

Efficient learning with Nyström projections

Lorenzo Rosasco

MaLGa, Università degli Studi di Genova, MIT, IIT

Joint with:

Daniele Calandriello, Raffaello Camoriano, Luigi Carratino, Giacomo Meanti (MaLGa),
Alessandro Rudi (INRIA Paris), Bharath Sriperumbudur, Nick Sterge (UPenn),
Alessandro Lazaric (Facebook), Michal Valko (DeepMind)

Data



Computations



The quest for **provably efficient** ML algorithms

Outline

Data size matters

$$\underbrace{\hat{X}}_{n \times d}$$

In many modern applications, space is the real constraint.

Think n, d large!

Dimensionality reduction

$$\underbrace{\hat{X}_M}_{n \times M} = \underbrace{\hat{X}}_{n \times d} \underbrace{S}_{d \times M}$$

Dimensionality reduction

$$\underbrace{\hat{X}_M}_{n \times M} = \underbrace{\hat{X}}_{n \times d} \underbrace{S}_{d \times M}$$

Classic example (not efficient): PCA

$$\hat{X} = \hat{U} \hat{\Lambda} \hat{V}^T \Rightarrow S = \hat{V}_M.$$

Dimensionality reduction

$$\underbrace{\hat{X}_M}_{n \times M} = \underbrace{\hat{X}}_{n \times d} \underbrace{S}_{d \times M}$$

Classic example (not efficient): PCA

$$\hat{X} = \hat{U} \hat{\Lambda} \hat{V}^T \Rightarrow S = \hat{V}_M.$$

Other example (efficient): random sketches

$$S_{ij} \sim \mathcal{N}(0, 1).$$

Nyström projections

$$\hat{X}_M = \hat{X} \bar{X}_M^T$$

Random subsampling¹ (efficient)

$$\bar{X}_M = \{\bar{x}_1, \dots, \bar{x}_M\} \subset \{x_1, \dots, x_n\} = \hat{X}.$$

Nyström projections

$$\hat{X}_M = \hat{X} \bar{X}_M^T$$

Random subsampling¹ (efficient)

$$\bar{X}_M = \{\bar{x}_1, \dots, \bar{x}_M\} \subset \{x_1, \dots, x_n\} = \hat{X}.$$

Computing \hat{X}_M is efficient!

[Williams, Seeger, Smola, Schölkopf, Bach, Muscoso, Clarkson, Mahoney, Woodruff, Avron, Drineas, Tropp, ...]

Nyström projections illustrated: least squares

From

$$\min_{w \in \mathbb{R}^d} \|\hat{X}w - \hat{y}\|^2, \quad \hat{X} \in \mathbb{R}^{n,d}$$

Nyström projections illustrated: least squares

From

$$\min_{w \in \mathbb{R}^d} \|\hat{X}w - \hat{y}\|^2, \quad \hat{X} \in \mathbb{R}^{n,d}$$

to

$$\min_{c \in \mathbb{R}^M} \|\hat{X}_M c - \hat{y}\|^2, \quad \hat{X}_M \in \mathbb{R}^{n,M}$$

Nyström projections illustrated: least squares

From

$$\min_{w \in \mathbb{R}^d} \|\hat{X}w - \hat{y}\|^2, \quad \hat{X} \in \mathbb{R}^{n,d}$$

to

$$\min_{c \in \mathbb{R}^M} \|\hat{X}_M c - \hat{y}\|^2, \quad \hat{X}_M \in \mathbb{R}^{n,M}$$

The latter problem is equivalent to

$$\min_{\substack{w = \hat{X}_M^\top c, \\ c \in \mathbb{R}^M}} \|\hat{X}w - y\|^2,$$

that is least squares projected on a random subspace.

[Engl, Hanke, Neubauer '96]

Nyström least squares: computations

(think n huge d ginormous)

From

$$w = \hat{X}^\top c, \quad c = (\underbrace{\hat{X}\hat{X}^\top}_{\hat{K} \in \mathbb{R}^{n,n}})^{-1}\hat{y} \in \mathbb{R}^n$$

Nyström least squares: computations

(think n huge d ginormous)

From

$$w = \hat{X}^\top c, \quad c = (\underbrace{\hat{X}\hat{X}^\top}_{\hat{K} \in \mathbb{R}^{n,n}})^{-1}\hat{y} \in \mathbb{R}^n$$

to

$$c = (\hat{X}_M^\top \underbrace{\hat{X}_M}_{\hat{K}_{n,M} \in \mathbb{R}^{n,M}})^{-1}\hat{X}_M^\top \hat{y} \in \mathbb{R}^M$$

Nyström least squares: computations

(think n huge d ginormous)

From

$$w = \hat{X}^\top c, \quad c = (\underbrace{\hat{X}\hat{X}^\top}_{\hat{K} \in \mathbb{R}^{n,n}})^{-1}\hat{y} \in \mathbb{R}^n$$

to

$$c = (\hat{X}_M^\top \underbrace{\hat{X}_M}_{\hat{K}_{n,M} \in \mathbb{R}^{n,M}})^{-1}\hat{X}_M^\top \hat{y} \in \mathbb{R}^M$$

From $O(n^2d + n^3)$ time/ $O(nd + n^2)$ space to $O(nM^2 + M^3)$ time/ $O(nM + M^2)$ space.

This matters for kernel methods

$$x^\top x' \quad \mapsto \quad k(x, x'), \quad \text{e.g.} \quad k(x, x') = e^{-\|x-x'\|^2\gamma}$$

This matters for kernel methods

$$x^\top x' \quad \mapsto \quad k(x, x'), \quad \text{e.g.} \quad k(x, x') = e^{-\|x-x'\|^2\gamma}$$

$$x^\top w = x^\top \hat{X}^\top c \quad \mapsto \quad f(x) = \sum_{i=1}^n k(x, x_i) c_i$$

This matters for kernel methods

$$x^\top x' \quad \mapsto \quad k(x, x'), \quad \text{e.g.} \quad k(x, x') = e^{-\|x-x'\|^2\gamma}$$

$$x^\top w = x^\top \hat{X}^\top c \quad \mapsto \quad f(x) = \sum_{i=1}^n k(x, x_i) c_i$$

$$\hat{X}w = \hat{X}\hat{X}^\top c = \hat{y} \quad \mapsto \quad \hat{K}c = \hat{y}$$

$$\hat{K}_{i,j} = k(x_i, x_j)$$

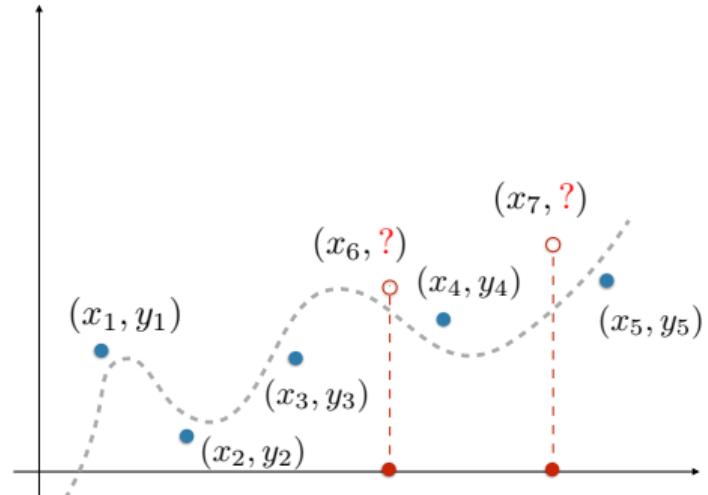
This matters for kernel methods

$$x^\top x' \quad \mapsto \quad k(x, x'), \quad \text{e.g.} \quad k(x, x') = e^{-\|x-x'\|^2\gamma}$$

$$x^\top w = x^\top \hat{X}^\top c \quad \mapsto \quad f(x) = \sum_{i=1}^n k(x, x_i) c_i$$

$$\hat{X}w = \hat{X}\hat{X}^\top c = \hat{y} \quad \mapsto \quad \hat{K}c = \hat{y}$$

$$\hat{K}_{i,j} = k(x_i, x_j)$$

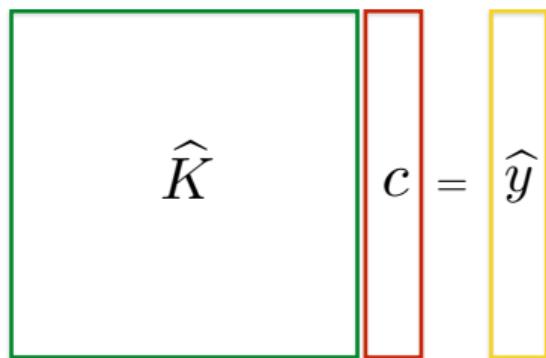


Nyström projections with kernels, aka column subsampling

$$\hat{X}w = \hat{X}\hat{X}^\top c = \hat{y}$$



$$f(x) = \sum_{i=1}^n k(x, x_i) c_i \quad \boxed{\hat{K}c = \hat{y}}$$



Nyström projections with kernels, aka column subsampling

$$\hat{X}_W = \hat{X}\hat{X}^\top c = \hat{y}$$

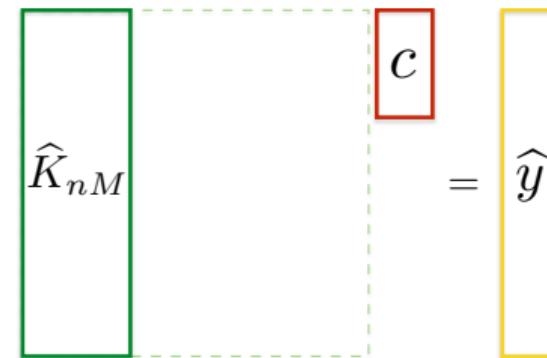
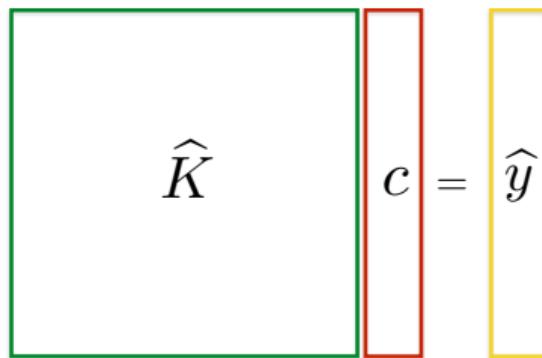


$$f(x) = \sum_{i=1}^n k(x, x_i) c_i \quad \boxed{\hat{K}c = \hat{y}}$$

$$\hat{X}_M c = \hat{y}$$



$$f(x) = \sum_{i=1}^n k(x, \bar{x}_i) c_i \quad \boxed{\hat{K}_{nM}c = \hat{y}}$$



From $O(n^3)$ time/ $O(n^2)$ space to $O(nM^2 + M^3)$ time/ $O(nM)$ space.

[Williams, Seeger, Smola Scholkopf, ... Mahoney, Drineas, ...]

Why Nyström?

Nyström approximation for integral equations

For all x

$$\int k(x, x') c(x') dx' = y(x) \quad \mapsto \quad \sum_{j=1}^M k(x, \bar{x}_j) c(\bar{x}_j) = y(x).$$

From operators to matrices

For all $i = 1, \dots, n$

$$\sum_{j=1}^n k(x_i, x_j) c_j = y_i \quad \mapsto \quad \sum_{j=1}^M k(x_i, \bar{x}_j) c_j = y_i.$$

[Kress '89]

Nyström approximation and subsampling

For all $i = 1, \dots, n$

$$\sum_{j=1}^n k(x_i, x_j) c_j = y_j \quad \mapsto \quad \sum_{j=1}^M k(x_i, \bar{x}_j) c_j = y_j.$$

[Williams, Seeger '00]

The above formulation highlights the connection to columns sampling,

$$\hat{K}c = \hat{y} \quad \mapsto \quad \hat{K}_{nM}c = \hat{y}.$$

So far

Nyström projections and connection to

- ▶ Sketching
- ▶ Projected least squares
- ▶ Column subsampling
- ▶ Nyström approximation

$$\hat{X}_M = \hat{X}\bar{X}_M^\top$$

Dimensionality reduction improves efficiency, but what about learning accuracy?

Outline

Supervised statistical Learning

Let $(x, y) \sim \rho$, $x \in X \subseteq \mathbb{R}^d$, $y \in Y \subseteq \mathbb{R}$

Solve

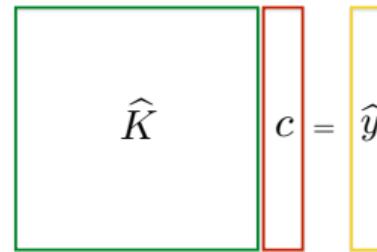
$$\min_{f \in \mathcal{H}} L(f), \quad L(f) = \mathbb{E}_{x,y} (y - f(x))^2$$

given $(x_i, y_i)_{i=1}^n \sim \rho^n$.

Kernel Ridge Regression (KRR)

aka Gaussian Process (GP) regression

$$\hat{f}_\lambda(x) = \sum_{i=1}^n k(x_i, x)c_i,$$
$$(\hat{K} + \lambda n I)c = \hat{y}$$



Theorem (Caponetto, De Vito '05)

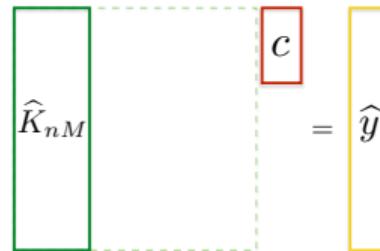
Let $\mathcal{H} = \text{span}\{k(x, \cdot) \mid x \in X\}$, if $\lambda = 1/\sqrt{n}$ then

$$\mathbb{E}L(\hat{f}_\lambda) - \min_{f \in \mathcal{H}} L(f) \lesssim \frac{1}{\sqrt{n}}$$

Nyström KRR

$$\hat{f}_{\lambda, M}(x) = \sum_{i=1}^M K(\tilde{x}_i, x) c_i$$

$$(\hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM}) c = \hat{K}_{nM}^\top \hat{y}$$



Theorem (Rudi, Camoriano, R. '15)

Let $(\tilde{x}_i)_{i=1}^M \subseteq (x_i)_{i=1}^n$ picked uniformly at random, if $\lambda = 1/\sqrt{n}$ and $M \geq \sqrt{n}$ then

$$\mathbb{E} L(\hat{f}_{\lambda, M}) - \min_{f \in \mathcal{H}} L(f) \lesssim \frac{1}{\sqrt{n}}$$

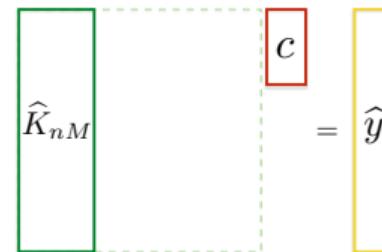
Iterative solvers and preconditioning

Consider an iterative solver, e.g. conjugate gradient (CG), on a **preconditioned** system

$$P^T (\hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM}) P \beta = P \hat{K}_{nM}^\top \hat{y}$$

...ideally

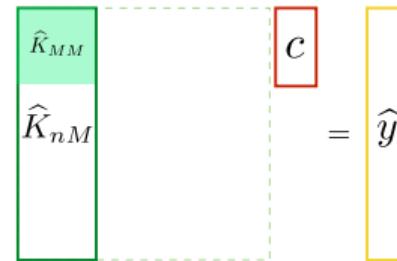
$$P P^T = (\hat{K}_{nM}^\top \hat{K}_{nM} + \lambda n \hat{K}_{MM})^{-1}$$



FALKON

$\widehat{f}_{\lambda, M, t}$ CG iteration with preconditioner

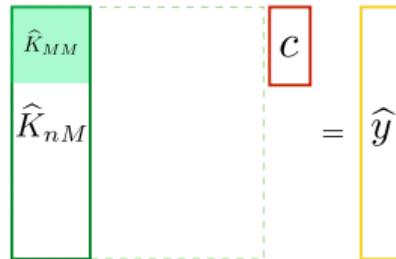
$$P P^\top = \left(\frac{n}{M} \widehat{K}_{MM}^2 + \lambda n \widehat{K}_{MM} \right)^{-1}$$



FALKON

$\widehat{f}_{\lambda, M, t}$ CG iteration with preconditioner

$$PP^\top = \left(\frac{n}{M} \widehat{\mathbf{K}}_{MM}^2 + \lambda n \widehat{\mathbf{K}}_{MM} \right)^{-1}$$



Theorem (Rudi, Carratino, Rosasco '17)

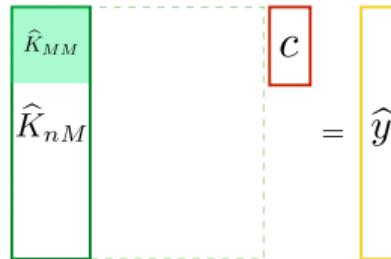
Let $(\tilde{x}_i)_{i=1}^M \subseteq (x_i)_{i=1}^n$ uniformly at random, then if $\lambda = 1/\sqrt{n}$, $M \geq \sqrt{n}$ and $t \geq \log(n)$

$$\mathbb{E} L(\widehat{f}_{\lambda, M, t}) - \min_{f \in \mathcal{H}} L(f) \lesssim \frac{1}{\sqrt{n}}$$

FALKON

$\hat{f}_{\lambda, M, t}$ CG iteration with preconditioner

$$PP^\top = \left(\frac{n}{M} \hat{K}_{MM}^2 + \lambda n \hat{K}_{MM} \right)^{-1}$$



Theorem (Rudi, Carratino, Rosasco '17)

Let $(\tilde{x}_i)_{i=1}^M \subseteq (x_i)_{i=1}^n$ uniformly at random, then if $\lambda = 1/\sqrt{n}$, $M \geq \sqrt{n}$ and $t \geq \log(n)$

$$\mathbb{E} L(\hat{f}_{\lambda, M, t}) - \min_{f \in \mathcal{H}} L(f) \lesssim \frac{1}{\sqrt{n}}$$

KRR: **Space** $O(n^2)$ / **Time** $O(n^3)$ **vs** **FALKON:** **Space** $O(n)$ / **Time** $O(n\sqrt{n}\log(n))$

Some experiments

	MillionSongs ($n \sim 10^6$)			YELP ($n \sim 10^6$)		TIMIT ($n \sim 10^6$)	
	MSE	Relative error	Time(s)	RMSE	Time(m)	c-err	Time(h)
FALKON	80.30	4.51×10^{-3}	55	0.833	20	32.3%	1.5
Prec. KRR	-	4.58×10^{-3}	289 [†]	-	-	-	-
Hierarchical	-	4.56×10^{-3}	293*	-	-	-	-
D&C	80.35	-	737*	-	-	-	-
Rand. Feat.	80.93	-	772*	-	-	-	-
Nyström	80.38	-	876*	-	-	-	-
ADMM R. F.	-	5.01×10^{-3}	958 [†]	-	-	-	-
BCD R. F.	-	-	-	0.949	42 [‡]	34.0%	1.7 [‡]
BCD Nyström	-	-	-	0.861	60 [‡]	33.7%	1.7 [‡]
KRR	-	4.55×10^{-3}	-	0.854	500 [‡]	33.5%	8.3 [‡]
EigenPro	-	-	-	-	-	32.6%	3.9 [‡]
Deep NN	-	-	-	-	-	32.4%	-
Sparse Kernels	-	-	-	-	-	30.9%	-
Ensemble	-	-	-	-	-	33.5%	-

Table: MillionSongs, YELP and TIMIT Datasets. Times obtained on: \ddagger = cluster of 128 EC2 r3.2xlarge machines, \dagger = cluster of 8 EC2 r3.8xlarge machines, \wr = single machine with two Intel Xeon E5-2620, one Nvidia GTX Titan X GPU and 128GB of RAM, $*$ = cluster with 512 GB of RAM and IBM POWER8 12-core processor, $*$ = unknown platform.

Some more experiments

	SUSY ($n \sim 10^6$)			HIGGS ($n \sim 10^7$)		IMAGENET ($n \sim 10^6$)	
	c-err	AUC	Time(m)	AUC	Time(h)	c-err	Time(h)
FALKON	19.6%	0.877	4	0.833	3	20.7%	4
EigenPro	19.8%	-	6 [‡]	-	-	-	-
Hierarchical	20.1%	-	40 [†]	-	-	-	-
Boosted Decision Tree	-	0.863	-	0.810	-	-	-
Neural Network	-	0.875	-	0.816	-	-	-
Deep Neural Network	-	0.879	4680 [‡]	0.885	78 [‡]	-	-
Inception-V4	-	-	-	-	-	20.0%	-

Table: Architectures: \dagger = cluster with IBM POWER8 12-core cpu, 512 GB RAM, \ddagger = single machine with two Intel Xeon E5-2620, one Nvidia GTX Titan X GPU, 128GB RAM, $\ddot{\dagger}$ = single machine.

Outline

Bandit Optimization

Given a set (of arms) $\mathcal{A} = \{x_1, \dots, x_A\} \subset \mathbb{R}^d$, let $f : \mathcal{A} \rightarrow \mathbb{R}$ unknown, $(\eta_t)_t$ random, and

$$x_* = \operatorname{argmax}_{x \in \mathcal{A}} f(x)$$

Bandit Optimization

Given a set (of arms) $\mathcal{A} = \{x_1, \dots, x_A\} \subset \mathbb{R}^d$, let $f : \mathcal{A} \rightarrow \mathbb{R}$ unknown, $(\eta_t)_t$ random, and

$$x_* = \operatorname{argmax}_{x \in \mathcal{A}} f(x)$$

For $t = 1, \dots, T$:

- (1) Estimate \hat{u}_t (ideally $\hat{u}_t \approx f$)

Bandit Optimization

Given a set (of arms) $\mathcal{A} = \{x_1, \dots, x_A\} \subset \mathbb{R}^d$, let $f : \mathcal{A} \rightarrow \mathbb{R}$ unknown, $(\eta_t)_t$ random, and

$$x_* = \operatorname{argmax}_{x \in \mathcal{A}} f(x)$$

For $t = 1, \dots, T$:

- (1) Estimate \hat{u}_t (ideally $\hat{u}_t \approx f$)
- (2) Select x_{t+1}

Bandit Optimization

Given a set (of arms) $\mathcal{A} = \{x_1, \dots, x_A\} \subset \mathbb{R}^d$, let $f : \mathcal{A} \rightarrow \mathbb{R}$ unknown, $(\eta_t)_t$ random, and

$$x_* = \operatorname{argmax}_{x \in \mathcal{A}} f(x)$$

For $t = 1, \dots, T$:

- (1) Estimate \hat{u}_t (ideally $\hat{u}_t \approx f$)
- (2) Select x_{t+1}
- (3) Receive noisy feedback $y_{t+1} = f(x_{t+1}) + \eta_{t+1}$

Bandit Optimization

Given a set (of arms) $\mathcal{A} = \{x_1, \dots, x_A\} \subset \mathbb{R}^d$, let $f : \mathcal{A} \rightarrow \mathbb{R}$ unknown, $(\eta_t)_t$ random, and

$$x_* = \operatorname{argmax}_{x \in \mathcal{A}} f(x)$$

For $t = 1, \dots, T$:

- (1) Estimate \hat{u}_t (ideally $\hat{u}_t \approx f$)
- (2) Select x_{t+1}
- (3) Receive noisy feedback $y_{t+1} = f(x_{t+1}) + \eta_{t+1}$

Goal: minimize cumulative **regret**

$$R_T = \sum_{t=1}^T f(x_*) - f(x_t)$$

Gaussian processes

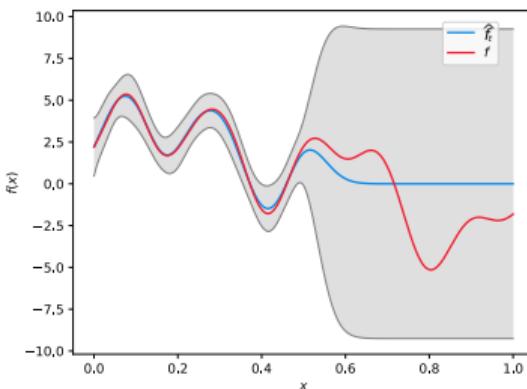
aka kernel ridge regression

$$\begin{aligned}\widehat{\mathbf{K}}_t &\in \mathbb{R}^{t,t} \text{ s.t. } (\widehat{\mathbf{K}})_{i,j} = k(x_i, x_j), i, j = 1, \dots, t \\ \widehat{\mathbf{k}}_t(x) &= (k(x_1, x), \dots, k(x_t, x)) \in \mathbb{R}^t\end{aligned}$$

$$\widehat{\mathbf{y}}_t = (y_1, \dots, y_t)$$

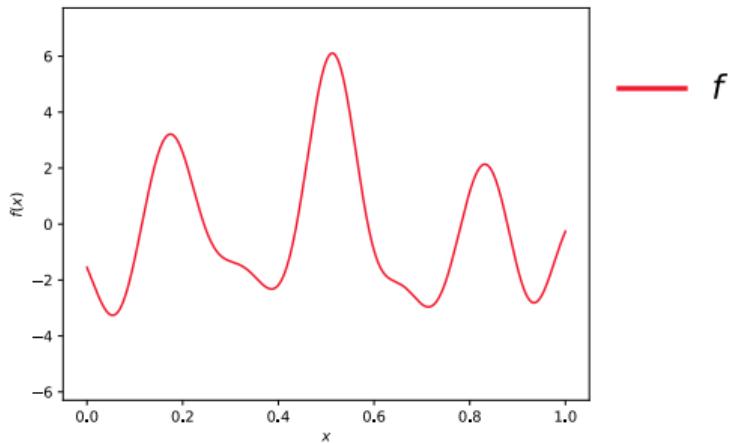
$$\widehat{f}_t(x) = \widehat{\mathbf{k}}_t(x)^\top (\widehat{\mathbf{K}}_t + \lambda I)^{-1} \widehat{\mathbf{y}}_t$$

$$\underbrace{\sigma_t^2(x)}_{\text{variance}} = k(x, x) - \widehat{\mathbf{k}}_t(x)^\top (\widehat{\mathbf{K}}_t + \lambda I)^{-1} \widehat{\mathbf{k}}_t(x)$$



GP-UCB

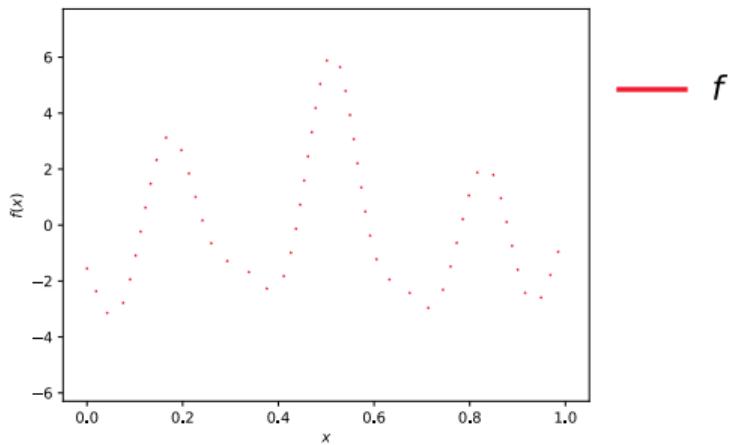
$f : \mathcal{A} \rightarrow \mathbb{R}$ unknown



(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB

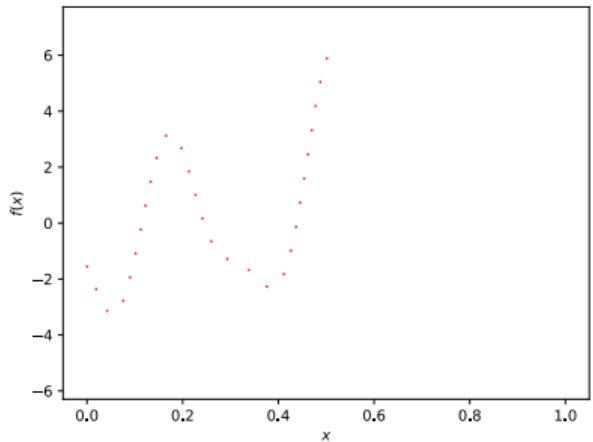
Arms $\mathcal{A} = \{x_i\}_{i=1}^A$



(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB

At time t , collected $(x_i, y_i)_{i=1}^t$



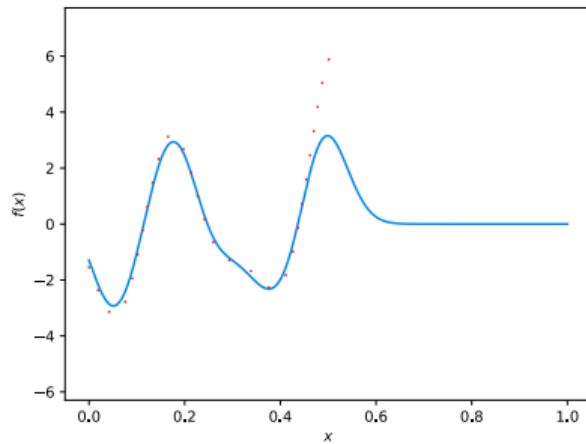
for $t = \{1, \dots, T - 1\}$ **do**

end

(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB

At time t , collected $(x_i, y_i)_{i=1}^t$



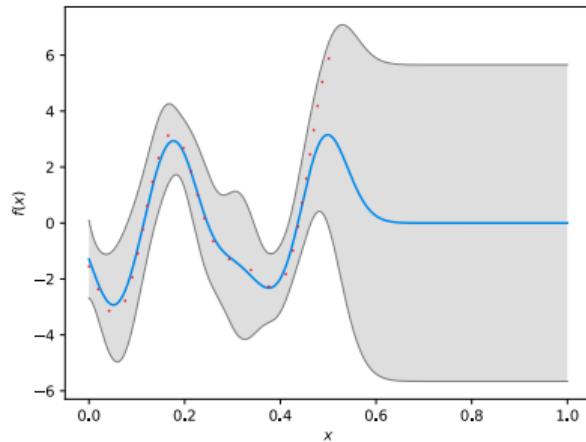
for $t = \{1, \dots, T - 1\}$ do
|
| end

$$\hat{f}_t(x) = \hat{k}_t(x)^\top (\hat{K}_t + \lambda I)^{-1} \hat{y}_t$$

(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB

At time t , collected $(x_i, y_i)_{i=1}^t$



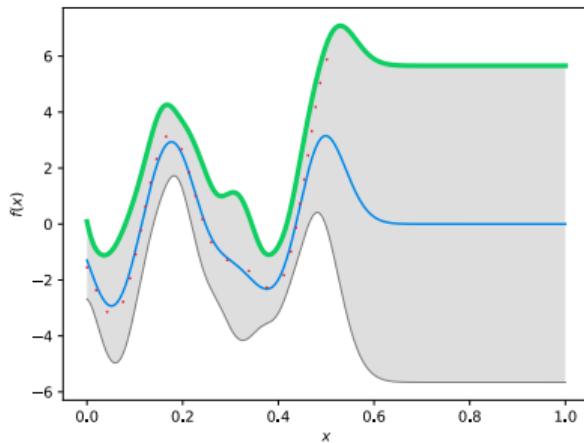
```
for t = {1, ..., T - 1} do
    |
    end
```

$$\hat{f}_t(x) = \hat{k}_t(x)^\top (\hat{K}_t + \lambda I)^{-1} \hat{y}_t \quad \sigma_t^2(x) = k(x, x) - \hat{k}_t(x)^\top (\hat{K}_t + \lambda I)^{-1} \hat{k}_t(x)$$

(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB

At time t , collected $(x_i, y_i)_{i=1}^t$



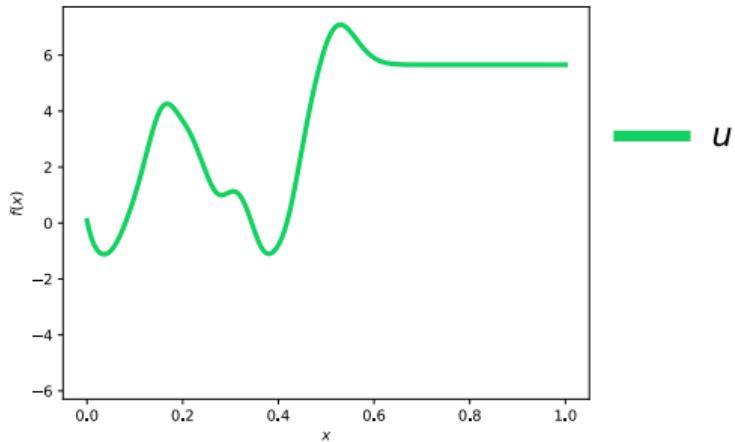
```
for t = {1, ..., T - 1} do
    for i = {1, ..., A} do
         $u_t(x_i) = \hat{f}_t(x_i) + \beta_t \sigma_t^2(x_i);$ 
    end
end
```

$$u_t(x) = \hat{f}_t(x) + \beta_t \sigma_t(x)$$

(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB

At time t , collected $(x_i, y_i)_{i=1}^t$



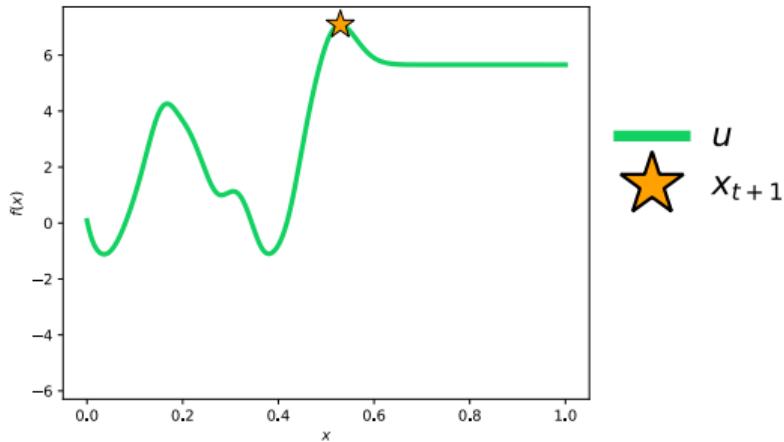
```
for t = {1, ..., T - 1} do
    for i = {1, ..., A} do
        u_t(x_i) =  $\hat{f}_t(x_i)$  +  $\beta_t \sigma_t^2(x_i)$ ;
    end
end
```

$$u_t(x) = \hat{f}_t(x) + \beta_t \sigma_t(x)$$

(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB

At time t , collected $(x_i, y_i)_{i=1}^t$

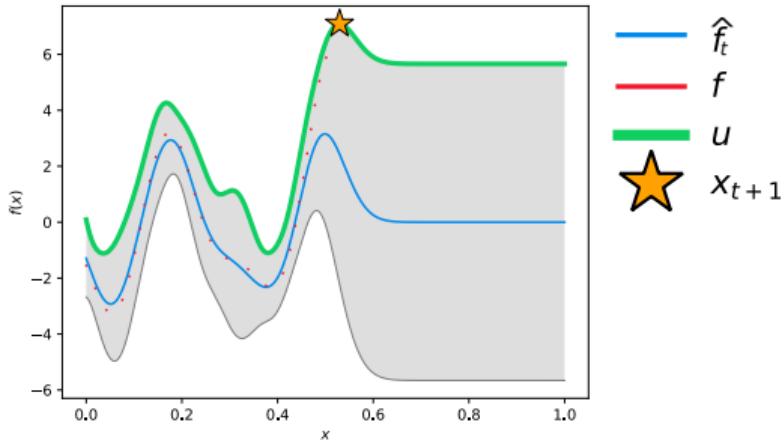


```
for t = {1, ..., T - 1} do
    for i = {1, ..., A} do
         $u_t(x_i) = \hat{f}_t(x_i) + \beta_t \sigma_t^2(x_i);$ 
    end
    Select  $x_{t+1} \leftarrow \operatorname{argmax}_{x_i \in \mathcal{A}} u_t(x_i);$ 
end
```

$$u_t(x) = \hat{f}_t(x) + \beta_t \sigma_t(x) \quad \rightarrow \quad x_{t+1} = \operatorname{argmax}_{x \in \mathcal{A}} u_t(x)$$

(Srinivas, Krause, Kakade, Seeger '10)

GP-UCB: Regret



Computations: Time $O(AT^3)$

Theorem (Srinivas, Krause, Kakade, Seeger '10)

For the proper β_t

$$R_T \leq \sqrt{T}$$

Nyström projection again

aka Sparse GP

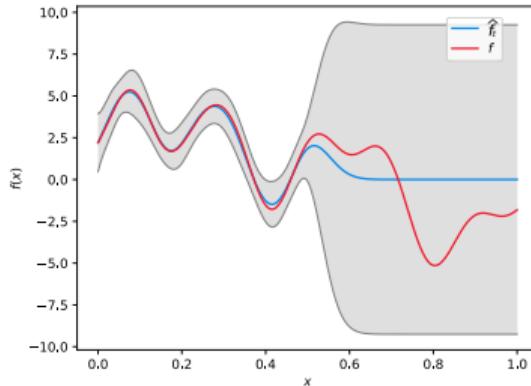
Equivalent formulation:

Given $S = (\bar{x}_i)_{i=1}^M \subseteq (x_i)_{i=1}^t \rightarrow \tilde{k}(x, x') = \tilde{k}_S(x)^\top \hat{K}_S^\dagger \tilde{k}_S(x')$,

with $\hat{K}_S \in \mathbb{R}^{M,M}$ s.t. $(\hat{K}_S)_{i,j} = k(\bar{x}_i, \bar{x}_j)$ and $\tilde{k}_S(x) = (k(\bar{x}_1, x), \dots, k(\bar{x}_M, x))$

$$\tilde{f}_t(x) = \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \hat{y}_t$$

$$\tilde{\sigma}_t^2(x) = \frac{1}{\lambda} \left(k(x, x) - \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \tilde{k}_t(x) \right)$$



Nyström projection again

aka Sparse GP

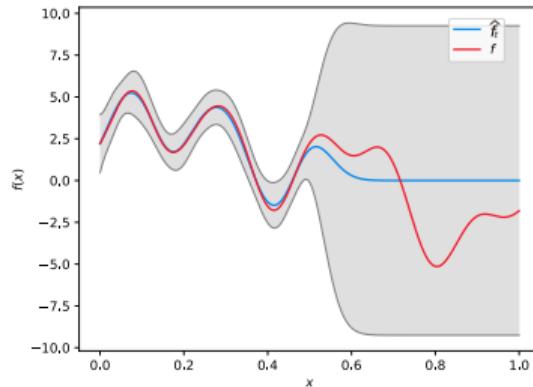
Equivalent formulation:

Given $S = (\bar{x}_i)_{i=1}^M \subseteq (x_i)_{i=1}^t \rightarrow \tilde{k}(x, x') = \tilde{k}_S(x)^\top \hat{K}_S^\dagger \tilde{k}_S(x')$,

with $\hat{K}_S \in \mathbb{R}^{M,M}$ s.t. $(\hat{K}_S)_{i,j} = k(\bar{x}_i, \bar{x}_j)$ and $\tilde{k}_S(x) = (k(\bar{x}_1, x), \dots, k(\bar{x}_M, x))$

$$\tilde{f}_t(x) = \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \hat{y}_t$$

$$\tilde{\sigma}_t^2(x) = \frac{1}{\lambda} \left(k(x, x) - \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \tilde{k}_t(x) \right)$$



Computations: Time $O(AM^2T)$

...but...

NO guarantees on regret (overconfident when S is "bad")

BKB: Regret

BKB: Regret

- $(x_i)_{i=1}^t$ changes with time $\rightarrow S_t$ must change with t

BKB: Regret

- ▶ $(x_i)_{i=1}^t$ changes with time $\rightarrow S_t$ must change with t
- ▶ $\sigma_t^2(\cdot)$ captures informative arms \rightarrow include x_i in S_t when $\sigma_t^2(x_i)$ is large

BKB: Regret

- $(x_i)_{i=1}^t$ changes with time $\rightarrow S_t$ must change with t
- $\sigma_t^2(\cdot)$ captures informative arms \rightarrow include x_i in S_t when $\sigma_t^2(x_i)$ is large

$$\tilde{f}_t(x) = \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \hat{y}$$

$$\tilde{\sigma}_t^2(x) = \frac{1}{\lambda} \left(k(x, x) - \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \tilde{k}_t(x) \right)$$

```
for t = {1, ..., T - 1} do
    for i = {1, ..., A} do
        |    $\tilde{u}_t(x_i) = \tilde{f}_t(x_i) + \tilde{\beta}_t \tilde{\sigma}_t^2(x_i);$ 
    end
    Select  $x_{t+1} \leftarrow \operatorname{argmax}_{x_i \in \mathcal{A}} \tilde{u}_t(x_i);$ 
    Set  $\tilde{p}_{t+1} \propto [\tilde{\sigma}_t^2(x_1), \dots, \tilde{\sigma}_t^2(x_{t+1})];$ 
    Sample  $S_{t+1} \sim \tilde{p}_{t+1};$ 
end
```

BKB: Regret

- $(x_i)_{i=1}^t$ changes with time $\rightarrow S_t$ must change with t
- $\sigma_t^2(\cdot)$ captures informative arms \rightarrow include x_i in S_t when $\sigma_t^2(x_i)$ is large

$$\tilde{f}_t(x) = \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \hat{y}$$

$$\tilde{\sigma}_t^2(x) = \frac{1}{\lambda} \left(k(x, x) - \tilde{k}_t(x)^\top (\tilde{K}_t + \lambda I)^{-1} \tilde{k}_t(x) \right)$$

```
for t = {1, ..., T - 1} do
    for i = {1, ..., A} do
        |    $\tilde{u}_t(x_i) = \tilde{f}_t(x_i) + \tilde{\beta}_t \tilde{\sigma}_t^2(x_i);$ 
    end
    Select  $x_{t+1} \leftarrow \operatorname{argmax}_{x_i \in \mathcal{A}} \tilde{u}_t(x_i);$ 
    Set  $\tilde{p}_{t+1} \propto [\tilde{\sigma}_t^2(x_1), \dots, \tilde{\sigma}_t^2(x_{t+1})];$ 
    Sample  $S_{t+1} \sim \tilde{p}_{t+1};$ 
end
```

Computations:

Time $O(A d_{\text{eff}}^2 T)$

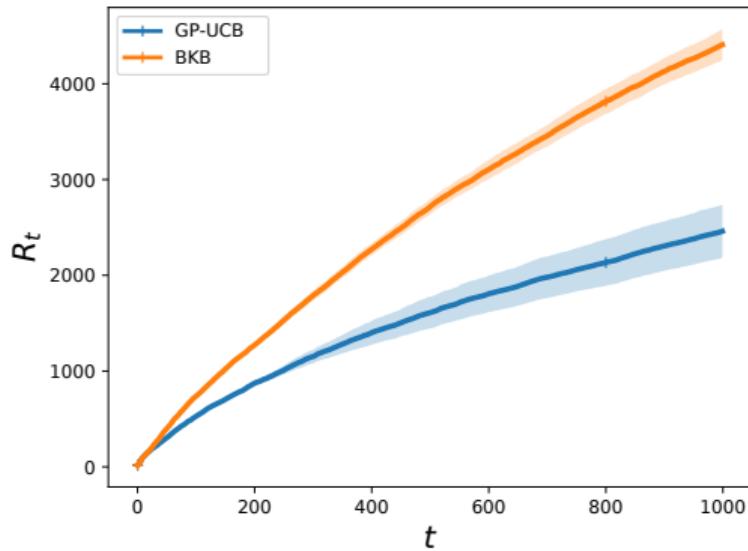
Theorem (Calandriello, Carratino, Lazaric, Valko, R. '19)

For the proper (and cheap to compute!) $\tilde{\beta}_t$, with $|S_T| \leq d_{\text{eff}}$ with $d_{\text{eff}} \ll T$

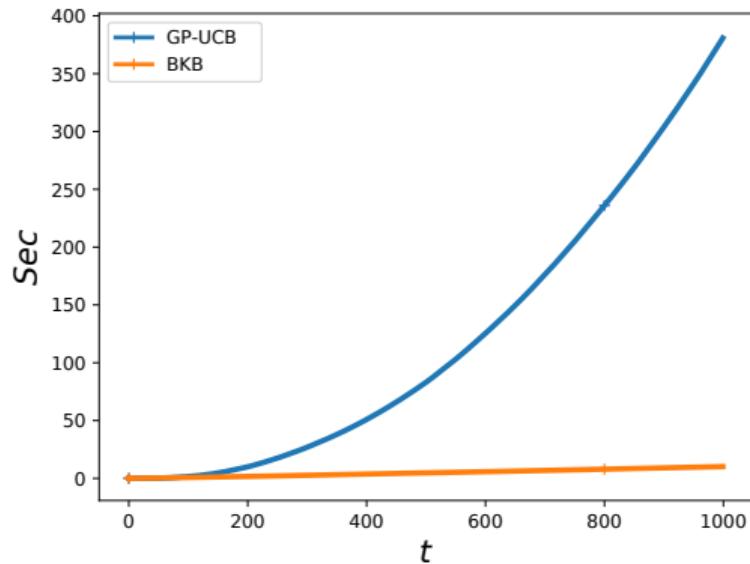
$$R_T \leq \sqrt{T}$$

In practice

Cumulative regret R_t



Times



Sublinear regret in a fraction of the time

Recent improvement using batching [Calandriello, Carratino, Lazaric, Valko, R. '20].

Outline

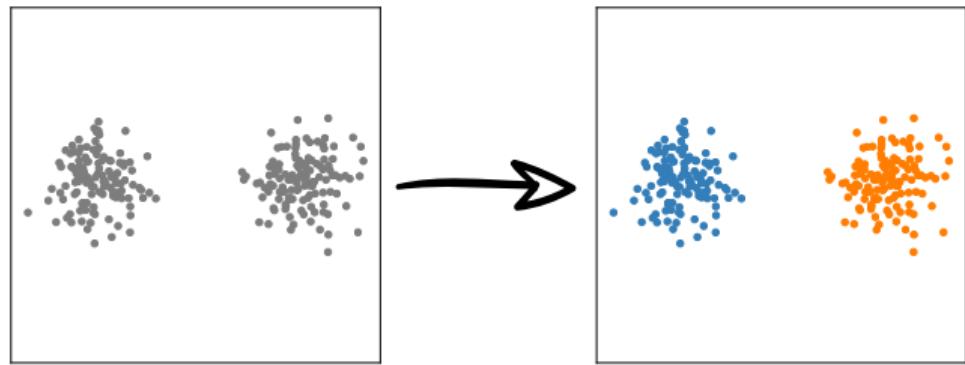
Nyström projection for unsupervised learning?

- ▶ Kernel K-means
- ▶ Kernel PCA

K-means

Partition n points into k clusters.

$$\hat{C}_K = \min_{[c_1, \dots, c_j]} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, K} \|x_i - c_j\|^2$$



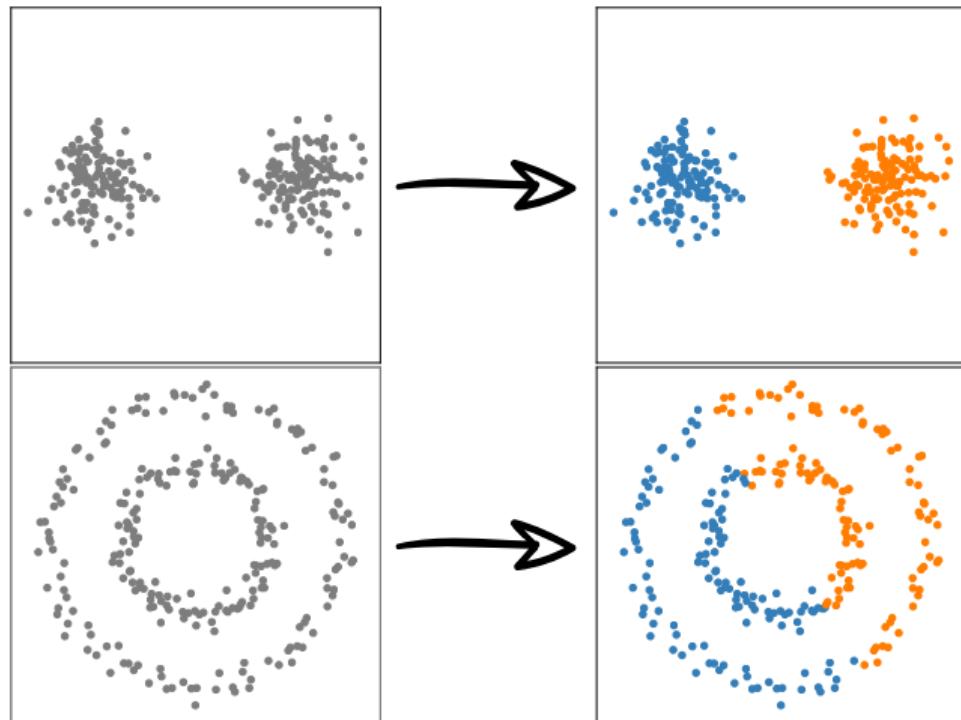
Only linear separations

K-means

Partition n points into k clusters.

$$\hat{C}_K = \min_{[c_1, \dots, c_j]} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, K} \|x_i - c_j\|^2$$

Only linear separations



From K-means to kernel K-means

Partition n points into K clusters.

$$\hat{C}_K = \min_{[c_1, \dots, c_j]} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, K} \|x_i - c_j\|^2$$

note that: $\|x - \bar{x}\|^2 = x^\top x + \bar{x}^\top \bar{x} - 2x^\top \bar{x}$

From K-means to kernel K-means

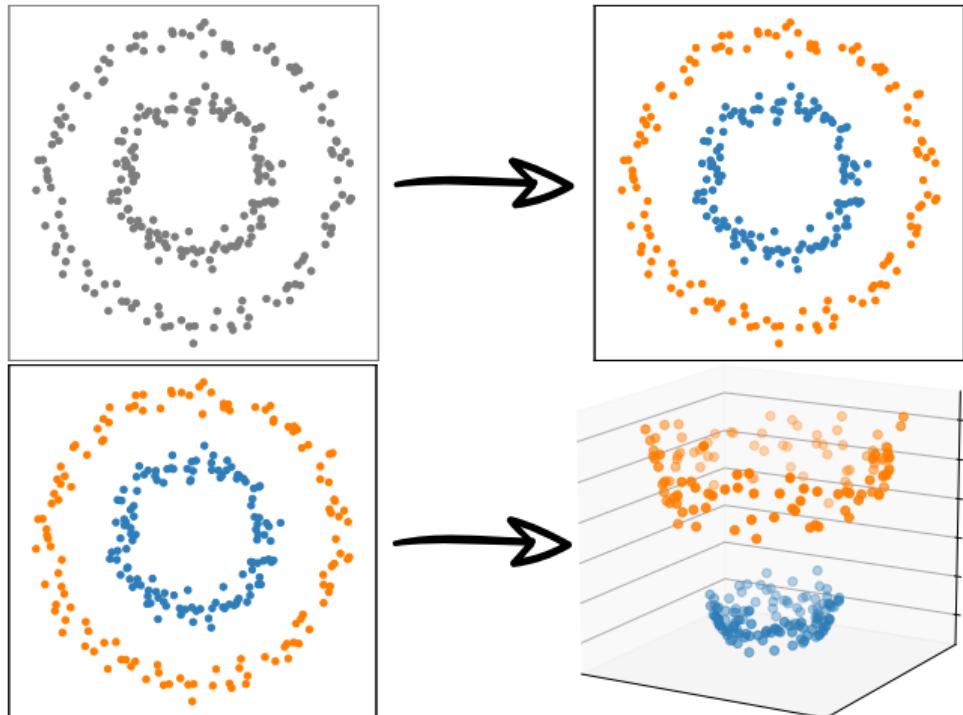
Partition n points into K clusters.

$$\hat{C}_K = \min_{[c_1, \dots, c_j]} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, K} \|x_i - c_j\|^2$$

note that: $\|x - \bar{x}\|^2 = x^\top x + \bar{x}^\top \bar{x} - 2x^\top \bar{x}$

Kernel to rescue!

$$x^\top \bar{x} \mapsto k(x, \bar{x})$$



From K-means to kernel K-means

Partition n points into K clusters.

$$\hat{C}_K = \min_{[c_1, \dots, c_j]} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, K} \|x_i - c_j\|^2$$

note that: $\|x - \bar{x}\|^2 = x^\top x + \bar{x}^\top \bar{x} - 2x^\top \bar{x}$

\widehat{K}

Kernel to rescue!

$$x^\top \bar{x} \quad \mapsto \quad k(x, \bar{x})$$

From K-means to Nyström K-means

Partition n points into K clusters.

$$\hat{C}_K = \min_{[c_1, \dots, c_j]} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, K} \|x_i - c_j\|^2$$

note that: $\|x - \bar{x}\|^2 = x^\top x + \bar{x}^\top \bar{x} - 2x^\top \bar{x}$

From K-means to Nyström K-means

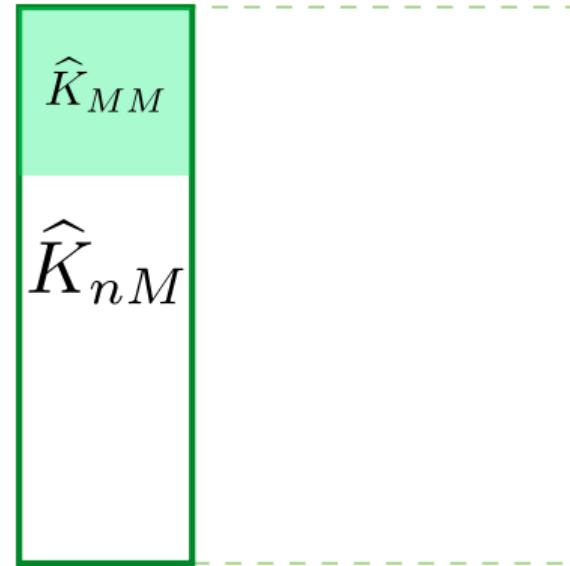
Partition n points into K clusters.

$$\hat{C}_K = \min_{[c_1, \dots, c_j]} \frac{1}{n} \sum_{i=1}^n \min_{j=1, \dots, K} \|x_i - c_j\|^2$$

note that: $\|x - \bar{x}\|^2 = x^\top x + \bar{x}^\top \bar{x} - 2x^\top \bar{x}$

Nyström to rescue!

$$k(x, x) \mapsto \tilde{k}(x, x) = \tilde{k}_M(x)^\top \hat{K}_M^\dagger \tilde{k}_M(x')$$



Guarantees for Nyström K-means

Assume $(x_i)_{i=1}^n \sim \rho^n$, \hat{C}_K the Nyström K-means solution and ²

$$L(\hat{C}_K) = \mathbb{E}_x [\min_{j=1,\dots,K} \|x - c_j\|^2].$$

Theorem (Calandriello, R. '19)

Assume $\|x\| \leq 1$, and the Nyström centers chosen uniformly at random, then

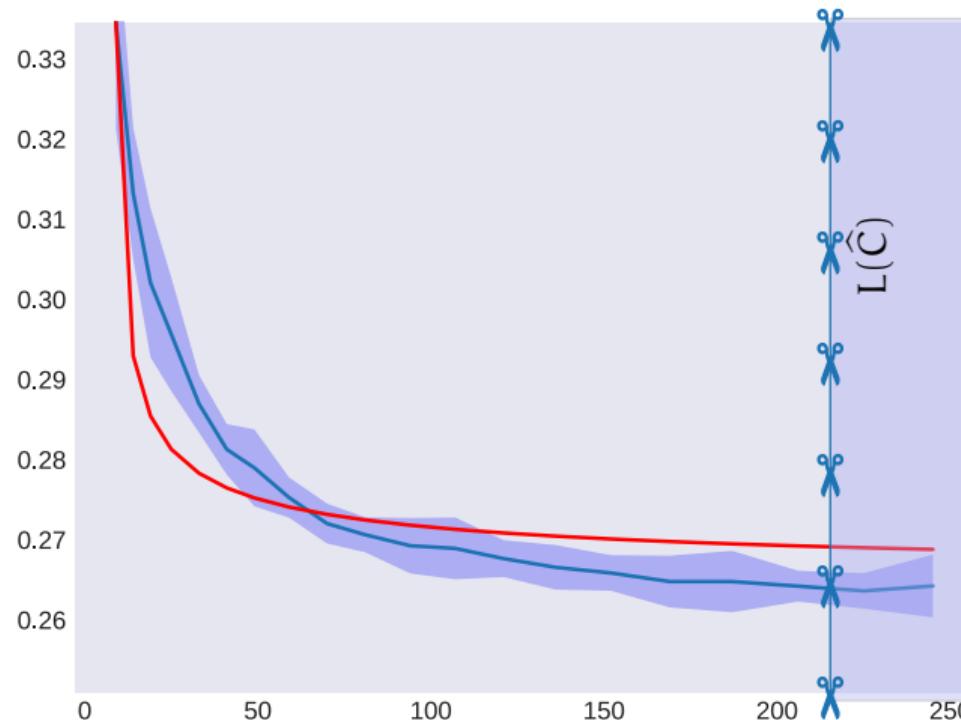
$$\mathbb{E}[L(\hat{C}_K)] \lesssim \frac{K}{\sqrt{n}} + \frac{K}{M},$$

so that if $M = \sqrt{n}$ then

$$\mathbb{E}[L(\hat{C}_K)] \lesssim \frac{K}{\sqrt{n}}.$$

The above bound matches that of exact kernel k-means.

MNIST-60k: expected loss vs projection size M



Nyström projections for unsupervised learning?

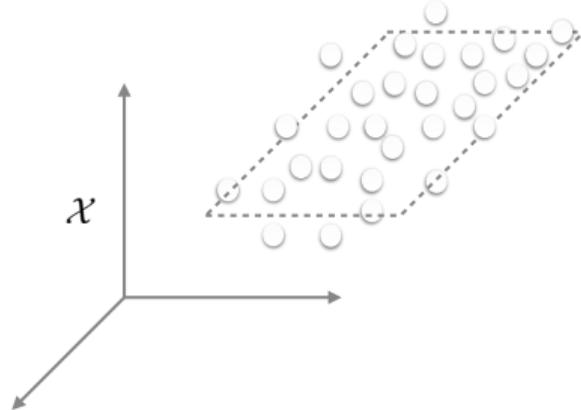
- ▶ Kernel K-means
- ▶ **Kernel PCA**

PCA

$$\widehat{\Sigma} = \frac{1}{n} \widehat{X}^\top \widehat{X} = \widehat{V} \widehat{\Lambda}^2 \widehat{V}^\top$$

$$\widehat{\Lambda} = \text{diag}(\widehat{\lambda}_1^2, \dots, \widehat{\lambda}_n^2).$$

$$x \mapsto (x^\top v_1, \dots, x^\top v_\ell)$$



Project only on linear subspaces

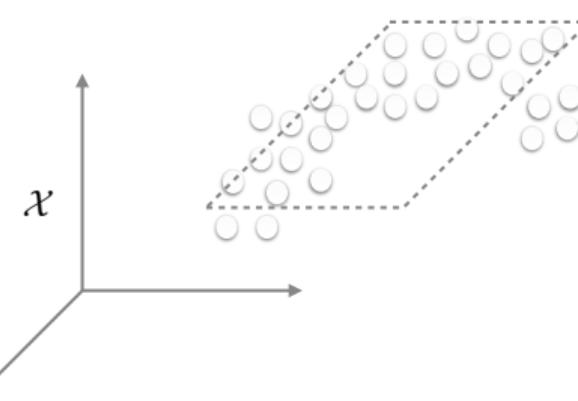
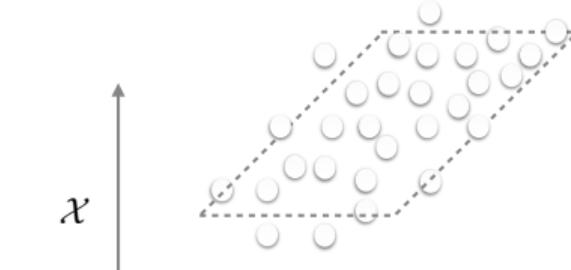
PCA

$$\widehat{\Sigma} = \frac{1}{n} \widehat{X}^\top \widehat{X} = \widehat{V} \widehat{\Lambda}^2 \widehat{V}^\top$$

$$\widehat{\Lambda} = \text{diag}(\widehat{\lambda}_1^2, \dots, \widehat{\lambda}_n^2).$$

$$x \mapsto (x^\top v_1, \dots, x^\top v_\ell)$$

Project only on linear subspaces



Kernel PCA

$$\hat{K} = \frac{1}{n} \hat{X} \hat{X}^\top = \hat{U} \hat{\Lambda}^2 \hat{U}$$

$$v_1 = \frac{1}{n\hat{\lambda}_1} \hat{X}^\top \hat{u}_1$$

Kernel PCA

$$\hat{K} = \frac{1}{n} \hat{X} \hat{X}^\top = \hat{U} \hat{\Lambda}^2 \hat{U}$$

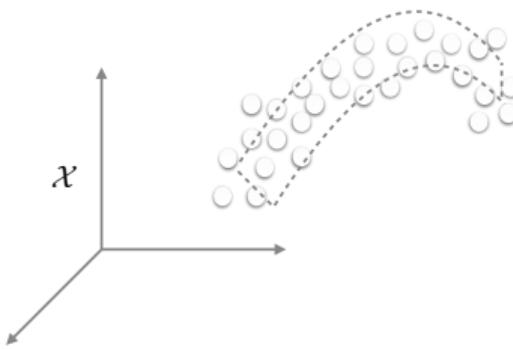
$$v_1 = \frac{1}{n\hat{\lambda}_1} \hat{X}^\top \hat{u}_1 \quad \Rightarrow \quad x^\top v_1 = \frac{1}{n\hat{\lambda}_1} \sum_{i=1}^n x^\top x_i (\hat{u}_1)_i$$

Kernel PCA

$$\hat{K} = \frac{1}{n} \hat{X} \hat{X}^\top = \hat{U} \hat{\Lambda}^2 \hat{U}$$

$$v_1 = \frac{1}{n\hat{\lambda}_1} \hat{X}^\top \hat{u}_1 \quad \Rightarrow \quad x^\top v_1 = \frac{1}{n\hat{\lambda}_1} \sum_{i=1}^n x^\top x_i (\hat{u}_1)_i$$

$$x^\top \bar{x} \quad \mapsto \quad k(x, \bar{x})$$



Project on non linear subspaces

Kernel PCA

$$\hat{K} = \frac{1}{n} \hat{X} \hat{X}^\top = \hat{U} \hat{\Lambda}^2 \hat{U}$$

$$v_1 = \frac{1}{n\hat{\lambda}_1} \hat{X}^\top \hat{u}_1 \quad \Rightarrow \quad x^\top v_1 = \frac{1}{n\hat{\lambda}_1} \sum_{i=1}^n x^\top x_i (\hat{u}_1)_i$$

$$x^\top \bar{x} \quad \mapsto \quad k(x, \bar{x})$$

Project on non linear subspaces

$$\widehat{K}$$

Nyström PCA

$$v_1 = \underset{\|w\|=1}{\operatorname{argmax}} w^\top \Sigma w \quad \Rightarrow \quad \tilde{v}_1 = \underset{w=\bar{X}_M^\top c : \|w\|=1}{\operatorname{argmax}} w^\top \Sigma w$$

Nyström PCA

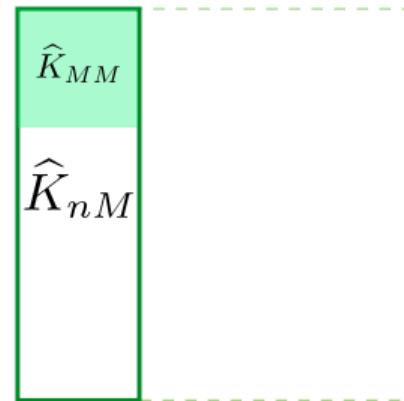
$$v_1 = \underset{\|w\|=1}{\operatorname{argmax}} w^\top \Sigma w \quad \Rightarrow \quad \tilde{v}_1 = \underset{w=\bar{X}_M^\top c : \|w\|=1}{\operatorname{argmax}} w^\top \Sigma w$$

...

$$x^\top v_1 \approx x^\top \tilde{v}_1 = \sum_{i=1}^M x^\top \bar{x}_i \hat{K}_M^{-1/2} (\tilde{u}_1)_i$$

with

$$\hat{K}_M^{-1/2} \hat{K}_{nM}^\top \hat{K}_{nM} \hat{K}_M^{-1/2} = \tilde{U} \tilde{\Lambda}^2 \tilde{U}^\top.$$



Guarantees for Nyström PCA

Assume $(x_i)_{i=1}^n \sim \rho^n$, \widehat{P}_ℓ the Nyström PCA projection and ³

$$L(\widehat{P}_\ell) = \mathbb{E}_x[\|x - P_\ell x\|^2].$$

Theorem (Serge, Sriperumbudur, R., Rudi '20)

Assume $\|x\| \leq 1$, and Nyström centers chosen uniformly at random. Let $\Sigma = \mathbb{E}_x[xx^\top]$ with

$$\lambda_j^2(\Sigma) \sim j^{-\alpha}, \quad \alpha > 1$$

Then for $\ell = n^{\frac{\theta}{\alpha}}$, $\theta < 1$ and $M \leq n^\theta \log n$

$$\mathbb{E}[L(\widehat{C})] \lesssim n^{-\theta(1-\frac{1}{\alpha})}.$$

The above bound matches that of exact KPCA.

Wrapping up

Contribution

- ▶ Nyström projections allow computational savings with no accuracy loss.
- ▶ Further results: adaptive sampling,
 - leverage scores [Calandriello, Rudi, Carratino, R. '18],
 - DPP sampling [Dereziński Calandriello, Valko '19]
- ▶ Related results: random features [Rudi, R. '16].

We are thinking about:

- ▶ More Nyström kernel [add your favorite].
- ▶ Interpolation regimes $n \ll 2^d$.
- ▶ Combine Nyström and multiscale approaches [Chen, Avron, Sindhwani '16].



PhD/Postdoc positions available!



Relevant stuff

Papers

Less is More: Nyström Computational Regularization

A. Rudi, R. Camoriano and L. Rosasco · NIPS15

FALKON: An Optimal Large Scale Kernel Method

A. Rudi, L. Carratino and L. Rosasco · NIPS17

Gaussian Process Optimization with Adaptive Sketching: Scalable and No Regret

D. Calandriello, L. Carratino, A. Lazaric, M. Valko and L. Rosasco · COLT19

Statistical and computational trade-offs in kernel k-means

D. Calandriello, L. Rosasco · NeurIPS18

Gain with no Pain: Efficient Kernel-PCA by Nyström Sampling

N. Sterge, B. Sriperumbur, L. Rosasco, A. Rudi · AISTATS20

Code

FALKON

G. Meanti, L. Carratino, L. Rosasco and A. Rudi · <http://lcsr.mit.edu>

BKB

D. Calandriello, L. Carratino, A. Lazaric, M. Valko and L. Rosasco · <http://lcsr.mit.edu>

More relevant stuff

Papers

Learning with SGD and Random Features

L. Carratino, A. Rudi and L. Rosasco · NeurIPS18

On Fast Leverage Score Sampling and Optimal Learning

A. Rudi, D. Calandriello, L. Carratino and L. Rosasco · NeurIPS18

Exact sampling of determinantal point processes with sublinear time preprocessing

M. Dereziński , D. Calandriello, M. Valko · NeurIPS19

Near-linear Time GP Optimization with Adaptive Batching and Resparsification

D. Calandriello, L. Carratino, A. Lazaric, M. Valko and L. Rosasco · Preprint 2020

Code

BLESS: leverage score sampling

A. Rudi, D. Calandriello, L. Carratino and L. Rosasco · <http://lcs.mit.edu>

DPP sampling

M. Dereziński , D. Calandriello, M. Valko · <http://lcs.mit.edu>