# Decentralized Online Nonparametric Learning

Alec Koppel[*], Santiago Paternain [†], Cédric Richard[§], Alejandro Ribeiro[†]
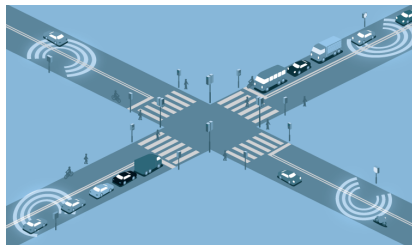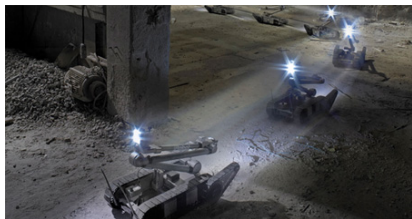[*]U.S. Army Research Laboratory, Adelphi, MD
[†]University of Pennsylvania
[§] Laboratoire Lagrange at the University of Nice Sophia-Antipolis.

Asilomar Conference, October 31, 2018, Pacific Grove, CA.

# Distributed Learning

- ▶ Network of agents $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ aims to make inferences from data
- ▶ Sensor Networks, multi-robot teams, internet of things
- ▶ For instance, distributed training of a classifier for some data set

- $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ is random pair $\Rightarrow$ training examples
- $\ell : \mathcal{W} \to \mathbb{R}$ convex loss ($\mathcal{W} \subset \mathbb{R}^p$), merit of statistical model
- Find parameters $\mathbf{w}^* \in \mathbb{R}^p$ that minimize expected risk $L(\mathbf{w})$

$$\mathbf{w}^* := \underset{\mathbf{w}}{\operatorname{argmin}}\, L(\mathbf{w}) := \underset{\mathbf{w}}{\operatorname{argmin}}\, \mathbb{E}_{\mathbf{x},\mathbf{y}}[\ell(\mathbf{w}^\top \mathbf{x}, \mathbf{y})]$$
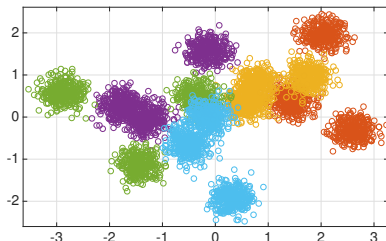
- Convex Optimization Problem for *linear statistical models*
  $\Rightarrow$ e.g., $y = \mathbf{w}^T \mathbf{x} \in \mathbb{R}$ or $y = \operatorname{sgn}(\mathbf{w}^T \mathbf{x}) \in \{-1, 1\}$
- Solve with favorite descent method $\Rightarrow$ Good Performance

# Easy to Implement over Networks

- ► Each agent $i$ has a local copy of the classifier $\mathbf{w}_i$ with $i = 1 \ldots |\mathcal{V}|$
  - ⇒ Observes some training examples ⇒ $(\mathbf{x}, \mathbf{y}) \in \mathcal{X}_i \times \mathcal{Y}_i$

$$\mathbf{w}^* := \underset{\mathbf{w} \in \mathbb{R}^{p|\mathcal{V}|}}{\operatorname{argmin}} \sum_{i=1}^{|\mathcal{V}|} \mathbb{E}_{\mathbf{x}_i, y_i} \left[ \ell(\mathbf{w}_i^\top \mathbf{x}_i, \mathbf{y}_i) \right]$$

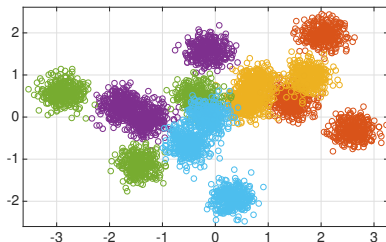$$s.t. \quad \mathbf{w}_i = \mathbf{w}_j \quad \text{for all} \quad j \in \mathcal{N}_i$$

- ► Convex Optimization Problem for *linear statistical models*
- ► Solve with saddle point algorithms or penalty methods
  - ⇒ Can be implemented in a distributed fashion

► The statistical model of complex data sets is nonlinear



► Neural Networks or Kernel Methods in centralized solution
► In this talk we focus on Distributed Kernel Methods
  ⇒ Contribution: each agent learns distinct kernel function
  ⇒ new penalty function that incentivizes coordination

▶ The statistical model of complex data sets is nonlinear

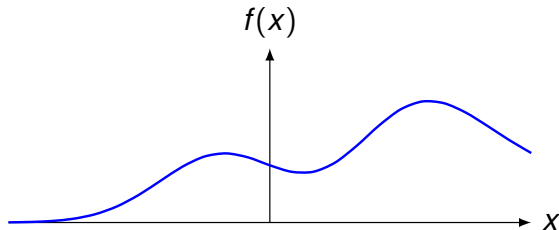

▶ Neural Networks or Kernel Methods in centralized solution
▶ In this talk we focus on Distributed Kernel Methods
    ⇒ Contribution: each agent learns distinct kernel function
    ⇒ new penalty function that incentivizes coordination

- Online consensus opt. for dist. learning (Tsitsiklis, Nedic, etc.)
  - ⇒ restrict statistical models to be linear (parameter vectors)

- Decentralized training of CNNs ⇒ non-convex consensus probs.
  - ⇒ Hong, Aldo, many others in past couple years
- Non-convexity precludes stable online model adaptation
  - ⇒ but good for stochastic algs. for large batch CNN training

- Focus on networked systems with nonlinear function approx.
  - ⇒ motivated by distributed intelligence w/ env. interaction
- Some prior works on distributed online kernel methods
  - ⇒ complexity reduction via fixing kernel matrix size, may diverge
- **Ours: globally convergent, sparse param. nonlinear funcs.**
  - ⇒ **in decentralized online setting**

- Equip $\mathcal{H}$ with a unique *kernel function*, $\kappa : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, such that:

$$(i) \; \langle f, \kappa(\mathbf{x}, \cdot) \rangle_{\mathcal{H}} = f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathcal{X} \, ,$$

$$(ii) \; \mathcal{H} = \overline{\text{span}\{\kappa(\mathbf{x}, \cdot)\}} \quad \text{for all } \mathbf{x} \in \mathcal{X} \, .$$



$f(x)$

$x$

- Property (i) $\Rightarrow$ Will allow us to compute derivatives
- Kernel examples:

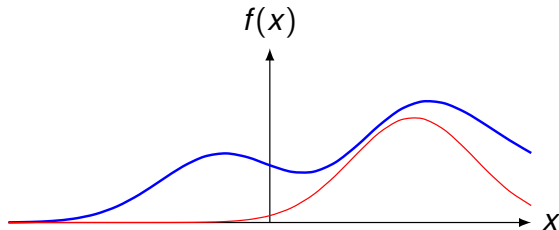$\Rightarrow$ Gaussian/RBF $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{ -\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2c^2} \right\}$

$\Rightarrow$ polynomial $\kappa(\mathbf{x}, \mathbf{x}') = \left( \mathbf{x}^T \mathbf{x}' + b \right)^c$

▶ Equip $\mathcal{H}$ with a unique *kernel function*, $\kappa : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, such that:

$$(i)\ \langle f, \kappa(\mathbf{x}, \cdot) \rangle_{\mathcal{H}} = f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathcal{X} ,$$

$$(ii)\ \mathcal{H} = \overline{\text{span}\{\kappa(\mathbf{x}, \cdot)\}} \quad \text{for all } \mathbf{x} \in \mathcal{X} .$$

$f(x)$



$x$

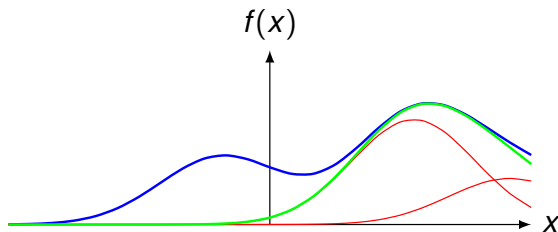▶ Property (i) $\Rightarrow$ Will allow us to compute derivatives

▶ Kernel examples:

$\Rightarrow$ Gaussian/RBF $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{ -\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2c^2} \right\}$

$\Rightarrow$ polynomial $\kappa(\mathbf{x}, \mathbf{x}') = \left( \mathbf{x}^T \mathbf{x}' + b \right)^c$

# Large-Scale Function Estimation

▶ Equip $\mathcal{H}$ with a unique *kernel function*, $\kappa : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, such that:

$$(i) \ \langle f, \kappa(\mathbf{x}, \cdot) \rangle_{\mathcal{H}} = f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathcal{X} ,$$

$$(ii) \ \mathcal{H} = \overline{\text{span}\{\kappa(\mathbf{x}, \cdot)\}} \quad \text{for all } \mathbf{x} \in \mathcal{X} .$$



▶ Property (i) $\Rightarrow$ Will allow us to compute derivatives

▶ Kernel examples:

$\Rightarrow$ Gaussian/RBF $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{ -\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2c^2} \right\}$

$\Rightarrow$ polynomial $\kappa(\mathbf{x}, \mathbf{x}') = \left(\mathbf{x}^T \mathbf{x}' + b\right)^c$

# Large-Scale Function Estimation

▶ Equip $\mathcal{H}$ with a unique *kernel function*, $\kappa : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, such that:

$$(i) \; \langle f, \kappa(\mathbf{x}, \cdot) \rangle_{\mathcal{H}} = f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathcal{X} ,$$

$$(ii) \; \mathcal{H} = \overline{\text{span}\{\kappa(\mathbf{x}, \cdot)\}} \quad \text{for all } \mathbf{x} \in \mathcal{X} .$$



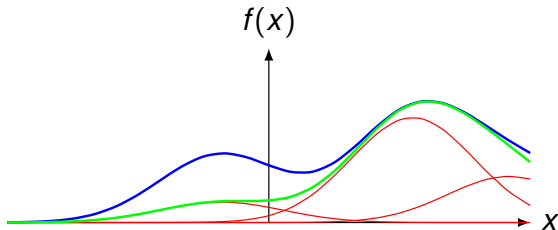▶ Property (i) $\Rightarrow$ Will allow us to compute derivatives
▶ Kernel examples:

$\Rightarrow$ Gaussian/RBF $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{ -\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2} \right\}$

$\Rightarrow$ polynomial $\kappa(\mathbf{x}, \mathbf{x}') = \left(\mathbf{x}^T\mathbf{x}' + b\right)^c$

► Equip $\mathcal{H}$ with a unique *kernel function*, $\kappa : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, such that:

$$(i) \ \langle f, \kappa(\mathbf{x}, \cdot) \rangle_{\mathcal{H}} = f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathcal{X} \ ,$$

$$(ii) \ \mathcal{H} = \overline{\text{span}\{\kappa(\mathbf{x}, \cdot)\}} \quad \text{for all } \mathbf{x} \in \mathcal{X} \ .$$



► Property (i) $\Rightarrow$ Will allow us to compute derivatives
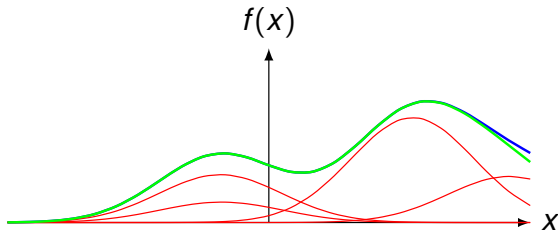► Kernel examples:

$\Rightarrow$ Gaussian/RBF $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{ -\dfrac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2c^2} \right\}$

$\Rightarrow$ polynomial $\kappa(\mathbf{x}, \mathbf{x}') = \left(\mathbf{x}^T \mathbf{x}' + b\right)^c$

# Large-Scale Function Estimation

▶ Equip $\mathcal{H}$ with a unique *kernel function*, $\kappa : \mathcal{X} \times \mathcal{X} \to \mathbb{R}$, such that:

$$(i) \; \langle f, \kappa(\mathbf{x}, \cdot) \rangle_{\mathcal{H}} = f(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \mathcal{X} \;,$$

$$(ii) \; \mathcal{H} = \overline{\text{span}\{\kappa(\mathbf{x}, \cdot)\}} \quad \text{for all } \mathbf{x} \in \mathcal{X} \;.$$



▶ Property (i) $\Rightarrow$ Will allow us to compute derivatives
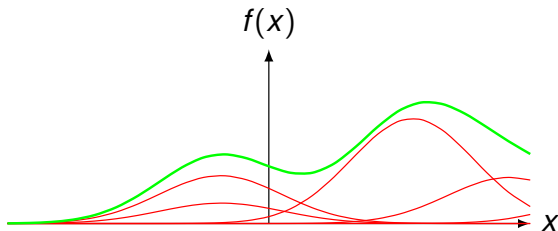▶ Kernel examples:
  $\Rightarrow$ Gaussian/RBF $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{ -\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2c^2} \right\}$
  $\Rightarrow$ polynomial $\kappa(\mathbf{x}, \mathbf{x}') = \left(\mathbf{x}^T \mathbf{x}' + b\right)^c$

# Function Representation

- ▶ Consider empirical risk minimization case: sample size $N < \infty$
- ▶ Representer Theorem:

$$f^* = \underset{f}{\arg\min} \frac{1}{N} \sum_{n=1}^{N} \ell(f(\mathbf{x}_n), y_n) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 = \sum_{m=1}^{N} w_m^* \, \kappa(\mathbf{x}_m, \mathbf{x}) \, .$$

- ▶ Representer Thm. into ERM $\Rightarrow$ opt. over $\mathcal{H}$ reduces to $\mathbf{w} \in \mathbb{R}^N$

$$f^* = \underset{\mathbf{w} \in \mathbb{R}^N}{\arg\min} \frac{1}{N} \sum_{n=1}^{N} \ell(\sum_{m=1}^{N} w_m \kappa(\mathbf{x}_m, \mathbf{x}_n), y_n) + \frac{\lambda}{2} \sum_{n,m=1}^{N} w_n w_m \kappa(\mathbf{x}_m, \mathbf{x}_n)$$

- ▶ Reduces to solve a convex optimization problem of dimension $N$.
- ▶ As $N \to \infty$ storage and computation issues are present
  - $\Rightarrow$ This is known as the Curse of Kernelization

# Distributed Function Estimation

- Each agents has a local copy $f_i \in \mathcal{H}$ with $i = 1 \ldots |\mathcal{V}|$
- Define the stacked function $f = [f_1, f_2, \ldots f_{|\mathcal{V}|}]^\top \in \mathcal{H}^{|\mathcal{V}|}$ and solve

$$p^* := \min_{f \in \mathcal{H}^{|\mathcal{V}|}} \sum_{i=1}^{|\mathcal{V}|} \mathbb{E}_{\mathbf{x}_i, y_i} \left[ \ell(f_i(\mathbf{x}_i), \mathbf{y}_i) \right] + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

$$s.t. \quad f_i = f_j \quad \text{for all} \quad i \in \mathcal{V} \quad \text{and} \quad j \in \mathcal{N}_i$$

- We solve it approximately using a penalty method

$$f_c^* = \operatorname*{argmin}_{f \in \mathcal{H}^{|\mathcal{V}|}} \psi_c(f) = \operatorname*{argmin}_{f \in \mathcal{H}^{|\mathcal{V}|}} \sum_{i \in \mathcal{V}} \mathbb{E}_{\mathbf{x}_i, \mathbf{y}_i} \left[ \ell_i(f_i(\mathbf{x}_i), y_i) \right] + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

$$+ \frac{c}{2} \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \mathbb{E}_{\mathbf{x}_i} \left[ \left( f_i(\mathbf{x}_i) - f_j(\mathbf{x}_i) \right)^2 \right]$$

▶ How far from consensus is the approximate solution?

## Proposition

*Let $f_c^* = \operatorname{argmin}_{f \in \mathcal{H}^{|\mathcal{V}|}} \psi_c(f)$ and let $p^*$ be the optimal cost of the distributed learning problem. Then for all penalties $c > 0$ we have that*

$$\frac{1}{2} \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{N}_i} \mathbb{E}_{\mathbf{x}_i} \left\{ \left[ f_{c,i}^*(\mathbf{x}_i) - f_{c,j}^*(\mathbf{x}_i) \right]^2 \right\} \leq \frac{p^*}{c}$$

▶ Expected disagreement arbitrarily small by increasing $c$

► Let $L(f)$ be the loss functional

$$L(f) = \sum_{i \in \mathcal{V}} \mathbb{E}_{\mathbf{x}_i, y_i}[\ell(f_i(\mathbf{x}_i), y_i)]$$

► Compute stochastic functional gradient of Ł$(f)$

$$\nabla_{f_i} \ell(f_i(\mathbf{x}_{i,t}), y_{i,t})(\cdot) = \frac{\partial \ell(f_i(\mathbf{x}_{i,t}), y_{i,t})}{\partial f_i(\mathbf{x}_{i,t})} \frac{\partial f_i(\mathbf{x}_{i,t})}{\partial f_i}(\cdot)$$

► Use reproducing property of kernel (i), differentiate both sides:

$$\frac{\partial f_i(\mathbf{x}_{i,t})}{\partial f_i}(\cdot) = \frac{\partial \langle f_i, \kappa(\mathbf{x}_{i,t}, \cdot) \rangle_{\mathcal{H}}}{\partial f_i} = \kappa(\mathbf{x}_{i,t}, \cdot)$$

- FDSGD applied to $\psi_c(f)$, given independent example $(\mathbf{x}_{i,t}, \mathbf{y}_{i,t})$:

$$f_{i,t+1} = f_{i,t} - \eta_t \hat{\nabla}_{f_i} \psi_c(f_{i,t}(\mathbf{x}_{i,t}), y_{i,t}) = (1 - \eta_t \lambda) f_{i,t} - \eta_t \omega_{i,t+1} \kappa(\mathbf{x}_{i,t}, \cdot)$$

$$\omega_{i,t+1} = \left( \ell'(f_i(\mathbf{x}_{i,t}), y_{i,t}) + c \sum_{j \in \mathcal{N}_i} \left( f_{i,t}(\mathbf{x}_{i,t}) - f_{j,t}(\mathbf{x}_{i,t}) \right) \right)$$

- Use the kernel expansion of $f_{i,t}$ to write

$$f_{i,t+1}(\mathbf{x}) = (1 - \eta_t \lambda) \sum_{n=1}^{t-1} w_{i,n} \kappa(\mathbf{x}_{i,n}, \mathbf{x}) - \eta_t \omega_{i,t+1} \kappa(\mathbf{x}_{i,t}, .)$$

- FDSGD: parametric updates on weights and dictionary

$$\mathbf{X}_{i,t+1} = [\mathbf{X}_{i,t}, \ \mathbf{x}_{i,t}], \quad \mathbf{w}_{i,t+1} = [(1 - \eta_t \lambda) \mathbf{w}_{i,t}, \ -\eta_t \omega_{i,t+1}],$$

- Note that model order $M_t = t - 1$ grows by one at each step

# Functional Distributed SGD

- FDSGD applied to $\psi_c(f)$, given independent example $(\mathbf{x}_{i,t}, \mathbf{y}_{i,t})$:

$$f_{i,t+1} = f_{i,t} - \eta_t \hat{\nabla}_{f_i} \psi_c(f_{i,t}(\mathbf{x}_{i,t}), y_{i,t}) = (1 - \eta_t \lambda) f_{i,t} - \eta_t \omega_{i,t+1} \kappa(\mathbf{x}_{i,t}, \cdot)$$

$$\omega_{i,t+1} = \left( \ell'(f_i(\mathbf{x}_{i,t}), y_{i,t}) + c \sum_{j \in \mathcal{N}_i} (f_{i,t}(\mathbf{x}_{i,t}) - f_{j,t}(\mathbf{x}_{i,t})) \right)$$

- Use the kernel expansion of $f_{i,t}$ to write

$$f_{i,t+1}(\mathbf{x}) = (1 - \eta_t \lambda) \sum_{n=1}^{t-1} w_{i,n} \kappa(\mathbf{x}_{i,n}, \mathbf{x}) - \eta_t \omega_{i,t+1} \kappa(\mathbf{x}_{i,t}, .)$$

**Consensus-term**

- FDSGD: parametric updates on weights and dictionary

$$\mathbf{X}_{i,t+1} = [\mathbf{X}_{i,t}, \ \mathbf{x}_{i,t}], \quad \mathbf{w}_{i,t+1} = [(1 - \eta_t \lambda) \mathbf{w}_{i,t}, \ -\eta_t \omega_{i,t+1}],$$

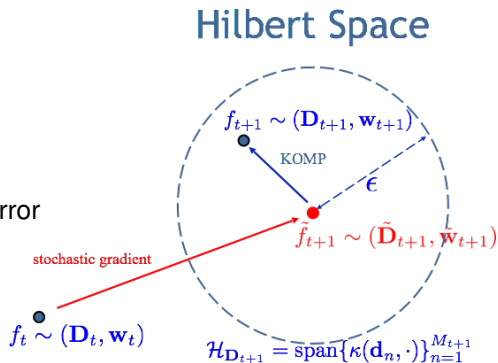- Note that model order $M_t = t - 1$ grows by one at each step

**Theorem**

Let $f_c^* := \mathrm{argmin}_{f \in \mathcal{H}} \psi_c(f)$, *under diminishing step-size rules*
$\sum_{t=1}^{\infty} \eta_t = \infty$ , $\sum_{t=1}^{\infty} \eta_t^2 < \infty$ , *with* $\eta_0 < 1/\lambda$,

$$\lim_{t \to \infty} \|f_t - f_c^*\|_{\mathcal{H}}^2 = 0 \qquad \textit{a.s.}$$

- ► Each agent learns $f^*_{c,i}$ in such a way that $M_{i,t} << \infty$ for each $f_{i,t}$

- ► Accomplished by fixing a error nbhd. around FDSGD iterates
  ⇒ Remove maximal no. kernel dict. elements while inside nbhd.

- ► We propose using KOMP ⇒ kernel orthogonal matching pursuit
  ⇒ a greedy compressive technique (Vincent & Bengio, 2002)

- Fix approximation error $\epsilon$
- $\tilde{f}_{t+1} = f_t - \eta \hat{\nabla}_f \psi_c(f_t)$
- Remove kernel element smallest error
- Project $\tilde{f}_{t+1}$ onto resulting RKHS
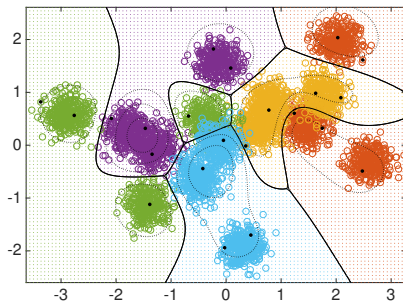- Repeat until error is larger than $\varepsilon$



**Hilbert Space**

$f_{t+1} \sim (\mathbf{D}_{t+1}, \mathbf{w}_{t+1})$

KOMP

$\boldsymbol{\epsilon}$

$\tilde{f}_{t+1} \sim (\tilde{\mathbf{D}}_{t+1}, \tilde{\mathbf{w}}_{t+1})$

stochastic gradient

$f_t \sim (\mathbf{D}_t, \mathbf{w}_t)$

$\mathcal{H}_{\mathbf{D}_{t+1}} = \overline{\text{span}}\{\kappa(\mathbf{d}_n, \cdot)\}_{n=1}^{M_{t+1}}$

### Theorem

Let $f_c^* := \text{argmin}_{f \in \mathcal{H}} \psi_c(f)$. Given regularizer $\lambda > 0$, constant algorithm step-size $\eta$ chosen such that $\eta < 1/\lambda$ and compression error $\epsilon = K\eta^{3/2} = \mathcal{O}(\eta^{3/2})$, where $K$ is a positive scalar,

$$\liminf_{t \to \infty} \|f_t - f_c^*\|_{\mathcal{H}} \leq \frac{\sqrt{\eta}}{\lambda}\left(K|\mathcal{V}| + \sqrt{K^2|\mathcal{V}|^2 + \lambda\sigma^2}\right) = \mathcal{O}(\sqrt{\eta}) \qquad a.s.$$

The model order of the function, $M_t$ is finite for all $t$

- ▶ Bias induced by sparsification asymptotically doesn't hurt too bad
- ▶ Constant step-size, approx. budget $\Rightarrow$ model order always finite

- 3 Gaussians per mixture, $C = 5$ classes total for this experiment
  - $\Rightarrow$ 15 total Gaussians generate data
- $\ell(\mathbf{f}(\mathbf{x}), y) = \max(0, 1 + f_r(\mathbf{x}) - f_y(\mathbf{x}))$, $r = \operatorname{argmax}_{c' \neq y} f_{c'}(\mathbf{x})$
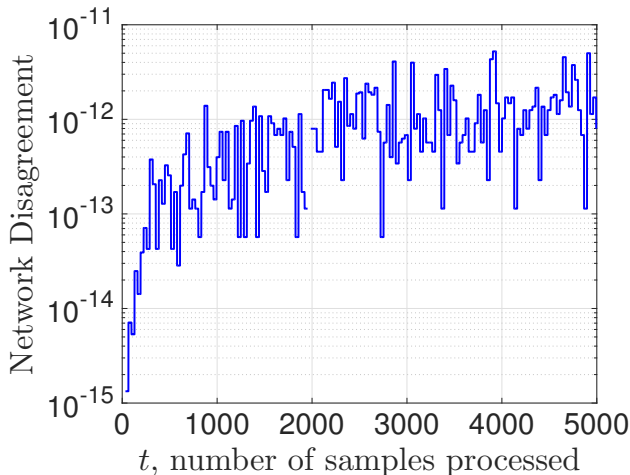


- Grid colors $\Rightarrow$ decision
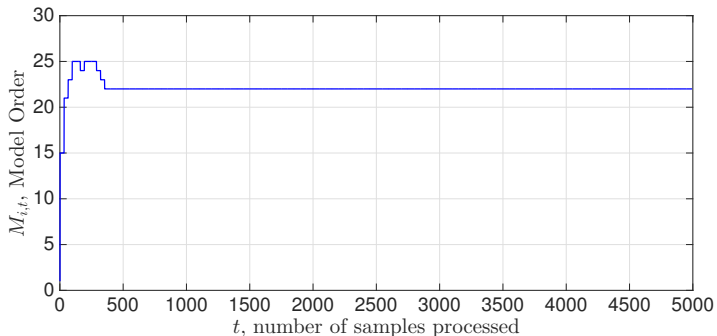- Black dots $\Rightarrow$ kernels
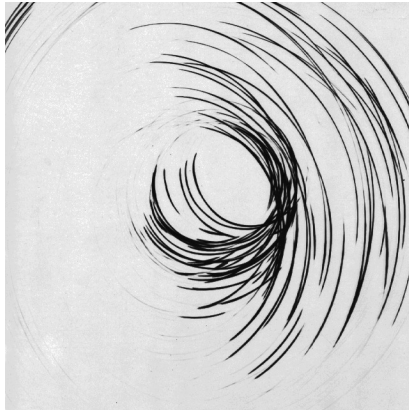- $\sim 95.7\%$ accuracy

► Convergence to optimal solution

- Consensus error remains small
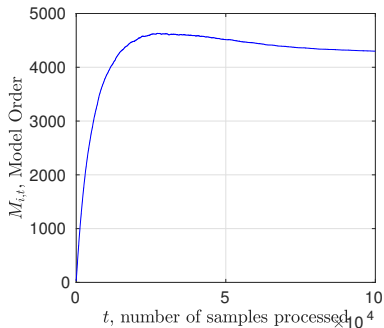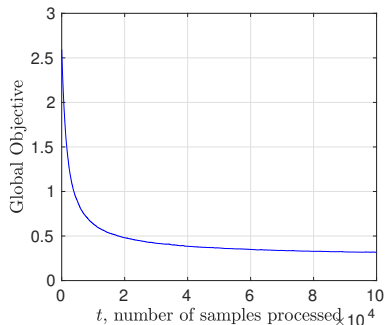
- ▶ Bounded model order

- Texture classification on Broadatz dataset via SVM

▶ We observe convergence and finite Model Order



▶ Accuracy of 93.5% comparable to centralized case (95.6%)

- ► We need to go beyond linear statistical models to do Learning
- ► Kernels and Neural Networks are the common tools to do so
  - ⇒ Kernel methods yield convex optimization problems
- ► We presented a distributed Learning algorithm (FDSGD)
  - ⇒ Converges to a neighborhood of the optimal function
  - ⇒ while ensuring a bound on the model order for all times
- ► Future directions: apply to, e.g., SLAM, exploration, navigation
  - ⇒ reduce communication overhead
  - ⇒ each agent learns kernel function w/ distinct bandwidth