



### Scalable Low-Rank Nonlinear Subspace Tracking

#### F. Sheikholeslami, D. Berberidis, and G. B. Giannakis

Dept. of ECE and DTC, University of Minnesota



Acknowledgement: NSF grants 1500713, 1514056, and NIH grant 1R01GM104975-01

# Learning from "Big Data"

- Challenges
  - > Big size ( $D \gg and/or N \gg$ )
  - Fast streaming
  - Incomplete
  - Noise and outliers

- Opportunities in key tasks
  - Dimensionality reduction
  - Online and robust regression, classification and clustering
  - Denoising and imputation





## Outline

- Scalable kernel-based learning
  - Sparsity-aware low-rank approximation of lifted data
  - Online subspace tracking
  - Affordable memory and computation
  - Theoretical guarantees and test results
  - Future directions

### Linear or nonlinear functions for learning?

**D** Regression or classification: Given  $\{y_n, \mathbf{x}_n\}_{n=1}^N$ , find  $\hat{f}: \mathbf{x} \to y = f(\mathbf{x}) + v$ 

 $\Box \text{ Lift via nonlinear map } \mathbf{x} \to \boldsymbol{\phi}(\mathbf{x}) \in \mathbb{R}^{\bar{D}} \text{ to linear } y_n = \boldsymbol{\phi}^\top(\mathbf{x}_n)\bar{\boldsymbol{\theta}} + v_n$ 

Pre-select kernel (inner product) function

$$\mathbf{x}_i^{\top} \mathbf{x}_j \rightarrow \frac{k(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\phi}^{\top}(\mathbf{x}_i) \boldsymbol{\phi}(\mathbf{x}_j)}{2\sigma^2}$$
e.g.,  $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2})$ 

RKHS basis expansion

 $\hat{f}(\mathbf{x}) = \sum_{n=1}^{N} \alpha_n k(\mathbf{x}, \mathbf{x}_n)$ 



 $\begin{array}{ccc} \Psi : \Lambda & \to \Lambda \\ (x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2}) x_1 x_2, x_2^2) \end{array}$ 

❑ Kernel-based nonparametric ridge regression

$$\mathbf{y} = \mathbf{K} \boldsymbol{\alpha} + \mathbf{v} \rightarrow \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{y}$$

> Memory requirement  $\mathcal{O}(N^2)$ , and complexity  $\mathcal{O}(N^3)$ 

### **Prior art**

#### Functional gradient descent

- Tractable gradient descent (batch form in primal domain)
  - Pegasus [Shaleve-Schwart et al.'11]; Norma [Kivinen et al.'04]; [Wang et al. '12, Lu et al. '16]

#### Kernel matrix approximation

- Tractable by leveraging low-rank attribute (in dual domain)
  - Nystrom approximation [Williamson-Seeger'00, Si et al. '14, Wang et al. '14]
  - Incomplete Cholesky decomposition [Fine-Scheinberg'01]

#### Random feature approximation

- Kernel approximation via feature approximation
  - Random sampling and kernel matrix factorization [Rahimi-Recht'08, Zhang et al.'11]

#### Reduced-dimensionality features

- Tractability through online extraction of low-dimensional features
  - Online (kernel) PCA [Honnein'012]; linear subspace tracking [Mardani et al.'15]

## Low-rank approximation of mapped features

lacksquare Nonlinear mapping  $\mathbf{x} o oldsymbol{\phi}(\mathbf{x})$ 

 $\Box \text{ Low-rank approximation } \phi(\mathbf{x}) \simeq \hat{\phi}(\mathbf{x}) := \bar{\mathbf{L}} \mathbf{q} \qquad \bar{\mathbf{L}} \in \mathbb{R}^{\bar{D} \times r} \qquad \mathbf{q} \in \mathbb{R}^{r \times 1}$ 

$$\min_{\bar{\mathbf{L}},\{\mathbf{q}_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N \|\boldsymbol{\phi}(\mathbf{x}_i) - \bar{\mathbf{L}}\mathbf{q}_i\|_{\mathcal{H}}^2 + \frac{\lambda}{N} \sum_{i=1}^N \|\mathbf{q}_i\|_2^2$$

□ Alternating minimization

**S1.** Update virtual subspace  $\bar{\mathbf{L}}$  while fixing feature vectors  $\{\mathbf{q}_i\}_{i=1}^N$ 

$$ar{\mathbf{L}} = oldsymbol{\Phi}_N \mathbf{Q}_N^\dagger := oldsymbol{\Phi}_N \mathbf{A} \qquad oldsymbol{\Phi}_N := [oldsymbol{\phi}(\mathbf{x}_1), \dots, oldsymbol{\phi}(\mathbf{x}_N)] \ \mathbf{Q}_N := [oldsymbol{q}_1, \dots, oldsymbol{q}_N]$$

**S2.** Update features  $\{\mathbf{q}_i\}_{i=1}^N$  while fixing subspace  $\bar{\mathbf{L}}$ 

$$\mathbf{q}_i = (\mathbf{A}^\top \mathbf{K} \mathbf{A} + \lambda \mathbf{I}_r)^{-1} \mathbf{A}^\top \mathbf{k}_i$$

### Sparsity-aware low-rank approximation

lacksquare Row-sparsity of  $\mathbf{A}$   $\Rightarrow$  fewer "basis vectors" for  $ar{\mathbf{L}}$ 

$$\min_{\mathbf{A}, \{\mathbf{q}_i\}_{i=1}^N} \frac{1}{N} \sum_{i=1}^N \|\boldsymbol{\phi}(\mathbf{x}_i) - \boldsymbol{\Phi}_N \mathbf{A} \mathbf{q}_i\|_2^2 + \frac{\lambda}{N} \sum_{i=1}^N \|\mathbf{q}_i\|_2^2 + \eta \sum_{i=1}^N \|\mathbf{a}_i\|_2$$

#### 

**S1.** Update  $\{\mathbf{q}_i\}_{i=1}^N$  with  $ar{\mathbf{L}}$  fixed

$$\mathbf{q}_i = (\mathbf{A}^\top \mathbf{K} \mathbf{A} + \lambda \mathbf{I}_r)^{-1} \mathbf{A}^\top \mathbf{k}_i$$

S2. Update 
$$\mathbf{A}$$
 with  $\{\mathbf{q}_i\}_{i=1}^N$  fixed

$$\mathbf{a}_{i}[k+1] = \arg\min_{\mathbf{a}_{i}} \frac{1}{2} \mathbf{a}_{i}^{\top} \mathbf{H}_{i} \mathbf{a}_{i} + \mathbf{p}_{i}^{\top} \mathbf{a}_{i} + \eta \|\mathbf{a}_{i}\|_{2}$$

Group shrinkage

- Solvers
  - i. BCD-Exact: Exact minimization- more expensive, larger descent
  - ii. BCD-PG: Proximal gradient-based update, inexpensive, smaller descent

F. Sheikholeslami and G. B. Giannakis, "Scalable Kernel-based Learning via Low-rank Approximation of Lifted Data," *Proc. of 55th Allerton Conf. on Comm., Control, and Computing*, Oct. 4-6, 2017.

### Convergence and generalization bounds

**Proposition 1** Sequences  $\{\mathbf{Q}[k], \mathbf{A}[k]\}_{k=1}^{\infty}$  generated by BCD-Ex and BCD-PG

iterations converge to a stationary point.

Generalization: performance guarantees on unseen data

Equivalent optimization  

$$\min_{\mathbf{A}\in\mathcal{A}} \frac{1}{N} \sum_{i=1}^{N} \min_{\mathbf{q}_i\in\mathcal{Q}} \|\phi(\mathbf{x}_i) - \Phi_N \mathbf{A}\mathbf{q}_i\|_{\mathcal{H}}^2 \qquad \mathcal{Q} := \{\mathbf{q} \mid \mathbf{q} \in \mathbb{R}^r, \|\mathbf{q}\|_2 < B\}$$

$$\mathcal{A} := \{\mathbf{A} \mid \sum_{i=1}^{N} \|\mathbf{a}_i\|_2 \leq \tau\}$$

Define 
$$\|\mathcal{L}\|_{\mathcal{Q}} := \sup_{\mathbf{L}\in\mathcal{L}} \|\mathbf{\bar{L}}\|_{\mathcal{Q}} = \sup_{\mathbf{L}\in\mathcal{L}} \sup_{\mathbf{q}\in\mathcal{Q}} \|\mathbf{\bar{L}}\mathbf{q}\|_{\mathcal{H}} \qquad \mathcal{L} := \{\mathbf{\bar{L}} \mid \mathbf{\bar{L}} = \Phi_N \mathbf{A}, \mathbf{A} \in \mathcal{A}\}$$

Proposition 2. If 
$$\|\mathcal{L}\|_{\mathcal{Q}} \geq 1$$
, and  $\kappa(\mathbf{x}, \mathbf{x}) \leq \kappa$ , then for  $\delta > 0$ , it holds wp  $> 1 - \delta$ 
 $\mathbb{E}[\ell(\mathbf{x}; \mathbf{\bar{L}})] - 1/N \sum_{i=1}^{N} \ell(\mathbf{x}_i; \mathbf{\bar{L}}) \leq \frac{d}{\sqrt{N}} \left(14\kappa B\tau + \frac{\kappa}{2}\sqrt{\ln(16N\kappa^2B^2\tau^2)}\right) + \kappa\sqrt{\frac{\ln(1/\delta)}{2N}} \quad \forall \mathbf{\bar{L}} \in \mathbb{R}$ 

F. Sheikholeslami and G. B. Giannakis, "Scalable Kernel-based Learning via Low-rank Approximation of 8 Lifted Data," Proc. Allerton, 2017.

 $\boldsymbol{Z}$ 

 $\sqrt{N}$ 

i=1

### Linearized kernel regression and classification



❑ Kernel matrix approximation

Proposition 3. If  $e_i := \|\phi(\mathbf{x}_i) - \hat{\phi}(\mathbf{x}_i)\|_{\mathcal{H}}^2$  iid with  $\bar{e} := \mathbb{E}[e_i]$  kernel matrix  $\mathbf{K} = \mathbf{\Phi}^{\top} \mathbf{\Phi}$ can be approximated as  $\hat{\mathbf{K}} = \hat{\mathbf{\Phi}}^{\top} \hat{\mathbf{\Phi}}$ , and wp >  $1 - 2e^{-2N\tau^2}$  $\frac{1}{N} \|\mathbf{K} - \hat{\mathbf{K}}\|_F \le \sqrt{\bar{e} + \tau} (\sqrt{\bar{e} + \tau} + 2)$ 

Bounds also on support vector machines for regression and classification

F. Sheikholeslami, D. K. Berberidis, and G. B. Giannakis, "Kernel-based Low-rank Feature Extraction on a Budget for Big Data Streams," arxiv:1601.07947.

### Simulation tests



F. Sheikholeslami and G. B. Giannakis, "Scalable Kernel-based Learning via Low-rank Approximation of Lifted Data," *Proc. Allerton Conf.* 2017. 10 Datasets available at UCI repository: http://archive.ics.uci.edu/ml/datasets.html

### Online function approximation on a budget

Low-rank (r) subspace tracking [Mardani-Mateos-GG'14] here on lifted data

$$\min_{\mathbf{A},\{\mathbf{q}_{\nu}\}_{\nu=1}^{n}} \frac{1}{n} \sum_{\nu=1}^{n} \|\boldsymbol{\phi}(\mathbf{x}_{\nu}) - \boldsymbol{\Phi}_{n} \mathbf{A} \mathbf{q}_{\nu}\|_{\mathcal{H}}^{2} + \frac{\lambda}{2n} \left( \mathbf{\Phi}_{n} \mathbf{A}^{1} \mathbf{Q}_{n} \mathbf{A}^{1} \|_{*F}^{2} + \|\mathbf{Q}_{n}\|_{F}^{2} \right)$$
$$\stackrel{\ell_{n}(\mathbf{x}_{\nu}; \mathbf{A}, \mathbf{q}_{\nu}; \mathbf{x}_{1:n}) :=}{k(\mathbf{x}_{\nu}, \mathbf{x}_{\nu}) - 2\mathbf{k}_{\nu}^{\top} \mathbf{A} \mathbf{q}_{\nu} + \mathbf{q}_{\nu}^{\top} \mathbf{A}^{\top} \mathbf{K}_{n} \mathbf{A} \mathbf{q}_{\nu}} \quad \text{Rank surrogate}$$

Censoring can mitigate curse of dimensionality as  $n \uparrow$ 

Iteration n-1:  $\mathbf{K}_{n-1}$  and  $\mathbf{A}[n-1]$  available

**S1.** Find projection coefficients

$$\begin{array}{c} 0.2 \\ 0.1 \\ 0.2 \\ -0.2 \\ -0.1 \\ -0.1 \\ -0.1 \\ -0.1 \\ 0.2 \\ -0.1 \\ -0.1 \\ 0.2 \\ -0.1 \\ -0.2 \\ -0.2 \\ -0.1 \\ -0.2 \\$$

**Budget** 

maintenance

11

$$\mathbf{q}_n = (\mathbf{A}^\top [n-1] \mathbf{K}_{n-1} \mathbf{A} [n-1] + \lambda \mathbf{I})^{-1} \mathbf{k}_{n-1} (\mathbf{x}_n)$$

**S2.** Update subspace factor (via stochastic gradient descent)

$$\mathbf{A}[n] = \min_{\mathbf{A}} \frac{1}{n} \sum_{\nu=1}^{n} \ell_n(\mathbf{x}_{\nu}; \mathbf{A}, \mathbf{q}_{\nu}; \mathbf{x}_{1:n}) + \frac{\lambda}{2n} \operatorname{Tr}\{\mathbf{A}^{\top} \mathbf{K}_n \mathbf{A}\}$$

 $\Box$  If budget exceeded, remove the row of **A** with minimum  $\ell_2$ -norm

F. Sheikholeslami, D. K. Berberidis, and G. B. Giannakis, "Kernel-based Low-rank Feature Extraction on a Budget for Big Data Streams," arxiv:1601.07947.

### **OK-FEB** with linear classification and regression

Year dataset (regression)

N=463,700 , d= 90, r=10, B=15

- □ Slice dataset (regression)
- □ N=53,500 , d= 384, r=10, B=15



✓ OK-FEB LSVM outperforms budgeted K-SVM/SVR variants in classification/regression

### Tracking dynamic subspaces

### Recency-aware removal rule

- Recency factor
- Associate factor  $r_i$  to the i-th SV in the budget
- Decay  $r_i$  with inclusion of a new SV,
- $r_i = \beta r_i \ 0 < \beta \le 1$

Elipsoid 1 Elipsoid 2

Removal rule

$$\hat{i}_* = \arg\min_{i=1,2,\dots,B+1} r_i \|\mathbf{a}_i[n]\|_2$$

- Parameter  $\beta$  trades off tracking for precision
- **Synthetic dataset**  $\{\mathbf{x}_t\}_{t=1}^{1000}$  and  $\{\mathbf{x}_t\}_{t=1000}^{2000}$  drawn from two subspaces



Successful use of affordable budget for nonlinear subspace tracking

Smaller values of  $\beta$  provide **faster tracking**, while larger values increase **fitting precision** 

### Online physical activity tracking

- Subjects asked to do various physical activities
- d =13 quantities measured from chest, ankle and wrist via wireless MUs
- > Gaussian kernel; and (r, B) = (10, 15)





#### Average LS-fit

Code	Activity	$\beta = 1$	$\beta = 0.9$	FIFO Bud.
0.3	Walking	0.099	0.074	0.074
	_	$\pm 0.016$	$\pm 0.012$	$\pm 0.012$
0.4	Running	0.227	0.187	0.187
		$\pm 0.025$	$\pm 0.022$	$\pm 0.022$
0.5	Cycling	0.058	0.028	0.028
		$\pm 0.027$	$\pm 0.012$	$\pm 0.12$
0.6	Nordic	0.130	0.103	0.103
	Walking	$\pm 0.020$	$\pm 0.016$	$\pm 0.016$
0.7	Ascending	0.079	0.063	0.063
	Stairs	$\pm 0.022$	$\pm 0.018$	$\pm 0.018$
0.8	Descending	0.094	0.066	0.065
	Stairs	$\pm 0.021$	$\pm 0.016$	$\pm 0.016$
0.9	Vacuum	0.045	0.029	0.029
	cleaning	$\pm 0.013$	$\pm 0.008$	$\pm 0.008$
1.0	Rope	0.272	0.238	0.238
	jumping	$\pm 0.063$	$\pm 0.057$	$\pm 0.057$

F. Sheikholeslami, D. K. Berberidis, and G. B. Giannakis, "Kernel-based Low-rank Feature Extraction on a Budget for Big Data Streams," arxiv:1601.07947.

### Summary

- Kernel-based learning
  - Sparsity-aware low-rank approximation of lifted data
  - Nuclear norm regularization for online approximation
  - Budget enforcement for affordable memory and computation
  - Theoretical guarantees and test results
- Future directions
  - Nonlinear feature extraction for canonical correlation analysis
  - Kernel-based feature extraction over 2-D signals on networks

Thank you!