

Online and Decentralized Statistical Learning

Alec Koppel
University of Pennsylvania, Philadelphia, PA

Phd Committee:
Alejandro Ribeiro (Advisor), Vijay Kumar (Chair),
Brian M. Sadler, Jonathan Fink

Phd Proposal
Philadelphia, PA, Dec. 1, 2016

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ Consider $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ as random pair \Rightarrow training examples
- ▶ Examples: classification $\mathbf{y} \in \{1, \dots, c\}$ or regression $\mathbf{y} \in \mathbb{R}$
 \Rightarrow Perform prediction \Rightarrow optimize statistical inference accuracy

$$\hat{\mathbf{y}}^*(\mathbf{x}) := \operatorname{argmin}_{\hat{\mathbf{y}}(\mathbf{x})} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\mathbb{I}\{\hat{\mathbf{y}}(\mathbf{x}) \neq \mathbf{y}\}]$$

- ▶ In general intractable \Rightarrow distribution $\mathbb{P}(\mathbf{x}, \mathbf{y})$ is unknown
 \Rightarrow but gives clear merit for choosing estimator $\hat{\mathbf{y}}$
 $\Rightarrow \hat{\mathbf{y}}^*(\mathbf{x})$ achieves optimal Bayes Risk

- ▶ Since optimizing statistical error rate directly is intractable
 - ⇒ Replace 0 – 1 loss ⇒ convex loss $\ell : \mathcal{W} \rightarrow \mathbb{R}$, $\mathcal{W} \subset \mathbb{R}^p$,
 - ⇒ Estimator $\hat{\mathbf{y}}(\mathbf{x}) = \hat{\mathbf{y}}(\mathbf{w}, \mathbf{x})$ depends on **model parameters \mathbf{w}**

$$\mathbf{w}^* := \underset{\mathbf{w}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\mathbf{w}; \mathbf{x}, \mathbf{y})] = \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}; \mathbf{x}_n, \mathbf{y}_n)$$

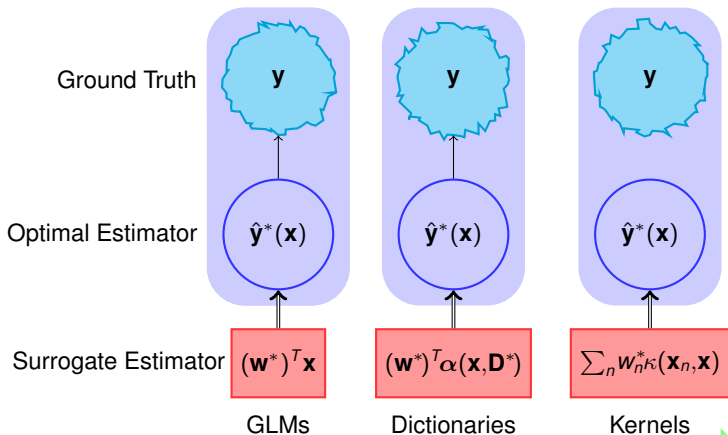
- ▶ Can solve optimization problem, define $L(\mathbf{w}) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\mathbf{w}; \mathbf{x}, \mathbf{y})]$
 - ⇒ but unclear how to choose $\hat{\mathbf{y}}(\mathbf{w}, \mathbf{x})$ such that $\hat{\mathbf{y}}(\mathbf{w}^*, \mathbf{x}) \approx \hat{\mathbf{y}}^*$

- ▶ Since optimizing statistical error rate directly is intractable
 - ⇒ Replace 0 – 1 loss ⇒ convex loss $\ell : \mathcal{W} \rightarrow \mathbb{R}$, $\mathcal{W} \subset \mathbb{R}^p$,
 - ⇒ Estimator $\hat{\mathbf{y}}(\mathbf{x}) = \hat{\mathbf{y}}(\mathbf{w}, \mathbf{x})$ depends on **model parameters \mathbf{w}**

$$\mathbf{w}^* := \underset{\mathbf{w}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\mathbf{w}; \mathbf{x}, \mathbf{y})] = \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}; \mathbf{x}_n, \mathbf{y}_n)$$

- ▶ Can solve optimization problem, define $L(\mathbf{w}) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\mathbf{w}; \mathbf{x}, \mathbf{y})]$
 - ⇒ but unclear how to choose $\hat{\mathbf{y}}(\mathbf{w}, \mathbf{x})$ such that $\hat{\mathbf{y}}(\mathbf{w}^*, \mathbf{x}) \approx \hat{\mathbf{y}}^*$
- ▶ We'll present approaches to indirectly optimizing accuracy
 - ⇒ Make use of surrogate losses for tractable model training
 - ⇒ Higher complexity $\hat{\mathbf{y}}(\mathbf{w}^*, \mathbf{x})$ ⇒ harder training, closer to $\hat{\mathbf{y}}^*$?
 - ⇒ Therefore, progressively increase complexity of $\hat{\mathbf{y}}(\mathbf{w}, \mathbf{x})$
- ▶ **Focus:** streaming data settings, multi-agent systems

Statistical Estimators



Increasing complexity of statistical model

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ Recall ERM prob. \Rightarrow loss ℓ , parameter vector $\mathbf{w} \in \mathbb{R}^p$

$$\mathbf{w}^* := \operatorname{argmin}_{\mathbf{w}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\mathbf{w}; \mathbf{x}, \mathbf{y})] = \operatorname{argmin}_{\mathbf{w}} \sum_{n=1}^N [\ell(\mathbf{w}; \mathbf{x}_n, \mathbf{y}_n)]$$

\Rightarrow Alternative measure of performance is **regret**

$$\mathbf{Reg}_N = \sum_{n=1}^N \ell(\mathbf{w}_n; \mathbf{x}_n, \mathbf{y}_n) - \sum_{n=1}^N \ell(\mathbf{w}^*; \mathbf{x}_n, \mathbf{y}_n)$$

- ▶ **Goal: asymptotic no-regret $\mathbf{Reg}_N/N \rightarrow 0$**
 - \Rightarrow benefit of this perspective is dropping i.i.d. assumption
 - $\Rightarrow \mathbf{Reg}_N/N \approx$ time-average sub-optimality

- ▶ Recall ERM prob. \Rightarrow loss ℓ , parameter vector $\mathbf{w} \in \mathbb{R}^p$
 \Rightarrow Alternative measure of performance is **regret**

$$\mathbf{Reg}_N = \sum_{n=1}^N \ell(\mathbf{w}_n; \mathbf{x}_n, \mathbf{y}_n) - \sum_{n=1}^N \ell(\mathbf{w}^*; \mathbf{x}_n, \mathbf{y}_n)$$

- ▶ **Goal: asymptotic no-regret $\mathbf{Reg}_N/N \rightarrow 0$**
 \Rightarrow benefit of this perspective is dropping i.i.d. assumption
 \Rightarrow **$\mathbf{Reg}_N/N \approx$ time-average sub-optimality**
- ▶ N large/data arrives online \Rightarrow can't compute full gradient of $L(\mathbf{w})$
 \Rightarrow Classically solved with stochastic (online) gradient

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla_{\mathbf{w}} \ell(\mathbf{w}_t; \mathbf{x}_t, \mathbf{y}_t), \quad \ell_t(\mathbf{w}) := \ell(\mathbf{w}; \mathbf{x}_t, \mathbf{y}_t)$$

- ▶ Descend w/ stoch. grad. rather than grad. \Rightarrow one sample/time
- ▶ One can establish that **$\mathbf{w}_t \rightarrow \mathbf{w}^*$ a.s. and $\mathbf{Reg}_N/N \rightarrow 0$**

Introduction

Generalized Linear Models

- Networked Regret Minimization
- Online Learning in Complex Networks

Dictionary Learning

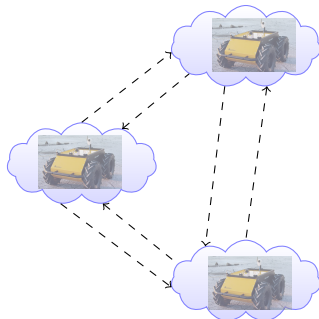
- Robot Path-Planning
- Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ Network $\mathcal{G} = (\mathcal{V}, \mathcal{E})$
 - ⇒ $|\mathcal{V}| = V, |\mathcal{E}| = E$
- ▶ Neighborhood of agent i
 - ⇒ $n_i = \{j : (j, i) \in \mathcal{E}\}$
- ▶ Repeated game at agent i , time t
 - ⇒ params. $\tilde{\mathbf{w}}_{i,t}$ ⇒ local loss $l_{i,t}$
- ▶ Minimize only local loss
 - ⇒ decoupled local learning
- ▶ Instead, each agent i aims to
 - ⇒ minimize global loss $l_t(\tilde{\mathbf{w}}_t) = \sum_{i=1}^V l_{i,t}(\tilde{\mathbf{w}}_t)$
 - ⇒ only observes local loss $l_{i,t}$ ⇒ collaborate with other agents



- ▶ **Local Regret** of **node i** of distributed online algorithm

$$\mathbf{Reg}_T^i = \sum_{t=1}^T \sum_{j=1}^V \ell_{j,t}(\tilde{\mathbf{w}}_{i,t}) - \sum_{t=1}^T \sum_{j=1}^V \ell_{j,t}(\mathbf{w}^*).$$

- ⇒ $\mathbf{w}^* = \operatorname{argmin}_{\mathbf{x}} \sum_{t=1}^T \sum_{j=1}^V \ell_{j,t}(\mathbf{w})$ is the global batch solution
- ⇒ Quality of node i 's prediction at others' losses

- ▶ **Global Networked Regret**

$$\mathbf{Reg}_T := \frac{1}{N} \sum_{i=1}^V \mathbf{Reg}_T^i = \frac{1}{N} \sum_{t=1}^T \sum_{i,j=1}^V \ell_{j,t}(\tilde{\mathbf{w}}_{i,t}) - \sum_{t=1}^T \sum_{j=1}^N \ell_{j,t}(\mathbf{w}^*),$$

- ▶ Networked online learning goal: $\mathbf{Reg}_T^i/T, \mathbf{Reg}_T/T \rightarrow 0$ as $T \uparrow$
 - ⇒ Measures how well agents learn global information

- ▶ **Convexity** of $\ell_{i,t} \implies$ **unique** globally optimal offline strategy \mathbf{w}^*
 - \implies Node predictions should coincide at optimality.
- ▶ At each t we want to enforce $\mathbf{w}_{i,t} = \mathbf{w}_{j,t}$ for $j \in n_i$, or $\mathbf{C}\mathbf{w}_t = \mathbf{0}$.
 - $\implies \mathbf{C}$ is node-edge incidence matrix of \mathcal{G} .
 - $\implies \mathbf{w}_t$ is stacked version of $\mathbf{w}_{i,t}$.
- ▶ Constraint enforcement requires global coordination
 - \implies Lagrangian relaxation allows distributed computation
- ▶ Online Lagrangian for networked learning problem:

$$\mathcal{L}_t(\mathbf{w}_t, \boldsymbol{\lambda}_t) = \sum_{i=1}^V \ell_{i,t}(\mathbf{w}_{i,t}) + \boldsymbol{\lambda}_t^\top \mathbf{C}\mathbf{w}_t$$

- ▶ Convex/concave function in the primal/dual variables
- ▶ $\mathcal{L}_t(\mathbf{w}_t, \boldsymbol{\lambda}_t)$ is stoch. approx. of Lagrangian of prob. (under i.i.d.):
 $\min_{\mathbf{w}} \sum_{i=1}^V \mathbb{E}_{\mathbf{x}_i, \mathbf{y}_i} [\ell(\mathbf{w}_i; \mathbf{x}_i, \mathbf{y}_i)]$ such that $\mathbf{w}_i = \mathbf{w}_j$

- ▶ Arrow-Hurwicz saddle pt. to online Lagrangian $\mathcal{L}_t(\mathbf{w}_t, \lambda_t)$
 - ⇒ Primal Lagrangian subgradient descent
 - ⇒ Dual Lagrangian subgradient ascent

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathcal{P}_X[\mathbf{w}_t - \epsilon \nabla_{\mathbf{w}} \mathcal{L}_t(\mathbf{w}_t, \lambda_t)] \\ \lambda_{t+1} &= \mathcal{P}_\Lambda[\lambda_t + \epsilon \nabla_{\lambda} \mathcal{L}_t(\mathbf{w}_t, \lambda_t)]\end{aligned}$$

- ▶ Initialize $\lambda_1 = \mathbf{0}$ for λ_t to remain in the image of \mathbf{C}
 - ⇒ Required for bounded dual subgradients and constraint slacks
- ▶ Yields a decentralized algorithm:

$$\mathbf{w}_{i,t+1} = \mathcal{P}_{\mathcal{W}} \left[\mathbf{w}_{i,t} - \epsilon \left(\nabla_{\mathbf{w}_i} \ell_{i,t}(\mathbf{w}_{i,t}) + \sum_{j \in n_i} \lambda_{ij,t} - \lambda_{ji,t} \right) \right]$$

- ▶ Dual step at each edge (i, j) ⇒ increase price of disagreement

$$\lambda_{ij,t+1} = \mathcal{P}_{\Lambda_{ij}} \left[\lambda_{ij,t} + \epsilon (\mathbf{w}_{i,t} - \mathbf{w}_{j,t}) \right]$$

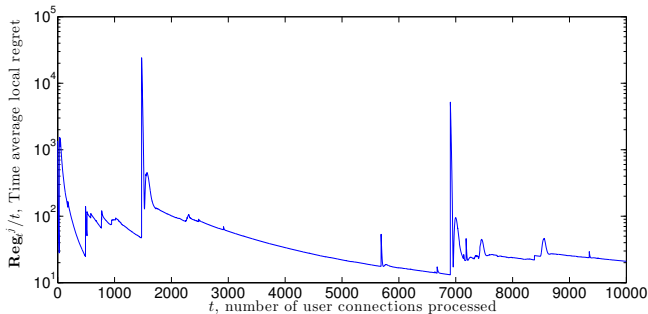
Theorem

The saddle point alg. sequence with initialization $\lambda_1 = \mathbf{0}$ and constant step size $\epsilon = 1/\sqrt{T}$ achieves the *Global networked regret bound*

$$\mathbf{Reg}_T \leq \frac{\sqrt{T}}{2} (\|\mathbf{w}_1 - \mathbf{w}^*\|^2 + EC_\lambda^2 + G_{\mathbf{w}}^2 + G_\lambda^2) = \mathcal{O}(\sqrt{T}).$$

- ▶ $\mathbf{Reg}_T/T \rightarrow 0$ with $T \uparrow$, learning constant depends on . . .
 - ⇒ Network size V and diameter D
 - ⇒ Initialization, Lipschitz constant K , gradient bounds $G_{\mathbf{w}}$, G_λ
- ▶ C_λ must satisfy $C_\lambda \geq DVK + 1 \Rightarrow$ dual set projection
- ▶ Comparable to centralized regret of online gradient descent

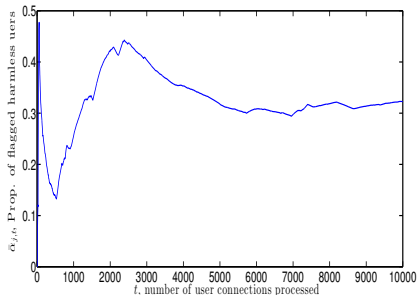
$\zeta = \log(41)$, $\epsilon = 1/\sqrt{T}$, $N = 50$ cycle network, Arbitrarily chosen $j \in V$



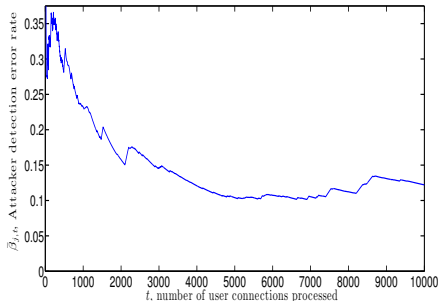
Local regret \mathbf{Reg}_t^j/t vs. no. user connections t

- ▶ Identify attackers in LAN
- ▶ KDDCup 99 data set:
 - ⇒ $VT \approx 5 \times 10^5$ data pts.
 - ⇒ $p = 41$ features
- ▶ $\mathbf{Reg}_T^j/T \rightarrow 0$
 - ⇒ Spikes ⇒ misclassifications
 - ⇒ Recover quickly after mistakes
 - ⇒ No raw info exchange

- ▶ Attacker detection protocol on fixed test set of size $T = 10^4$



(a) Avg. false alarm rate vs. no. of user connections t



(b) Avg. error rate vs. no. of user connections t

- ▶ $\bar{\alpha}_{j,t}$: False alarm rate
- ▶ Friendly users flagged
⇒ within $[\cdot33, \cdot30]$

- ▶ $\bar{\beta}_{j,t}$: error rate
- ▶ Undetected attacks
⇒ within $[\cdot13, \cdot17]$

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ Before, we solved the problem

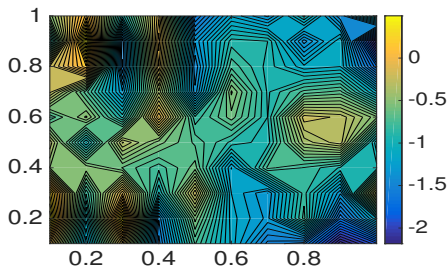
$$\min_{\mathbf{w}} \sum_{i=1}^V \mathbb{E}_{\mathbf{x}_i, \mathbf{y}_i} [\ell(\mathbf{w}_i, \mathbf{x}_i, \mathbf{y}_i)] \quad \text{such that } \mathbf{w}_i = \mathbf{w}_j$$

- ⇒ implicitly assumes nodes seek independent, common params
- ▶ If there is latent correlation among variables
 - ⇒ equality constraint will harm estimation accuracy
- ▶ Introduce proximity function $h(\mathbf{w}_i, \mathbf{w}_j)$
 - ⇒ couples variables of i and j according to a prior γ_{ij} on $\rho(\mathbf{w}_i, \mathbf{w}_j)$

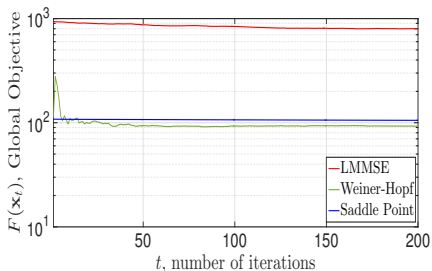
$$\min_{\mathbf{w}} \sum_{i=1}^V \mathbb{E}_{\mathbf{x}_i, \mathbf{y}_i} [\ell(\mathbf{w}_i, \mathbf{x}_i, \mathbf{y}_i)] \quad \text{such that } h(\mathbf{w}_i, \mathbf{w}_j) \leq \gamma_{ij}$$

- ⇒ can solve this problem with comparable saddle pt. techniques

- ▶ $N = 100$ grid sensor network
⇒ deployed in 200 sq. m. region
- ▶ Linear estimation w/ corr. obs.
⇒ distance corr. $\rho_{ij} = e^{-\|l_i - l_j\|}$
- ▶ Constant step-size $\epsilon = 10^{-2.75}$
⇒ Prox. func. $\|\mathbf{w}_i - \mathbf{w}_j\|^2 \leq \gamma_{ij}$
⇒ γ_{ij} ⇒ sample correlation
- ▶ Comparable performance to (recursive) Wiener-Hopf estimator
⇒ via proximity constraints

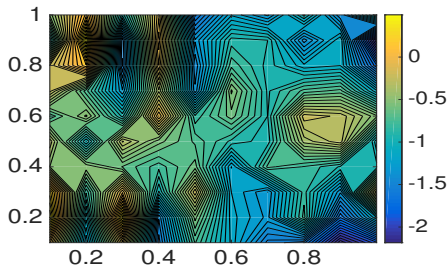


(a) Snapshot of random field

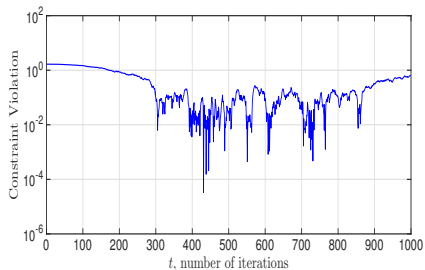


(b) Objective over iteration t

- ▶ $N = 100$ grid sensor network
⇒ deployed in 200 sq. m. region
- ▶ Linear estimation w/ corr. obs.
⇒ distance corr. $\rho_{ij} = e^{-\|l_i - l_j\|}$
- ▶ Constant step-size $\epsilon = 10^{-2.75}$
⇒ Prox. func. $\|\mathbf{w}_i - \mathbf{w}_j\|^2 \leq \gamma_{ij}$
⇒ γ_{ij} ⇒ sample correlation
- ▶ Comparable performance to (recursive) Wiener-Hopf estimator
⇒ via proximity constraints



(a) Snapshot of random field



(b) Constraint Violation over iteration t

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ In distributed convex setting \Rightarrow we achieve optimum of ERM
 - \Rightarrow even when data is scattered across a multi-agent network
 - \Rightarrow mediocre accuracy \Rightarrow complexity of relating \mathbf{x} and \mathbf{y}
- ▶ Need better estimator \Rightarrow alternative data representations
 - \Rightarrow Signal encoding \Rightarrow Fourier, wavelets, PCA, or data-driven
 - \Rightarrow Task-driven: tailor dictionary to inference (Mairal '12)

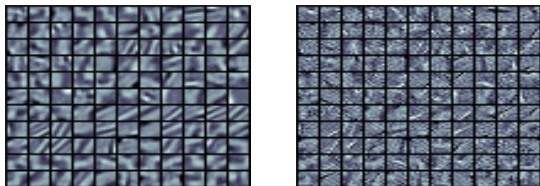


Figure: Initialized (left) and learned (right) dictionary for small image patches.

- ▶ Represent signals \mathbf{x}_t as combos. of k basis elements $\{\mathbf{d}_l\}_{l=1}^k$
 - ⇒ learn dictionary $\mathbf{D} \in \mathbb{R}^{m \times k}$ from data
 - ⇒ Denote the coding of \mathbf{x}_t as $\alpha_t \in \mathbb{R}^k$
- ▶ **Representation loss** $g(\alpha_t, \mathbf{D}; \mathbf{x}_t) \Rightarrow$ small if $\mathbf{D}\alpha_t$ and \mathbf{x}_t close
 - ⇒ $\mathbf{D}\alpha_t$ is representation of \mathbf{x}_t w.r.t dictionary \mathbf{D}
- ▶ Formulate the **coding problem**

$$\alpha^*(\mathbf{D}; \mathbf{x}_t) := \underset{\alpha_t \in \mathbb{R}^k}{\operatorname{argmin}} g(\alpha_t, \mathbf{D}; \mathbf{x}_t) .$$

- ▶ Dictionary learning
 - ⇒ seek \mathbf{D} such that signals \mathbf{x}_t well-represented by $\mathbf{D}\alpha^*(\mathbf{D}; \mathbf{x}_t)$

- ▶ Tailor dictionary to **discriminative modeling task**
- ▶ Use coding $\alpha^*(\mathbf{D}; \mathbf{x}_t)$ as representation of signal \mathbf{x}_t
- ▶ Decision variable $\mathbf{w} \Rightarrow$ predict the label/vector \mathbf{y}_t given $\alpha^*(\mathbf{D}; \mathbf{x}_t)$.
- ▶ **Loss function** $\ell(\mathbf{D}, \mathbf{w}; (\mathbf{x}_t, \mathbf{y}_t)) = \ell(\alpha^*(\mathbf{D}; \mathbf{x}_t), \mathbf{D}, \mathbf{w}; (\mathbf{x}_t, \mathbf{y}_t))$
 \Rightarrow predictive quality of \mathbf{w} for var. \mathbf{y}_t given coding $\alpha^*(\mathbf{D}; \mathbf{x}_t)$

- ▶ Discriminative dictionary learning

$$(\mathbf{D}^*, \mathbf{w}^*) := \operatorname{argmin}_{\mathbf{D} \in \mathcal{D}, \mathbf{w} \in \mathcal{W}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\ell(\mathbf{D}, \mathbf{w}; (\mathbf{x}, \mathbf{y}))].$$

- \Rightarrow Learn jointly regression weights \mathbf{w} and dictionary \mathbf{D}
- \Rightarrow **Non-convex** stochastic program

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ One role for **learning** in robotics: learn model uncertainty
- ▶ Simplified physics models used for control due to complexity
 - ⇒ Models are available. Not perfect but not useless either
 - ⇒ Replace mechanical models with learned models
- ▶ Learn mismatch between model and reality when
 - ⇒ This mismatch has **variability across different terrains**
- ▶ Use **sensory input to learn uncertainty** in execution of controller

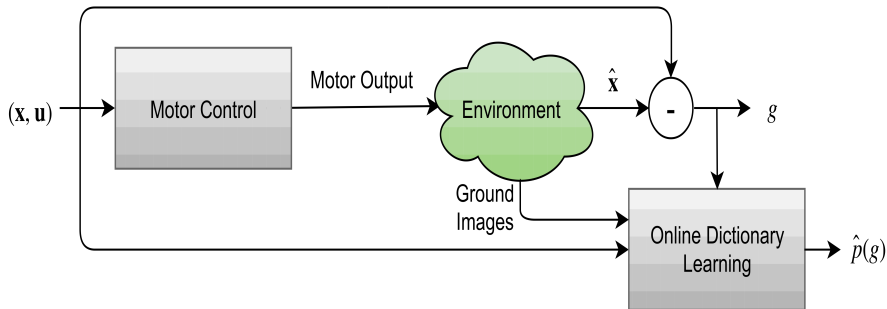
- ▶ Consider a discrete nonlinear state-space system of equations

$$\mathbf{x}_{k+1} = f(\mathbf{x}_k, \mathbf{u}_k) + g(\mathbf{a}_k) = f(\mathbf{x}_k, \mathbf{u}_k) + g(\mathbf{u}_k, \mathbf{z}_k)$$

- ▶ $\mathbf{x}_k \Rightarrow$ state vector, $\mathbf{u}_k \Rightarrow$ control input, $\mathbf{z}_k \Rightarrow$ sensory input
- ▶ Kinematic model $f(\mathbf{x}_k, \mathbf{u}_k)$ not exact \Rightarrow add mismatch term $g(\mathbf{a}_k)$
 \Rightarrow Want to learn $g(\mathbf{a}_k)$ to use as input to robust control block
- ▶ Measure estimate $\hat{\mathbf{x}}_k$ of state \mathbf{x}_k (with on-board IMU, for instance)

$$\hat{g}(\mathbf{a}_{k-1}) = \hat{\mathbf{x}}_k - f(\mathbf{x}_{k-1}, \mathbf{u}_{k-1}).$$

- ▶ Learning $\hat{g}(\mathbf{a}_{k-1})$ is challenging
 - \Rightarrow Captures difficult-to-model physics we typically ignore
 - \Rightarrow Dictionary leaning approach



- ▶ Platform's state \mathbf{x} , control \mathbf{u} intended by a kinematic planner
⇒ differ from measured ground truth $\hat{\mathbf{x}}$ by disturbance g
- ▶ This difference, as well as state, control, and visual features
⇒ Fed into dictionary learning method ⇒ disturbance pred. \hat{g}
⇒ Dictionary is a statistical model using sparse approximation

- ▶ Model disturbance $\hat{g}(\mathbf{a})$ as Gaussian conditional on features \mathbf{a}

$$\mathbb{P}[\hat{g}(\mathbf{a}) \mid \mathbf{a}] = \frac{1}{\sqrt{2\pi\sigma^2(\mathbf{a})}} \exp \left[-\frac{(\hat{g}(\mathbf{a}) - \mu(\mathbf{a}))^2}{2\sigma^2(\mathbf{a})} \right].$$

- ▶ Distribution parameterized by unknown mean $\mu(\mathbf{a})$, var. $\sigma^2(\mathbf{a})$
⇒ which depend on control \mathbf{u}_k and sensory input \mathbf{z}_k
- ▶ Realizations of $(\mathbf{a}, \hat{g}(\mathbf{a}))$ **available online**
- ▶ Sequentially obtained while exploring feature space
- ▶ Utilize to learn parametric representation of the distribution

- ▶ Learn online mean, variance \Rightarrow introduce **regressors** $\mathbf{w}_1, \mathbf{w}_2$
 \Rightarrow predict first, second-order stats. $\mu(\mathbf{a})$ and $\sigma^2(\mathbf{a})$, given \mathbf{a}

$$\hat{\mu}(\mathbf{a}) = \mathbf{w}_1^T \mathbf{a}, \quad \hat{\sigma}^2(\mathbf{a}) = \sigma_{\min}^2 + \left(\mathbf{w}_2^T \mathbf{a} + \sigma_{\text{init}}^2 \right)^2$$

- ▶ Rather than use \mathbf{a} directly, use a **sparse code** $\alpha^*(\mathbf{D}; \mathbf{a})$

$$\hat{\mu}(\mathbf{a}) = \mathbf{w}_1^T \alpha^*(\mathbf{D}; \mathbf{a}), \quad \hat{\sigma}^2(\mathbf{a}) = \sigma_{\min}^2 + \left(\mathbf{w}_2^T \alpha^*(\mathbf{D}; \mathbf{a}) + \sigma_{\text{init}}^2 \right)^2$$

\Rightarrow Regress on sparse approximation $\alpha^*(\mathbf{D}; \mathbf{a})$ w.r.t. dictionary \mathbf{D}

- ▶ Motivation for using sparse code $\alpha^*(\mathbf{D}; \mathbf{a})$, learning dictionary \mathbf{D} :
 - $\Rightarrow g(\cdot)$ relates robotic sensory perception, unexpected dynamics
 - \Rightarrow Relationship between $(\mathbf{a}, \hat{g}(\mathbf{a}))$ highly nonlinear
 - \Rightarrow Estimation accuracy \Rightarrow boosted via **alternative encoding**

- ▶ Dictionary $\mathbf{D} = \{\mathbf{d}_l\}_{l=1}^m$, $\mathbf{d}_l \in \mathbb{R}^k$ composed of m basis elements
- ▶ Estimate $\hat{\mathbf{a}}_k = \mathbf{D}\alpha_k$ as linear combo. of dictionary elements
- ▶ Select coefficients that yield a **sparse code (elastic net)**

$$\alpha^*(\mathbf{D}; \mathbf{a}_k) := \operatorname{argmin}_{\alpha \in \mathbb{R}^k} \|\mathbf{a}_k - \mathbf{D}\alpha\|^2 + \lambda \|\alpha\|_1 + \nu \|\alpha\|_2^2$$

- ▶ Jointly learn dictionary and regressors \mathbf{w}_1 and \mathbf{w}_2

$$(\mathbf{D}^*, \mathbf{w}_1^*, \mathbf{w}_2^*) := \operatorname{argmin}_{\mathbf{D} \in \mathcal{D}, \mathbf{w}_1, \mathbf{w}_2} \mathbb{E}_{\mathbf{a}, \hat{g}(\mathbf{a})} \left(-\log \mathbb{P}[\hat{g}(\mathbf{a}) | \mathbf{a}, \mathbf{D}, \mathbf{w}_1, \mathbf{w}_2] \right).$$

- ▶ **Non-convex** but convex w.r.t. \mathbf{D} and \mathbf{w}_1 and \mathbf{w}_2 separately
- ▶ Objective is an expectation over dataset \Rightarrow use stochastic grad.

- ▶ Observe signals \mathbf{z}_k , use past control \mathbf{u}_k to compute coding

$$\alpha_k^* := \underset{\alpha_k \in \mathbb{R}^s}{\operatorname{argmin}} (1/2) \|\mathbf{a}_k - \mathbf{D}\alpha_k\|_2^2 + \lambda \|\alpha_k\|_1 + \nu \|\alpha_k\|_2,$$

- ⇒ Update dictionary using stoch. grad. step w.r.t. dictionary

$$\mathbf{D}_{k+1} = \mathbf{D}_k - \epsilon_k (\nabla_{\mathbf{D}} \log \mathbb{P}[\hat{g}(\mathbf{a}_k) | \mathbf{a}_k, \mathbf{D}_k, \mathbf{w}_{1,k}, \mathbf{w}_{2,k}])$$

- ▶ Update regressors along regressor gradient of loss function

$$\mathbf{w}_{1,k+1} = \mathbf{w}_{1,k} + \epsilon_k (\nabla_{\mathbf{w}_1} \log \mathbb{P}[\hat{g}(\mathbf{a}_k) | \mathbf{a}_k, \mathbf{D}_k, \mathbf{w}_{1,k}, \mathbf{w}_{2,k}]) ,$$

$$\mathbf{w}_{2,k+1} = \mathbf{w}_{2,k} + \epsilon_k (\nabla_{\mathbf{w}_2} \log \mathbb{P}[\hat{g}(\mathbf{a}_k) | \mathbf{a}_k, \mathbf{D}_k, \mathbf{w}_{1,k}, \mathbf{w}_{2,k}]) ,$$

- ▶ Converges to locally optimal dictionary and regressors

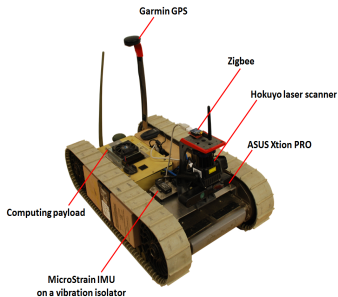


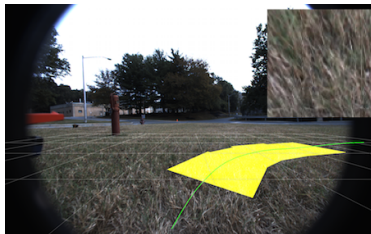
Figure: An iRobot *Packbot* was used in our experiments. It was additionally configured with a high-resolution camera.

- ▶ We consider a differential drive model of a skid-steer robot

$$f(\mathbf{x}_k, \mathbf{u}_k) = \begin{bmatrix} \dot{x}_k \\ \dot{y}_k \\ \dot{\theta}_k \end{bmatrix} = \underbrace{\begin{bmatrix} \cos(\theta_k) & -\sin(\theta_k) & 0 \\ \sin(\theta_k) & \cos(\theta_k) & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{A(\theta)} \begin{bmatrix} \nu_k \\ \omega_k \end{bmatrix}$$

- ▶ Disturbance \Rightarrow **commanded & actual angular velocity** difference

- ▶ Visual patch $\mathbf{z}_k \Rightarrow$ associated with the portion of ground
 \Rightarrow Collect images over time horizon of planned robot trajectory



- ▶ From the raw patch we construct statistical visual features \mathbf{c}_k
 \Rightarrow mean, variance, skewness, kurtosis of RGB color channels
 \Rightarrow Textures \mathbf{h}_k via texton histogram (Leung '99)
- ▶ Concatenated with average linear, angular velocity in $[k, k + 1]$

- ▶ Model fit of disturbance to task driven dictionary learnt dist.
- ▶ Compared to (windowed) recursive average of mean, variance

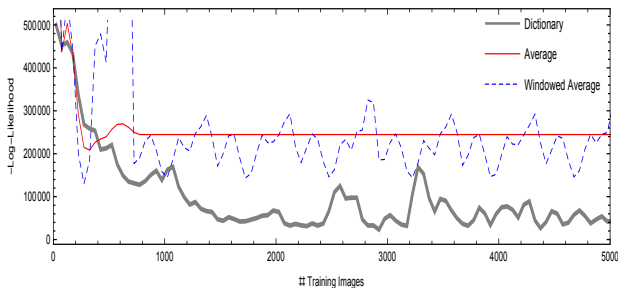
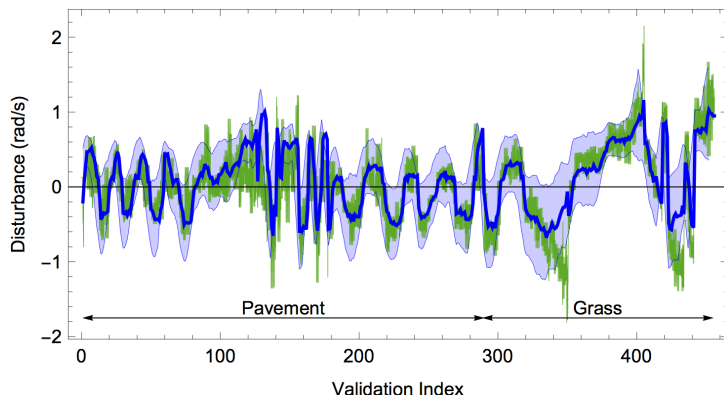


Figure: Comparison of dictionary learning vs. classical alternatives.

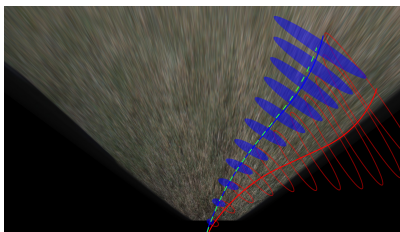
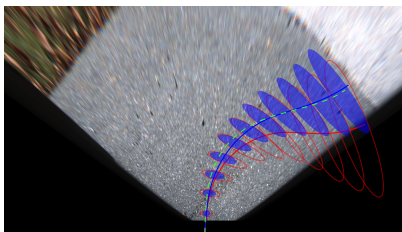
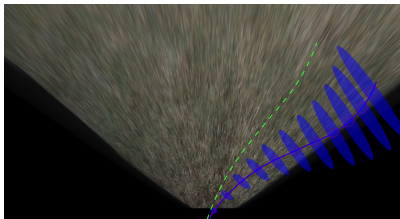
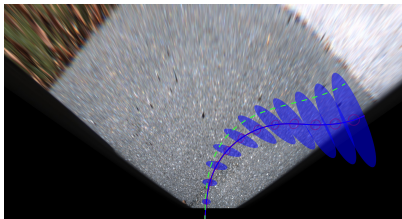
- ▶ Superior model fit to Gaussian approximation of disturbance
- ▶ Exploits sensory input to identify terrain type
 - ⇒ Pavement or grass essentially. But more granular than that

- ▶ Test trajectory with predicted & actual disturbance stats. overlaid
- ▶ Measured dist. (green) and predicted dist. (blue) for trajectory
⇒ Predicted mean and $\pm 2\sigma$ envelope shown



- ▶ Observations are mostly contained within confidence envelopes

- ▶ Actual trajectory not contained within cones for initial dictionary
⇒ But contained within cone after dictionary is properly learnt



Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ Dictionary-based estimations better than GLM
⇒ but what if data is scattered across a network?
- ▶ Incentivize agreement via constraint $\mathbf{D}_i = \mathbf{D}_j, \mathbf{w}_i = \mathbf{w}_j$ for all $j \in n_i$
- ▶ Decentralized task-driven dictionary learning problem

$$\{\mathbf{D}_i^*, \mathbf{w}_i^*\}_{i=1}^V := \underset{\mathbf{D}_i \in \mathcal{D}, \mathbf{w}_i \in \mathcal{W}}{\operatorname{argmin}} \sum_{i=1}^V \mathbb{E}_{\mathbf{x}_i, \mathbf{y}_i} [\ell(\mathbf{D}_i, \mathbf{w}_i; (\mathbf{x}_i, \mathbf{y}_i))] .$$

such that $\mathbf{D}_i = \mathbf{D}_j, \mathbf{w}_i = \mathbf{w}_j$ for all $j \in n_i$

- ▶ Similar to networked regret minimization
- ▶ Enforcing agreement constraint would require global coordination
⇒ Define stochastic Lagrangian ⇒ distributed alg.

$$\hat{\mathcal{L}}_t(\mathbf{D}, \mathbf{w}, \Lambda, \nu) = \sum_{i=1}^V [\ell(\mathbf{D}_i, \mathbf{w}_i; (\mathbf{x}_{i,t}, \mathbf{y}_{i,t}))] + \operatorname{tr}(\Lambda^T \mathbf{C}_D \mathbf{D}) + \nu^T \mathbf{C}_w \mathbf{w}$$

- ▶ Decentralized online dict. learning \Rightarrow Block saddle point alg.
- ▶ Stochastic approximation: $\mathcal{L}(\mathbf{D}, \mathbf{w}, \Lambda, \nu) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\hat{\mathcal{L}}_t(\mathbf{D}, \mathbf{w}, \Lambda, \nu)]$
 \Rightarrow **primal** stochastic gradient descent

$$\mathbf{D}_{t+1} = \mathbf{D}_t - \epsilon_t \nabla_{\mathbf{D}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t),$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \epsilon_t \nabla_{\mathbf{w}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t).$$

\Rightarrow **dual** stochastic gradient ascent

$$\Lambda_{t+1} = \Lambda_t + \epsilon_t \nabla_{\Lambda} \hat{\mathcal{L}}_t(\mathbf{D}_{t+1}, \mathbf{w}_{t+1}, \Lambda_t, \nu_t),$$

$$\nu_{t+1} = \nu_t + \epsilon_t \nabla_{\nu} \hat{\mathcal{L}}_t(\mathbf{D}_{t+1}, \mathbf{w}_{t+1}, \Lambda_t, \nu_t).$$

- ▶ $\nabla_{\mathbf{D}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t) \Rightarrow$ Projected stoch. grad. w.r.t. \mathbf{D}
 \Rightarrow gradient approximated with current signals $\{\mathbf{x}_{i,t}, \mathbf{y}_{i,t}\}_{i=1}^V$

- ▶ At agent i , time t , observe $(\mathbf{x}_{i,t}, \mathbf{y}_{i,t})$,
- ▶ Compute coding $\alpha_{i,t+1}^* = \operatorname{argmin}_{\alpha \in \mathbb{R}^k} g(\alpha, \mathbf{D}_{i,t}; \mathbf{x}_{i,t})$
⇒ In practice chosen as *sparse coding* via lasso or elastic-net
- ▶ Update primal variables at agent i

$$\mathbf{D}_{i,t+1} = \mathbf{D}_{i,t} - \epsilon_t \left(\nabla_{\mathbf{D}_i} \ell_i(\mathbf{D}_{i,t}, \mathbf{w}_{i,t}; (\mathbf{x}_{i,t}, \mathbf{y}_{i,t})) + \sum_{j \in n_i} (\Lambda_{ij,t} - \Lambda_{ji,t}) \right),$$

$$\mathbf{w}_{i,t+1} = \mathbf{w}_{i,t} - \epsilon_t \left(\nabla_{\mathbf{w}_i} \ell_i(\mathbf{D}_{i,t}, \mathbf{w}_{i,t}; (\mathbf{x}_{i,t}, \mathbf{y}_{i,t})) + \sum_{j \in n_i} (\nu_{ij,t} - \nu_{ji,t}) \right),$$

- ▶ Update dual variables at network communication link (i, j)

$$\Lambda_{ij,t+1} = \Lambda_{ij,t} + \epsilon_t (\mathbf{D}_{i,t} - \mathbf{D}_{j,t})$$

$$\nu_{ij,t+1} = \nu_{ij,t} + \epsilon_t (\mathbf{w}_{i,t} - \mathbf{w}_{j,t})$$

Theorem

Saddle pt. seq. $(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t)$ converges to stationarity in expectation:

$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\mathbf{D}} \mathcal{L}(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t)\|] = 0,$$

$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t)\|] = 0$$

Asymptotic feasibility condition achieved in expectation:

$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\Lambda} \mathcal{L}(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t)\|] = 0$$

$$\lim_{t \rightarrow \infty} \mathbb{E}[\|\nabla_{\nu} \mathcal{L}(\mathbf{D}_t, \mathbf{w}_t, \Lambda_t, \nu_t)\|] = 0$$

- ▶ Performance guarantee for D4L
 - ⇒ convergence in non-convex stochastic opt.
 - ⇒ sensitive to data distribution, step-size, network structure

- ▶ Texture database classification problem \Rightarrow Brodatz textures
 - \Rightarrow Insight into **dynamic image processing problems**
 - \Rightarrow **Toy model of real-time navigability analysis** in robotic teams
- ▶ Real-time image data \Rightarrow train multi-class logistic reg. weights
- ▶ Decentralized dynamic texture classification
 - \Rightarrow Subset of textures: {grass, bark, straw, herringbone_weave}

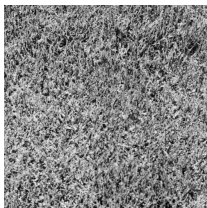


Figure: Sample images from Brodatz textures.

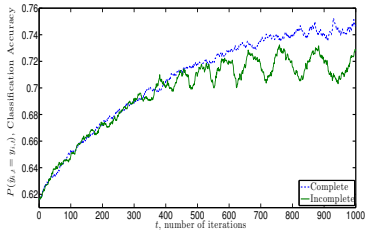
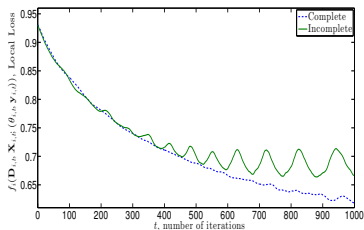


Figure: Local loss (left) and classification accuracy (right) versus iteration t .

- ▶ $V = 10$ node random network, results shown for random $j \in \mathcal{V}$
- ▶ Agents observe **random incomplete subsets** of feature space
- ▶ Still learn **global information** and reach consensus
- ▶ Moderate classifier performance
 - ⇒ due to small step-size required for convergence
 - ⇒ Small step-sizes required for convergence

- ▶ Each robot in the network sequentially observes images
 - ⇒ partitions them into small patches
 - ⇒ classify patches with multi-class logistic regression.
- ▶ Terrain classification has been used as a layer in robust control
 - ⇒ Classes are terrains of varying traversability
- ▶ Data from Lejeune Robotics Test Facility ⇒ Thanks to ARL!

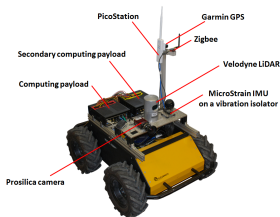


Figure: Sample image (left) from a $N = 3$ robot network of Huskies (right).

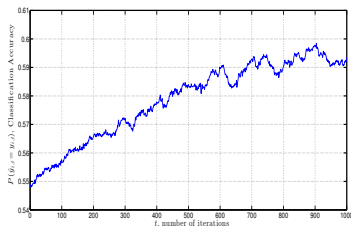
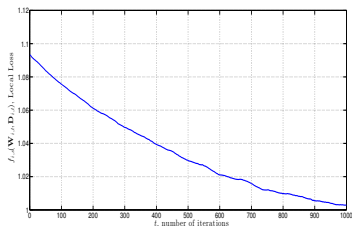


Figure: Local loss (left) and classification accuracy (right) versus iteration t . $V = 3$ complete graph w/ complete sampling.

- ▶ Main takeaways: training is slow, but convergent
 - ⇒ due to necessity of small step-sizes
 - ⇒ good for centralized learning, slow in distributed case
- ▶ Non-convexity makes training/optimization challenging
 - ⇒ nonlinear classifiers + convex training? ⇒ flexible dictionary

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ Learning nonlinear statistical models \Rightarrow function estimation
- ▶ Want to find $f^* \in \mathcal{H}$ to minimize regularized expected risk $R(f)$

$$f^* = \underset{f \in \mathcal{H}}{\operatorname{argmin}} R(f) := \underset{f \in \mathcal{H}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(f(\mathbf{x}), \mathbf{y})] + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

- \Rightarrow Loss $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ penalize deviations between $f(\mathbf{x})$, \mathbf{y}
- ▶ Generally intractable \Rightarrow infinite dimensional data-dependent opt.

- ▶ Learning nonlinear statistical models \Rightarrow function estimation
- ▶ Want to find $f^* \in \mathcal{H}$ to minimize regularized expected risk $R(f)$

$$f^* = \operatorname{argmin}_{f \in \mathcal{H}} R(f) := \operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(f(\mathbf{x}), y)] + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2$$

\Rightarrow Loss $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ penalize deviations between $f(\mathbf{x})$, \mathbf{y}

- ▶ Generally intractable \Rightarrow infinite dimensional data-dependent opt.
- ▶ **Reproducing kernels** \Rightarrow framework to make this task possible!
 $\Rightarrow \mathcal{H}$ is equipped with unique *kernel function*, $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$

$$f^* := \operatorname{argmin}_{f \in \mathcal{H}} \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(\sum_{n \in \mathcal{I}} w_n \kappa(\mathbf{x}_n, \mathbf{x}), y)] + \frac{\lambda}{2} \left\| \sum_{n, m \in \mathcal{I}} w_n w_m \kappa(\mathbf{x}_n, \mathbf{x}_m) \right\|_{\mathcal{H}}^2$$

\Rightarrow via use of Representer Theorem, dating back to Riesz Thm.

- ▶ Kernel methods learn function $f(\mathbf{x}) = \sum_{n \in \mathcal{I}} w_n \kappa(\mathbf{x}_n, \mathbf{x})$
 - ⇒ comes from Representer Theorem, dating back to Riesz Thm.
 - ⇒ \mathcal{I} is a countably infinite indexing set
- ▶ Maintain **convexity** while learning nonlinear statistical model
 - ⇒ complicated representation: $|\mathcal{I}| = \infty$ vs. k basis elements
 - ⇒ flexible dictionary, model conditional density of inference
- ▶ Train with functional stochastic gradient descent?

$$f_{t+1}(\cdot) = (1 - \eta_t \lambda) f_t - \eta_t \ell'(f_t(\mathbf{x}_t), \mathbf{y}_t) \kappa(\mathbf{x}_t, \cdot)$$

- ▶ **Problem: training complexity with SGD is cubic in iteration index**
 - ⇒ Sparsify the solution? Need to sparsify **training**

- ▶ We implement functional generalization of SGD

$$\tilde{f}_{t+1}(\cdot) = (1 - \eta_t \lambda) f_t - \eta_t \ell'(f_t(\mathbf{x}_t), \mathbf{y}_t) \kappa(\mathbf{x}_t, \cdot)$$

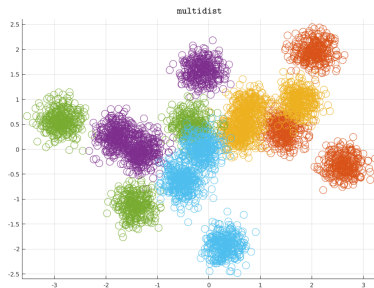
$$(f_{t+1}, \mathbf{D}_{t+1}, \mathbf{w}_{t+1}) = \mathbf{KOMP}(\tilde{f}_{t+1}, \tilde{\mathbf{D}}_{t+1}, \tilde{\mathbf{w}}_{t+1}, \epsilon_t)$$

- ⇒ operating in tandem with projection step onto subspaces of \mathcal{H}
- ⇒ subspaces greedily constructed via matching pursuit

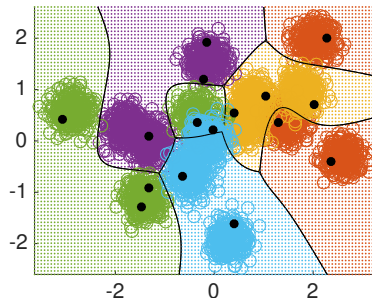
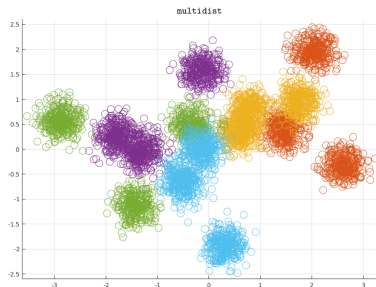
$$f_{t+1} = \operatorname{argmin}_{f \in \mathcal{H}_{\mathbf{D}_{t+1}}} \left\| f - \left((1 - \eta_t \lambda) f_t - \eta_t \nabla_f \ell(f_t(\mathbf{x}_t), \mathbf{y}_t) \right) \right\|_{\mathcal{H}}^2$$

- ▶ Proposed work: online training of sparsified kernel classifiers
 - ⇒ Main problem: sparsification bias may ruin stochastic descent

- ▶ Case where training examples for a fixed class
⇒ drawn from a distinct Gaussian mixture
- ▶ 3 Gaussians per mixture, $C = 5$ classes total for this experiment
⇒ 15 total Gaussians generate data



- ▶ Case where training examples for a fixed class
⇒ drawn from a distinct Gaussian mixture
- ▶ 3 Gaussians per mixture, $C = 5$ classes total for this experiment
⇒ 15 total Gaussians generate data



- ▶ Grid colors ⇒ decision, bold black dots ⇒ kernel dict. elements
- ▶ ~ 97% accuracy

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

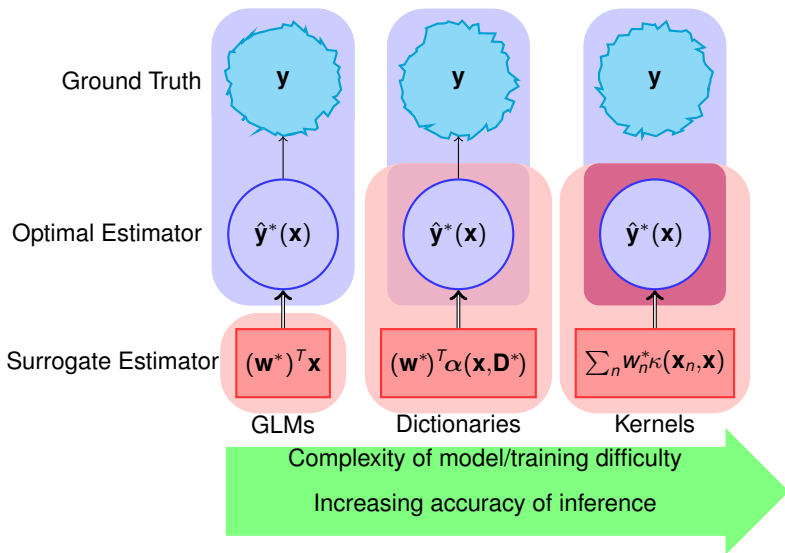
Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix



- ▶ Can establish strong convergence guarantees for GLMs
 - ⇒ In centralized and distributed online cases
 - ⇒ But we are solving the wrong problem ⇒ too far from $\hat{y}^*(\mathbf{x})$
- ▶ Dictionary methods ⇒ richer estimators, nonlinear classifiers
 - ⇒ promising results on robotics application
 - ⇒ non-convexity ⇒ challenge to train in distributed setting
- ▶ Kernel methods allow learning nonlinear classifiers
 - ⇒ preserve convexity, directly model conditional density
 - ⇒ but training has prohibitive complexity in online setting
- ▶ Proposal: learn kernel classifiers online with low complexity

Journals

- ▶ A. Koppel, F. Jakubeic, and A. Ribeiro, "A saddle point algorithm for networked online convex optimization," IEEE Trans. Signal Process., vol.PP, no.99, June. 2015.
- ▶ A. Koppel, B. Sadler, and A. Ribeiro, "Proximity without Consensus in Online Multi-Agent Optimization," in IEEE Trans. Signal Proc. (submitted), June 2016.
- ▶ A. Koppel, G. Warnell, E. Stump, and A. Ribeiro, "D4L: Decentralized Dynamic Discriminative Dictionary Learning," in IEEE Trans. Signal Info. Process over Networks (submitted)., June. 2016.

Conferences

- ▶ A. Koppel, F. Jakubeic and A. Ribeiro, "Regret Bounds of a distributed saddle point algorithm," in Proc. Int. Conf. Acoustics Speech Signal Process., Brisbane, Australia, Apr 19-24 2015.
- ▶ A. Koppel, F. Y. Jakubiec, and A. Ribeiro, "A saddle point algorithm for networked online convex optimization." in 39th Proc. Int. Conf. Acoust. Speech Signal Process., Florence, Italy, May 4-9 2014, pp. 8292 - 8296.
- ▶ A. Koppel, B. M. Sadler and A. Ribeiro, "Proximity without consensus in online multi-agent optimization," in Proc. Int. Conf. Acoustics Speech Signal Process., Shanghai, China, Mar. 20-25 2016.
- ▶ A. Koppel, B. M. Sadler, and A. Ribeiro, "Decentralized Online Optimization with Heterogeneous Data Sources", IEEE Global Conference on Signal and Information Processing (to appear), Washington, DC, Dec. 7-9, 2016.
- ▶ A. Koppel, J. Fink, G. Warnell, E. Stump, and A. Ribeiro, "Online Learning for Characterizing Unknown Environments in Ground Robotic Vehicle Models," in Proc. Int. Conf. Intelligent Robotics and Systems, Daejeon, Korea, Oct9-Oct14 2016
- ▶ A. Koppel, G. Warnell, and E. Stump. "Task-Driven Dictionary Learning in Distrubted Online Settings." in Proc. Asilomar Conf. on Signals Systems Computers, Pacific Grove, CA, November 8-11 2015.
- ▶ A. Koppel, G. Warnell, E. Stump, and A. Ribeiro, "D4L: Decentralized Dynamic Discriminative Dictionary Learning," in Proc. Int. Conf. Intelligent Robotics and Systems, Hamburg, Germany, Sep 28-Oct2 2015.

Introduction

Generalized Linear Models

Networked Regret Minimization

Online Learning in Complex Networks

Dictionary Learning

Robot Path-Planning

Decentralized Dynamic Discriminative Dictionaries

Nonparametric Regression

Conclusion

Appendix

- ▶ The network \mathcal{G} is symmetric and connected with diameter D .
- ▶ Loss function gradients for any \mathbf{w} bounded by constant G , i.e.

$$\|\nabla \ell_t(\mathbf{w})\|_2 \leq G.$$

- ▶ Losses $\ell_{i,t}(\mathbf{x})$ are Lipschitz continuous with modulus $K_{i,t} \leq K$

$$\|\ell_{i,t}(\mathbf{w}) - \ell_{i,t}(\mathbf{v})\|_2 \leq K_{i,t} \|\mathbf{w} - \mathbf{v}\|_2 \leq K \|\mathbf{w} - \mathbf{v}\|_2$$

- ▶ Network $\mathcal{G} \Rightarrow$ symmetric and connected with diameter D .
- ▶ Diminishing step-size rules: $\sum_{t=0}^{\infty} \epsilon_t = \infty$ and $\sum_{t=0}^{\infty} \epsilon_t^2 < \infty$
- ▶ Mean and variance conditions of Lagrangian stochastic gradients

$$\mathbb{E}[\|\delta_{\mathbf{D},t}\| \mid \mathcal{F}_t] \leq A\epsilon_t ,$$

$$\mathbb{E}[\|\nabla_{\mathbf{D}} \hat{\mathcal{L}}_t(\mathbf{D}_t, \mathbf{x}_t, \boldsymbol{\Lambda}_t, \nu_t)\|^2 \mid \mathcal{F}_t] \leq \sigma^2 .$$

- ▶ Feasible dictionary set is those with unit column-norms

$$\mathcal{D} = \{\mathbf{D} \in \mathbb{R}^{m \times k} : \|\mathbf{d}_j\| \leq 1, j = 1 \dots k\} .$$

- ▶ Instantaneous loss \Rightarrow hinge loss at current data point

$$\ell_{i,t}(\mathbf{w}) = \frac{\zeta}{2} \|\mathbf{w}\|_2^2 + \max\left(0, 1 - y_{i,t} \mathbf{w}^T \mathbf{x}_{i,t}\right)$$

- ▶ \mathbf{Reg}_T measures price of distributed causal classifier training
- ▶ **Algorithm Formulation**

$$\mathbf{w}_{i,t+1} = \mathcal{P}_W \left[\mathbf{w}_{i,t} - \epsilon \left(\zeta \mathbf{w}_{i,t} - y_{i,t} \mathbf{x}_{i,t} \mathbb{I}(y_{i,t} \mathbf{w}_{i,t}^T \mathbf{x}_{i,t} < 1) + \sum_{j \in n_i} (\lambda_{ij,t} - \lambda_{ji,t}) \right) \right],$$

$$\lambda_{ij,t+1} = \mathcal{P}_{\Lambda_{ij}} \left[\lambda_{ij,t} + \epsilon (\mathbf{w}_{i,t} - \mathbf{w}_{j,t}) \right].$$

- ▶ Limit classifier complexity to set $W = \{\mathbf{w}_i \in \mathbb{R}^p : \|\mathbf{w}_i\|_2 \leq \zeta\}$
- ▶ $\mathbb{I}(y_{i,t} \mathbf{w}^T \mathbf{x}_{i,t} < 1) = 1$ if $y_{i,t} \mathbf{w}^T \mathbf{x}_{i,t} < 1$, $\mathbb{I}(y_{i,t} \mathbf{w}^T \mathbf{x}_{i,t} < 1) = 0$ else
- ▶ Projected Perceptron with **dual correction (neighbor info)**

- ▶ Multi-class logistic reg. prob. \Rightarrow Agent i receives signals $\mathbf{x}_{i,t}$
 \Rightarrow output a **decision variable** $\mathbf{y}_{i,t} \in \{0, 1\}^C \Rightarrow C$ no. of classes
- ▶ $[\mathbf{y}_{i,t}]_c \Rightarrow$ binary indicator of whether signal falls in class c .
- ▶ Local loss $\ell_i \Rightarrow$ negative log-likelihood of prob. model

$$\ell_i(\mathbf{D}_i, \mathbf{W}_i; (\boldsymbol{\theta}_i, \mathbf{y}_i)) = \log \left(\sum_{c=1}^C e^{\mathbf{w}_{i,c}^T \boldsymbol{\alpha}_i^* + w_{i,c}^0} \right) - \sum_{c=1}^C \left(y_{i,c} \mathbf{w}_{i,c}^T \boldsymbol{\alpha}_i^* + w_{i,c}^0 \right) + \xi \|\mathbf{W}_i\|_F^2,$$

- ▶ $\boldsymbol{\alpha}_i^* \Rightarrow$ sparse coding via elastic-net min. prob.
- ▶ $g_c(\boldsymbol{\alpha}_i^*) = e^{\mathbf{w}_{i,c}^T \boldsymbol{\alpha}_i^* + w_{i,c}^0}$ is activation function;
 $\Rightarrow g_c(\mathbf{z}_i) / \sum_{c'} g_{c'}(\mathbf{z}_i) \Rightarrow$ prob. \mathbf{z}_i in class c
 $\Rightarrow \mathbf{z}_i \Rightarrow$ average of image sub-patches
- ▶ Classification decision \Rightarrow **maximum likelihood class label**
 $\Rightarrow \tilde{c} = \operatorname{argmax}_c g_c(\mathbf{z}_i) / \sum_{c'} g_{c'}(\mathbf{z}_i); [\mathbf{y}_{i,t}]_c = 0$ for $c \neq \tilde{c}$

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

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

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

- ▶ Property (i) \Rightarrow source of “kernel trick:”
 - \Rightarrow define nonlinear map $\phi(\mathbf{x}) = \kappa(\mathbf{x}, \cdot)$ of feature vector \mathbf{x}
 - \Rightarrow accessed only via inner products $\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}} = \kappa(\mathbf{x}, \mathbf{x}')$
- ▶ Property (ii) \Rightarrow Representation Thms. from functional analysis
- ▶ Kernel examples:
 - \Rightarrow Gaussian/RBF $\kappa(\mathbf{x}, \mathbf{x}') = \exp \left\{ -\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2} \right\}$
 - \Rightarrow polynomial $\kappa(\mathbf{x}, \mathbf{x}') = (\mathbf{x}^T \mathbf{x}' + b)^c$

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

$$f^* = \operatorname{argmin}_f \frac{1}{N} \sum_{n=1}^N \ell(f(\mathbf{x}_n), \mathbf{y}_n) \text{ takes the form } f(\mathbf{x}) = \sum_{n=1}^N w_n \kappa(\mathbf{x}_n, \mathbf{x}).$$

$\Rightarrow \mathbf{x}_n$ are feature vectors, and w_n is a scalar weight.

$\Rightarrow f$ is a kernel expansion over training set

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

$$f^* = \operatorname{argmin}_f \frac{1}{N} \sum_{n=1}^N \ell(f(\mathbf{x}_n), \mathbf{y}_n) \text{ takes the form } f(\mathbf{x}) = \sum_{n=1}^N w_n \kappa(\mathbf{x}_n, \mathbf{x}).$$

$\Rightarrow \mathbf{x}_n$ are feature vectors, and w_n is a scalar weight.

$\Rightarrow f$ is a kernel expansion over training set

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

$$f^* = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^N} \frac{1}{N} \sum_{n=1}^N \ell\left(\sum_{m=1}^N w_m \kappa(\mathbf{x}_m, \mathbf{x}_n), \mathbf{y}_n\right) + \frac{\lambda}{2} \left\| \sum_{n,m=1}^N w_n w_m \kappa(\mathbf{x}_m, \mathbf{x}_n) \right\|_{\mathcal{H}}^2$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^N} \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{w}^T \kappa_{\mathbf{x}}(\mathbf{x}_n), \mathbf{y}_n) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{K}_{\mathbf{x}, \mathbf{x}} \mathbf{w}$$

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

$$f^* = \operatorname{argmin}_f \frac{1}{N} \sum_{n=1}^N \ell(f(\mathbf{x}_n), \mathbf{y}_n) \text{ takes the form } f(\mathbf{x}) = \sum_{n=1}^N w_n \kappa(\mathbf{x}_n, \mathbf{x}).$$

$\Rightarrow \mathbf{x}_n$ are feature vectors, and w_n is a scalar weight.

$\Rightarrow f$ is a kernel expansion over training set

- ▶ Example: kernel logistic regression $\mathbb{P}(y = 0 \mid \mathbf{x}) = \frac{\exp\{f(\mathbf{x})\}}{1 + \exp\{f(\mathbf{x})\}}$.

$$= \operatorname{argmin}_{f \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \left[\log(1 + \exp\{f(\mathbf{x}_n)\}) - \mathbb{I}(y_n = 1) - f(\mathbf{x}_n) \mathbb{I}(y_n = 0) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 \right]$$

$$= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^N} \frac{1}{N} \sum_{n=1}^N \left[\log(1 + \exp\{\mathbf{w}^T \kappa_{\mathbf{x}}(\mathbf{x}_n)\}) - \mathbb{I}(y_n = 1) - \mathbf{w}^T \kappa_{\mathbf{x}}(\mathbf{x}_n) \mathbb{I}(y_n = 0) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{K}_{\mathbf{x}, \mathbf{x}} \mathbf{w} \right]$$

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

$$f^* = \operatorname{argmin}_f \frac{1}{N} \sum_{n=1}^N \ell(f(\mathbf{x}_n), \mathbf{y}_n) \text{ takes the form } f(\mathbf{x}) = \sum_{n=1}^N w_n \kappa(\mathbf{x}_n, \mathbf{x}) .$$

⇒ \mathbf{x}_n are feature vectors, and w_n is a scalar weight.

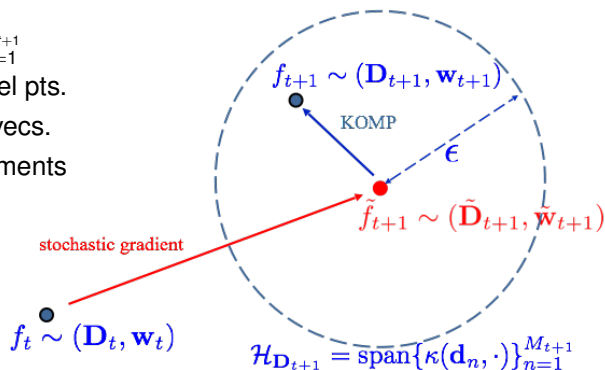
⇒ f is a kernel expansion over training set

- ▶ Unfortunately, as sample size $N \rightarrow \infty$
 - ⇒ kernel matrix $[\mathbf{K}_{\mathbf{X},\mathbf{X}}]_{m,n} := \kappa(\mathbf{x}_m, \mathbf{x}_n)$ becomes infinite too!
 - ⇒ $\mathbf{X} = [\mathbf{x}_1; \mathbf{x}_2; \dots]$ ⇒ kernel dictionary
 - ⇒ $\kappa_{\mathbf{X}}(\cdot) = [\kappa(\mathbf{x}_1, \cdot) \dots \kappa(\mathbf{x}_N, \cdot)]^T$ ⇒ empirical kernel map
 - ⇒ model order $M(= N) \rightarrow \infty$ ⇒ number of dictionary columns
- ▶ **Storage/representation issue** ⇒ “the curse of kernelization”

- ▶ Fix approx. error ϵ_t
- ▶ Define subspace

$$\mathcal{H}_{\mathbf{D}_{t+1}} = \text{span}\{\kappa(\mathbf{d}_n, \cdot)\}_{n=1}^{M_{t+1}}$$
- ▶ $\{\mathbf{d}_n\} \subset \{\mathbf{x}_u\}_{u \leq t} \Rightarrow$ model pts.
 \Rightarrow subset of past feat. vecs.
- ▶ Remove kernel dict. elements
- ▶ Stopping criterion:
 $\|\tilde{f}_{t+1} - f_{t+1}\|_{\mathcal{H}} \leq \epsilon_t$
- ▶ New model order:
 $M_{t+1} \leq M_t + 1$

Hilbert Space



- ▶ Define un-projected/unsparsified iterate at step $t + 1$

$$\tilde{f}_{t+1} = (1 - \eta_t \lambda) f_t - \eta_t \nabla_f \ell(f_t; \mathbf{x}_t, \mathbf{y}_t).$$

⇒ parameterized by dictionary and coefficients

$$\tilde{\mathbf{D}}_{t+1} = [\mathbf{D}_t, \mathbf{x}_t], \quad \tilde{\mathbf{w}}_{t+1} = [(1 - \eta_t \lambda) \mathbf{w}_t, -\eta_t \ell'(f_t(\mathbf{x}_t), y_t)].$$

- ▶ Our method: $(f_{t+1}, \mathbf{D}_{t+1}, \mathbf{w}_{t+1}) = \mathbf{KOMP}(\tilde{f}_{t+1}, \tilde{\mathbf{D}}_{t+1}, \tilde{\mathbf{w}}_{t+1}, \epsilon_t)$
- ▶ This amounts to a certain orthogonal subspace projection

$$\begin{aligned} f_{t+1} &= \operatorname{argmin}_{f \in \mathcal{H}_{\mathbf{D}_{t+1}}} \left\| f - \left((1 - \eta_t \lambda) f_t - \eta_t \nabla_f \ell(f_t(\mathbf{x}_t), y_t) \right) \right\|_{\mathcal{H}}^2 \\ &:= \mathcal{P}_{\mathcal{H}_{\mathbf{D}_{t+1}}} \left[(1 - \eta_t \lambda) f_t - \eta_t \nabla_f \ell(f_t(\mathbf{x}_t), y_t) \right]. \end{aligned}$$

- ▶ Recall the Hilbert subspace $\mathcal{H}_{\mathbf{D}_{t+1}} = \operatorname{span}\{\kappa(\mathbf{d}_n, \cdot)\}_{n=1}^{M_{t+1}}$
 ⇒ \mathbf{d}_n are model points ⇒ subset of past feature vectors $\{\mathbf{x}_u\}_{u \leq t}$

Require: function \tilde{f} defined by dict. $\tilde{\mathbf{D}} \in \mathbb{R}^{p \times \tilde{M}}$, coeffs. $\tilde{\mathbf{w}} \in \mathbb{R}^{\tilde{M}}$, approx. budget $\epsilon_t > 0$

Initialize $f = \tilde{f}$, dict. $\mathbf{D} = \tilde{\mathbf{D}}$ (indices \mathcal{I}), model order $M = \tilde{M}$, coeffs. $\mathbf{w} = \tilde{\mathbf{w}}$.

while candidate dictionary is non-empty $\mathcal{I} \neq \emptyset$ **do**

for $j = 1, \dots, \tilde{M}$ **do**

Find minimal approx. error with dict. element \mathbf{d}_j removed

$$\gamma_j = \min_{\mathbf{w}_{\mathcal{I} \setminus \{j\}} \in \mathbb{R}^{M-1}} \|\tilde{f}(\cdot) - \sum_{k \in \mathcal{I} \setminus \{j\}} \mathbf{w}_k \kappa(\mathbf{d}_k, \cdot)\|_{\mathcal{H}}.$$

end for

Find dictionary index minimizing approx. error: $j^* = \operatorname{argmin}_{j \in \mathcal{I}} \gamma_j$

if minimal approximation error exceeds threshold $\gamma_{j^*} > \epsilon_t$

stop

else

Prune dictionary $\mathbf{D} \leftarrow \mathbf{D}_{\mathcal{I} \setminus \{j^*\}}$, revise $\mathcal{I} \leftarrow \mathcal{I} \setminus \{j^*\}$

Revise model order $M \leftarrow M - 1$; compute \mathbf{w} defined by \mathbf{D}

$$\mathbf{w} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^M} \|\tilde{f}(\cdot) - \mathbf{w}^T \kappa_{\mathbf{D}}(\cdot)\|_{\mathcal{H}}$$

end

end while

return $f, \mathbf{D}, \mathbf{w}$ of model order $M \leq \tilde{M}$ such that $\|f - \tilde{f}\|_{\mathcal{H}} \leq \epsilon_t$

- ▶ Given feature vectors \mathbf{x} with labels $y \in \{1, \dots, C\}$
 - ⇒ Classifier ⇒ optimize a geometric criterion of separation
- ▶ Kernel multi-class SVM ⇒ for class c , function $f_c : \mathcal{X} \rightarrow \mathbb{R}$.
 - ⇒ Classify \mathbf{x} ⇒ max class-conditional prob. $y = \max_{y'} f_{y'}(\mathbf{x})$.
- ▶ Define vectorized function $\mathbf{f} = [f_1, \dots, f_C] \in \mathcal{H}^C$.
 - ⇒ Want to minimize λ -regularized multi-class hinge loss

$$\mathbf{f}^* = \operatorname{argmin}_{\mathbf{f}} \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{f}(\mathbf{x}_n), y_n) + \lambda \sum_{c'=1}^C \|f_{c'}\|_{\mathcal{H}}^2,$$

$$\Rightarrow \ell(\mathbf{f}(\mathbf{x}), y) = \max(0, 1 + f_r(\mathbf{x}) - f_y(\mathbf{x})), \quad r = \operatorname{argmax}_{c' \neq y} f_{c'}(\mathbf{x}).$$

- ▶ Parameter selection:
 - ⇒ Gaussian kernel with bandwidth $\tilde{\sigma}^2 = 0.6$
 - ⇒ regularizer $\lambda = 10^{-6}$, constant learning rate $\eta = 6.0$
 - ⇒ approximation budget $\epsilon = K\eta^{3/2}$
 - ⇒ parsimony constant $K = 0.04$
 - ⇒ Initialize kernel classifier as null, i.e., $f_0 = 0$