)}$.
- ▶ Each node carries m_i labelled data points.
- ▶ Node-wise ML fails if $m_i \ll d$ (overfitting).

Idea:

Use the neighbours $\mathcal{N}^{(i)} := \{i' : \{i, i'\} \in \mathcal{E}\}$ to regularize!

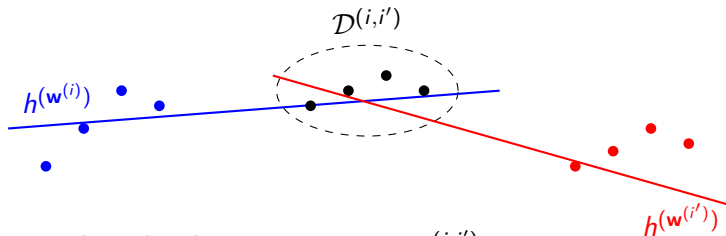
FL via Regularization (ctd.)

As for basic ML, regularization can be done either via

- ▶ **Data augmentation** using data from the neighbours.
- ▶ **Prune local models** by requiring them to agree across edges.
- ▶ **Add a penalty term** to the local loss function.

Building a Penalty Across Edges

- ▶ Consider two nodes i, i' with local datasets $\mathcal{D}^{(i)}, \mathcal{D}^{(i')}$.
- ▶ Assume there is a non-empty overlap $\mathcal{D}^{(i)} \cap \mathcal{D}^{(i')}$.



We penalize the disagreement on $\mathcal{D}^{(i,i')}$:

$$\begin{aligned} \sum_{\mathbf{x} \in \mathcal{D}^{(i,i')}} (h(\mathbf{w}^{(i)})(\mathbf{x}) - h(\mathbf{w}^{(i')})(\mathbf{x}))^2 &= \sum_{\mathbf{x} \in \mathcal{D}^{(i,i')}} (\mathbf{x}^T \mathbf{w}^{(i)} - \mathbf{x}^T \mathbf{w}^{(i')})^2 \\ &= (\mathbf{w}^{(i)} - \mathbf{w}^{(i')})^T \left[\sum_{\mathbf{x} \in \mathcal{D}^{(i,i')}} \mathbf{x}^T \mathbf{x} \right] (\mathbf{w}^{(i)} - \mathbf{w}^{(i')}). \end{aligned}$$

Generalized TV Minimization (GTVMin)

Learn model params. $\widehat{\mathbf{w}}^{(i)}$ by balancing local loss and GTV

$$\min_{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)} \in \mathbb{R}^d} \sum_{i=1}^n \left[L_i(\mathbf{w}^{(i)}) + \alpha \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} \phi(\mathbf{w}^{(i)} - \mathbf{w}^{(i')}) \right]$$

- ▶ Penalty function $\phi(\mathbf{u})$ is a design choice.
- ▶ Previous slide used $\phi(\mathbf{u}) = \mathbf{u}^T \mathbf{Q} \mathbf{u}$ with $\mathbf{Q} := \sum_{\mathbf{x} \in \mathcal{D}(i, i')} \mathbf{x} \mathbf{x}^T$.
- ▶ Our focus is on the choice $\phi(\mathbf{u}) := \|\mathbf{u}\|_2^2$.
- ▶ Another popular choice is $\phi(\mathbf{u}) := \|\mathbf{u}\|_1$.

¹Y. SarcheshmehPour, et.al, "Clustered Federated Learning via Generalized Total Variation Minimization," in IEEE Trans. Sig. Proc, 2023, doi: 10.1109/TSP.2023.3322848.

Model-Agnostic GTVMin

Replacing $\phi(\mathbf{w}^{(i)} - \mathbf{w}^{(i')})$ with the disagreement measure

$$D(h^{(i)}, h^{(i')}) := \sum_{\mathbf{x} \in \mathcal{D}(i, i')} (h^{(i)}(\mathbf{x}) - h^{(i')}(\mathbf{x}))^2$$

yields a model-agnostic generalization of GTVMin

$$\min_{\substack{h^{(i)} \in \mathcal{H}^{(i)} \\ i \in \mathcal{V}}} \sum_{i \in \mathcal{V}} \left[L_i(h^{(i)}) + \alpha \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} D(h^{(i)}, h^{(i')}) \right].$$

This allows for VERY heterogeneous FL networks, e.g., $\mathcal{H}^{(1)} = \text{lin.model}$, $\mathcal{H}^{(2)} = \text{LLM}$, $\mathcal{H}^{(3)} = \text{decision tree}$.

Table of Contents

Formulating FL as Optimization

Computational Aspects

Statistical Aspects

Interpretations

Computational Aspects

$$\min_{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)} \in \mathbb{R}^d} \sum_{i=1}^n \left[L_i(\mathbf{w}^{(i)}) + \alpha \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} \phi(\mathbf{w}^{(i)} - \mathbf{w}^{(i')}) \right]$$

- ▶ How can we solve it efficiently over an FL network?
- ▶ How much compute/comm. is needed at least?
- ▶ What is the effect of different choices for the edges \mathcal{E} , loss funcs. $L_i(\cdot)$, and GTV penalty ϕ ?

Computational Aspects - Smooth GTVmin

Consider a GTVMin instance

$$\min_{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)} \in \mathbb{R}^d} \sum_{i=1}^n \left[L_i(\mathbf{w}^{(i)}) + \alpha \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(i')} \right\|_2^2 \right]$$

with a smooth (differentiable) $L_i(\cdot)$.

If we use (distributed) gradient descent to solve GTVMin:

- ▶ How many iterations should we run?
- ▶ What is a good choice for the learning rate?
- ▶ How to communicate gradients over comm. links?

Characterizing GTVMin Solutions

- ▶ Consider GTVMin solution $\hat{\mathbf{w}}^{(i)} \in \mathbb{R}^d$, for $i = 1, \dots, n$.
- ▶ We stack them into a long vector

$$\hat{\mathbf{w}} := (\hat{\mathbf{w}}^{(1)}, \dots, \hat{\mathbf{w}}^{(n)})^T \in \mathbb{R}^{dn}.$$

- ▶ We characterize the solutions as a fixed-point of some \mathcal{F} ,

$$\hat{\mathbf{w}} \text{ solves GTVMin} \Leftrightarrow \hat{\mathbf{w}} = \mathcal{F}\hat{\mathbf{w}}$$

- ▶ The operator \mathcal{F} is not unique (design choice!).

Convex and Smooth GTVMin

Consider GTVMin with a smooth and convex $L_i(\cdot)$,

$$\min_{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)} \in \mathbb{R}^d} \sum_{i=1}^n \left[L_i(\mathbf{w}^{(i)}) + \alpha \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(i')} \right\|_2^2 \right] \quad (1)$$

$$\hat{\mathbf{w}} \text{ solves (1)} \Leftrightarrow \hat{\mathbf{w}} = \mathcal{F}^{(\eta)} \hat{\mathbf{w}}$$

$\mathcal{F}^{(\eta)}$ maps $\mathbf{u} = (\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(n)})^T$ to $\mathbf{v} = (\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)})^T$,

$$\mathbf{v}^{(i)} = \mathbf{u}^{(i)} - \eta \left[\nabla L_i(\mathbf{u}^{(i)}) + 2\alpha \sum_{i' \in \mathcal{N}^{(i)}} A_{i, i'} (\mathbf{u}^{(i)} - \mathbf{u}^{(i')}) \right].$$

Different choices for “step-size” $\eta > 0$ yield different \mathcal{F} .

Fixed-Point Iterations

Q: How to compute a fixed point $\widehat{\mathbf{w}}$ of \mathcal{F} ?

A: Start with initial guess $\widehat{\mathbf{w}}^{(0)}$ and iterate

$$\widehat{\mathbf{w}}^{(k)} = \mathcal{F}\widehat{\mathbf{w}}^{(k-1)}, \text{ for } k = 1, 2, \dots$$

If \mathcal{F} is **firmly non-expansive** $\lim_{k \rightarrow \infty} \widehat{\mathbf{w}}^{(k)} = \widehat{\mathbf{w}}$.²

If \mathcal{F} is even **contractive** with constant $\kappa < 1$,

$$\|\widehat{\mathbf{w}}^{(k)} - \widehat{\mathbf{w}}\|_2 \leq \kappa^k \|\widehat{\mathbf{w}}^{(0)} - \widehat{\mathbf{w}}\|_2.$$

²H. Bauschke, P. Combettes, "Convex Analysis and Monotone Operator Theory in Hilbert Spaces," Springer, 2017.

Gradient Descent as Fixed-Point Iteration

GD for smooth and convex objective function $f(\mathbf{w})$,

$$\mathbf{w}^{(k)} = \mathbf{w}^{(k-1)} - \eta \nabla f(\mathbf{w}^{(k-1)})$$

is a fixed-point iteration with $\mathcal{F}^{(\eta)} : \mathbf{w} \mapsto \mathbf{w} - \eta \nabla f(\mathbf{w})$.

- ▶ In general, $\mathcal{F}^{(\eta)}$ is neither firmly non-exp. nor contractive.
- ▶ Convergence can still be ensured if η is sufficiently small.
- ▶ E.g., using learning rate $\eta_k = 1/k$ for smooth $f(\mathbf{w})$.

Table of Contents

Formulating FL as Optimization

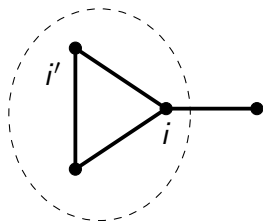
Computational Aspects

Statistical Aspects

Interpretations

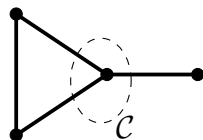
Statistical Aspects

- ▶ GTVMin solution yields model params. $\widehat{\mathbf{w}}^{(i)}$, $i = 1, \dots, n$
- ▶ How useful are these model params. ?
- ▶ The local loss $L_i(\widehat{\mathbf{w}}^{(i)})$ can be misleading (why?)
- ▶ Better to use aggregate $\sum_{i \in \mathcal{C}^{(i)}} L_i(\widehat{\mathbf{w}}^{(i)})$, with cluster

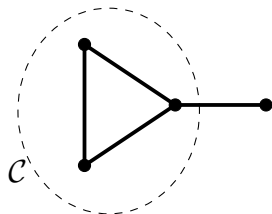


$$\mathcal{C}^{(i)} := \{i' : \widehat{\mathbf{w}}^{(i)} \approx \widehat{\mathbf{w}}^{(i')}\}.$$

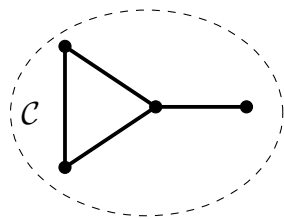
Clustering of GTVMin³



small α



moderate α

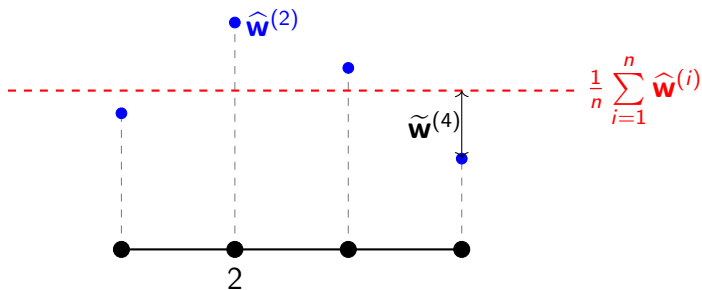


large α

³Y. SarcheshmehPour, Y. Tian, L. Zhang and A. Jung, "Clustered Federated Learning via Generalized Total Variation Minimization," in IEEE Transactions on Signal Processing, 2023,

Analysis of Clustering - Assumptions

- ▶ Consider a connected FL network \mathcal{G} with $\lambda_2 > 0$.
- ▶ Assume loss funcs. satisfy $\min_{\mathbf{v} \in \mathbb{R}^d} \sum_{i=1}^n L_i(\mathbf{v}) \leq \varepsilon$
- ▶ Use GTVMin to learn local params. $\widehat{\mathbf{w}}^{(i)}$.
- ▶ Define the variation $\widetilde{\mathbf{w}}^{(i)} := \widehat{\mathbf{w}}^{(i)} - \frac{1}{n} \sum_{i=1}^n \widehat{\mathbf{w}}^{(i)}$.



Analysis of Clustering - Upper Bound

The variation $\tilde{\mathbf{w}}^{(i)}$ is upper bounded as

$$\sum_{i=1}^n \|\tilde{\mathbf{w}}^{(i)}\|_2^2 \leq \frac{\varepsilon}{\alpha\lambda_2}.$$

This bound involves the

- ▶ connectivity of FL network (via λ_2),
- ▶ the properties of local loss functions (via ε), and
- ▶ the GTVMin parameter α .

A large $\alpha\lambda_2$ results in nearly identical local params. $\tilde{\mathbf{w}}^{(i)} \approx \mathbf{0}$.

Table of Contents

Formulating FL as Optimization

Computational Aspects

Statistical Aspects

Interpretations

Interpretations

We next discuss some interpretations of GTVMin

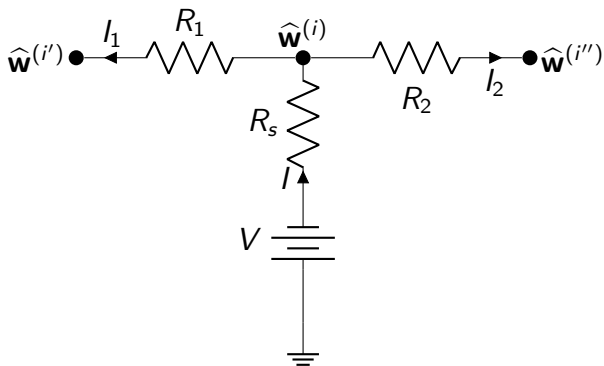
$$\min_{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)} \in \mathbb{R}^d} \sum_{i=1}^n \left[L_i(\mathbf{w}^{(i)}) + \alpha \sum_{\{i, i'\} \in \mathcal{E}} A_{i, i'} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(i')} \right\|_2^2 \right]$$

for some FL network with weighted undirected graph \mathcal{G} and smooth and convex loss func. $L_i(\mathbf{w}^{(i)})$.

We assume that there exists a solution $\hat{\mathbf{w}}^{(1)}, \dots, \hat{\mathbf{w}}^{(n)}$. (Do we really need to make this assumption?)

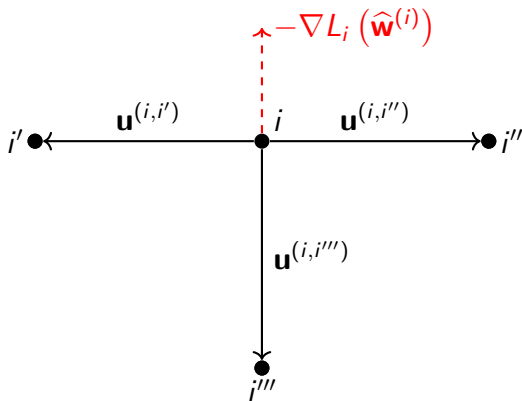
Electronic Circuit

Consider a node i with neighbours $\mathcal{N}^{(i)} = \{i', i''\}$.



$$\underbrace{-\nabla L_i}_{I}(\widehat{\mathbf{w}}^{(i)}) = \left[\underbrace{A_{i,i'}}_{I_1}(\widehat{\mathbf{w}}^{(i)} - \widehat{\mathbf{w}}^{(i')}) + \underbrace{A_{i,i''}}_{I_2}(\widehat{\mathbf{w}}^{(i)} - \widehat{\mathbf{w}}^{(i'')}) \right]$$

Vector-Valued Flows⁴

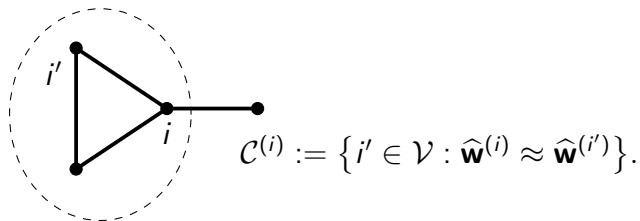


Vector-valued flow $\mathbf{u}^{(i,i')} := \nabla \phi(\mathbf{u}) \big|_{\mathbf{u}=\widehat{\mathbf{w}}^{(i)}-\widehat{\mathbf{w}}^{(i')}}.$

⁴AJ, "On the Duality Between Network Flows and Network Lasso," in IEEE Signal Processing Letters, 2020.

Locally Weighted Learning

GTVMin delivers local params. $\widehat{\mathbf{w}}^{(i)}$ that are clustered.

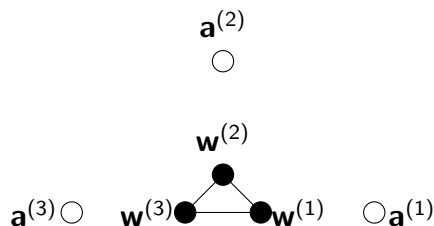


For node i , GTVmin is the same as locally weighted learning⁵

$$\min_{\mathbf{w}^{(i)} \in \mathbb{R}^d} \sum_{i'=1}^n L_{i'}(\mathbf{w}^{(i)}) \rho_{i'} \text{ with } \rho_{i'} = \begin{cases} 1 & \text{if } i' \in \mathcal{C}^{(i)} \\ 0 & \text{, otherwise.} \end{cases}$$

⁵C. G. Atkeson, S. A. Schaal and Andrew W. Moore, Locally Weighted Learning, AI Review, Volume 11, Pages 11-73 (Kluwer Publishers) 1997.

Generalized Convex Clustering⁶



$$\min_{\mathbf{w}^{(1)}, \dots, \mathbf{w}^{(n)}} \sum_{i=1}^n \left\| \mathbf{w}^{(i)} - \mathbf{a}^{(i)} \right\|_2^2 + \alpha \sum_{i, i' \in \mathcal{V}} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(i')} \right\|_2.$$

⁶D. Sun, et.al, Convex Clustering: Model, Theoretical Guarantee and Efficient Algorithm, JMLR, 2021.

What's Next?

The next module applies optimization methods to solve GTVMin.

We can implement these methods as message passing over the edges of an FL network.

Further Resources

- ▶ **YouTube:** @alexjung111
- ▶ **LinkedIn:** Alexander Jung
- ▶ **GitHub:** alexjungaalto

