

Learning in Heterogeneous Networks

Stefan Vlaski

work with Roula Nassif, Mert Kayaalp, Virginia Bordignon, Cédric Richard and Ali H. Sayed

Imperial College London, United Kingdom

Optimization by Quantum and Machine Learning Workshop

**Imperial College
London**

Multi-Agent Systems

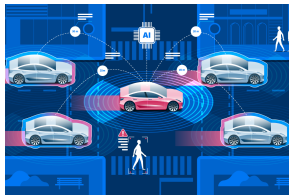


Figure: Autonomous vehicles
[smartcitiesworld.net]

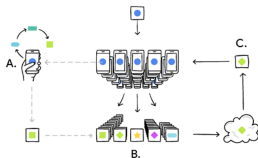


Figure: Phones [ai.googleblog.com]



Figure: Social network [medium.com]

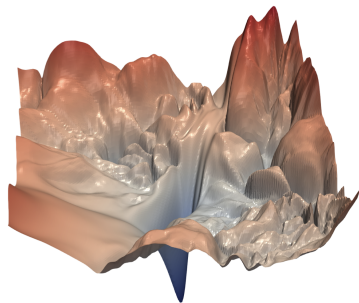
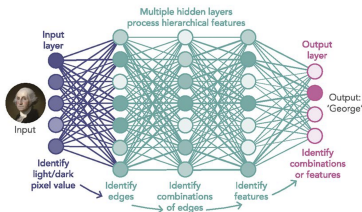


Figure: Drone swarms [ft.com]

Optimisation for Learning

- Learning is building models from data.
- We want the model that fits **best**.
- Optimisation enables “good” learning.

0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9



Towards Provable and Efficient Learning over Networks

- Performance metrics

- ▶ Error probability
- ▶ Excess risk
- ▶ Mean-square deviation
- ▶ Convergence rate
- ▶ First-order stationarity
- ▶ Second-order stationarity
- ▶ ...

- Limitations of distributed systems:

- ▶ Network topology
- ▶ Limited, streaming data
- ▶ Unreliable participation
- ▶ Noisy links
- ▶ Quantisation
- ▶ Privacy
- ▶ Heterogeneity

Aim

Relate performance and limitations in a unified manner to inform design of distributed systems.

Learning Problems as Optimisation Problems

- Consider a linear regression (channel estimation) problem:

$$\gamma_{k,i} = w_k^{o\top} \mathbf{h}_{k,i} + \mathbf{v}_{k,i} \quad (1)$$

- We can pursue the solution via least mean-squares:

$$\arg \min_{w_k} \mathbb{E} \|\gamma_{k,i} - w_k^\top \mathbf{h}_{k,i}\|^2 \quad (2)$$

- If the relation is non-linear, we may use a deep neural net:

$$\arg \min_{w_k} \mathbb{E} \|\gamma_{k,i} - \sigma(W_{k,L} \cdot \sigma(W_{k,L-1} \cdot \sigma(W_{k,1} \mathbf{h}_{k,i})))\|^2 \quad (3)$$

- Unified formulation:**

$$\arg \min_{w_k} \mathbb{E} Q(w_k, \mathbf{x}_k) \quad (4)$$

Learning Paradigms

- **Non-cooperative learning:**

$$w_k^o = \arg \min_{w_k} \mathbb{E}Q(w_k, \mathbf{x}_k) \quad (5)$$

- **Single-task learning** (i.e., consensus optimization):

$$w^o = \arg \min_w \sum_{k=1}^K p_k \mathbb{E}Q(w, \mathbf{x}_k) \quad (6)$$

- **Multi-task learning** over aggregate tasks $\mathcal{W} = \text{col} \{w_k\}$:

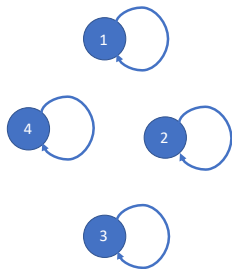
$$w^o = \arg \min_{\mathcal{W}} \sum_{k=1}^K \mathbb{E}Q(w_k, \mathbf{x}_k) + \frac{\eta}{2} \mathcal{R}(\mathcal{W}) \quad (7)$$

$$\text{subject to } \mathcal{W} \in \Omega \quad (8)$$

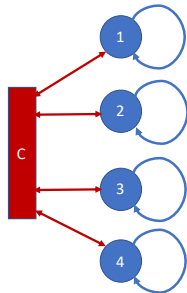
► $\mathcal{R}(\mathcal{W})$ and Ω encode priors on w_k .

Communication Paradigms

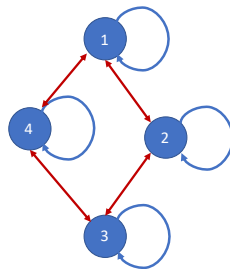
- Restrictions on the flow of information.



Non-cooperative



Centralized/federated



Decentralized

- Communication efficiency,
- robustness to failure,
- privacy.

Single- vs. Multi-Task Learning

Single-Task Objective:

$$\sum_{k=1}^K p_k \mathbb{E} Q(w, \mathbf{x}_k)$$

Example algorithm [Chen and Sayed '12]:

$$\begin{aligned}\psi_{k,i} &= \mathbf{w}_{k,i-1} - \mu \widehat{\nabla} J_k(\mathbf{w}_{k,i-1}) \\ \mathbf{w}_{k,i} &= \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \psi_{\ell,i}\end{aligned}$$

Multi-Task Objective:

$$\begin{aligned}\sum_{k=1}^K \mathbb{E} Q(w_k, \mathbf{x}_k) + \frac{\eta}{2} \mathcal{R}(w) \\ \text{subject to } w \in \Omega\end{aligned}$$

Generic framework:

$$\begin{aligned}\psi_{k,i} &= \mathbf{w}_{k,i-1} - \mu \widehat{\nabla} J_k(\mathbf{w}_{k,i-1}) \\ \mathbf{w}_{k,i} &= \text{Agg} \left(\{ \psi_{\ell,i} \}_{\ell \in \mathcal{N}_k} \right)\end{aligned}$$

Example 1 – Regularized Multitask Learning

- Relationship prior through smoothness regularization:

$$\mathbf{w}_\eta^o = \arg \min_{\mathcal{W}} \sum_{k=1}^K \mathbb{E} Q(\mathbf{w}_k; \mathbf{x}_k) + \frac{\eta}{2} \mathbf{w}^\top (L \otimes I) \mathbf{w} \quad (9)$$

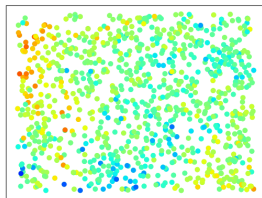
- Example:**

$$\mathbf{w}_\eta^o = \arg \min_{\mathcal{W}} \frac{1}{2\sigma_v^2} \sum_{k=1}^K \mathbb{E} \|\gamma_k - \mathbf{h}_k^\top \mathbf{w}_k\|^2 + \frac{\eta}{2} \mathbf{w}^\top (L \otimes I) \mathbf{w}$$

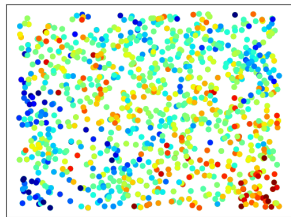
denotes the maximum a posteriori estimate for a linear model with GMRF prior:

$$\gamma_k = \mathbf{h}_k^\top \mathbf{w}_k^o + \mathbf{v}_k$$

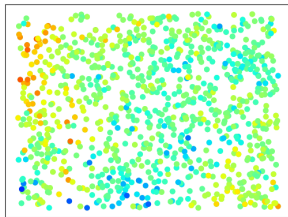
$$f(\mathcal{W}) \triangleq (2\pi)^{-M(K-1)/2} (|\mathcal{L}|^*)^{1/2} e^{-\eta \frac{1}{2} \mathbf{w}^\top \mathcal{L} \mathbf{w}}$$



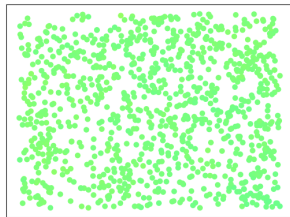
Visualizing GMRFs



$$\eta = 0.001$$



$$\eta = 0.005$$



$$\eta = 1$$

- Graph chosen according to Euclidean distance, η controls level of smoothness.

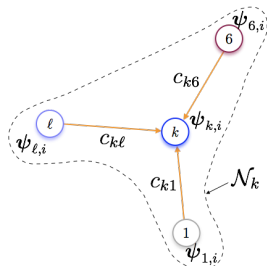
Multitask learning under smoothness

- Network global cost:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{k=1}^N J_k(\mathbf{w}_k) + \frac{\eta}{2} \mathbf{w}^\top (L \otimes I) \mathbf{w}$$

- Multitask learning algorithm:

$$\begin{cases} \psi_{k,i} = \mathbf{w}_{k,i-1} - \mu \widehat{\nabla J_k}(\mathbf{w}_{k,i-1}) & (\text{self-learning}) \\ \mathbf{w}_{k,i} = \left(1 - \mu\eta \sum_{\ell \in \mathcal{N}_k} c_{k\ell}\right) \psi_{k,i} + \mu\eta \sum_{\ell \in \mathcal{N}_k} c_{k\ell} \psi_{\ell,i} & (\text{social learning}) \end{cases}$$

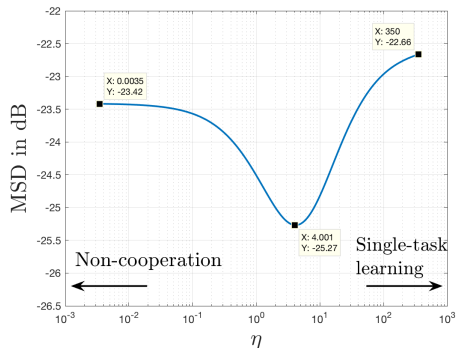


Analytical Performance Guarantee

$$\text{MSD} \approx \underbrace{\overline{\text{MSD}}}_{O(\mu), \eta} + \underbrace{\|w_\eta^o - w^o\|^2}_{\text{smoothness}, \eta}$$

By increasing η :

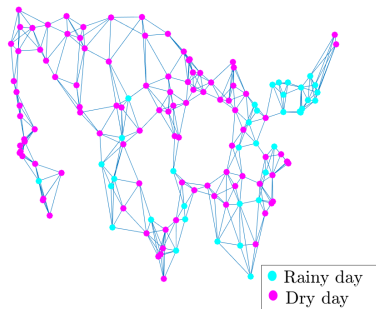
- First term is more likely to decrease
- Second term is more likely to increase and the size of this increase depends on the smoothness of w^o



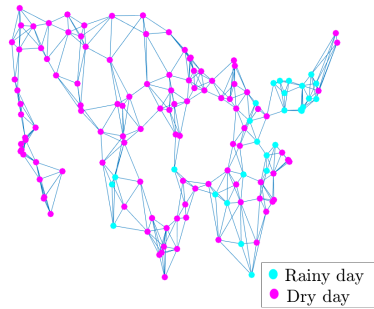
Application: Weather forecasting

Graph constructed based on geodesical distance (4-nearest neighbors). Training set (2004 – 2012), test set (2013 – 2017)

η	0	10	45	100	1000	μ^{-1}
prediction error	0.309	0.232	0.225	0.226	0.228	0.232



Occurrence of rain on July 30, 2015



Prediction of rain occurrence ($\eta = 45$)

Example 2 – Subspace Prior

- An alternative model-based setting may be one where tasks are not necessarily smooth over the graph, but instead linearly related, i.e., $w \in \text{Range}(\mathcal{U})$ for some \mathcal{U} .

$$w_{\mathcal{U}}^o = \arg \min_w J(w) \triangleq \sum_{k=1}^K J_k(w_k), \quad (10)$$

subject to $w \in \text{Range}(\mathcal{U})$,

where $\text{Range}(\cdot)$ denotes the range space operator, and \mathcal{U} is an $KM \times P$ full-column rank matrix with $P \ll KM$.

- Appears naturally in many machine learning and signal processing applications.

Problem Formulation

- Network global cost (\mathcal{U} full column-rank) [Nassif et al., 2020, Di Lorenzo et al., 2020]:

$$\begin{aligned} \mathcal{W}^* = \quad & \arg \min_{\mathcal{W}} \sum_{k=1}^N J_k(\mathbf{w}_k) \\ & \text{subject to } \mathcal{W} \in \text{Range}(\mathcal{U}) \end{aligned}$$

- Centralized stochastic gradient projection approach (\mathcal{P}_u : orthogonal projector onto $\text{Range}(\mathcal{U})$)

$$\mathbf{w}_i = \mathcal{P}_u \left(\mathbf{w}_{i-1} - \mu \text{col} \left\{ \widehat{\nabla_{\mathbf{w}_k} J_k(\mathbf{w}_{k,i-1})} \right\}_{k=1}^N \right)$$

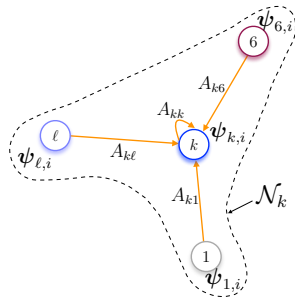
Algorithm

- To get a decentralized algorithm, replace \mathcal{P}_u by an \mathcal{A} such that:

$$\lim_{i \rightarrow \infty} \mathcal{A}^i = \mathcal{P}_u, \quad A_{k\ell} = [\mathcal{A}]_{k\ell} = 0 \text{ if } \ell \notin \mathcal{N}_k$$

- Multitask learning algorithm [Nassif et al., 2020]:

$$\begin{cases} \psi_{k,i} = \mathbf{w}_{k,i-1} - \mu \widehat{\nabla J_k}(\mathbf{w}_{k,i-1}) & (\text{self-learning}) \\ \mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} A_{k\ell} \psi_{\ell,i} & (\text{social learning}) \end{cases}$$



Performance

- Network global cost:

$$\begin{aligned} \mathcal{W}^* = & \arg \min_{\mathcal{W}} \sum_{k=1}^N J_k(w_k) \\ & \text{subject to } \mathcal{W} \in \text{Range}(\mathcal{U}) \end{aligned}$$

- Data characteristics: $H_k = \nabla^2 J_k(w_k^*)$, $R_k = \mathbb{E}[\mathbf{s}_{k,i}(w_k^*) \mathbf{s}_{k,i}^\top(w_k^*)]$

$$\mathcal{H} = \text{diag}\{H_1, \dots, H_N\}, \quad \mathcal{S} = \text{diag}\{R_1, \dots, R_N\}$$

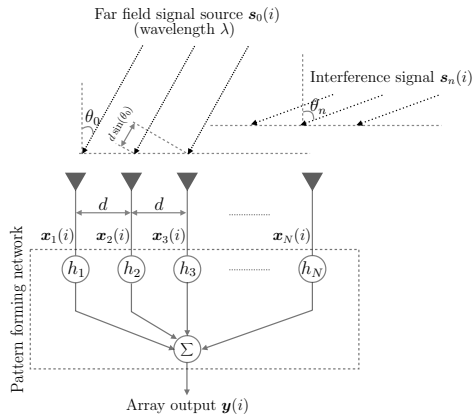
Mean-square-error w.r.t. \mathcal{W}^* (in the small adaptation regime) [Nassif et al., 2020]

$$\text{MSD} = \lim_{i \rightarrow \infty} \frac{1}{N} \mathbb{E} \|\mathcal{W}^* - \mathcal{W}_i\|^2 \approx \frac{\mu}{2N} \text{Tr} \left(\left(\mathcal{U}^\top \mathcal{H} \mathcal{U} \right)^{-1} \left(\mathcal{U}^\top \mathcal{S} \mathcal{U} \right) \right)$$

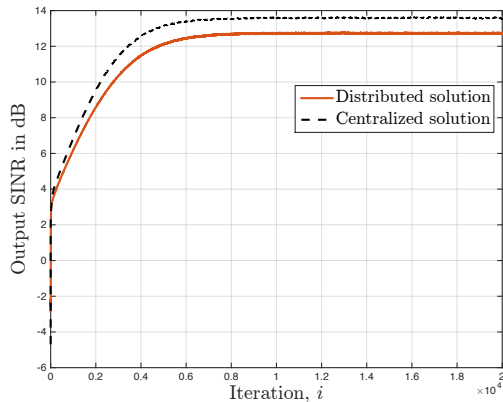
- For sufficiently small step-sizes, the decentralized strategy attains the same MSD performance as the centralized one

Application in Beamforming

- The framework is applied to an approximate linearly-constrained minimum-variance (LCMV) beamforming problem [Nassif et al., 2020, 2022]



Uniform linear array of N antennas



Comparison of output SINR

Example 3 – MAML for Multi-Agent Systems

- Instead of directly modeling the relationship between tasks w_k^o and w_ℓ^o , in model-agnostic meta-learning one assumes that both one or several (stochastic) gradient step away from a common launch-model:

$$w_k^o \approx w^o - \mu \nabla Q(w^o; \mathbf{x}_k) \quad (11)$$

- One then optimizes:

$$w^o \triangleq \arg \min_w \frac{1}{K} \sum_{k=1}^K \mathbb{E} Q(w - \mu \nabla Q(w; \mathbf{x}_k^1); \mathbf{x}_k^2) \quad (12)$$

to determine a common launch model w^o , which adapts quickly to other tasks w_k^o via one or several (stochastic) gradient steps.

- See [Smith et al. 2017] for centralized MAML and [Fallah, Mokhtari, and Ozdaglar 2020] for federated implementation and analysis.

Diffusion-MAML: Decentralization

- If we denote:

$$\overline{Q}(w; \mathbf{x}_k^1, \mathbf{x}_k^2) \triangleq Q(w - \mu \nabla Q(w; \mathbf{x}_k^1); \mathbf{x}_k^2) \quad (13)$$

then the optimization problem

$$w^o \triangleq \arg \min_w \frac{1}{K} \sum_{k=1}^K \mathbb{E} \overline{Q}(w; \mathbf{x}_k^1, \mathbf{x}_k^2) \quad (14)$$

is a **single-task** problem over the common launch model w , and can be pursued via diffusion in a decentralized manner:

$$\phi_{k,i} = \mathbf{w}_{k,i-1} - \mu \nabla \overline{Q}(\mathbf{w}_{k,i-1}; \mathbf{x}_{k,i}^1, \mathbf{x}_{k,i}^2) \quad (15)$$

$$= \mathbf{w}_{k,i-1} - \mu \nabla Q(\mathbf{w}_{k,i-1} - \mu \nabla Q(\mathbf{w}_{k,i-1}; \mathbf{x}_{k,i}^1); \mathbf{x}_{k,i}^2) \quad (16)$$

$$\mathbf{w}_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \phi_{\ell,i} \quad (17)$$

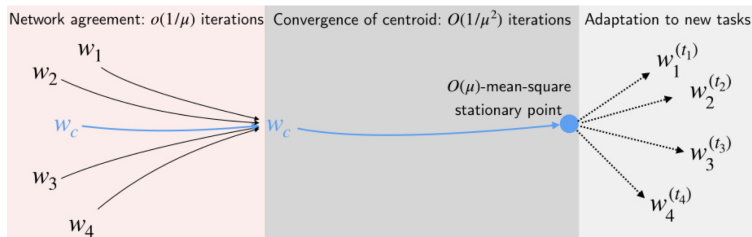
Diffusion-MAML: Performance

Convergence Guarantee [Kayaalp, Vlaski, and Sayed 2022]

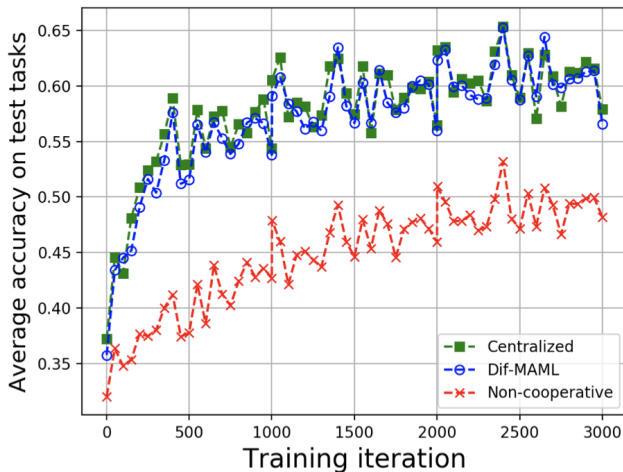
All agents agree on a common launch model in $o(1/\mu)$ iterations, and together find an approximately first-order stationary point of the aggregate objective:

$$\min_w \frac{1}{K} \sum_{k=1}^K \mathbb{E} Q(w - \mu \nabla Q(w; \mathbf{x}_k^1); \mathbf{x}_k^2) \quad (18)$$

in $O(1/\mu^2)$ iterations.

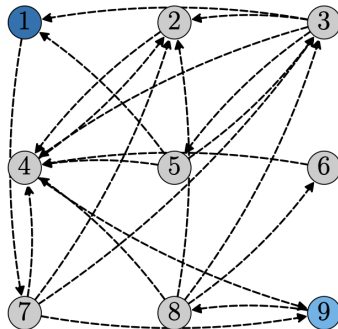
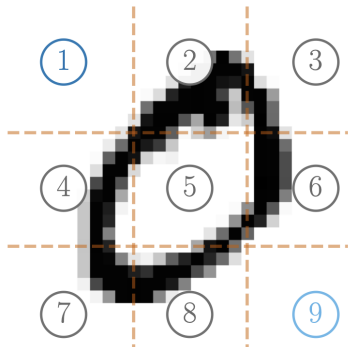


Diffusion-MAML: Few-Shot Image Recognition

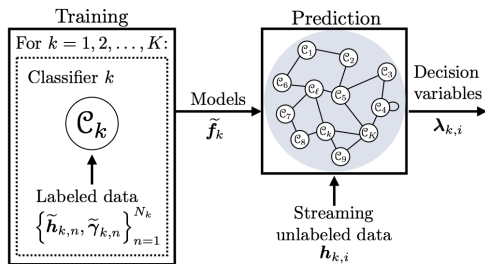


Example 4: Social Machine Learning

- So far we either explicitly induced a prior (multitask learning) or assumed an unknown relationship (model-agnostic meta-learning).
- What if there is no relationship between models, but rather their decisions?



Algorithm: Independent Training followed by Cooperative Inference



Social inference using consensus protocol on the decision [Bordignon et al, 2023]:

$$\delta_{k,i} = \lambda_{k,i-1} + \frac{\hat{L}_k^{\tilde{\mathbf{f}}_k}(\mathbf{h}_{k,i} | \gamma = +1)}{\hat{L}_k^{\tilde{\mathbf{f}}_k}(\mathbf{h}_{k,i} | \gamma = -1)}$$

$$\lambda_{k,i} = \sum_{\ell \in \mathcal{N}_k} a_{\ell k} \delta_{\ell,i}$$

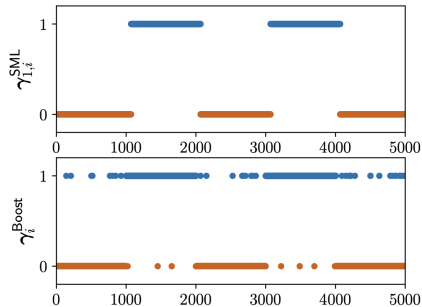
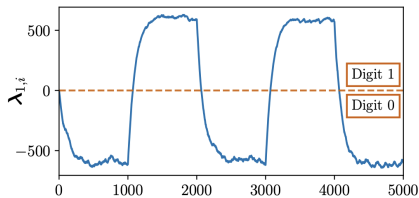
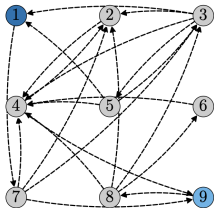
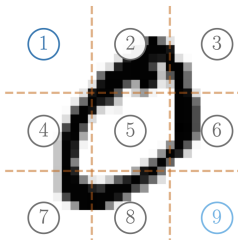
Consistency of Social Machine Learning [Bordignon et al, 2023]

The social machine learning strategy enables consistent learning with probability:

$$P_c \geq 1 - 2 \exp \left\{ -\frac{8N_{\max}}{\alpha^2 \beta^2} (\epsilon - \rho)^2 \right\} \quad (19)$$

where $\rho \geq \epsilon$, ϵ quantifies the complexity of the of the classification problem and ρ quantifies the Rademacher complexity of the classifiers.

Application: Partial MNIST



Take-Aways

- Ordinary averaging for distributed learning yields improvement only in sufficiently homogeneous environments.
- **Multitask learning:** Induce priors through regularization and constraints.
- **Model-agnostic meta-learning:** Absence of task priors.
- **Social machine learning:** Fully heterogeneous classifiers observing a common state.

References

- R. Nassif, S. Vlaski, C. Richard and A. H. Sayed, "Learning Over Multitask Graphs—Part I: Stability Analysis," in IEEE Open Journal of Signal Processing, vol. 1, pp. 28-45, 2020.
- R. Nassif, S. Vlaski and A. H. Sayed, "Adaptation and Learning Over Networks Under Subspace Constraints—Part I: Stability Analysis," in IEEE Transactions on Signal Processing, vol. 68, pp. 1346-1360, 2020.
- M. Kayaalp, S. Vlaski and A. H. Sayed, "Dif-MAML: Decentralized Multi-Agent Meta-Learning," in IEEE Open Journal of Signal Processing, vol. 3, pp. 71-93, 2022.
- V. Bordignon, S. Vlaski, V. Matta and A. H. Sayed, "Learning from Heterogeneous Data Based on Social Interactions over Graphs," in IEEE Transactions on Information Theory, to appear.