

Projected Stochastic Primal-Dual Method for Constrained Online Learning with Kernels

Kaiqing Zhang[§], Tame Başar[§], Hao Zhu[†], Alec Koppel^{*}
[§]University of Illinois Urbana-Champaign [†]UT Austin
^{*}CISD, U.S. Army Research Laboratory

Large-Scale Distributed Optimization and Computation on Graphs
IEEE Conference on Decision and Control
Miami, FL, Dec. 18, 2018

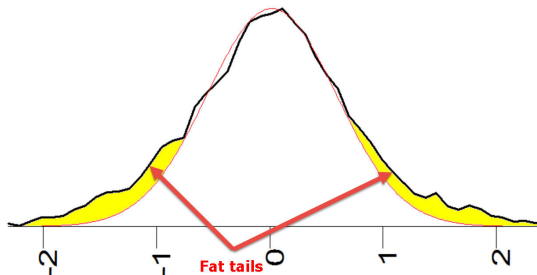
- ▶ Want to find $f^* \in \mathcal{H}$ to minimize some expected cost $R(f)$

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

such that $\mathbf{G}(f) \leq \mathbf{0}$

- ⇒ Loss $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ penalize deviations between $f(\mathbf{x})$, \mathbf{y}
- ⇒ interpret $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^p$ as features/state variables
- ⇒ $y \in \mathcal{Y} \Rightarrow$ targets, e.g., reference trajectory or binary labels
- ⇒ expected risk $L(f) := \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\ell(f(\mathbf{x}), \mathbf{y})]$
- ▶ Examples:
 - ⇒ learning with risk constraints
 - ⇒ trajectory planning with obstacle avoidance
 - ⇒ wireless utility maximization with quality of service guarantees
 - ⇒ wireless beamforming with robustness constraints

- ▶ Learning with nonlinear statistical models and risk constraints
 - ⇒ when distribution $\mathbb{P}(\mathbf{x}, \mathbf{y})$ has heavy tails
 - ⇒ then learning $f(\mathbf{x})$ by minimizing **average** loss will “**overfit**”



- ▶ Impose risk constraint, such as CVaR (Rockafeller '2000)

$$\begin{aligned}
 G(f) &= \text{CVaR}_\alpha(f) - \gamma \\
 &= \min_{z \in \mathbb{R}} \left\{ z + \frac{1}{1-\alpha} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \{ [\ell(f(\mathbf{x}), \mathbf{y}) - z]_+ \} \right\} - \gamma
 \end{aligned}$$

- ▶ Trajectory planning with obstacle avoidance
⇒ obstacles define region of state space to be avoided



- ▶ Let $g(f(x)) > 0$ represent safe area in \mathbb{R}^p
⇒ upper bound $\mathbb{P}(g(f(x)) > 0) \leq \gamma$ for a given γ
- ▶ Convexified chance constraint:

$$\inf_{\lambda > 0} [\Psi(f, \lambda) - \lambda\gamma] \leq 0,$$

here $\Psi(f, \lambda) = \lambda \mathbb{E}_x[\phi(\lambda^{-1} g(f(x)))] \Rightarrow \phi(\cdot)$: MGF of $\mathbb{P}(x)$

- ▶ Want to find $f^* \in \mathcal{H}$ to minimize some expected cost $R(f)$

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

such that $\mathbf{G}(f) \leq \mathbf{0}$

- ▶ Classically address using calculus of variations (Hamilton 1800s)
 - ⇒ w/o special structure, can't solve Euler-Lagrange equations
 - ⇒ with special structure on distribution ⇒ variational inference
- ▶ Without hypotheses on distribution, “learning” approaches
 - ⇒ parameterize f , then estimate parameters via samples of \mathbf{x}, y
- ▶ Choose parameterization s.t. solution close to original problem
 - ⇒ “universal parameterizations:” Bayesian/**nonparametric**/DNN

- ▶ Want to find $f^* \in \mathcal{H}$ to minimize some expected cost $R(f)$

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

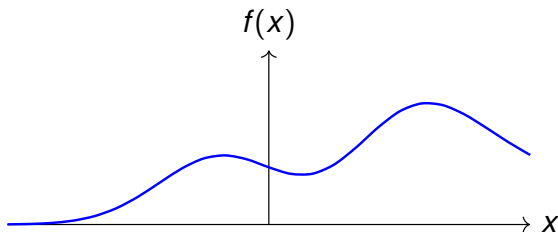
such that $\mathbf{G}(f) \leq \mathbf{0}$

- ▶ We adopt a **nonparametric** parameterization of f
 - ⇒ using a reproducing kernel Hilbert space (RKHS)
 - ⇒ motivated by the fact this param. **preserves convexity**
 - ⇒ therefore Lagrange duality applies
- ▶ Extend Representer Theorem to constrained settings
 - ⇒ for certain constraints, using augmented Lagrangian
- ▶ Propose a projected stochastic primal-dual method
 - ⇒ custom projection trades off convergence and complexity
 - ⇒ generalize existing convergence rates from vector case

- ▶ 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 Will allow us to compute derivatives
- ▶ Kernel examples:

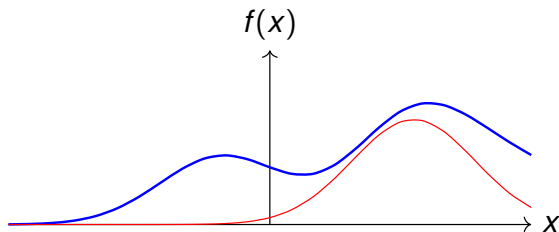
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

- ▶ Kernel examples:

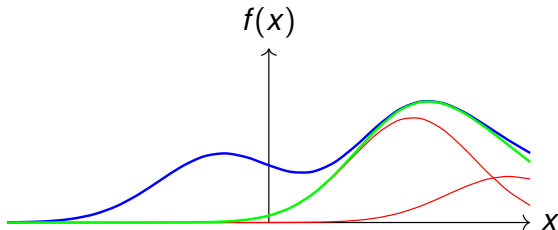
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

- ▶ Kernel examples:

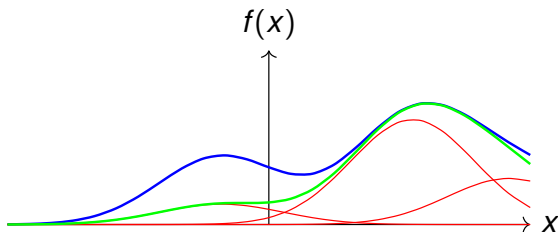
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

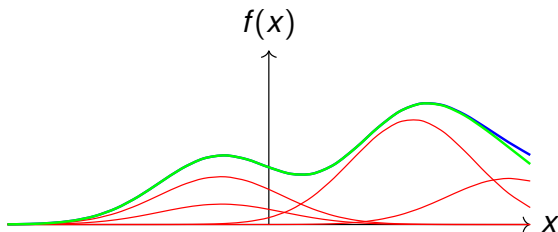
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

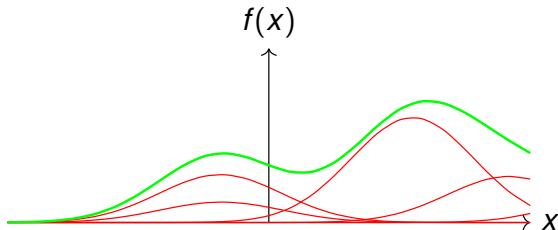
$$\Rightarrow \text{Gaussian/RBF } \kappa(\mathbf{x}, \mathbf{x}') = \exp\left\{-\frac{\|\mathbf{x}-\mathbf{x}'\|_2^2}{2c^2}\right\}$$

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

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

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

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



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

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

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

- ▶ Consider empirical risk minimization case: sample size $N < \infty$
- ▶ **Classic 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
 - ⇒ dates to Riesz & Weiner, to ML by Scholkopf/Smola
- ▶ Does not apply to constrained settings . . .

- ▶ Consider Lagrangian of constrained problem

$$\mathcal{L}^o(f, \boldsymbol{\mu}) = L(f) + \boldsymbol{\mu}^\top \mathbf{G}(f) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2,$$

Theorem

Suppose constraint takes form $\mathbf{G}(f) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\mathbf{g}(f(\mathbf{x}), \mathbf{y})]$

Define saddle pt. prob: $(\check{f}^*, \check{\boldsymbol{\mu}}^*) = \arg \max_{\boldsymbol{\mu} \in \mathbb{R}_+^m} \min_{f \in \mathcal{H}} \mathcal{L}^o(f, \boldsymbol{\mu}; \mathcal{S})$,

\Rightarrow Consider sample avg. approx. of Lagrangian w/ $\mathcal{S}_N = \{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$:

$$\mathcal{L}^o(f, \boldsymbol{\mu}; \mathcal{S}) := \frac{1}{N} \sum_{n=1}^N \left[\ell(f(\mathbf{x}_n), \mathbf{y}_n) + \sum_{j=1}^m \mu_j g_j(f(\mathbf{x}_n), \mathbf{y}_n) \right] + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2.$$

Then primal minimizer of takes form $f(\mathbf{x}) = \sum_{n=1}^N w_n \kappa(\mathbf{x}_n, \mathbf{x})$.

- Formulate augmented Lagrangian of constrained prob:

$$\mathcal{L}(f, \boldsymbol{\mu}) = L(f) + \boldsymbol{\mu}^\top \mathbf{G}(f) + \frac{\lambda}{2} \|f\|_{\mathcal{H}}^2 - \frac{\delta\eta}{2} \|\boldsymbol{\mu}\|^2.$$

⇒ $\boldsymbol{\mu}$ is Lagrange multiplier, δ is a regularization parameter

⇒ Define stoch. approx. based on sample $(\mathbf{x}_t, \mathbf{y}_t)$ as $\hat{\mathcal{L}}_t(f, \boldsymbol{\mu})$

- Set $\tilde{\ell}(f(\mathbf{x}), \mathbf{y}, \boldsymbol{\mu}) = \ell(f(\mathbf{x}), \mathbf{y}) + \sum_{j=1}^m \mu_j g_j(f(\mathbf{x}), \mathbf{y})$.

⇒ Then let's compute primal stochastic gradient:

$$\begin{aligned} \nabla_f \tilde{\ell}(f(\mathbf{x}_t), \mathbf{y}_t, \boldsymbol{\mu})(\cdot) &= \frac{\partial \tilde{\ell}(f(\mathbf{x}_t), \mathbf{y}_t, \boldsymbol{\mu})}{\partial f(\mathbf{x}_t)} \frac{\partial f(\mathbf{x}_t)}{\partial f}(\cdot) \\ &= \tilde{\ell}'(f(\mathbf{x}_t), \mathbf{y}_t, \boldsymbol{\mu}) \kappa(\mathbf{x}_t, \cdot) \end{aligned}$$

- ▶ Stochastic primal/dual descent/ascent steps:

$$\left\{ \begin{array}{l} f_{t+1} = (1 - \eta\lambda)f_t - \eta \left[\ell'(f_t(\mathbf{x}_t), \mathbf{y}_t) + \sum_{j=1}^m \mu_j g'_j(f_t(\mathbf{x}_t), \mathbf{y}_t) \right] \kappa(\mathbf{x}_t, \cdot), \\ \boldsymbol{\mu}_{t+1} = [(1 - \eta^2\delta)\boldsymbol{\mu}_t + \eta \mathbf{g}(f_t(\mathbf{x}_t), \mathbf{y}_t)]_+, \end{array} \right.$$

- ▶ Via induction, can show $f_t(\mathbf{x}) = \sum_{t=1}^{t-1} \mathbf{w}_t \kappa(\mathbf{x}_t, \mathbf{x}) = \mathbf{w}_t^\top \boldsymbol{\kappa}_{\mathbf{x}_t}(\mathbf{x})$
 \Rightarrow hence f_t is parameterized by a growing matrix, weight vec.:

$$\mathbf{X}_{t+1} = [\mathbf{X}_t, \mathbf{x}_t],$$

$$\mathbf{w}_{t+1} = \left[(1 - \eta\lambda)\mathbf{w}_t, -\eta \ell'(f_t(\mathbf{x}_t), \mathbf{y}_t) - \eta \sum_{j=1}^m \mu_j g'_j(f_t(\mathbf{x}_t), \mathbf{y}_t) \right].$$

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

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

⇒ parameterized by dictionary and coefficients

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

- ▶ 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

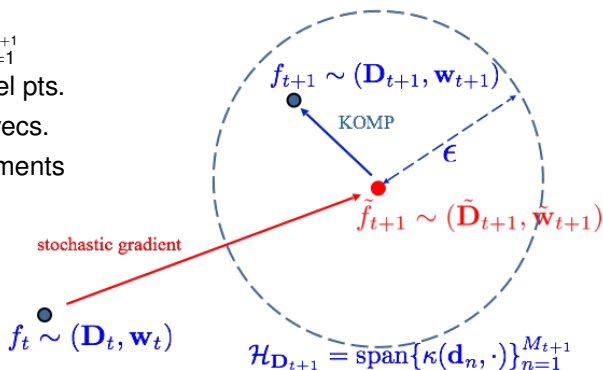
$$\begin{aligned} f_{t+1} &= \operatorname{argmin}_{f \in \mathcal{H}_{\mathbf{D}_{t+1}}} \left\| f - \left((1 - \eta_t \lambda) f_t - \eta_t \nabla_f \tilde{\ell}(f_t; \mathbf{x}_t, \mathbf{y}_t; \boldsymbol{\mu}_t) \right) \right\|_{\mathcal{H}}^2 \\ &:= \mathcal{P}_{\mathcal{H}_{\mathbf{D}_{t+1}}}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

- ▶ where we define 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 $\{\mathbf{x}_u\}_{u \leq t}$

- ▶ 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



- ▶ The feature space $\mathcal{X} \subset \mathbb{R}^p$, target domain $\mathcal{Y} \subset \mathbb{R}$ are compact, and kernel is bounded $\sup_{\mathbf{x} \in \mathcal{X}} \sqrt{\kappa(\mathbf{x}, \mathbf{x})} = X < \infty$
- ▶ Instantaneous loss $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is uniformly C_1 -Lipschitz continuous, and the constraints $g_i : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is C_2 -Lipschitz, for all $z \in \mathbb{R}$ for a fixed $y \in \mathcal{Y}$.
- ▶ The primal loss $\ell(f(\mathbf{x}), y)$ is convex and differentiable with respect to its scalar argument $f(\mathbf{x})$, as are the constraints $g_i(f(\mathbf{x}), y)$ on \mathbb{R} for all $\mathbf{x} \in \mathcal{X}$ and $y \in \mathcal{Y}$.
- ▶ There exists a strictly feasible pt: some $f \in \mathcal{H}$ satisfies $\mathbf{G}(f) < \mathbf{0}$.
- ▶ The output f_{t+1} of the KOMP update has Hilbert norm bounded by $R_B < \infty$, and the optimal f^* lies in the ball \mathcal{B} with radius R_B

Theorem

Denote the projected stochastic primal-dual sequence as (f_t, μ_t) . After T iterations with a constant step-size selected as $\eta = 1/\sqrt{T}$ and the approximation budget $\epsilon_t = \epsilon = P\eta^2$, where $P > 0$ is a fixed constant, we have

$$\sum_{t=1}^T \mathbb{E}[R(f_t) - R(f^*)] = \mathcal{O}(\sqrt{T})$$

Moreover, the time aggregation of the expected constraint violation of the algorithm grows sublinearly in T as

$$\sum_{j=1}^m \mathbb{E} \left[\sum_{t=1}^T G_j(f_t) \right]_+ \leq \mathcal{O}(T^{3/4}).$$

Corollary

For $\bar{f}_T = \sum_{t=1}^T f_t / T$ as the functional formed by averaging the primal iterates f_t over time $t = 1, \dots, T$, its objective function satisfies

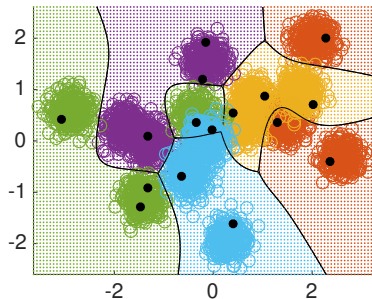
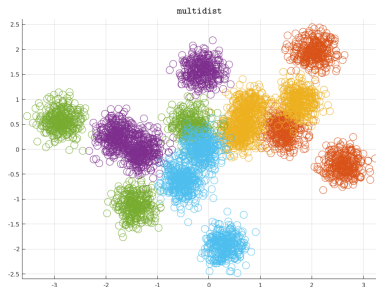
$$\mathbb{E}[R(\bar{f}_T) - R(f^*)] \leq \mathcal{O}(1/\sqrt{T}).$$

In addition, the constraint violation evaluated at \bar{f}_T satisfies

$$\sum_{j=1}^m \mathbb{E} \left[(G_j(\bar{f}_T)) \right]_+ \leq \mathcal{O}(T^{-1/4}).$$

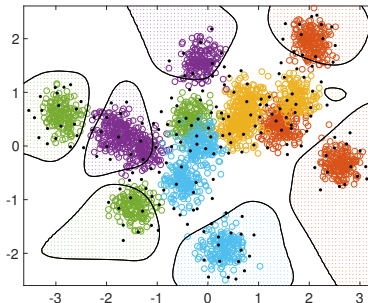
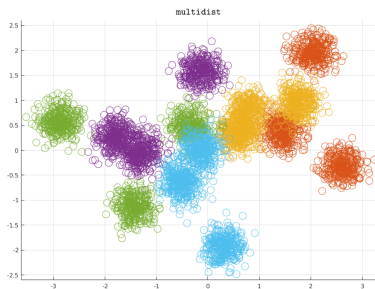
- ▶ Comparable to existing results for vector-valued case
 - ⇒ both in terms of primal sub-optimality and constraint violation
- ▶ Using a constant step-size and compression budget
 - ⇒ yields fact that model complexity of RKHS function is finite
- ▶ Complexity depends on
 - ⇒ parsimony constant P
 - ⇒ kernel choice
 - ⇒ data domain radius X

- ▶ 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

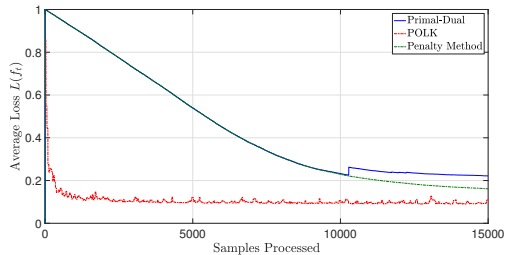


- ▶ Grid colors ⇒ decision, bold black dots ⇒ kernel dict. elements
- ▶ ~ 96% accuracy

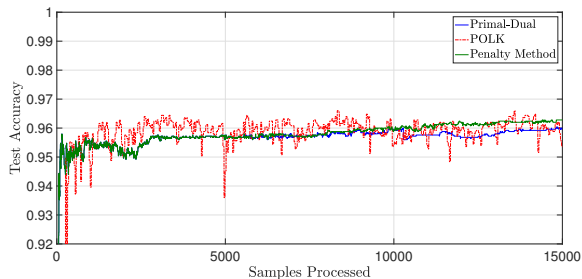
- ▶ Case where training examples for a fixed class
⇒ drawn from a distinct Gaussian mixture
- ▶ 3 Gaussians per mixture, $C = 5$ classes total for this experiment
⇒ 15 total Gaussians generate data



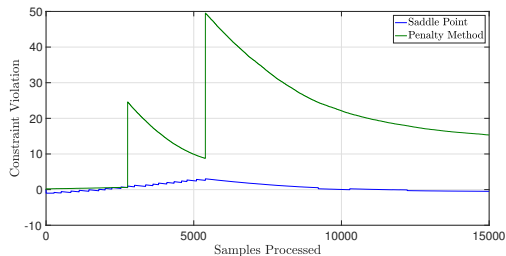
- ▶ Grid colors ⇒ decision, bold black dots ⇒ kernel dict. elements
- ▶ risk constraint prevents confidence in areas of class overlap



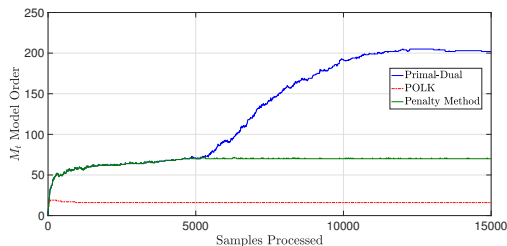
- ▶ Comparison with unconstrained (POLK) and penalty approach
⇒ Constrained optimizer has large objective than unconstrained



- Comparable in terms of test accuracy



- Primal-dual most effective for maintaining feasibility



- ▶ Constrained optimizer higher complexity than unconstrained
 - ⇒ control decision uncertainty ⇒ higher order data moments
 - ⇒ define more complicated subspace than mean loss

- ▶ Focus on stochastic nonlinear interpolation with constraints
 - ⇒ MPC with obstacle avoidance, risk-aware learning, etc.
 - ⇒ parameterized with RKHS: “universal,” preserves convexity
- ▶ Extended Representer Theorem to constrained settings
 - ⇒ via use of empirical Lagrangian
- ▶ Proposed stochastic primal-dual method to solve it
 - ⇒ operates in parallel with subspace projection scheme
- ▶ Generalized convergence results for vector-valued case
 - ⇒ sub-optimality $\mathcal{O}(\sqrt{T})$; constraint violation as $\mathcal{O}(T^{3/4})$
- ▶ Online kernel multi-class SVM example
 - ⇒ demonstrates effect of incorporating risk into decisions
 - ⇒ confidence regions repelled by class overlap