



Scalable Kernel-based Learning via Low-rank Approximation of Lifted Data

Fatemeh Sheikholeslami

Acknowledgement

Prof. Georgios B. Giannakis

NSF grants 1343248, 1442686, 1500713, 1514056, and NIH grant 1R01GM104975-01

DTC and Dept. of ECE, University of Minnesota

Oct. 4th 2017

55th Allerton Conf. on Comm., Control, and Computing

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

- Non-linear (kernel-based) learning
- Low rank approximation of lifted data
- Sparsity regularized approximation
- Theoretical guarantees
- Simulation tests

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 Lift via nonlinear map $\mathbf{x} \to \boldsymbol{\phi}(\mathbf{x}) \in \mathbb{R}^{\bar{D}}$ 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})}{e.g., \ k(\mathbf{x}_{i}, \mathbf{x}_{j}) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \xrightarrow{\mathbf{x} \times \mathbf{x}_{i}} \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i}}{\mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i} \times \mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} - \mathbf{x}_{i}\|^{2}}{2\sigma^{2}}) \left(\frac{\mathbf{x}_{i} \times \mathbf{x}_{i} \times \mathbf{x}_{i} \times \mathbf{x}_{i}} \right) = \exp(\frac{-\|\mathbf{x}_{i} \times \mathbf{x}_{i} \times \mathbf{x}_{i} \times \mathbf{x}_{i}} \right)$$

RKHS basis expansion

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



□ 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

igcup Recall the nonlinear mapping $\mathbf{x} o oldsymbol{\phi}(\mathbf{x})$

$$\label{eq:Low-rank approximation} \begin{split} \square \text{ Low-rank approximation} \phi(\mathbf{x}) \simeq \hat{\phi}(\mathbf{x}) := \bar{\mathbf{L}} \mathbf{q} \ \text{ where } \ \bar{\mathbf{L}} \in \mathbb{R}^{\bar{D} \times r} \\ \mathbf{q} \in \mathbb{R}^{r \times 1} \end{split}$$

Regularized least-squares minimization

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

□ Alternating minimization

• Update virtual subspace $ar{\mathbf{L}}$ while fixing feature vectors $\{\mathbf{q}_i\}_{i=1}^N$

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

• Update vector of extracted features $\{\mathbf{q}_i\}_{i=1}^N$ while fixing subspace $\bar{\mathbf{L}}$

$$\Rightarrow$$
 $\mathbf{q}_i = (\mathbf{A}^\top \mathbf{K} \mathbf{A} + \lambda \mathbf{I}_r)^{-1} \mathbf{A}^\top \mathbf{k}_i$ where $\mathbf{K}_{N \times N} := \mathbf{\Phi}_N^\top \mathbf{\Phi}_N$

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.



Low-rank appr. with row-sparsity regularization

 \Box Row-sparsity on factor $\, A \, \Rightarrow$ fewer "basis" for the virtual subspace $\, \bar{L} \,$

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

Block coordinate descent

- Update vector of extracted features $\{\mathbf{q}_i\}_{i=1}^N$ while fixing subspace $\bar{\mathbf{L}}$ $\Rightarrow \mathbf{q}_i = (\mathbf{A}^\top \mathbf{K} \mathbf{A} + \lambda \mathbf{I}_r)^{-1} \mathbf{A}^\top \mathbf{k}_i$
- Update rows of subspace factor \mathbf{A} while fixing feature vectors $\{\mathbf{q}_i\}_{i=1}^N$

$$\Rightarrow \mathbf{a}_i[k+1] = \arg\min_{\mathbf{a}_i} \frac{1}{2} \mathbf{a}_i^\top \mathbf{H}_i \mathbf{a}_i + \mathbf{p}_i^\top \mathbf{a}_i + \eta \|\mathbf{a}_i\|_2$$

Group shrinkage

Group shrinkage- exact solver

□ Target minimization $\mathbf{a}_i[k+1] = \arg\min_{\mathbf{a}_i} F[k] + \eta \|\mathbf{a}_i\|_2$ where $F[k] := \frac{1}{2} \mathbf{a}_i^\top \mathbf{H}_i \mathbf{a}_i + \mathbf{p}_i^\top \mathbf{a}_i$ $\frac{1}{\text{Algorithm BCD-Ex}}$

Exact solver

$$\mathbf{a}_{i}[k+1] = \begin{cases} \mathbf{0} & \text{if } \|\mathbf{p}_{i}\|_{2} \leq \eta \\ \left(\mathbf{H}_{i} + \frac{\eta}{\Delta}\mathbf{I}_{r}\right)^{-1}\mathbf{p}_{i} & \text{o.w.} \end{cases}$$

where $\Delta := \|\mathbf{a}_i[k+1]\|$

- Not closed form!
- Root-finding iterations required

Input: data $\{x_i\}_{i=1}^N$, I_{\max} , or, ϵ Initialize $\mathbf{a}_1 = \mathbf{1}_{r \times 1}$, $\mathbf{a}_i = \mathbf{0}_{r \times 1}$ $\forall i \neq 1$ repeat for i = 1 to N do $\mathbf{q}_i[k+1] = (\mathbf{A}^\top \mathbf{K} \mathbf{A} + \lambda \mathbf{I}_r)^{-1} \mathbf{A}^\top \mathbf{k}_i$ end for for i = 1 to N do Set \mathbf{H}_i and \mathbf{p}_i . if $\|\mathbf{p}_i\| \leq \eta$ then $a_i[k+1] = 0$ else $\{\gamma_j, \mathbf{u}_j\}_{j=1}^r = \operatorname{svd}(\mathbf{H}_i)$ Using root finding algs., find Δ s.t. $\|\mathbf{y}_i(\Delta)\|^2 := \sum_k \frac{(\mathbf{u}_k^\top \mathbf{p}_i)^2}{(\gamma_k \Delta + n)^2} = 1$ $\mathbf{a}_i[k+1] = (\mathbf{H}_i + \frac{\eta}{\Lambda} \mathbf{I})^{-1} \mathbf{p}_i$ end if end for k = k + 1until $k = I_{\max}$ or $\frac{\|\mathbf{Q}[k+1] - \mathbf{Q}[k]\|_F}{\|\mathbf{Q}[k]\|_F} < \epsilon$

8

Z. Qin, K. Scheinberg, and D. Goldfarb, "Efficient block-coordinate descent algorithms for the group lasso," Mathematical Programming Computation, vol. 5, no. 2, pp. 143–169, 2013.

Group shrinkage- inexact solver

 $\Box \text{ Majorized minimization } \mathbf{a}_i[k+1] = \arg\min_{\mathbf{a}_i} \hat{F}_i(\mathbf{a}_i; \mathbf{a}_i[k]) + \eta \|\mathbf{a}_i\|_2$ $\hat{F}_i(\mathbf{a}_i; \mathbf{a}_i[k]) := F_i(\mathbf{a}_i[k]) + \langle \mathbf{a}_i - \mathbf{a}_i[k], \nabla F_i(\mathbf{a}_i[k]) \rangle + \frac{1}{2\mu_k} \|\mathbf{a}_i - \mathbf{a}_i[k]\|_2^2 \qquad \mu_k \le \|\nabla^2 F_i(\mathbf{a}_i)\|^{-1}$

- Proximal gradient
- Forward step

$$\hat{\mathbf{a}}_{i}[k+1] := \mathbf{a}_{i}[k] - \mu_{k} \nabla^{\top} \hat{F}_{i}(\mathbf{a}_{i};\mathbf{a}_{i}[k]) \Big|_{\mathbf{a}_{i}=\mathbf{a}_{i}[k]}$$
$$= \mathbf{a}_{i}[k] - \mu_{k} \Big(\mathbf{H}_{i} \mathbf{a}_{i}[k] + \mathbf{p}_{i} \Big)$$

• Backward step

$$\mathbf{a}_{i}[k+1] = \arg\min_{\mathbf{a}_{i}} \{ \frac{1}{2\mu_{k}} \| \mathbf{a}_{i} - \hat{\mathbf{a}}_{i}[k+1] \|_{2}^{2} + \eta \| \mathbf{a}_{i} \|_{2} \}$$

$$\mathbf{a}_{i}[k+1] = \begin{cases} \mathbf{0}, & \text{if } \|\hat{\mathbf{a}}_{i}[k+1]\| \leq \mu_{k}\eta \\ \left(1 + \frac{\mu_{k}\eta}{\|\hat{\mathbf{a}}_{i}[k+1]\| - \mu_{k}\eta}\right)^{-1} \hat{\mathbf{a}}_{i}[k+1], \text{o.w.} \end{cases}$$

AlgorithmBCD-PGInput:data $\{x_i\}_{i=1}^N$, I_{max} , or, ϵ Initialize $\mathbf{a}_1 = \mathbf{1}_{r \times 1}$, $\mathbf{a}_i = \mathbf{0}_{r \times 1}$ $\forall i \neq 1$

repeat

for i = 1 to N do $\mathbf{q}_i[k+1] = (\mathbf{A}^\top \mathbf{K} \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^\top \mathbf{k}_i$ end for for i = 1 to N do set \mathbf{H}_i and \mathbf{p}_i . $\hat{\mathbf{a}}_i[k+1] = \mathbf{a}_i[k] - \mu_k \left(\mathbf{H}_i \mathbf{a}_i[k] + \mathbf{p}_i\right)$ if $\|\hat{\mathbf{a}}_i[k+1]\| \le \mu_k \eta$ then $\mathbf{a}_i[k+1] = \mathbf{0}$ else $\Delta = \|\hat{\mathbf{a}}_i[k+1]\| - \mu_k \eta$ $\mathbf{a}_i[k+1] = (1 + \frac{\eta\mu_k}{\Delta})^{-1}\hat{\mathbf{a}}_i[k+1]$ end if end for k = k + 1until $k = I_{\text{max}}$ or $\frac{\|\mathbf{Q}[k+1] - \mathbf{Q}[k]\|_F}{\|\mathbf{Q}[k]\|_F} < \epsilon$

Convergence and generalization bounds

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

PG algorithms converge to a stationary point of the minimization.

Generalization: performance guarantees on unseen data

• Consider the equivalent optimization $\min_{A \in A}$

where
$$\begin{aligned} \mathcal{Q} &:= \{ \mathbf{q} \, | \, \mathbf{q} \in \mathbb{R}^r \, , \| \mathbf{q} \|_2 < B \} \\ \mathcal{A} &:= \{ \mathbf{A} | \sum_{i=1}^N \| \mathbf{a}_i \|_2 \leq \tau \} \end{aligned}$$

$$\min_{\mathbf{A}\in\mathcal{A}}\frac{1}{N}\sum_{i=1}^{N}\min_{\mathbf{q}_{i}\in\mathcal{Q}}\|\boldsymbol{\phi}(\mathbf{x}_{i})-\boldsymbol{\Phi}_{N}\mathbf{A}\mathbf{q}_{i}\|_{\mathcal{H}}^{2}$$
$$=:\ell(\mathbf{x},\bar{\mathbf{L}})$$

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

Proposition 2. Assume that $\|\mathcal{L}\|_{\mathcal{Q}} \geq 1$, and $\kappa(\mathbf{x}, \mathbf{x}) \leq \kappa, \forall \mathbf{x}$. Then, for a fixed $\delta > 0$, with probability at least $1 - \delta$ we have $\forall \bar{\mathbf{L}} \in \mathcal{L}$ $\mathbb{E}[\ell(\mathbf{x}; \bar{\mathbf{L}})] - 1/N \sum_{i=1}^{N} \ell(\mathbf{x}_i; \bar{\mathbf{L}}) \leq \frac{d}{\sqrt{N}} \left(14\kappa B\tau + \frac{\kappa}{2} \sqrt{\ln(16N\kappa^2 B^2 \tau^2)} \right) + \kappa \sqrt{\frac{\ln(1/\delta)}{2N}}$

Linearized kernel regression and classification



Bounds also on support vector machines for regression and classification





IJCNN Dataset N = 140KBCD-Ex (r=50) BCD-PG (r=50) BCD-Ex (r=100) BCD-Ex (r=200) - BCD-PG (r=200) cost 0 10² 10-2 10⁰ 10⁴ CPU time(sec) -BCD-Ex BCD-PG --- Nyst K missmatch 🔶 Prb-Nyst -X-ImpNyst 10 SVD 10² 50 100 150 0 200 Rank (r) 10⁴ Run time (sec) 10² 10⁰ 10⁻² 0 50 100 150 200 Rank (r)

Slice Dataset N = 53K



10⁻¹

10

20

30

Rank (r)

40

50

Summary

- Kernel-based learning
 - Low-rank approximation of lifted data
 - Sparsity regularized approximation
 - Theoretical guarantees and test results
- Future work
 - Utilization of other regularization techniques
 - Exploitation of censoring to mitigate number of updates
 - Budgeted techniques for bounded memory and complexity growth

