



Fast Kernel Methods for Generic Lipschitz Losses via p -Sparsified Sketches

CAp 2023

Tamim El Ahmad^{*}, Pierre Laforgue[†], Florence d'Alché-Buc^{*}

^{*} LTCI, Télécom Paris, Institut Polytechnique de Paris

[†] Università degli Studi di Milano

July 5, 2023

Supervised Scalar Regression

We have:

- i.i.d. training sample $(x_i, y_i)_{i=1}^n \in (\mathcal{X}, \mathbb{R})^n \sim P$
- loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$

Goal: Approach $f^* = \arg \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_{(X, Y) \sim P} [\ell(f(X), Y)]$ (ERM).

Supervised Scalar Regression

We have:

- i.i.d. training sample $(x_i, y_i)_{i=1}^n \in (\mathcal{X}, \mathbb{R})^n \sim P$
- loss function $\ell : \mathbb{R} \times \mathbb{R} \rightarrow [0, \infty)$

Goal: Approach $f^* = \arg \inf_{f: \mathcal{X} \rightarrow \mathbb{R}} \mathbb{E}_{(X, Y) \sim P} [\ell(f(X), Y)]$ (ERM).

$\mathcal{F} = \{f : \mathcal{X} \rightarrow \mathbb{R}\}$ is too large: which hypothesis space?

Reminder: positive definite kernels and Reproducing Kernel Hilbert Space

Positive definite kernel: $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that

- for all $(x, x') \in \mathcal{X}^2$, $k(x, x') = k(x', x)$
- for all $n \in \mathbb{N}$ and any $(x_i, \alpha_i)_{i=1}^n \in (\mathcal{X} \times \mathbb{R})^n$, $\sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) \geq 0$

RKHS (Aronszajn, 1950): k is uniquely associated to a Hilbert space \mathcal{H} of functions $f : \mathcal{X} \rightarrow \mathbb{R}$ s. t. for all $f \in \mathcal{H}$ and $x \in \mathcal{X}$

1. $x' \mapsto k(x, x') \in \mathcal{H}$,
2. $\langle f, k(\cdot, x) \rangle_{\mathcal{H}} = f(x)$ (**reproducing property**).

Kernel-Based Regression

Given k and its associated RKHS \mathcal{H} , $\lambda_n > 0$

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) + \frac{\lambda_n}{2} \|f\|_{\mathcal{H}}^2.$$

Kernel-Based Regression

Given k and its associated RKHS \mathcal{H} , $\lambda_n > 0$

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) + \frac{\lambda_n}{2} \|f\|_{\mathcal{H}}^2.$$

Representer Theorem: $\hat{f} = \sum_{j=1}^n k(\cdot, x_j) \hat{\alpha}_j$, where

$$(\hat{\alpha}_1, \dots, \hat{\alpha}_n)^\top = \hat{\alpha} = \arg \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n \ell([K\alpha]_i^\top, y_i) + \frac{\lambda_n}{2} \alpha^\top K \alpha.$$

Optimisation problem on n parameters: can we reduce n ?

Table of contents

1. Sketched Kernel Machines
2. p -Sparsified Sketches
3. Experiments
4. Conclusion

Sketched Kernel Machines

First Idea: Sub-Sampling, i.e. Nyström Approximation

$\tilde{f} = \sum_{j=1}^s k(\cdot, x_j) \tilde{\gamma}_j$, where

$$(\tilde{\gamma}_1, \dots, \tilde{\gamma}_s)^\top = \tilde{\gamma} = \arg \min_{\gamma \in \mathbb{R}^s} \frac{1}{n} \sum_{i=1}^n \ell \left(\left[\begin{array}{c} \underbrace{K_{ns}}_{n \times s} \gamma \\ \vdots \end{array} \right]^\top, y_i \right) + \frac{\lambda_n}{2} \gamma^\top \underbrace{K_{ss}}_{s \times s} \gamma.$$

First Idea: Sub-Sampling, i.e. Nyström Approximation

$\tilde{f} = \sum_{j=1}^s k(\cdot, x_j) \tilde{\gamma}_j$, where

$$(\tilde{\gamma}_1, \dots, \tilde{\gamma}_s)^\top = \tilde{\gamma} = \arg \min_{\gamma \in \mathbb{R}^s} \frac{1}{n} \sum_{i=1}^n \ell \left(\left[\begin{array}{c} \underbrace{K_{ns}}_{n \times s} \gamma \\ \vdots \end{array} \right]^\top, y_i \right) + \frac{\lambda_n}{2} \gamma^\top \underbrace{K_{ss}}_{s \times s} \gamma.$$

Sampling the wrong data can lead to poor results \implies
data-dependent sampling schemes (e.g. leverage scores) (Alaoui and Mahoney, 2015; Musco and Musco, 2017; Rudi et al., 2018; Chen and Yang, 2021b)

Sub-Sampling is Random Projection

Let $n = 5$, $X = \{x_1, \dots, x_5\}$, $k_X^x = (k(x, x_1), \dots, k(x, x_5))$, $s = 2$ and

$$S = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

Sub-Sampling is Random Projection

Let $n = 5$, $X = \{x_1, \dots, x_5\}$, $k_X^x = (k(x, x_1), \dots, k(x, x_5))$, $s = 2$ and

$$S = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$K_{Sn} = \begin{pmatrix} k_X^{x_1} \\ k_X^{x_4} \end{pmatrix} = SK \quad \text{and} \quad K_{Ss} = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_4) \\ k(x_4, x_1) & k(x_4, x_4) \end{pmatrix} = SKS^T$$

Sub-Sampling is Random Projection

Let $n = 5$, $X = \{x_1, \dots, x_5\}$, $k_X^x = (k(x, x_1), \dots, k(x, x_5))$, $s = 2$ and

$$S = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$K_{Sn} = \begin{pmatrix} k_X^{x_1} \\ k_X^{x_4} \end{pmatrix} = SK \quad \text{and} \quad K_{Ss} = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_4) \\ k(x_4, x_1) & k(x_4, x_4) \end{pmatrix} = SKS^T$$

$\tilde{f} = \sum_{j=1}^s k(\cdot, x_{i_j}) \tilde{\gamma}_j = \sum_{j=1}^n k(\cdot, x_{i_j}) [S^T \tilde{\gamma}]_j$, where

$$(\tilde{\gamma}_1, \dots, \tilde{\gamma}_s)^T = \tilde{\gamma} = \arg \min_{\gamma \in \mathbb{R}^s} \frac{1}{n} \sum_{i=1}^n \ell \left([KS^T \gamma]_i^T, y_i \right) + \frac{\lambda_n}{2} \gamma^T SKS^T \gamma.$$

Sub-Sampling is Random Projection

Let $n = 5$, $X = \{x_1, \dots, x_5\}$, $k_X^x = (k(x, x_1), \dots, k(x, x_5))$, $s = 2$ and

$$S = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$K_{Sn} = \begin{pmatrix} k_X^{x_1} \\ k_X^{x_4} \end{pmatrix} = SK \quad \text{and} \quad K_{Ss} = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_4) \\ k(x_4, x_1) & k(x_4, x_4) \end{pmatrix} = SKS^T$$

$\tilde{f} = \sum_{j=1}^s k(\cdot, x_{i_j}) \tilde{\gamma}_j = \sum_{j=1}^n k(\cdot, x_{i_j}) [S^T \tilde{\gamma}]_j$, where

$$(\tilde{\gamma}_1, \dots, \tilde{\gamma}_s)^T = \tilde{\gamma} = \arg \min_{\gamma \in \mathbb{R}^s} \frac{1}{n} \sum_{i=1}^n \ell \left([KS^T \gamma]_i^T, y_i \right) + \frac{\lambda_n}{2} \gamma^T SKS^T \gamma.$$

Could we use other random matrix distributions?

Johnson-Linderstrauss Lemma

Lemma

Given $0 < \varepsilon < 1$, a set \mathcal{S} of n points in \mathbb{R}^D , and an integer $d > 8(\log n)/\varepsilon^2$, there is a linear map $h : \mathbb{R}^D \rightarrow \mathbb{R}^d$ such that

$$(1 - \varepsilon) \|u - v\|^2 \leq \|h(u) - h(v)\|^2 \leq (1 + \varepsilon) \|u - v\|^2,$$

for all $u, v \in \mathcal{S}$.

Johnson-Linderstrauss Lemma

Lemma

Given $0 < \varepsilon < 1$, a set \mathcal{S} of n points in \mathbb{R}^D , and an integer $d > 8(\log n)/\varepsilon^2$, there is a linear map $h : \mathbb{R}^D \rightarrow \mathbb{R}^d$ such that

$$(1 - \varepsilon) \|u - v\|^2 \leq \|h(u) - h(v)\|^2 \leq (1 + \varepsilon) \|u - v\|^2,$$

for all $u, v \in \mathcal{S}$.

Most famous proof:

1. take $h = \frac{1}{\sqrt{d}}S \in \mathbb{R}^{d \times D}$, where $S_{ij} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1) \implies$ **Gaussian sketching**
2. prove the above equation with high probability

Gaussian sketching then?

$$(\tilde{\gamma}_1, \dots, \tilde{\gamma}_s)^\top = \hat{\gamma} = \arg \min_{\gamma \in \mathbb{R}^s} \frac{1}{n} \sum_{i=1}^n \ell([KS^\top \gamma]_i, y_i) + \frac{\lambda_n}{2} \gamma^\top SKS^\top \gamma.$$

Gaussian sketching then?

$$(\tilde{\gamma}_1, \dots, \tilde{\gamma}_s)^\top = \hat{\gamma} = \arg \min_{\gamma \in \mathbb{R}^s} \frac{1}{n} \sum_{i=1}^n \ell([KS^\top \gamma]_i, y_i) + \frac{\lambda_n}{2} \gamma^\top SKS^\top \gamma.$$

Problems:

1. computing SK : $\mathcal{O}(n^2s)$ time complexity \rightarrow **still high complexity**
2. storing K : $\mathcal{O}(n^2)$ space complexity \rightarrow **space complexity does not change**

Which property should sketching distributions satisfy?

- $K/n = UDU^T$
- $D = \text{diag}(\mu_1, \dots, \mu_n)$ where $\mu_1 \geq \dots \geq \mu_n$
- δ_n^2 the lowest value s. t. $\psi(\delta_n) = \left(\frac{1}{n} \sum_{i=1}^n \min(\delta_n^2, \mu_i)\right)^{1/2} \leq \delta_n^2$
(Bartlett et al., 2005)
- $d_n = \min \{j \in \{1, \dots, n\} : \mu_j \leq \delta_n^2\}$

Which property should sketching distributions satisfy?

- $K/n = UDU^T$
- $D = \text{diag}(\mu_1, \dots, \mu_n)$ where $\mu_1 \geq \dots \geq \mu_n$
- δ_n^2 the lowest value s. t. $\psi(\delta_n) = \left(\frac{1}{n} \sum_{i=1}^n \min(\delta_n^2, \mu_i)\right)^{1/2} \leq \delta_n^2$
(Bartlett et al., 2005)
- $d_n = \min \{j \in \{1, \dots, n\} : \mu_j \leq \delta_n^2\}$

Definition (K -satisfiability (Yang et al., 2017))

Let $c > 0$ independent of n . Let $U_1 \in \mathbb{R}^{n \times d_n}$ and $U_2 \in \mathbb{R}^{n \times (n-d_n)}$ be the left and right blocks of matrix U previously defined, and $D_2 = \text{diag}(\mu_{d_n+1}, \dots, \mu_n)$. A sketch matrix S is said to be K -satisfiable for c if S is such that

$$\left\| (SU_1)^T SU_1 - I_{d_n} \right\|_{\text{op}} \leq 1/2, \quad \text{and} \quad \left\| SU_2 D_2^{1/2} \right\|_{\text{op}} \leq c\delta_n.$$

Intuition: S is K -satisfiable \implies isometry on the largest eigenvectors of K/n and small operator norm on the smallest eigenvectors

p -Sparsified Sketches

Definition

Let $s < n$, and $p \in (0, 1]$. A p -sparsified sketch $S \in \mathbb{R}^{s \times n}$ is composed of i.i.d. entries

$$S_{ij} = \frac{1}{\sqrt{sp}} B_{ij} R_{ij},$$

where $B_{ij} \stackrel{\text{i.i.d.}}{\sim} \text{Ber}(p)$ and $R_{ij} \stackrel{\text{i.i.d.}}{\sim} \text{Rad}(\frac{1}{2})$ (p -SR) or $\mathcal{N}(0, 1)$ (p -SG).

K -Satisfiability of p -Sparsified Sketches

Theorem

Let S be a p -sparsified sketch. Then, there are some universal constants $C_0, C_1 > 0$ and a constant $c(p)$, increasing with p , such that for $s \geq \max(C_0 d_n / p^2, \delta_n^2 n)$ and with a probability at least $1 - C_1 e^{-sc(p)}$, the sketch S is K -satisfiable for $c = \frac{2}{\sqrt{p}} \left(1 + \sqrt{\log(5)}\right) + 1$.

Intuitive behavior of p :

- $p = 1$: we recover Yang et al. (2017)'s result for Gaussian sketching
- the larger it is, the denser S is, and the more likely S is K -satisfiable
- the smaller it is, the larger s is needed

Computational Property: *Decomposition trick*

Let $s' = \sum_{j=1}^n \mathbb{I}\{S_{:j} \neq 0_S\}$,

$$S = S_{SG} S_{SS},$$

where

- $S_{SG} \in \mathbb{R}^{s \times s'}$: **sparse sub-gaussian sketch** obtained by deleting the null columns from S
- $S_{SS} \in \mathbb{R}^{s' \times n}$: **sub-sampling sketch** obtained by sampling the rows of I_n corresponding to the indices of non-zero columns of S

Example:

$$\begin{pmatrix} 1 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 0 & -1 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$s' \sim \text{Binom}(n, 1 - (1 - p)^S) \implies \mathbb{E}[s'] = n(1 - (1 - p)^S) \underset{\rho \rightarrow 0}{\sim} nsp$$

Time and Space Complexities

Let C_k = cost of computing $k(x, x')$, complexities of **Gaussian** vs **p -sparsified** sketch:

Time: $\mathcal{O}(C_k n^2 + n^2 s)$ vs $\mathcal{O}(C_k n^2 s p + n^2 s^2 p)$

Space: $\mathcal{O}(n^2)$ vs $\mathcal{O}(n^2 s p)$

Time and Space Complexities

Let C_k = cost of computing $k(x, x')$, complexities of Gaussian vs p -sparsified sketch:

Time: $\mathcal{O}(C_k n^2 + n^2 s)$ vs $\mathcal{O}(C_k n^2 s p + n^2 s^2 p)$

Space: $\mathcal{O}(n^2)$ vs $\mathcal{O}(n^2 s p)$

p -sparsified sketch's goal \rightarrow best of both worlds:

1. computational efficiency of sub-sampling sketch
2. statistical accuracy of Rademacher or Gaussian sketch

Time and Space Complexities

Let C_k = cost of computing $k(x, x')$, complexities of **Gaussian** vs **p -sparsified** sketch:

Time: $\mathcal{O}(C_k n^2 + n^2 s)$ vs $\mathcal{O}(C_k n^2 s p + n^2 s^2 p)$

Space: $\mathcal{O}(n^2)$ vs $\mathcal{O}(n^2 s p)$

p -sparsified sketch's goal \rightarrow best of both worlds:

1. computational efficiency of sub-sampling sketch
2. statistical accuracy of Rademacher or Gaussian sketch

Related works:

1. sub-sampling sketch with data-dependent sampling schemes (e.g. leverage scores) (Alaoui and Mahoney, 2015; Musco and Musco, 2017; Rudi et al., 2018; Chen and Yang, 2021b)
2. accumulation sketch (Chen and Yang, 2021a): sum of sub-sampling sketches

Experiments

Scalar regression with synthetic dataset: settings

1) $n = 10,000$, $(x_i, y_i) \in \mathbb{R}^{10} \times \mathbb{R}$

2) Inhomogeneous input data distribution

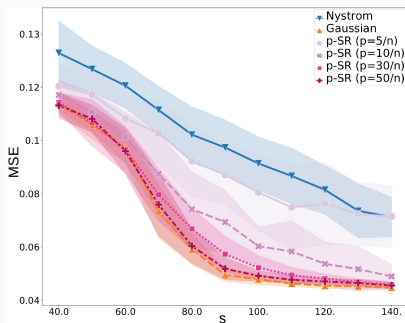
$$x_i \sim \begin{cases} \mathcal{U}([0_{10}, \mathbb{1}_{10}]) , & \text{if } i = 1, \dots, 9,900 , \\ \mathcal{N}(1.5\mathbb{1}_{10}, 0.25I_{10}) , & \text{if } i = 9,901, \dots, 10,000 , \end{cases}$$

3) $y = f^*(x) + \epsilon$, where $\epsilon \sim \mathcal{N}(0, 1)$ and

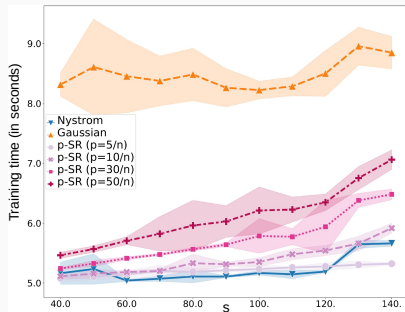
$$f^*(x) = 0.1 \exp(4x_1) + \frac{4}{1 + \exp(-20(x_2 - 0.5))} + 3x_3 + 2x_4 + x_5 .$$

4) loss: κ -Huber

Scalar regression with synthetic dataset



(a) Test relative MSE w.r.t. sketch size s



(b) Training time (sec) w.r.t. sketch size s

Scalar regression with synthetic dataset

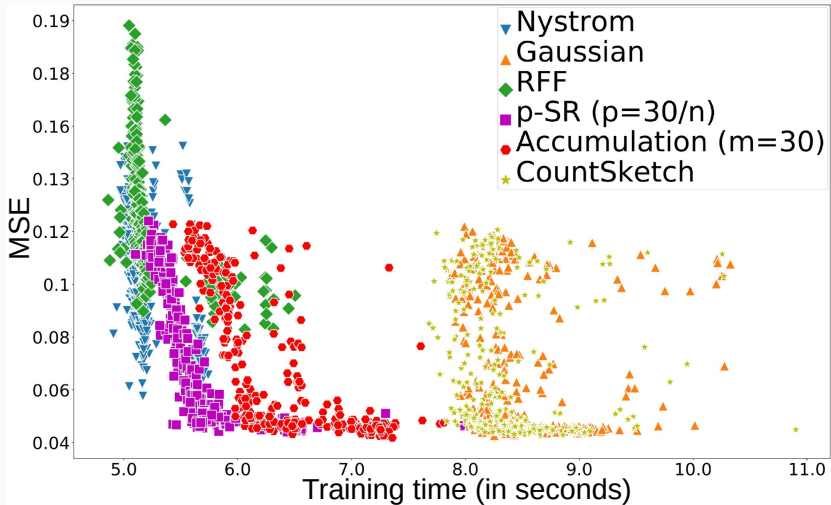


Figure 2: Test relative MSE w.r.t. training times

Conclusion

Conclusion

- Extend scalar regression framework to **multi-output** regression thanks to **decomposable matrix-valued kernels**

Conclusion

- Extend scalar regression framework to **multi-output** regression thanks to **decomposable matrix-valued kernels**
- Extend previous results and provide **excess risk bounds** for the **multiple output setting** and with any **generic Lipschitz loss** thanks to decomposable matrix-valued kernels

Conclusion

- Extend scalar regression framework to **multi-output** regression thanks to **decomposable matrix-valued kernels**
- Extend previous results and provide **excess risk bounds** for the **multiple output setting** and with any **generic Lipschitz loss** thanks to decomposable matrix-valued kernels
- Provide new K -satisfiable sketching distribution – **p -sparsified** – well-suited to kernel methods thanks to the *decomposition trick*

Conclusion

- Extend scalar regression framework to **multi-output** regression thanks to **decomposable matrix-valued kernels**
- Extend previous results and provide **excess risk bounds** for the **multiple output setting** and with any **generic Lipschitz loss** thanks to decomposable matrix-valued kernels
- Provide new K -satisfiable sketching distribution – **p -sparsified** – well-suited to kernel methods thanks to the *decomposition trick*
- When the input data distribution shows some inhomogeneity, p -sparsified sketches
 1. **outperform Nyström approximation and RFFs**
 2. **compete with statistically accurate sketches (Gaussian, CountSketch, Accumulation) while being faster**

Conclusion

- Extend scalar regression framework to **multi-output** regression thanks to **decomposable matrix-valued kernels**
- Extend previous results and provide **excess risk bounds** for the **multiple output setting** and with any **generic Lipschitz loss** thanks to decomposable matrix-valued kernels
- Provide new K -satisfiable sketching distribution – **p -sparsified** – well-suited to kernel methods thanks to the *decomposition trick*
- When the input data distribution shows some inhomogeneity, p -sparsified sketches
 1. **outperform Nyström approximation and RFFs**
 2. **compete with statistically accurate sketches (Gaussian, CountSketch, Accumulation) while being faster**
- Sketched kernel algorithms show **similar performances** – and even outperform in some cases – non-sketched kernel algorithms, while **being significantly faster**

References

- Alaoui, A. and Mahoney, M. W. (2015). Fast randomized kernel ridge regression with statistical guarantees. In Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., and Garnett, R., editors, *Advances in Neural Information Processing Systems (NeurIPS)*, volume 28.
- Aronszajn, N. (1950). Theory of reproducing kernels. *Transactions of the American Mathematical Society*, pages 337–404.
- Bartlett, P. L., Bousquet, O., and Mendelson, S. (2005). Local rademacher complexities. *Ann. Statist.*, 33(4):1497–1537.

- Chen, Y. and Yang, Y. (2021a). Accumulations of projections—a unified framework for random sketches in kernel ridge regression. In *International Conference on Artificial Intelligence and Statistics*, pages 2953–2961. PMLR.
- Chen, Y. and Yang, Y. (2021b). Fast statistical leverage score approximation in kernel ridge regression. In *International Conference on Artificial Intelligence and Statistics*, pages 2935–2943. PMLR.
- Harrison Jr, D. and Rubinfeld, D. L. (1978). Hedonic housing prices and the demand for clean air. *Journal of environmental economics and management*, 5(1):81–102.
- Koenker, R. (2005). *Quantile regression*. Cambridge university press.

References iii

- Moen, E., Handegard, N. O., Allken, V., Albert, O. T., Harbitz, A., and Malde, K. (2018). Automatic interpretation of otoliths using deep learning. *PLoS One*, 13(12):e0204713.
- Musco, C. and Musco, C. (2017). Recursive sampling for the nyström method. *Advances in Neural Information Processing Systems*, 2017:3834–3846.
- Ordoñez, A., Eikvil, L., Salberg, A.-B., Harbitz, A., Murray, S. M., and Kampffmeyer, M. C. (2020). Explaining decisions of deep neural networks used for fish age prediction. *PloS one*, 15(6):e0235013.
- Rudi, A., Calandriello, D., Carratino, L., and Rosasco, L. (2018). On fast leverage score sampling and optimal learning. In *NeurIPS*.
- Sangnier, M., Fercoq, O., and d’Alché Buc, F. (2016). Joint quantile regression in vector-valued RKHSs. In *Advances in Neural Information Processing Systems (NeurIPS)*, Barcelona, France.

- Spyromitros-Xioufis, E., Tsoumakas, G., Groves, W., and Vlahavas, I. (2016). Multi-target regression via input space expansion: treating targets as inputs. *Machine Learning*, 104(1):55–98.
- Yang, Y., Pilanci, M., Wainwright, M. J., et al. (2017). Randomized sketches for kernels: Fast and optimal nonparametric regression. *The Annals of Statistics*, 45(3):991–1023.

Lipschitz Losses

$\ell(y, y') = g(y - y')$, where g is:

- **For κ -Huber:** For $\kappa > 0$:

$$\forall y \in \mathcal{Y}, g(y) = \begin{cases} \frac{1}{2} \|y\|_{\mathcal{Y}}^2 & \text{if } \|y\|_{\mathcal{Y}} \leq \kappa \\ \kappa (\|y\|_{\mathcal{Y}} - \frac{\kappa}{2}) & \text{otherwise} \end{cases} .$$

- **The pinball loss (Koenker, 2005) for joint quantile regression:**
For d quantile levels, $\tau_1 < \tau_2 < \dots < \tau_d$ with $\tau_i \in (0, 1)$, we define:

$$\ell_{\tau}(f(x), y) = L_{\tau}(f(x) - y \mathbb{1}_d),$$

with the following definition for L_{τ} the extension of pinball loss to \mathbb{R}^d (Sangnier et al., 2016):

For $r \in \mathbb{R}^d$:

$$L_{\tau}(r) = \sum_{j=1}^d \begin{cases} \tau_j r_j & \text{if } r_j \geq 0, \\ (\tau_j - 1) r_j & \text{if } r_j < 0. \end{cases}$$

Example: Kernel Ridge Multi-Output Regression

With $\mathcal{K} = kl_d$

- Without sketching: $\hat{A} = (K + n\lambda I_n)^{-1} Y \implies$ inversion of $n \times n$ matrix
- With sketching: $\tilde{\Gamma} = (SK^2S^\top + n\lambda SKS^\top)^{-1} SKY \implies$ inversion of $s \times s$ matrix

Previous work

Settings in Yang et al. (2017):

- $d = 1 \implies$ scalar regression only
- $\ell(y, y') = (y - y')^2 \implies$ KRR only
- $y_i = f^*(x_i) + \sigma\omega_i$, where ω_i s i.i.d. standard Gaussian variates
- Focus on the squared $L^2(\mathbb{P}_n)$ error, i.e.,
$$\left\| \tilde{f}_s - f^* \right\|_n^2 = \frac{1}{n} \sum_{i=1}^n \left(\tilde{f}_s(x_i) - f^*(x_i) \right)^2 \implies$$
 not excess risk in expectation

Previous work

Settings in Yang et al. (2017):

- $d = 1 \implies$ scalar regression only
- $\ell(y, y') = (y - y')^2 \implies$ KRR only
- $y_i = f^*(x_i) + \sigma\omega_i$, where ω_i s i.i.d. standard Gaussian variates
- Focus on the squared $L^2(\mathbb{P}_n)$ error, i.e.,
$$\left\| \tilde{f}_s - f^* \right\|_n^2 = \frac{1}{n} \sum_{i=1}^n \left(\tilde{f}_s(x_i) - f^*(x_i) \right)^2 \implies$$
 not excess risk in expectation

Yang et al. (2017, Theorem 2): If $f^* \in \mathcal{H}$, then for any $\lambda \geq 2\delta_n^2$, with a probability greater than $1 - c_1 e^{-c_2 n \delta_n^2}$

$$\left\| \tilde{f}_s - f^* \right\|_n^2 \leq c_u (\lambda + \delta_n^2), \quad (1)$$

where c_u only depends on $\|f^*\|_{\mathcal{H}}$.

Theoretical Guarantees

Assumptions

A. 1: Expected risk is minimized over \mathcal{H} at $f_{\mathcal{H}} = \operatorname{arginf}_{f \in \mathcal{H}} \mathbb{E}[\ell(f(X), Y)]$.

Assumptions

A. 1: Expected risk is minimized over \mathcal{H} at $f_{\mathcal{H}} = \operatorname{arginf}_{f \in \mathcal{H}} \mathbb{E}[\ell(f(X), Y)]$.

A. 2: The hypothesis set considered is the unit ball $\mathcal{B}(\mathcal{H})$ of \mathcal{H} .

Assumptions

A. 1: Expected risk is minimized over \mathcal{H} at $f_{\mathcal{H}} = \operatorname{arginf}_{f \in \mathcal{H}} \mathbb{E}[\ell(f(X), Y)]$.

A. 2: The hypothesis set considered is the unit ball $\mathcal{B}(\mathcal{H})$ of \mathcal{H} .

A. 3: $\forall y \in \mathbb{R}^d, z \mapsto \ell(z, y)$ is L -Lipschitz over $\mathcal{H}(\mathcal{X}) = \{f(x) : f \in \mathcal{H}, x \in \mathcal{X}\}$.

Assumptions

A. 1: Expected risk is minimized over \mathcal{H} at

$$f_{\mathcal{H}} = \operatorname{arginf}_{f \in \mathcal{H}} \mathbb{E}[\ell(f(X), Y)].$$

A. 2: The hypothesis set considered is the unit ball $\mathcal{B}(\mathcal{H})$ of \mathcal{H} .

A. 3: $\forall y \in \mathbb{R}^d, z \mapsto \ell(z, y)$ is L -Lipschitz over

$$\mathcal{H}(\mathcal{X}) = \{f(x) : f \in \mathcal{H}, x \in \mathcal{X}\}.$$

A. 4: $\exists \kappa > 0$ s. t. $k(x, x) \leq \kappa, \forall x \in \mathcal{X}$ and M is non-singular.

Assumptions

A. 1: Expected risk is minimized over \mathcal{H} at

$$f_{\mathcal{H}} = \operatorname{arginf}_{f \in \mathcal{H}} \mathbb{E}[\ell(f(X), Y)].$$

A. 2: The hypothesis set considered is the unit ball $\mathcal{B}(\mathcal{H})$ of \mathcal{H} .

A. 3: $\forall y \in \mathbb{R}^d, z \mapsto \ell(z, y)$ is L -Lipschitz over

$$\mathcal{H}(\mathcal{X}) = \{f(x) : f \in \mathcal{H}, x \in \mathcal{X}\}.$$

A. 4: $\exists \kappa > 0$ s. t. $k(x, x) \leq \kappa, \forall x \in \mathcal{X}$ and M is non-singular.

A. 5: The sketch S is K -satisfiable for a $c > 0$ independent of n .

Excess Risk Bound

Theorem

Under **A. 1, 2, 3, 4 and 5**, let $C = 1 + \sqrt{6}c$, for any $\delta \in (0, 1)$, then with probability at least $1 - \delta$,

$$\begin{aligned}\mathbb{E} [\ell_{\tilde{f}}] &\leq \mathbb{E} [\ell_{f_{\mathcal{H}}}] + LC\sqrt{\lambda_n + \|M\|_{\text{op}} \delta_n^2} + \frac{\lambda_n}{2} \\ &\quad + 8L\sqrt{\frac{\kappa \text{Tr}(M)}{n}} + 2\sqrt{\frac{8 \log(4/\delta)}{n}}.\end{aligned}$$

If $\ell(z, y) = \|z - y\|_2^2 / 2$ and $\mathcal{Y} \subset \mathcal{B}(\mathbb{R}^d)$, then with probability at least $1 - \delta$,

$$\begin{aligned}\mathbb{E} [\ell_{\tilde{f}}] &\leq \mathbb{E} [\ell_{f_{\mathcal{H}}}] + \left(C^2 + \frac{1}{2}\right) \lambda_n + C^2 \|M\|_{\text{op}} \delta_n^2 \\ &\quad + 8 \text{Tr}(M)^{1/2} \frac{\kappa \|M\|_{\text{op}}^{1/2} + \kappa^{1/2}}{\sqrt{n}} + 2\sqrt{\frac{8 \log(4/\delta)}{n}}.\end{aligned}$$

Sketch of proof: Decomposition Error

$$\begin{aligned}\mathbb{E}[\ell_{\tilde{f}_s}] - \mathbb{E}[\ell_{f_{\mathcal{H}_k}}] &= \mathbb{E}_{(X,Y) \sim P}[\ell(\tilde{f}_s(X), Y)] - \frac{1}{n} \sum_{i=1}^n \ell(\tilde{f}_s(x_i), y_i) \leftarrow \text{gen. error} \\ &+ \frac{1}{n} \sum_{i=1}^n \ell(\tilde{f}_s(x_i), y_i) - \frac{1}{n} \sum_{i=1}^n \ell(f_{\mathcal{H}_k}(x_i), y_i) \leftarrow \text{approx. error} \\ &+ \frac{1}{n} \sum_{i=1}^n \ell(f_{\mathcal{H}_k}(x_i), y_i) - \mathbb{E}_{(X,Y) \sim P}[\ell(f_{\mathcal{H}_k}(X), Y)] \leftarrow \text{gen. error}\end{aligned}$$

Sketch of proof: Approximation Error

$$\text{Let } \mathcal{H}_S = \left\{ f = \sum_{i=1}^n k(\cdot, x_i) M \left[S^T \tilde{\Gamma} \right]_i \mid \gamma \in \mathbb{R}^{s \times d} \right\}$$

$$\frac{1}{n} \sum_{i=1}^n \ell(\tilde{f}_S(x_i), y_i) - \frac{1}{n} \sum_{i=1}^n \ell(f_{\mathcal{H}_K}(x_i), y_i)$$

$$\leq \inf_{\substack{f \in \mathcal{H}_S \\ \|f\|_{\mathcal{H}_K} \leq 1}} \frac{L}{n} \sum_{i=1}^n \|f(x_i) - f_{\mathcal{H}_K}(x_i)\|_2 \leftarrow \text{A. 2}$$

$$\leq L \inf_{\substack{f \in \mathcal{H}_S \\ \|f\|_{\mathcal{H}_K} \leq 1}} \sqrt{\frac{1}{n} \sum_{i=1}^n \|f(x_i) - f_{\mathcal{H}_K}(x_i)\|_2^2} \leftarrow \text{Jensen}$$

Which sketching distribution to use for kernels?

With $\mathcal{K} = kl_d$

- Without sketching: $\hat{A} = (K + n\lambda I_n)^{-1} Y \implies$ inversion of $n \times n$ matrix
- With sketching: $\hat{\Gamma} = (SK^2S^\top + n\lambda SKS^\top)^{-1} SKY \implies$ inversion of $s \times s$ matrix

Which sketching distribution to use for kernels?

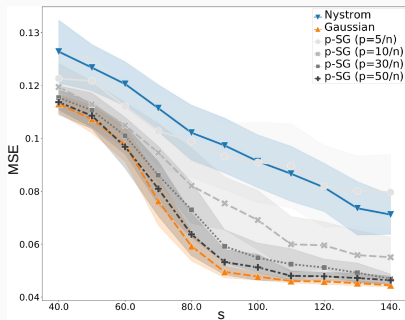
With $\mathcal{K} = kl_d$

- Without sketching: $\hat{A} = (K + n\lambda I_n)^{-1} Y \implies$ inversion of $n \times n$ matrix
- With sketching: $\hat{\Gamma} = (SK^2S^T + n\lambda SKS^T)^{-1} SKY \implies$ inversion of $s \times s$ matrix

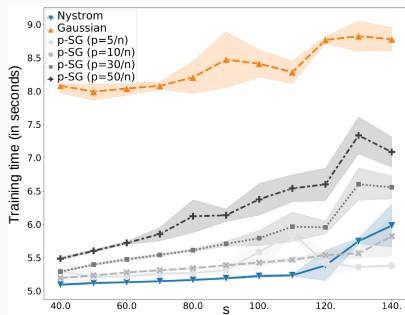
Problems:

1. computing SK : $\mathcal{O}(n^2s)$ time complexity \rightarrow still high complexity
2. storing K : $\mathcal{O}(n^2)$ space complexity \rightarrow space complexity does not change

Scalar regression with synthetic dataset



(a) Test relative MSE w.r.t. sketch size s



(b) Training time (sec) w.r.t. sketch size s

Scalar regression with synthetic dataset

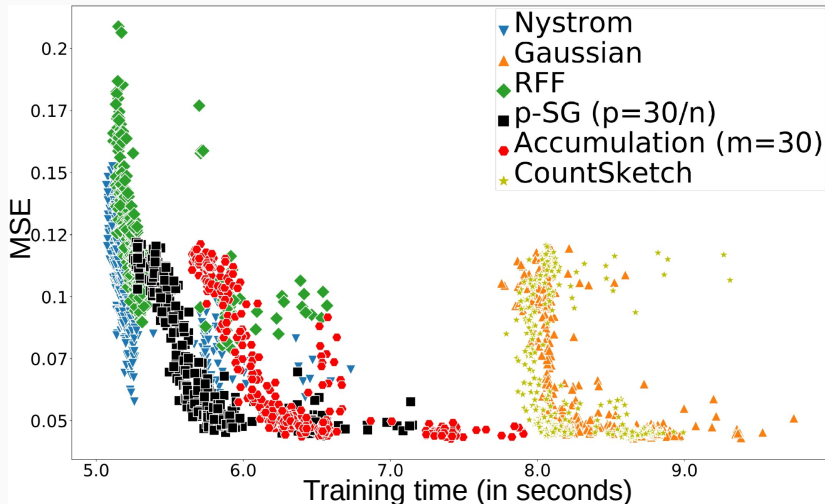


Figure 4: Test relative MSE w.r.t. training times with κ -Huber

Joint Quantile Regression on real data

- Boston dataset (Harrison Jr and Rubinfeld, 1978): house price prediction, $n = 506$
- Otoliths dataset (Moen et al., 2018; Ordoñez et al., 2020): fish age prediction, $n = 3780$

Quantile levels to predict: (0.1, 0.3, 0.5, 0.7, 0.9)

Table 1: Empirical test pinball and crossing loss and training times (in sec) without sketching and with sketching ($s = 50$).

Dataset	Metrics	w/o Sketch	$20/n_{tr}$ -SR	$20/n_{tr}$ -SG	Acc. $m = 20$
Boston	Pinball loss	51.28 ± 0.67	54.75 ± 0.74	54.78 ± 0.72	54.73 ± 0.75
	Crossing loss	0.34 ± 0.13	0.26 ± 0.08	0.11 ± 0.07	0.15 ± 0.07
	Training time	6.97 ± 0.25	1.43 ± 0.07	1.38 ± 0.08	1.48 ± 0.05
otoliths	Pinball loss	2.78	2.66 ± 0.02	2.64 ± 0.02	2.67 ± 0.03
	Crossing loss	5.18	5.46 ± 0.06	5.43 ± 0.05	5.46 ± 0.06
	Training time	606.8	20.4 ± 0.5	20.0 ± 0.3	22.1 ± 0.4

Multi-target Regression on real data

- rf1 and rf2 datasets (Spyromitros-Xioufis et al., 2016): river network flows prediction, $n = 4108, 4108$
- scm1d and scm20d datasets (Spyromitros-Xioufis et al., 2016): products price prediction, $n = 8145, 7463$

Table 2: ARRMSSE and training times (in sec) with square loss and $s = 100$ when using Sketching.

Dataset	Metrics	w/o Sketch	$20/n_{tr}$ -SR	$20/n_{tr}$ -SG	Acc. $m = 20$
rf1	ARRMSE	0.575	0.584 ± 0.003	0.583 ± 0.003	0.592 ± 0.001
	Training time	1.73	0.22 ± 0.025	0.25 ± 0.005	0.60 ± 0.0004
rf2	ARRMSE	0.578	0.671 ± 0.009	0.656 ± 0.006	0.796 ± 0.006
	Training time	1.77	0.28 ± 0.003	0.27 ± 0.003	0.82 ± 0.003
scm1d	ARRMSE	0.418	0.422 ± 0.002	0.423 ± 0.001	0.423 ± 0.001
	Training time	9.36	0.45 ± 0.022	0.45 ± 0.019	0.86 ± 0.006
scm20d	ARRMSE	0.755	0.754 ± 0.003	0.754 ± 0.003	0.753 ± 0.001
	Training time	6.16	0.38 ± 0.016	0.38 ± 0.017	0.70 ± 0.032