

Exact Nonparametric Decentralized Online Optimization

Hrusikesh Pradhan[†], Amrit Singh Bedi[†], Alec Koppel[§], Ketan Rajawat[†]

[†] Dept. of EE, India Institute of Technology Kanpur,

[§] CISD, U.S. Army Research Laboratory

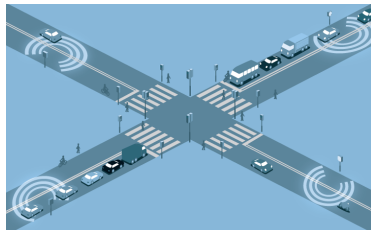
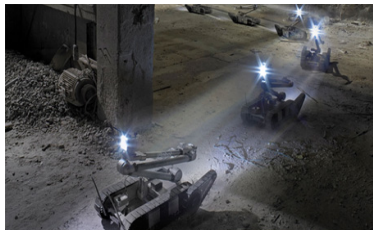
IEEE Global Conference on Signal and Information Processing
Anaheim, California, Nov. 28, 2018

A diagram consisting of three light green circles with a slight gradient and a drop shadow. The circles are arranged in a triangle: one on the left, one on the right, and one centered below the other two. Each circle contains text in a black, sans-serif font.

Consensus

Linear
Statistical
Models

Offline



Consensus

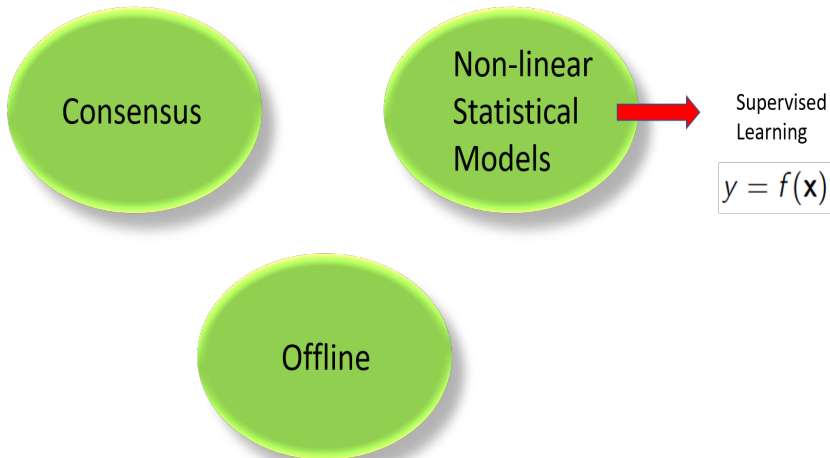
Linear
Statistical
Models

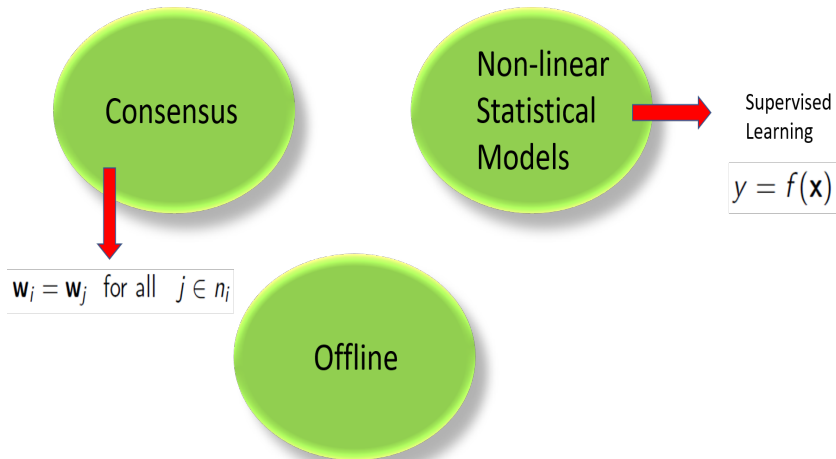


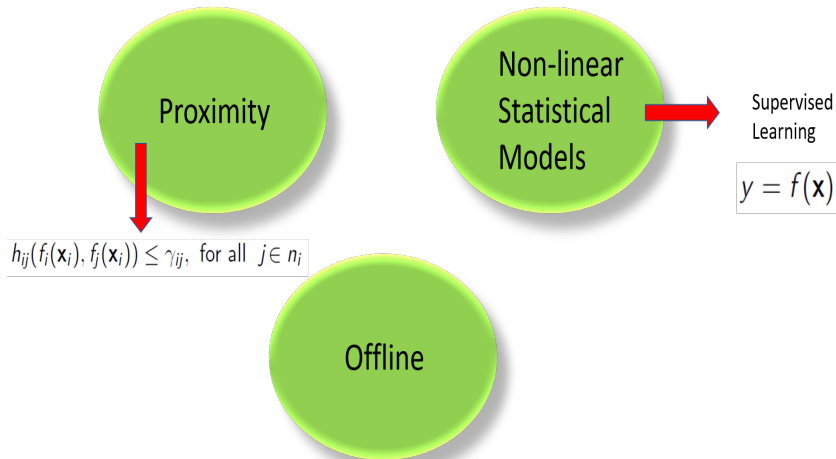
Supervised
Learning

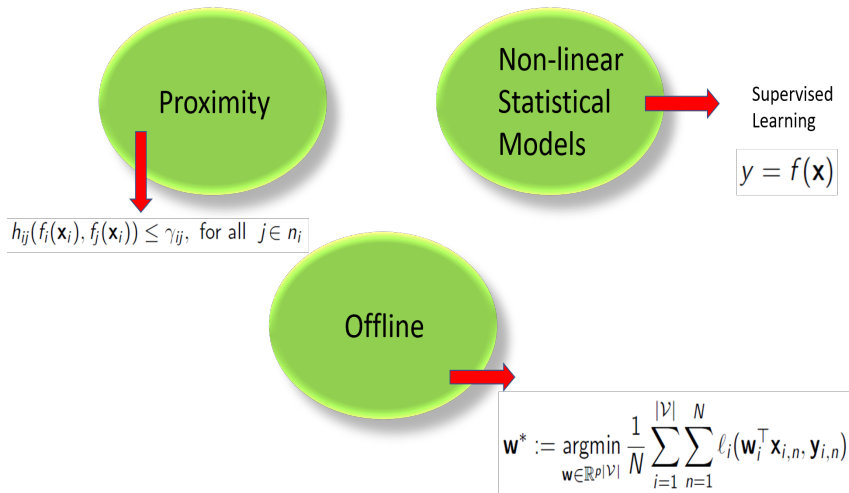
$$y_i = \mathbf{w}_i^T \mathbf{x}_i$$

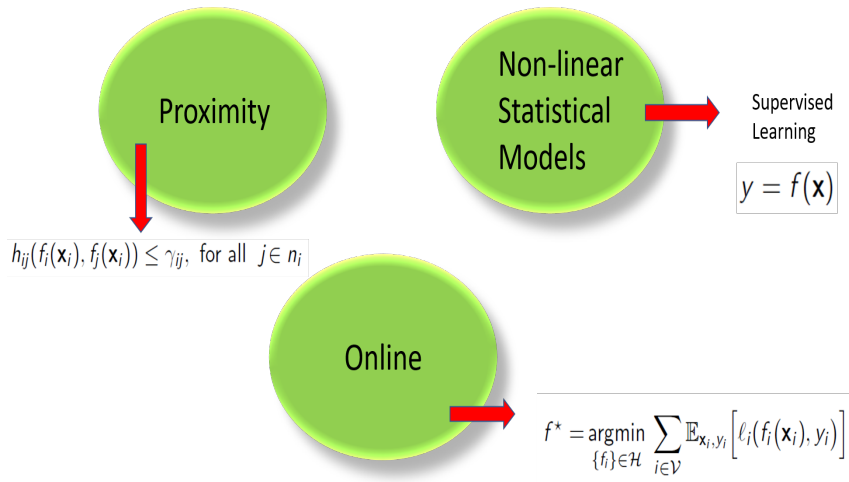
Offline












Distributed Online Learning	Linear	Nonlinear
Homogeneous		Koppel ,
Heterogeneous	Sayed, Chen, Nassif, Lee, Sadler	Our Work



Haoran Sun, Mingyi Hong,
Nikos D. Sidiropoulos

-
- A. Nedic and A. Ozdaglar, Distributed subgradient methods for multiagent optimization, IEEE Transactions on Automatic Control
- Mokhtari, Aryan, and Alejandro Ribeiro. "DSA: Decentralized double stochastic averaging gradient algorithm." JMLR
- Sirb, B., Ye, X. (2018). Decentralized consensus algorithm with delayed and stochastic gradients. SIAM Journal on Optimization.
- Chen, T., Ling, Q., Giannakis, G. B. (2017). An online convex optimization approach to proactive network resource allocation. IEEE Transactions on Signal Processing, 65(24), 6350-6364.
- Sun Haoran, and Mingyi Hong. "Distributed non-convex first-order optimization and information processing: Lower complexity bounds and rate optimal algorithms." arXiv preprint arXiv:1804.02729 (2018)
- Sun Haoran, Xiangyi Chen, Qingjiang Shi, Mingyi Hong, Xiao Fu, and Nikos D. Sidiropoulos. "Learning to optimize: Training deep neural networks for wireless resource management." In SPAWC, 2017 IEEE 18th International Workshop on, pp. 1-6. IEEE, 2017.
- Chen, J., Richard, C. and Sayed, A. H. (2014). Multitask diffusion adaptation over networks. IEEE Transactions on Signal Processing
- Nassif, R., Richard, C., Ferrari, A. and Sayed, A. H. Distributed learning over multitask networks with linearly related tasks.
- Lee, S., Zavlanos, M. M. (2017). On the sublinear regret of distributed primal-dual algorithms for online constrained optimization.
- A. Koppel, S. Paternain, C. Richard, and A. Ribeiro, "Decentralized Online Learning with Kernels", in IEEE Trans. Signal Process
- A. Koppel, B. Sadler, and A. Ribeiro, "Proximity without Consensus in Online Multi-Agent Optimization," in IEEE Trans. Signal

- ▶ **Approach:**
 - ⇒ Hypothesized non-linear function in kernel Hilbert space
 - ⇒ Form stochastic lagrangian
 - ⇒ Apply stochastic primal dual method
 - ⇒ Take subspace projection (to handle memory growth)
- ▶ Sublinear convergence
 - ⇒ $\mathcal{O}(\sqrt{T})$ for primal optimality
 - ⇒ $\mathcal{O}(T^{3/4})$ for constraint violation
- ▶ Generalizes existing rate results for primal-dual method
 - ⇒ to case of non-linear statistical models
- ▶ **Application:** Decentralized correlated random field estimation
 - ⇒ online non-linear observation model.

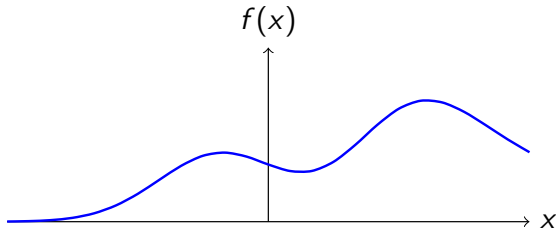
- ▶ Symmetric, connected and directed network of agents $\mathcal{G} = (\mathcal{V}, \mathcal{E})$
- ▶ Learning nonlinear statistical models is equivalent to finding
 $\Rightarrow f : \mathcal{X} \rightarrow \mathcal{Y}$, such that $y = f(\mathbf{x})$
- ▶ Loss $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ penalize deviations between $f(\mathbf{x})$, \mathbf{y}
- ▶ Encoded by a convex proximity function $h_{ij}(f_i(\mathbf{x}_i), f_j(\mathbf{x}_i))$
 \Rightarrow incentivizes nearby agents to make similar decisions
- ▶ Yields the constrained functional stochastic program:

$$f^* = \underset{\{f_i\} \in \mathcal{H}}{\operatorname{argmin}} \sum_{i \in \mathcal{V}} \left(\mathbb{E}_{\mathbf{x}_i, y_i} \left[\ell_i(f_i(\mathbf{x}_i), y_i) \right] + \frac{\lambda}{2} \|f_i\|_{\mathcal{H}}^2 \right)$$
$$\text{s.t. } \mathbb{E}_{\mathbf{x}_i} \left[h_{ij}(f_i(\mathbf{x}_i), f_j(\mathbf{x}_i)) \right] \leq \gamma_{ij}, \text{ for all } j \in n_i. \quad (1)$$

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

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

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



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

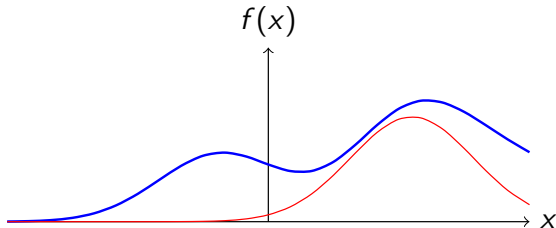
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

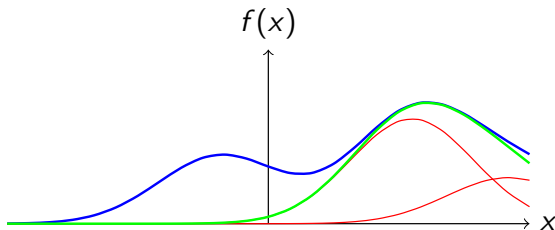
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

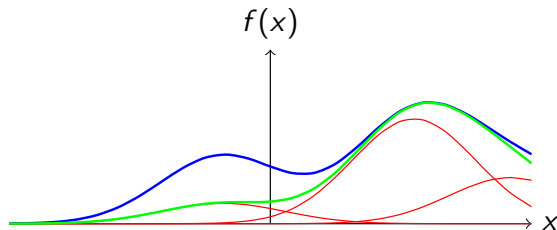
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

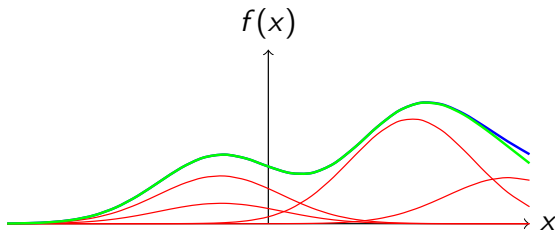
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

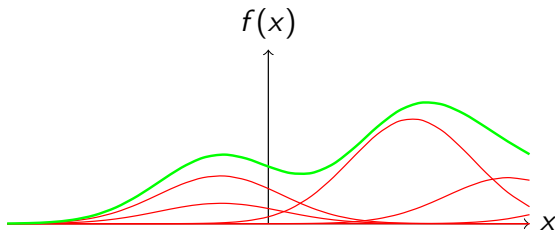
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

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

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

- ▶ Stochastic augmented Lagrangian function of (1) at time t

$$\hat{\mathcal{L}}_t(f, \boldsymbol{\mu}) := \sum_{i \in \mathcal{V}} \left[\ell_i(f_i(\mathbf{x}_{i,t}), y_{i,t}) + \frac{\lambda}{2} \|f_i\|_{\mathcal{H}}^2 + \sum_{j \in n_i} \left\{ \left[\mu_{ij} (h_{ij}(f_i(\mathbf{x}_{i,t}), f_j(\mathbf{x}_{i,t})) - \gamma_{ij}) \right] - \frac{\delta \eta}{2} \mu_{ij}^2 \right\} \right] \quad (2)$$

where $\boldsymbol{\mu}$ is a lagrange multiplier, with μ_{ij} defined for each $(i, j) \in \mathcal{E}$.

Corollary

Consider the sample average approximation of (1), and its associated Lagrangian relaxation. The each i th component of the solution to the resulting saddle-point problem can be expressed as

$$f_i^* = \sum_{t=1}^T w_{i,t} \kappa(\mathbf{x}_{i,t}, \cdot) \quad (3)$$

where $w_{i,t}$ are real-valued coefficients.

- Functional stochastic gradient of local loss in (2):

$$\ell'_i(f_i(\mathbf{x}_{i,t}), y_{i,t}) := \nabla_{f_i} \ell_i(f_i(\mathbf{x}_{i,t}), y_{i,t})(\cdot) = \frac{\partial \ell_i(f_i(\mathbf{x}_{i,t}), y_{i,t})}{\partial f_i(\mathbf{x}_{i,t})} \frac{\partial f_i(\mathbf{x}_{i,t})}{\partial f_i}(\cdot)$$

- Using the reproducing property of the kernel we obtain

$$\frac{\partial f_i(\mathbf{x}_{i,t})}{\partial f_i} = \frac{\partial \langle f_i, \kappa(\mathbf{x}_{i,t}, \cdot) \rangle_{\mathcal{H}}}{\partial f_i} = \kappa(\mathbf{x}_{i,t}, \cdot) \quad (4)$$

- Now the full gradient result can be written as

$$\begin{aligned} \nabla_{f_i} \hat{\mathcal{L}}_t(f_t, \mu_t) &= \ell'_i(f_i(\mathbf{x}_{i,t}), y_{i,t}) \kappa(\mathbf{x}_{i,t}, \cdot) + \lambda f_i \\ &+ \sum_{j \in n_i} \mu_{ij} h'_{ij}(f_i(\mathbf{x}_{i,t}), f_j(\mathbf{x}_{i,t})) \kappa(\mathbf{x}_{i,t}, \cdot) \end{aligned} \quad (5)$$

- ▶ **loop in parallel** for agent $i \in \mathcal{V}$
- ▶ Observe local training example realization $(\mathbf{x}_{i,t}, y_{i,t})$
- ▶ Send $\mathbf{x}_{i,t}$ to the neighboring nodes, $j \in n_i$ and receive $f_{j,t}(\mathbf{x}_{i,t})$
- ▶ Receive $\mathbf{x}_{j,t}$ from the neighbouring nodes, $j \in n_i$ and send $f_{i,t}(\mathbf{x}_{j,t})$
- ▶ Stochastic primal descent step on Lagrangian:

$$f_{i,t+1} = f_{i,t}(1 - \eta\lambda) - \eta \left[\ell'_i(f_{i,t}(\mathbf{x}_{i,t}), y_{i,t}) + \sum_{j \in n_i} \mu_{ij} h'_{ij}(f_{i,t}(\mathbf{x}_{i,t}), f_{j,t}(\mathbf{x}_{i,t})) \right] \kappa(\mathbf{x}_{i,t}, \cdot) \quad (6)$$

- ▶ Stochastic dual ascent step on Lagrangian:

$$\mu_{ij,t+1} = \left[\mu_{ij,t}(1 - \delta\eta^2) + \eta \left(h_{ij}(f_{i,t}(\mathbf{x}_{i,t}), f_{j,t}(\mathbf{x}_{i,t})) - \gamma_{ij} \right) \right]_+ \quad (7)$$

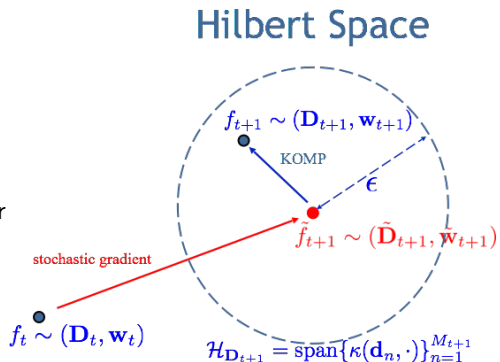
- ▶ Using $f_{i,t}(\mathbf{x}) = \sum_{n=1}^{t-1} w_{i,n} \kappa(\mathbf{x}_{i,n}, \mathbf{x}) = \mathbf{w}_{i,t}^T \boldsymbol{\kappa}_{\mathbf{x}_{i,t}}(\mathbf{x})$, V parallel parametric updates on both kernel dictionaries \mathbf{X}_i and \mathbf{w}_i are

$$\begin{aligned} \mathbf{X}_{i,t+1} &= [\mathbf{X}_{i,t}, \mathbf{x}_{i,t}], \\ [\mathbf{w}_{i,t+1}]_u &= \begin{cases} (1 - \eta\lambda)[\mathbf{w}_{i,t}]_u, & 0 \leq u \leq t-1 \\ -\eta \left(\ell'_i(f_{i,t}(\mathbf{x}_{i,t}), y_{i,t}) \right. \\ \left. + \sum_{j \in n_i} \mu_{ij} h'_{ij}(f_{i,t}(\mathbf{x}_{i,t}), f_{j,t}(\mathbf{x}_{i,t})) \right), & u = t \end{cases} \end{aligned} \quad (8)$$

- ▶ Data points $M_{i,t}$ grows by one each time (**curse of kernelization**).
- ▶ **Proj. Funct. Update:** Onto $\mathcal{H}_{\mathbf{D}_{i,t+1}} = \text{span}\{\kappa(\mathbf{d}_{i,n}, \cdot)\}_{n=1}^{M_{t+1}} \subset \mathcal{H}$

$$f_{i,t+1} := \mathcal{P}_{\mathcal{H}_{\mathbf{D}_{i,t+1}}} \left[f_{i,t} - \eta \nabla_{f_i} \hat{\mathcal{L}}_t(f_t, \mu_t) \right]. \quad (9)$$

- ▶ Fix approximation error ϵ
- ▶ $\tilde{f}_{t+1} = f_t - \eta \nabla_{f_t} \hat{\mathcal{L}}_t(f_t, \mu_t)$
- ▶ Remove kernel element smallest error
- ▶ Project \tilde{f}_{t+1} onto resulting RKHS
- ▶ Repeat until error is larger than ϵ



Theorem

Let $S(f_t) := \sum_{i \in \mathcal{V}} \mathbb{E}[\ell_i(f_{i,t}(\mathbf{x}_{i,t}), y_{i,t})] + \frac{\lambda}{2} \|f_{i,t}\|_{\mathcal{H}}^2$ as the objective with f^* defined in (1). Then with constant step-size $\eta = 1/\sqrt{T}$ and approximation budget $\epsilon_t = \epsilon = P\eta^2$, where $P > 0$ is termed as the parsimony constant.

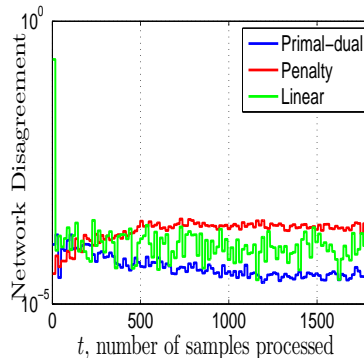
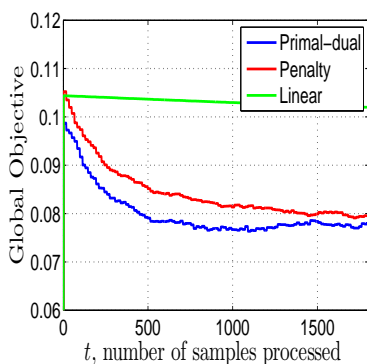
$$\sum_{t=1}^T \mathbb{E}[S(f_t) - S(f^*)] \leq \mathcal{O}(\sqrt{T}). \quad (10)$$

Furthermore, the time-aggregation of the expected constraint violation grows sub-linearly with iteration T as

$$\sum_{(i,j) \in \mathcal{E}} \mathbb{E} \left[\sum_{t=1}^T (h_{ij}(f_i(\mathbf{x}_{i,t}), f_j(\mathbf{x}_{i,t})) - \gamma_{ij}) \right]_+ \leq \mathcal{O}(T^{3/4}). \quad (11)$$

- ▶ A planar field is a random function of spatial components
- ▶ Random field is parameterized by the correlation matrix \mathbf{R}_x
 - ⇒ $[\mathbf{R}_x]_{ij}$ is assumed to have a structure of the form $e^{-\|i-j\|}$,
- ▶ Observation model: $y_{i,t} = h_i \|\mathbf{x}_{i,t}\|^2 + n_{i,t}$, $n_{i,t} \sim \mathcal{N}(0, \sigma^2)$
- ▶ 10 nodes distributed in a 10×10 meter square area.
- ▶ Instantaneous observation, $\mathbf{x}_t = \pi + \mathbf{C}^T \mathbf{v}_t$
 - ⇒ $\pi = \{1/V, 2/V, \dots, 1\}$
 - ⇒ \mathbf{C} : Cholesky factorization of the correlation matrix \mathbf{R}_x
 - ⇒ $\mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$

- Convergence of global objective and reduced network disagreement



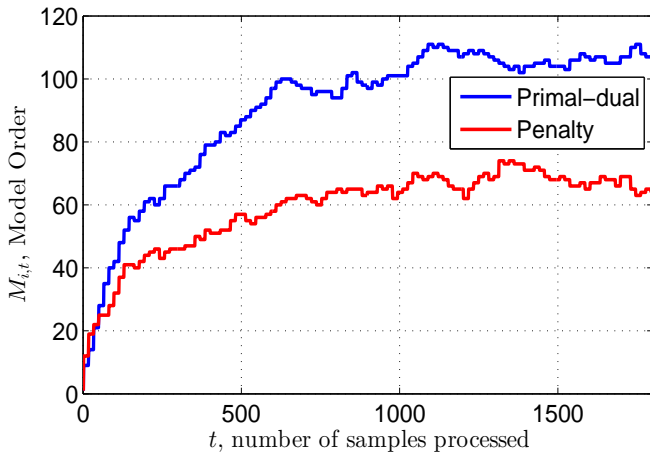


Figure: Model Order $M_{i,t}$ v.s. Samples

- ▶ Focused on **online** learning
 - ⇒ Decentralized **heterogeneous** networks
 - ⇒ **Non-linear** statistical models
- ▶ Proposed new variant of projected stochastic primal dual method
 - ⇒ **Convergence** to the optimum
 - ⇒ **Finite** growth of **model order**
 - ⇒ Observed good empirical performance
- ▶ **Future Work:**
 - ⇒ Asynchrony
 - ⇒ Reduce complexity of projections

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

$$\mathbf{w}^* := \underset{\mathbf{w}}{\operatorname{argmin}} L(\mathbf{w}) := \underset{\mathbf{w}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\mathbf{w}^T \mathbf{x}, \mathbf{y})]$$

- ▶ **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

- ▶ Symmetric, connected and directed network of agents $\mathcal{G} = (\mathcal{V}, \mathcal{E})$
- ▶ The nodes aims to make inferences from local data
- ▶ $|\mathcal{V}| = V$ nodes, $|\mathcal{E}| = M$ edges, and $n_i := \{j : (i, j) \in \mathcal{E}\}$
- ▶ Agent $i \in \mathcal{V}$ has a local copy of the classifier \mathbf{w}_i
 - ⇒ Observes some training examples ⇒ $(\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{X}_i \times \mathcal{Y}_i$

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

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

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

- ▶ Project (6) onto a lower dimensional subspace $\mathcal{H}_{\mathbf{D}} \subseteq \mathcal{H}$
- ▶ $\mathcal{H}_{\mathbf{D}}$ is represented by a dictionary $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_M] \in \mathbb{R}^{p \times M}$.
- ▶ $\mathcal{H}_{\mathbf{D}} = \{f : f(\cdot) = \sum_{n=1}^M w_n \kappa(\mathbf{d}_n, \cdot) = \mathbf{w}^T \boldsymbol{\kappa}_{\mathbf{D}}(\cdot)\}, \{\mathbf{d}_n\} \subset \{\mathbf{x}_u\}_{u \leq t}$.
- ▶ We denote the un-projected functional update as

$$\tilde{f}_{i,t+1} = f_{i,t} - \eta \nabla_{f_i} \hat{\mathcal{L}}_t(f_t, \mu_t). \quad (12)$$

where $\nabla_{f_i} \hat{\mathcal{L}}_t(f_t, \mu_t) :=$

$$\lambda f_{i,t} + \left[\ell'_i(f_{i,t}(\mathbf{x}_{i,t}), y_{i,t}) + \sum_{j \in n_i} \mu_{ij} h'_{ij}(f_{i,t}(\mathbf{x}_{i,t}), f_{j,t}(\mathbf{x}_{i,t})) \right] \kappa(\mathbf{x}_{i,t}, \cdot).$$

- ▶ $\tilde{f}_{i,t+1}$ in form of dictionary and coefficient vector:

$$\tilde{\mathbf{D}}_{i,t+1} = [\mathbf{D}_{i,t}, \mathbf{x}_{i,t}],$$

$$[\tilde{\mathbf{w}}_{i,t+1}]_{u=0} = \begin{cases} (1 - \eta\lambda)[\mathbf{w}_{i,t}]_u, & \text{for } 0 \leq u \leq t-1 \\ -\eta \left(\ell'_i(f_{i,t}(\mathbf{x}_{i,t}), y_{i,t}) + \sum_{j \in n_i} \mu_{ij} h'_{ij}(f_{i,t}(\mathbf{x}_{i,t}), f_{j,t}(\mathbf{x}_{i,t})) \right), & \text{for } u = t \end{cases}$$