

Sketch In, Sketch Out: Accelerating both Learning and Inference for Structured Prediction with Kernels

Journée de Statistique 2024

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

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

[†] L2S, CentraleSupélec

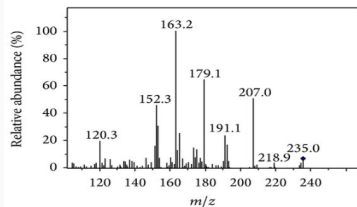
[‡] Università degli Studi di Milano

May 28, 2024

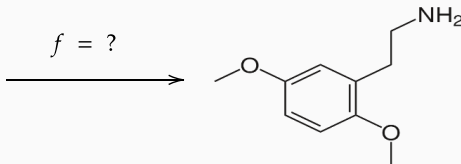
Structured Prediction

Goal: learn a mapping $f: \mathcal{X} \rightarrow \mathcal{Y}$ with \mathcal{Y} a space of structured objects (graphs, rankings, sequences, binary vectors, etc.).

MS/MS spectra



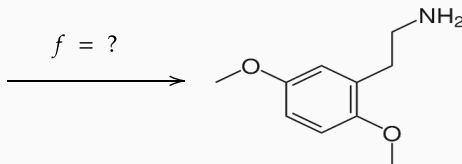
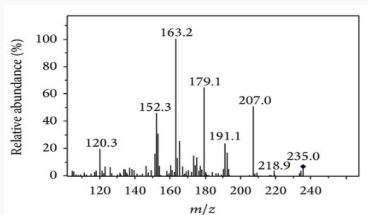
$f = ?$



Structured Prediction

Goal: learn a mapping $f: \mathcal{X} \rightarrow \mathcal{Y}$ with \mathcal{Y} a space of structured objects (graphs, rankings, sequences, binary vectors, etc.).

MS/MS spectra



Existing works: Energy-based models (Lafferty et al., 2001; Taskar et al., 2003; Tsochantaridis et al., 2004; LeCun et al., 2007; Belanger and McCallum, 2016):

$$f(x) = \arg \min_{y \in \mathcal{Y}} E(x, y) \quad (1)$$

Table of contents

1. Input Output Kernel Regression
2. Sketched Input Sketched Output Kernel Regression
3. Theoretical Analysis
4. Experiments
5. Conclusion

Input Output Kernel Regression

Output Kernel Regression

Given a p.d. kernel $k_{\mathcal{Y}} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ defining a relevant similarity measure and $\psi_{\mathcal{Y}} : y \in \mathcal{Y} \mapsto k_{\mathcal{Y}}(\cdot, y) \in \mathcal{H}_{\mathcal{Y}}$,

we define $\Delta(y, y') = \|\psi_{\mathcal{Y}}(y) - \psi_{\mathcal{Y}}(y')\|_{\mathcal{H}_{\mathcal{Y}}}^2 = k_{\mathcal{Y}}(y, y) - 2k_{\mathcal{Y}}(y, y') + k_{\mathcal{Y}}(y', y')$ (Weston et al., 2003; Cortes et al., 2005), and solve

$$\min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(X, Y) \sim \rho} [\|\psi_{\mathcal{Y}}(f(X)) - \psi_{\mathcal{Y}}(Y)\|_{\mathcal{H}_{\mathcal{Y}}}^2] \quad (2)$$

Output Kernel Regression

Given a p.d. kernel $k_{\mathcal{Y}} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ defining a relevant similarity measure and $\psi_{\mathcal{Y}} : y \in \mathcal{Y} \mapsto k_{\mathcal{Y}}(\cdot, y) \in \mathcal{H}_{\mathcal{Y}}$,

we define $\Delta(y, y') = \|\psi_{\mathcal{Y}}(y) - \psi_{\mathcal{Y}}(y')\|_{\mathcal{H}_{\mathcal{Y}}}^2 = k_{\mathcal{Y}}(y, y) - 2k_{\mathcal{Y}}(y, y') + k_{\mathcal{Y}}(y', y')$ (Weston et al., 2003; Cortes et al., 2005), and solve

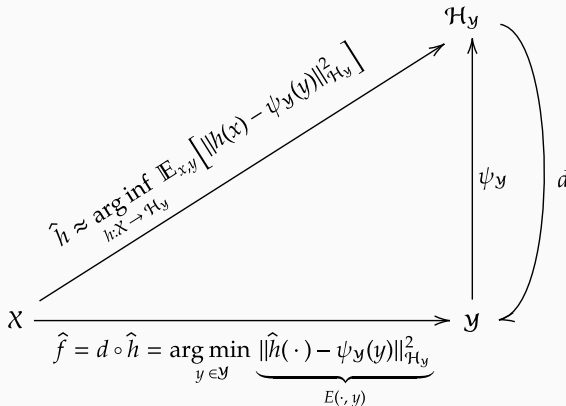
$$\min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(X, Y) \sim \rho} [\|\psi_{\mathcal{Y}}(f(X)) - \psi_{\mathcal{Y}}(Y)\|_{\mathcal{H}_{\mathcal{Y}}}^2] \quad (3)$$

How to learn f through $\psi_{\mathcal{Y}}$?

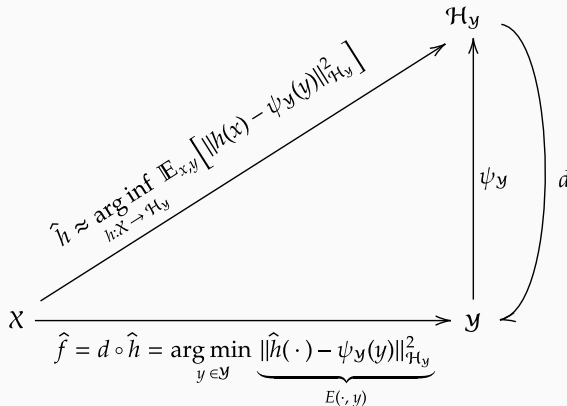
Output Kernel Regression

$$\mathcal{X} \xrightarrow{\quad} \mathcal{Y}$$
$$\hat{f} \approx \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{x,y} \left[\|\psi_{\mathcal{Y}}(f(x)) - \psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}}^2 \right]$$

Output Kernel Regression



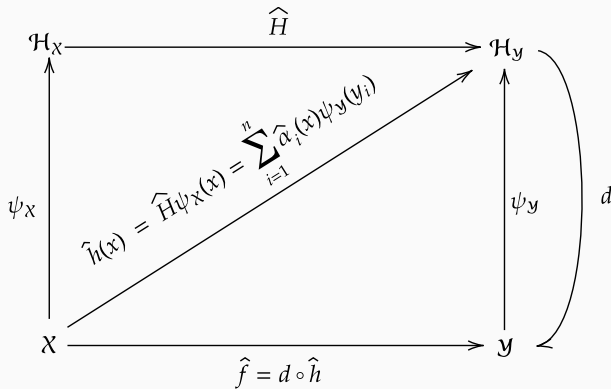
Output Kernel Regression



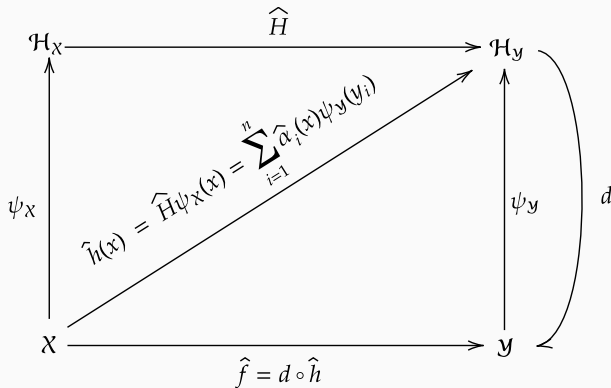
Which hypothesis space for \hat{h} ?

How to deal with infinite-dimensional output feature space \mathcal{H}_Y ?

Input Output Kernel Regression

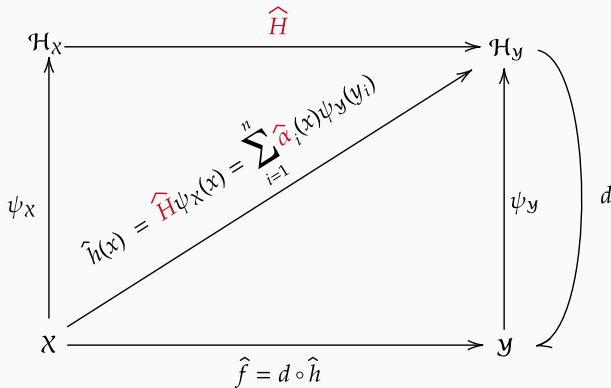


Input Output Kernel Regression



$$\hat{\alpha}(x) = \underbrace{(K_X + n\lambda I_n)^{-1}}_{n \times n} k_X^x = \hat{\Omega} k_X^x \quad \text{where } n = \text{number of training data}$$

Input Output Kernel Regression



$$\hat{\alpha}(x) = \underbrace{(K_X + n\lambda I_n)^{-1}}_{n \times n} k_X^x = \hat{\Omega} k_X^x \quad \text{where } n = \text{number of training data}$$

Training complexity: $\mathcal{O}(n^3)$

Input Output Kernel Regression: Inference

$$\begin{aligned}\hat{f}(x) &= d(\hat{h}(x)) = \arg \min_{y \in \mathcal{Y}} \left\| \hat{h}(x) - \psi_{\mathcal{Y}}(y) \right\|_{\mathcal{H}_{\mathcal{Y}}}^2 = \\ &\arg \min_{y \in \mathcal{Y}} k_{\mathcal{Y}}(y, y) - 2k_x^T \hat{\Omega} k_y^y\end{aligned}$$

Input Output Kernel Regression: Inference

$$\hat{f}(x) = d(\hat{h}(x)) = \arg \min_{y \in \mathcal{Y}} \left\| \hat{h}(x) - \psi_{\mathcal{Y}}(y) \right\|_{\mathcal{H}_{\mathcal{Y}}}^2 =$$
$$\arg \min_{y \in \mathcal{Y}} k_{\mathcal{Y}}(y, y) - 2k_X^{xT} \hat{\Omega} k_Y^y$$

- **Test set:** X_{te} of size n_{te}
- **Candidate set:** $\mathcal{Y}_c \subseteq \mathcal{Y}$ of size n_c

$$\underbrace{K_X^{te, tr}}_{n_{te} \times n} \underbrace{\hat{\Omega}}_{n \times n} \underbrace{K_Y^{tr, c}}_{n \times n_c} \quad (4)$$

Input Output Kernel Regression: Inference

$$\hat{f}(x) = d(\hat{h}(x)) = \arg \min_{y \in \mathcal{Y}} \left\| \hat{h}(x) - \psi_{\mathcal{Y}}(y) \right\|_{\mathcal{H}_{\mathcal{Y}}}^2 =$$
$$\arg \min_{y \in \mathcal{Y}} k_{\mathcal{Y}}(y, y) - 2k_X^{xT} \hat{\Omega} k_Y^y$$

- **Test set:** X_{te} of size n_{te}
- **Candidate set:** $\mathcal{Y}_c \subseteq \mathcal{Y}$ of size n_c

$$\underbrace{K_X^{te, tr}}_{n_{te} \times n} \underbrace{\hat{\Omega}}_{n \times n} \underbrace{K_Y^{tr, c}}_{n \times n_c} \quad (5)$$

Inference complexity: $\mathcal{O}(n_{te} n n_c)$ if $n_{te} < n \leq n_c$

Fisher consistency and excess risk bound

Lemma 1 and Theorem 3 from Ciliberto et al. (2020). Let \mathcal{Y} be compact, $k_{\mathcal{Y}} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a p.d. kernel and $\psi_{\mathcal{Y}} : y \mapsto k_{\mathcal{Y}}(\cdot, y)$ s.t. $\|\psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}} = 1, \forall y \in \mathcal{Y}$, and

$$f^* = \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathcal{R}(f) = \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(x,y) \sim \rho} [\|\psi_{\mathcal{Y}}(f(x)) - \psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}}^2].$$

Then,

$$\hat{f}^*(x) = \arg \min_{y \in \mathcal{Y}} \|h^*(x) - \psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}}^2 = d \circ h^*(x), \quad h^*(x) = \mathbb{E}_y[\psi_{\mathcal{Y}}(y)|x],$$

almost surely with respect to $\rho_{\mathcal{X}}$.

Fisher consistency and excess risk bound

Lemma 1 and Theorem 3 from Ciliberto et al. (2020). Let \mathcal{Y} be compact, $k_{\mathcal{Y}} : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ be a p.d. kernel and $\psi_{\mathcal{Y}} : y \mapsto k_{\mathcal{Y}}(\cdot, y)$ s.t. $\|\psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}} = 1, \forall y \in \mathcal{Y}$, and

$$f^* = \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathcal{R}(f) = \arg \inf_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(x,y) \sim \rho} [\|\psi_{\mathcal{Y}}(f(x)) - \psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}}^2].$$

Then,

$$f^*(x) = \arg \min_{y \in \mathcal{Y}} \|h^*(x) - \psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}}^2 = d \circ h^*(x), \quad h^*(x) = \mathbb{E}_y[\psi_{\mathcal{Y}}(y)|x],$$

almost surely with respect to $\rho_{\mathcal{X}}$.

Moreover, let $h : \mathcal{X} \rightarrow \mathcal{H}_{\mathcal{Y}}$ be measurable and $f : \mathcal{X} \rightarrow \mathcal{Y}$ such that, for any $x \in \mathcal{X}$,

$$f(x) = \arg \min_{y \in \mathcal{Y}} \|h(x) - \psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}}^2 = d \circ h(x).$$

Then,

$$\mathcal{R}(f) - \mathcal{R}(f^*) \leq 12\sqrt{\mathcal{E}(h) - \mathcal{E}(h^*)},$$

where $\mathcal{E}(h) = \mathbb{E}_{(x,y) \sim \rho} [\|h(x) - \psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}}^2]$.

1) **Strong theoretical grounding:** derived from the operator-valued kernel and surrogate methods literature.

- 1) **Strong theoretical grounding:** derived from the operator-valued kernel and surrogate methods literature.
- 2) **Very general algorithm for structured prediction:** ability to tackle many different tasks through an appropriate choice of the output kernel

- 1) **Strong theoretical grounding:** derived from the operator-valued kernel and surrogate methods literature.
- 2) **Very general algorithm for structured prediction:** ability to tackle many different tasks through an appropriate choice of the output kernel
- 3) **Closed-form solution of kernel Ridge regression:** no need for any optimization algorithm to be solved, unlike deep models (Belanger and McCallum, 2016; Belanger et al., 2017; Gygli et al., 2017)

Research question

Can we scale IOKR up to large datasets at both the training and inference phases, especially since they employ not only an input but also an output kernel, while keeping good empirical and theoretical statistical guarantees?

Sketched Input Sketched Output Kernel Regression

Motivation: build a **low-rank** approximation \tilde{h} thanks to **input and output** random projectors \tilde{P}_X and \tilde{P}_Y to obtain a **scalable** predictor \tilde{f} together with an excess risk bound

Some notations

- For \mathcal{Z} a Polish space, $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ a p.d. kernel,
 $\psi_{\mathcal{Z}}(z) := k_{\mathcal{Z}}(\cdot, z)$, $\mathcal{H}_{\mathcal{Z}}$ its RKHS

Some notations

- For \mathcal{Z} a Polish space, $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ a p.d. kernel, $\psi_{\mathcal{Z}}(z) := k_{\mathcal{Z}}(\cdot, z)$, $\mathcal{H}_{\mathcal{Z}}$ its RKHS
- For an i.i.d. sample $(z_i)_{i=1}^n \in \mathcal{Z}^n \sim \rho_{\mathcal{Z}}$

Some notations

- For \mathcal{Z} a Polish space, $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ a p.d. kernel, $\psi_{\mathcal{Z}}(z) := k_{\mathcal{Z}}(\cdot, z)$, $\mathcal{H}_{\mathcal{Z}}$ its RKHS
- For an i.i.d. sample $(z_i)_{i=1}^n \in \mathcal{Z}^n \sim \rho_{\mathcal{Z}}$
- $S_{\mathcal{Z}} : f \in \mathcal{H}_{\mathcal{Z}} \mapsto (1/\sqrt{n})(f(z_1), \dots, f(z_n))^{\top} \in \mathbb{R}^n$ sampling operator

Some notations

- For \mathcal{Z} a Polish space, $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ a p.d. kernel, $\psi_{\mathcal{Z}}(z) := k_{\mathcal{Z}}(\cdot, z)$, $\mathcal{H}_{\mathcal{Z}}$ its RKHS
- For an i.i.d. sample $(z_i)_{i=1}^n \in \mathcal{Z}^n \sim \rho_{\mathcal{Z}}$
- $S_{\mathcal{Z}} : f \in \mathcal{H}_{\mathcal{Z}} \mapsto (1/\sqrt{n})(f(z_1), \dots, f(z_n))^{\top} \in \mathbb{R}^n$ sampling operator
- $S_{\mathcal{Z}}^{\#} : \alpha \in \mathbb{R}^n \mapsto (1/\sqrt{n}) \sum_{i=1}^n \alpha_i \psi_{\mathcal{Z}}(z_i) \in \mathcal{H}_{\mathcal{Z}}$ its adjoint

Some notations

- For \mathcal{Z} a Polish space, $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ a p.d. kernel, $\psi_{\mathcal{Z}}(z) := k_{\mathcal{Z}}(\cdot, z)$, $\mathcal{H}_{\mathcal{Z}}$ its RKHS
- For an i.i.d. sample $(z_i)_{i=1}^n \in \mathcal{Z}^n \sim \rho_{\mathcal{Z}}$
- $S_{\mathcal{Z}} : f \in \mathcal{H}_{\mathcal{Z}} \mapsto (1/\sqrt{n})(f(z_1), \dots, f(z_n))^{\top} \in \mathbb{R}^n$ sampling operator
- $S_{\mathcal{Z}}^{\#} : \alpha \in \mathbb{R}^n \mapsto (1/\sqrt{n}) \sum_{i=1}^n \alpha_i \psi_{\mathcal{Z}}(z_i) \in \mathcal{H}_{\mathcal{Z}}$ its adjoint
- $K_{\mathcal{Z}} = (k_{\mathcal{Z}}(z_i, z_j))_{1 \leq i, j \leq n} = n S_{\mathcal{Z}} S_{\mathcal{Z}}^{\#}$

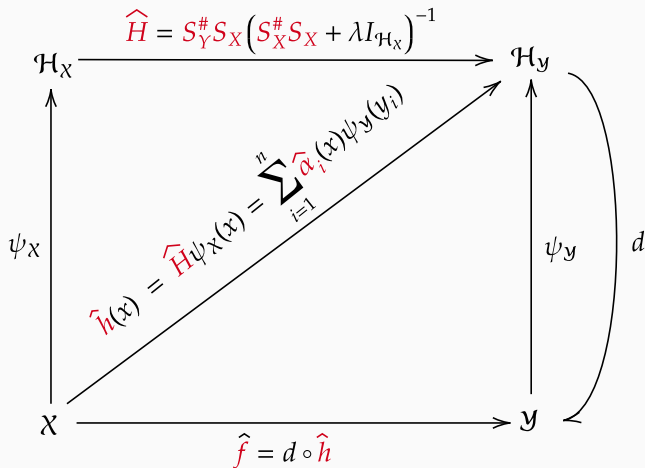
Some notations

- For \mathcal{Z} a Polish space, $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ a p.d. kernel, $\psi_{\mathcal{Z}}(z) := k_{\mathcal{Z}}(\cdot, z)$, $\mathcal{H}_{\mathcal{Z}}$ its RKHS
- For an i.i.d. sample $(z_i)_{i=1}^n \in \mathcal{Z}^n \sim \rho_{\mathcal{Z}}$
- $S_{\mathcal{Z}} : f \in \mathcal{H}_{\mathcal{Z}} \mapsto (1/\sqrt{n})(f(z_1), \dots, f(z_n))^{\top} \in \mathbb{R}^n$ sampling operator
- $S_{\mathcal{Z}}^{\#} : \alpha \in \mathbb{R}^n \mapsto (1/\sqrt{n}) \sum_{i=1}^n \alpha_i \psi_{\mathcal{Z}}(z_i) \in \mathcal{H}_{\mathcal{Z}}$ its adjoint
- $K_{\mathcal{Z}} = (k_{\mathcal{Z}}(z_i, z_j))_{1 \leq i, j \leq n} = n S_{\mathcal{Z}} S_{\mathcal{Z}}^{\#}$
- $C_{\mathcal{Z}} = \mathbb{E}_{\mathcal{Z}}[\psi_{\mathcal{Z}}(z) \otimes \psi_{\mathcal{Z}}(z)]$ covariance operator

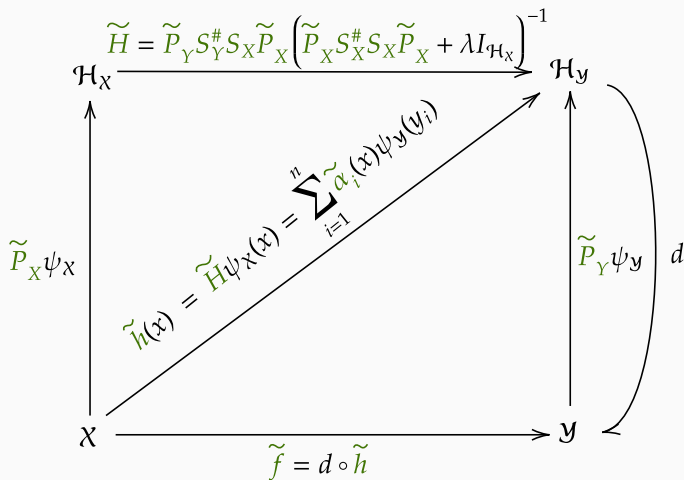
Some notations

- For \mathcal{Z} a Polish space, $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ a p.d. kernel, $\psi_{\mathcal{Z}}(z) := k_{\mathcal{Z}}(\cdot, z)$, $\mathcal{H}_{\mathcal{Z}}$ its RKHS
- For an i.i.d. sample $(z_i)_{i=1}^n \in \mathcal{Z}^n \sim \rho_{\mathcal{Z}}$
- $S_{\mathcal{Z}} : f \in \mathcal{H}_{\mathcal{Z}} \mapsto (1/\sqrt{n})(f(z_1), \dots, f(z_n))^{\top} \in \mathbb{R}^n$ sampling operator
- $S_{\mathcal{Z}}^{\#} : \alpha \in \mathbb{R}^n \mapsto (1/\sqrt{n}) \sum_{i=1}^n \alpha_i \psi_{\mathcal{Z}}(z_i) \in \mathcal{H}_{\mathcal{Z}}$ its adjoint
- $K_{\mathcal{Z}} = (k_{\mathcal{Z}}(z_i, z_j))_{1 \leq i, j \leq n} = n S_{\mathcal{Z}} S_{\mathcal{Z}}^{\#}$
- $C_{\mathcal{Z}} = \mathbb{E}_{\mathcal{Z}}[\psi_{\mathcal{Z}}(z) \otimes \psi_{\mathcal{Z}}(z)]$ covariance operator
- $\hat{C}_{\mathcal{Z}} = (1/n) \sum_{i=1}^n \psi_{\mathcal{Z}}(z_i) \otimes \psi_{\mathcal{Z}}(z_i) = S_{\mathcal{Z}}^{\#} S_{\mathcal{Z}}$ its empirical counterpart

Low-rank Estimator: from IOKR to SISOKR



Low-rank Estimator: from IOKR to SISOKR



Sketching: random linear projections

Sketching: random linear projections

Let $m_{\mathcal{Z}} \ll n$, $R_{\mathcal{Z}} \in \mathbb{R}^{m_{\mathcal{Z}} \times n}$ be a random matrix and n data $(z_i)_{i=1}^n \in \mathcal{Z}$

Sketching: random linear projections

Let $m_{\mathcal{Z}} \ll n$, $R_{\mathcal{Z}} \in \mathbb{R}^{m_{\mathcal{Z}} \times n}$ be a random matrix and n data $(z_i)_{i=1}^n \in \mathcal{Z}$

Basic idea: Sketching-based operator $\tilde{P}_{\mathcal{Z}}$ projects onto the following linear subspace of $\mathcal{H}_{\mathcal{Z}}$

$$\sum_{j=1}^n (R_{\mathcal{Z}})_{ij} \psi_{\mathcal{Z}}(z_j) \in \mathcal{H}_{\mathcal{Z}}, \quad i = 1, \dots, m_{\mathcal{Z}} \quad (6)$$

Construction of the orthogonal projector \tilde{P}_Z

- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z$

Construction of the orthogonal projector \tilde{P}_Z

- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z$
- $\tilde{K}_Z = R_Z K_Z R_Z^\top$, and $\left\{ \left(\sigma_i(\tilde{K}_Z), \tilde{\mathbf{v}}_i^Z \right), i \in [m_Z] \right\}$ its eigenpairs

Construction of the orthogonal projector \tilde{P}_Z

- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z$
- $\tilde{K}_Z = R_Z K_Z R_Z^\top$, and $\left\{ \left(\sigma_i(\tilde{K}_Z), \tilde{\mathbf{v}}_i^Z \right), i \in [m_Z] \right\}$ its eigenpairs
- $p_Z = \text{rank}(\tilde{K}_Z)$, and for all $1 \leq i \leq p_Z$, $\tilde{\mathbf{e}}_i^Z = \sqrt{\frac{n}{\sigma_i(\tilde{K}_Z)}} S_Z^\# R_Z^\top \tilde{\mathbf{v}}_i^Z \in \mathcal{H}_Z$

Construction of the orthogonal projector \tilde{P}_Z

- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z$
- $\tilde{K}_Z = R_Z K_Z R_Z^\top$, and $\left\{ \left(\sigma_i(\tilde{K}_Z), \tilde{\mathbf{v}}_i^Z \right), i \in [m_Z] \right\}$ its eigenpairs
- $p_Z = \text{rank}(\tilde{K}_Z)$, and for all $1 \leq i \leq p_Z$, $\tilde{\mathbf{e}}_i^Z = \sqrt{\frac{n}{\sigma_i(\tilde{K}_Z)}} S_Z^\# R_Z^\top \tilde{\mathbf{v}}_i^Z \in \mathcal{H}_Z$

Proposition

The $\tilde{\mathbf{e}}_i^Z$ s are the eigenfunctions, associated to the eigenvalues $\sigma_i(\tilde{K}_Z)/n$, of \tilde{C}_Z .

Construction of the orthogonal projector \tilde{P}_Z

- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z$
- $\tilde{K}_Z = R_Z K_Z R_Z^\top$, and $\left\{ \left(\sigma_i(\tilde{K}_Z), \tilde{\mathbf{v}}_i^Z \right), i \in [m_Z] \right\}$ its eigenpairs
- $p_Z = \text{rank}(\tilde{K}_Z)$, and for all $1 \leq i \leq p_Z$, $\tilde{\mathbf{e}}_i^Z = \sqrt{\frac{n}{\sigma_i(\tilde{K}_Z)}} S_Z^\# R_Z^\top \tilde{\mathbf{v}}_i^Z \in \mathcal{H}_Z$

Proposition

The $\tilde{\mathbf{e}}_i^Z$ s are the eigenfunctions, associated to the eigenvalues $\sigma_i(\tilde{K}_Z)/n$, of \tilde{C}_Z .

Furthermore, let $\tilde{\mathcal{H}}_Z = \text{span}(\tilde{\mathbf{e}}_1^Z, \dots, \tilde{\mathbf{e}}_{p_Z}^Z)$, the orthogonal projector \tilde{P}_Z onto $\tilde{\mathcal{H}}_Z$ writes as

$$\tilde{P}_Z = (R_Z S_Z)^\# (R_Z S_Z (R_Z S_Z)^\#)^\dagger R_Z S_Z. \quad (7)$$

Construction of the orthogonal projector \tilde{P}_Z

- $\tilde{C}_Z = S_Z^\# R_Z^\top R_Z S_Z$
- $\tilde{K}_Z = R_Z K_Z R_Z^\top$, and $\left\{ \left(\sigma_i(\tilde{K}_Z), \tilde{\mathbf{v}}_i^Z \right), i \in [m_Z] \right\}$ its eigenpairs
- $p_Z = \text{rank}(\tilde{K}_Z)$, and for all $1 \leq i \leq p_Z$, $\tilde{\mathbf{e}}_i^Z = \sqrt{\frac{n}{\sigma_i(\tilde{K}_Z)}} S_Z^\# R_Z^\top \tilde{\mathbf{v}}_i^Z \in \mathcal{H}_Z$

Proposition

The $\tilde{\mathbf{e}}_i^Z$ s are the eigenfunctions, associated to the eigenvalues $\sigma_i(\tilde{K}_Z)/n$, of \tilde{C}_Z .

Furthermore, let $\tilde{\mathcal{H}}_Z = \text{span}(\tilde{\mathbf{e}}_1^Z, \dots, \tilde{\mathbf{e}}_{p_Z}^Z)$, the orthogonal projector \tilde{P}_Z onto $\tilde{\mathcal{H}}_Z$ writes as

$$\tilde{P}_Z = (R_Z S_Z)^\# (R_Z S_Z (R_Z S_Z)^\#)^\dagger R_Z S_Z. \quad (7)$$

Related work on Nyström: Yang et al. (2012); Rudi et al. (2015)

Proposition (Expression of SISOKR)

$$\tilde{h}(x) = \sum_{i=1}^n \tilde{\alpha}_i(x) \psi_Y(y_i), \quad \text{where} \quad \tilde{\alpha}(x) = R_Y^\top \tilde{\Omega} R_X k_X^x, \quad (8)$$

with

$$\tilde{\Omega} = \underbrace{(R_Y K_Y R_Y^\top)^\dagger}_{m_Y \times m_Y} R_Y K_Y K_X R_X^\top \underbrace{(R_X K_X^2 R_X^\top + n\lambda R_X K_X R_X^\top)^\dagger}_{m_X \times m_X} \quad (9)$$

Proposition (Expression of SISOKR)

$$\tilde{h}(x) = \sum_{i=1}^n \tilde{\alpha}_i(x) \psi_{\mathcal{Y}}(y_i), \quad \text{where} \quad \tilde{\alpha}(x) = R_{\mathcal{Y}}^{\top} \tilde{\Omega} R_{\mathcal{X}} k_{\mathcal{X}}^x, \quad (8)$$

with

$$\tilde{\Omega} = \underbrace{(R_{\mathcal{Y}} K_{\mathcal{Y}} R_{\mathcal{Y}}^{\top})^{\dagger}}_{m_{\mathcal{Y}} \times m_{\mathcal{Y}}} R_{\mathcal{Y}} K_{\mathcal{Y}} K_{\mathcal{X}} R_{\mathcal{X}}^{\top} \underbrace{(R_{\mathcal{X}} K_{\mathcal{X}}^2 R_{\mathcal{X}}^{\top} + n \lambda R_{\mathcal{X}} K_{\mathcal{X}} R_{\mathcal{X}}^{\top})^{\dagger}}_{m_{\mathcal{X}} \times m_{\mathcal{X}}} \quad (9)$$

Inversion complexity: $\mathcal{O}(n^3) \rightarrow \mathcal{O}(\max(m_{\mathcal{X}}^3, m_{\mathcal{Y}}^3))$

Proposition (Expression of SISOKR)

$$\tilde{h}(x) = \sum_{i=1}^n \tilde{\alpha}_i(x) \psi_Y(y_i), \quad \text{where} \quad \tilde{\alpha}(x) = R_Y^\top \tilde{\Omega} R_X k_X^x, \quad (8)$$

with

$$\tilde{\Omega} = \underbrace{(R_Y K_Y R_Y^\top)^\dagger}_{m_Y \times m_Y} R_Y K_Y K_X R_X^\top \underbrace{(R_X K_X^2 R_X^\top + n\lambda R_X K_X R_X^\top)^\dagger}_{m_X \times m_X} \quad (9)$$

Inversion complexity: $\mathcal{O}(n^3) \rightarrow \mathcal{O}(\max(m_X^3, m_Y^3))$

Complexity of $R_Z K_Z$: depends on the sketch matrix, between $\mathcal{O}(nm_Z)$ and $\mathcal{O}(n^2 m_Z)$

Proposition (Expression of SISOKR)

$$\tilde{h}(x) = \sum_{i=1}^n \tilde{\alpha}_i(x) \psi_Y(y_i), \quad \text{where} \quad \tilde{\alpha}(x) = R_Y^\top \tilde{\Omega} R_X k_X^x, \quad (8)$$

with

$$\tilde{\Omega} = \underbrace{(R_Y K_Y R_Y^\top)^\dagger}_{m_Y \times m_Y} R_Y K_Y K_X R_X^\top \underbrace{(R_X K_X^2 R_X^\top + n\lambda R_X K_X R_X^\top)^\dagger}_{m_X \times m_X} \quad (9)$$

Inversion complexity: $\mathcal{O}(n^3) \rightarrow \mathcal{O}(\max(m_X^3, m_Y^3))$

Complexity of $R_Z K_Z$: depends on the sketch matrix, between $\mathcal{O}(nm_Z)$ and $\mathcal{O}(n^2 m_Z)$

\implies Training complexity reduced!

- Test set: X_{te} of size n_{te}
- Candidate set: $\mathcal{Y}_c \subseteq \mathcal{Y}$ of size n_c

$$\underbrace{K_X^{te,tr} R_{\mathcal{X}}^T}_{n_{te} \times m_{\mathcal{X}}} \underbrace{\tilde{\Omega}}_{m_{\mathcal{X}} \times m_{\mathcal{Y}}} \underbrace{R_{\mathcal{Y}} K_Y^{tr,c}}_{m_{\mathcal{Y}} \times n_c} \quad (10)$$

- Test set: X_{te} of size n_{te}
- Candidate set: $\mathcal{Y}_c \subseteq \mathcal{Y}$ of size n_c

$$\underbrace{K_X^{te,tr} R_{\mathcal{X}}^T}_{n_{te} \times m_{\mathcal{X}}} \underbrace{\tilde{\Omega}}_{m_{\mathcal{X}} \times m_{\mathcal{Y}}} \underbrace{R_{\mathcal{Y}} K_Y^{tr,c}}_{m_{\mathcal{Y}} \times n_c} \quad (10)$$

Decoding complexity: $\mathcal{O}(n_{te} n n_c) \rightarrow \mathcal{O}(n_{te} m_{\mathcal{Y}} n_c)$ if

$$n_{te} \leq m_{\mathcal{X}}, m_{\mathcal{Y}} < n \leq n_c$$

- Test set: X_{te} of size n_{te}
- Candidate set: $\mathcal{Y}_c \subseteq \mathcal{Y}$ of size n_c

$$\underbrace{K_X^{te,tr} R_{\mathcal{X}}^T}_{n_{te} \times m_{\mathcal{X}}} \underbrace{\tilde{\Omega}}_{m_{\mathcal{X}} \times m_{\mathcal{Y}}} \underbrace{R_{\mathcal{Y}} K_Y^{tr,c}}_{m_{\mathcal{Y}} \times n_c} \quad (10)$$

Decoding complexity: $\mathcal{O}(n_{te} n n_c) \rightarrow \mathcal{O}(n_{te} m_{\mathcal{Y}} n_c)$ if

$$n_{te} \leq m_{\mathcal{X}}, m_{\mathcal{Y}} < n \leq n_c$$

\Rightarrow Inference complexity reduced!

Theoretical Analysis

Assumptions

Asm. 1 (Attainability): Recall that $h^*(x) := \mathbb{E}_Y[\psi_Y(Y) \mid X = x]$. $h^* \in \mathcal{H}$,
i.e. there exists $H : \mathcal{H}_X \rightarrow \mathcal{H}_Y$ with $\|H\|_{\text{HS}} < +\infty$ such that

$$h^*(x) = H\psi_X(x) \quad \forall x \in \mathcal{X}. \quad (11)$$

Assumptions

Asm. 1 (Attainability): Recall that $h^*(x) := \mathbb{E}_Y[\psi_Y(Y) \mid X = x]$. $h^* \in \mathcal{H}$, i.e. there exists $H : \mathcal{H}_X \rightarrow \mathcal{H}_Y$ with $\|H\|_{\text{HS}} < +\infty$ such that

$$h^*(x) = H\psi_X(x) \quad \forall x \in \mathcal{X}. \quad (11)$$

Asm. 2 (Bounded kernel): there exists $\kappa_Z > 0$ such that

$$k_Z(z, z) \leq \kappa_Z^2 \quad \forall z \in \mathcal{Z}. \quad (12)$$

Assumptions

Asm. 1 (Attainability): Recall that $h^*(x) := \mathbb{E}_Y[\psi_Y(Y) \mid X = x]$. $h^* \in \mathcal{H}$, i.e. there exists $H : \mathcal{H}_X \rightarrow \mathcal{H}_Y$ with $\|H\|_{\text{HS}} < +\infty$ such that

$$h^*(x) = H\psi_X(x) \quad \forall x \in \mathcal{X}. \quad (11)$$

Asm. 2 (Bounded kernel): there exists $\kappa_Z > 0$ such that

$$k_Z(z, z) \leq \kappa_Z^2 \quad \forall z \in \mathcal{Z}. \quad (12)$$

Asm. 3 (Capacity condition): there exists $\gamma_Z \in [0, 1]$ such that

$$Q_Z := \text{Tr}(C_Z^{\gamma_Z}) < +\infty. \quad (13)$$

Assumptions

Asm. 1 (Attainability): Recall that $h^*(x) := \mathbb{E}_Y[\psi_Y(Y) \mid X = x]$. $h^* \in \mathcal{H}$, i.e. there exists $H : \mathcal{H}_X \rightarrow \mathcal{H}_Y$ with $\|H\|_{\text{HS}} < +\infty$ such that

$$h^*(x) = H\psi_X(x) \quad \forall x \in \mathcal{X}. \quad (11)$$

Asm. 2 (Bounded kernel): there exists $\kappa_Z > 0$ such that

$$k_Z(z, z) \leq \kappa_Z^2 \quad \forall z \in \mathcal{Z}. \quad (12)$$

Asm. 3 (Capacity condition): there exists $\gamma_Z \in [0, 1]$ such that

$$Q_Z := \text{Tr}(C_Z^{\gamma_Z}) < +\infty. \quad (13)$$

Asm. 4 (Embedding property): there exists $b_Z > 0$ and $\mu_Z \in [0, 1]$ such that almost surely

$$\psi_Z(z) \otimes \psi_Z(z) \preceq b_Z C_Z^{1-\mu_Z}. \quad (14)$$

Assumptions

Asm. 1 (Attainability): Recall that $h^*(x) := \mathbb{E}_Y[\psi_Y(Y) \mid X = x]$. $h^* \in \mathcal{H}$, i.e. there exists $H : \mathcal{H}_X \rightarrow \mathcal{H}_Y$ with $\|H\|_{\text{HS}} < +\infty$ such that

$$h^*(x) = H\psi_X(x) \quad \forall x \in \mathcal{X}. \quad (11)$$

Asm. 2 (Bounded kernel): there exists $\kappa_Z > 0$ such that

$$k_Z(z, z) \leq \kappa_Z^2 \quad \forall z \in \mathcal{Z}. \quad (12)$$

Asm. 3 (Capacity condition): there exists $\gamma_Z \in [0, 1]$ such that

$$Q_Z := \text{Tr}(C_Z^{\gamma_Z}) < +\infty. \quad (13)$$

Asm. 4 (Embedding property): there exists $b_Z > 0$ and $\mu_Z \in [0, 1]$ such that almost surely

$$\psi_Z(z) \otimes \psi_Z(z) \preceq b_Z C_Z^{1-\mu_Z}. \quad (14)$$

Asm. 5 (Sub-Gaussian sketches): $R_Z \in \mathbb{R}^{m_Z \times n}$ composed with i.i.d. entries s.t. (i) $\mathbb{E}[R_{Z_{ij}}] = 0$, (ii) $\mathbb{E}[R_{Z_{ij}}^2] = 1/m_Z$ and (iii)

$R_{Z_{ij}} \sim \frac{\nu_Z}{m_Z} - \text{sub-Gaussian with } \nu_Z \geq 1.$

SISOKR Learning Rates

Corollary (SISOKR learning rates)

Under **Asm. 1, 2, 3, 4 and 5**, if for all $y \in \mathcal{Y}$, $\|\psi_{\mathcal{Y}}(y)\|_{\mathcal{H}_{\mathcal{Y}}} = \kappa_{\mathcal{Y}}$, for $n \in \mathbb{N}$ sufficiently large such that $\frac{9}{n} \log(n/\delta) \leq n^{-\frac{1}{1+\gamma_{\mathcal{X}}}} \leq \|C_{\mathcal{X}}\|_{\text{op}}/2$, and $\frac{9}{n} \log(n/\delta) \leq n^{-\frac{1}{1+\gamma_{\mathcal{Y}}}} \leq \|C_{\mathcal{Y}}\|_{\text{op}}/2$, and for sketching size $m_{\mathcal{X}}, m_{\mathcal{Y}} \in \mathbb{N}$ such that

$$m_{\mathcal{X}} \gtrsim \max \left(\nu_{\mathcal{X}}^2 n^{\frac{\gamma_{\mathcal{X}} + \mu_{\mathcal{X}}}{1+\gamma_{\mathcal{X}}}}, \nu_{\mathcal{X}}^4 \log(1/\delta) \right), \quad (15)$$

$$m_{\mathcal{Y}} \gtrsim \max \left(\nu_{\mathcal{Y}}^2 n^{\frac{\gamma_{\mathcal{Y}} + \mu_{\mathcal{Y}}}{1+\gamma_{\mathcal{Y}}}}, \nu_{\mathcal{Y}}^4 \log(1/\delta) \right), \quad (16)$$

then with probability $1 - \delta$

$$\mathbb{E}[\|\tilde{h}(x) - h^*(x)\|_{\mathcal{H}_{\mathcal{Y}}}^2]^{\frac{1}{2}} \lesssim \log(4/\delta) n^{-\frac{1-\gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}}}{2(1+\gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}})}}, \quad (17)$$

and

$$\mathcal{R}(\tilde{f}) - \mathcal{R}(f^*) \lesssim \log(4/\delta) n^{-\frac{1-\gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}}}{2(1+\gamma_{\mathcal{X}} \vee \gamma_{\mathcal{Y}})}}. \quad (18)$$

Experiments

Multi-Label Classification: Statistical Performance

Table 1: F_1 score on tag prediction from text data.

Method	Bibtex	Bookmarks
SISOKR	44.1 ± 0.07	39.3 ± 0.61
ISOKR	44.8 ± 0.01	NA
SIOKR	44.7 ± 0.09	39.1 ± 0.04
IOKR	44.9	NA
LR	37.2	30.7
NN	38.9	33.8
SPEN	42.2	34.4
PRLR	44.2	34.9
DVN	44.7	37.1

Multi-Label Classification: Computational Performance

Table 2: Comparison of training/inference computation times (in seconds).

Method	Bibtex	Bookmarks
SISOKR	1.41 ± 0.03 / 0.46 ± 0.01	118 ± 1.5 / 20 ± 0.2
ISOKR	2.51 ± 0.06 / 0.58 ± 0.01	NA
SIOKR	1.99 ± 0.07 / 1.22 ± 0.03	354 ± 2.1 / 297 ± 2.1
IOKR	2.54 / 1.18	NA

Conclusion

Take-home messages

- Scale up surrogate kernel methods for structured prediction by leveraging **random projections**, in **both input and output** feature spaces, to **accelerate training and inference** phases

Take-home messages

- Scale up surrogate kernel methods for structured prediction by leveraging **random projections**, in **both input and output** feature spaces, to **accelerate training and inference** phases
- Derive **excess risk bounds** for the **sketched estimator**

Take-home messages

- Scale up surrogate kernel methods for structured prediction by leveraging **random projections**, in **both input and output** feature spaces, to **accelerate training and inference** phases
- Derive **excess risk bounds** for the **sketched estimator**
- Show that **sub-Gaussian** sketches are admissible sketches in the sense that they lead to **close to optimal learning rates** with **sketching sizes $m < n$**

Take-home messages

- Scale up surrogate kernel methods for structured prediction by leveraging **random projections**, in **both input and output** feature spaces, to **accelerate training and inference** phases
- Derive **excess risk bounds** for the **sketched estimator**
- Show that **sub-Gaussian** sketches are admissible sketches in the sense that they lead to **close to optimal learning rates** with **sketching sizes $m < n$**
- Provide structured prediction **experiments** on **real-world data sets** showing **similar performances** as IOKR while **being faster** in **both training and inference** phases.

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.
- Bach, F. (2013). Sharp analysis of low-rank kernel matrix approximations. In *Proc. of the 26th annual Conference on Learning Theory*, pages 185–209. PMLR.

- Belanger, D. and McCallum, A. (2016). Structured prediction energy networks. In *International Conference on Machine Learning*, pages 983–992.
- Belanger, D., Yang, B., and McCallum, A. (2017). End-to-end learning for structured prediction energy networks. In Precup, D. and Teh, Y. W., editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 429–439. PMLR.
- Brault, R., Heinonen, M., and Buc, F. (2016). Random fourier features for operator-valued kernels. In *Asian Conference on Machine Learning*, pages 110–125. PMLR.
- Carmeli, C., De Vito, E., and Toigo, A. (2006). Vector valued reproducing kernel hilbert spaces of integrable functions and mercer theorem. *Analysis and Applications*, 4(04):377–408.

- Carmeli, C., De Vito, E., Toigo, A., and Umanitá, V. (2010). Vector valued reproducing kernel hilbert spaces and universality. *Analysis and Applications*, 8(01):19–61.
- 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.
- Ciliberto, C., Rosasco, L., and Rudi, A. (2020). A general framework for consistent structured prediction with implicit loss embeddings. *J. Mach. Learn. Res.*, 21(98):1–67.

- Cortes, C., Mohri, M., and Weston, J. (2005). A general regression technique for learning transductions. In *International Conference on Machine Learning (ICML)*, pages 153–160.
- Drineas, P., Mahoney, M. W., and Cristianini, N. (2005). On the nyström method for approximating a gram matrix for improved kernel-based learning. *JMLR*, 6(12).
- El Ahmad, T., Laforgue, P., and d’Alché Buc, F. (2023). Fast kernel methods for generic lipschitz losses via p -sparsified sketches. *Transactions on Machine Learning Research*.
- Gazagnadou, N., Ibrahim, M., and Gower, R. M. (2021). ***RidgeSketch***: A fast sketching based solver for large scale ridge regression.

- Gygli, M., Norouzi, M., and Angelova, A. (2017). Deep value networks learn to evaluate and iteratively refine structured outputs. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, page 1341–1351. JMLR.org.
- Katakis, I., Tsoumakas, G., and Vlahavas, I. (2008). Multilabel text classification for automated tag suggestion. In *Proceedings of the ECML/PKDD*, volume 18, page 5. Citeseer.
- Kpotufe, S. and Sriperumbudur, B. K. (2020). Gaussian sketching yields a J-L lemma in RKHS. In Chiappa, S. and Calandra, R., editors, *AISTATS 2020*, volume 108 of *Proceedings of Machine Learning Research*, pages 3928–3937. PMLR.
- Lacotte, J. and Pilanci, M. (2020). Adaptive and oblivious randomized subspace methods for high-dimensional optimization: Sharp analysis and lower bounds. *arXiv preprint arXiv:2012.07054*.

- Lacotte, J., Pilanci, M., and Pavone, M. (2019). High-dimensional optimization in adaptive random subspaces. In *Proc. of the 33rd International Conference on Neural Information Processing Systems*, pages 10847–10857.
- Lafferty, J. D., McCallum, A., and Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *International Conference on Machine Learning*.
- LeCun, Y., Chopra, S., Ranzato, M., and Huang, F.-J. (2007). Energy-based models in document recognition and computer vision. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, volume 1, pages 337–341. IEEE.
- Li, Z., Ton, J.-F., Oglic, D., and Sejdinovic, D. (2021). Towards a unified analysis of random fourier features. *Journal of Machine Learning Research*, 22(108):1–51.

- Mahoney, M. W. et al. (2011). Randomized algorithms for matrices and data. *Foundations and Trends® in Machine Learning*, 3(2):123–224.
- Meanti, G., Carratino, L., Rosasco, L., and Rudi, A. (2020). Kernel methods through the roof: Handling billions of points efficiently. *Advances in Neural Information Processing Systems (NeurIPS)*, 33.
- Micchelli, C. A. and Pontil, M. (2005). On learning vector-valued functions. *Neural computation*, 17(1):177–204.
- Musco, C. and Musco, C. (2017). Recursive sampling for the nyström method. *Advances in Neural Information Processing Systems*, 2017:3834–3846.
- Rahimi, A. and Recht, B. (2007). Random features for large scale kernel machines. *NIPS*, 20:1177–1184.
- Rudi, A., Calandriello, D., Carratino, L., and Rosasco, L. (2018). On fast leverage score sampling and optimal learning. In *NeurIPS*.

- Rudi, A., Camoriano, R., and Rosasco, L. (2015). Less is more: Nyström computational regularization. *Advances in Neural Information Processing Systems*, 28.
- Rudi, A., Carratino, L., and Rosasco, L. (2017). Falkon: an optimal large scale kernel method. In *Proceedings of the 31st International Conference on Advances on Neural Information Processing Systems (NeurIPS)*, pages 3891–3901.
- Rudi, A. and Rosasco, L. (2017). Generalization properties of learning with random features. In *Advances on Neural Information Processing Systems (NeurIPS)*, pages 3215–3225.
- Senkene, E. and Tempel'man, A. (1973). Hilbert spaces of operator-valued functions. *Lithuanian Mathematical Journal*, 13(4):665–670.

- Sriperumbudur, B. K. and Szabó, Z. (2015). Optimal rates for random fourier features. In *NIPS*.
- Taskar, B., Guestrin, C., and Koller, D. (2003). Max-margin markov networks. In Thrun, S., Saul, L., and Schölkopf, B., editors, *Advances in Neural Information Processing Systems*, volume 16. MIT Press.
- Tsochantaridis, I., Hofmann, T., Joachims, T., and Altun, Y. (2004). Support vector machine learning for interdependent and structured output spaces. *Machine Learning*.
- Weston, J., Chapelle, O., Vapnik, V., Elisseeff, A., and Schölkopf, B. (2003). Kernel dependency estimation. In Becker, S., Thrun, S., and Obermayer, K., editors, *Advances in Neural Information Processing Systems 15*, pages 897–904. MIT Press.

- Williams, C. and Seeger, M. (2001). Using the nyström method to speed up kernel machines. In Leen, T., Dietterich, T., and Tresp, V., editors, *Advances in Neural Information Processing Systems*, volume 13, pages 682–688. MIT Press.
- Woodruff, D. P. (2014). Sketching as a tool for numerical linear algebra. *Found. Trends Theor. Comput. Sci.*, 10(1-2):1–157.
- Yang, T., Li, Y.-f., Mahdavi, M., Jin, R., and Zhou, Z.-H. (2012). Nyström method vs random fourier features: A theoretical and empirical comparison. In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc.
- 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.

Reminder: positive definite kernels and Reproducing Kernel Hilbert Space

Positive definite kernel: $k_{\mathcal{Z}} : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$ such that

- for all $(z, z') \in \mathcal{Z}^2$, $k_{\mathcal{Z}}(z, z') = k_{\mathcal{Z}}(z', z)^{\top}$
- for all $n \in \mathbb{N}$ and any $(z_i, \alpha_i)_{i=1}^n \in (\mathcal{Z} \times \mathbb{R})^n$,
 $\sum_{i,j=1}^n \alpha_i \alpha_j k_{\mathcal{Z}}(z_i, z_j) \geq 0$

RKHS (Aronszajn, 1950): Hilbert space $\mathcal{H}_{\mathcal{Z}}$ of functions $f : \mathcal{Z} \rightarrow \mathbb{R}$ s. t.
for all $f \in \mathcal{H}_{\mathcal{Z}}$ and $z \in \mathcal{Z}$

1. $z' \mapsto k_{\mathcal{Z}}(z, z') \in \mathcal{H}_{\mathcal{Z}}$,
2. $\langle f, k_{\mathcal{Z}}(\cdot, z) \rangle_{\mathcal{H}_{\mathcal{Z}}} = f(z)$ (reproducing property).

Vector-Valued Reproducing Kernel Hilbert Space

Operator-valued kernel (Senkene and Tempel'man, 1973; Micchelli and Pontil, 2005; Carmeli et al., 2006, 2010): $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathcal{L}(\mathcal{F})$, where \mathcal{F} is a Hilbert space, such that

- for all $(x, x') \in \mathcal{X}^2$, $\mathcal{K}(x, x') = \mathcal{K}(x', x)^\#$
- for all $n \in \mathbb{N}$ and any $(x_i, \varphi_i)_{i=1}^n \in (\mathcal{X} \times \mathcal{F})^n$,
$$\sum_{i,j=1}^n \langle \varphi_i, \mathcal{K}(x_i, x_j) \varphi_j \rangle_{\mathcal{F}} \geq 0$$

vv-RKHS: Hilbert space \mathcal{H} of functions $f : \mathcal{X} \rightarrow \mathcal{F}$ s. t. for all $f \in \mathcal{H}$, $\varphi \in \mathcal{F}$ and $x \in \mathcal{X}$

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

Background: Scalability to large datasets

1) Random Fourier Features (Rahimi and Recht, 2007; Rudi and Rosasco, 2017; Sriperumbudur and Szabó, 2015; Brault et al., 2016; Li et al., 2021)

Background: Scalability to large datasets

1) **Random Fourier Features** (Rahimi and Recht, 2007; Rudi and Rosasco, 2017; Sriperumbudur and Szabó, 2015; Brault et al., 2016; Li et al., 2021)

2) **Sketching** (Mahoney et al., 2011; Woodruff, 2014): dimension reduction approach based on random linear projections

- **Nyström approximation** (\iff sub-sampling sketch) (Williams and Seeger, 2001; Drineas et al., 2005; Bach, 2013; Rudi et al., 2017; Meanti et al., 2020)
- **Gaussian, Randomized Orthogonal Systems, sparse sketches** etc. (Yang et al., 2017; Lacotte et al., 2019; Kpotufe and Sriperumbudur, 2020; Lacotte and Pilanci, 2020; Chen and Yang, 2021a; Gazagnadou et al., 2021)

Example: Sketching for scalar Kernel Ridge Regression ($\mathcal{Y} = \mathbb{R}$)

Representer theorem: $\hat{f} = \sum_{i=1}^n \hat{\alpha}_i k_x(\cdot, x_i)$, where

$$\begin{aligned}\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_n)^\top &= \arg \min_{\alpha \in \mathbb{R}^n} \alpha^\top (K_X^2 + n\lambda K_X) \alpha - 2Y^\top K_X \alpha \\ &= \underbrace{(K_X + n\lambda I_n)}_{n \times n}^{-1} Y\end{aligned}$$

Example: Sketching for scalar Kernel Ridge Regression ($\mathcal{Y} = \mathbb{R}$)

Representer theorem: $\hat{f} = \sum_{i=1}^n \hat{\alpha}_i k_x(\cdot, x_i)$, where

$$\begin{aligned}\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_n)^\top &= \arg \min_{\alpha \in \mathbb{R}^n} \alpha^\top (K_X^2 + n\lambda K_X) \alpha - 2Y^\top K_X \alpha \\ &= \underbrace{(K_X + n\lambda I_n)}_{n \times n}^{-1} Y\end{aligned}$$

Let $m \ll n$, $R \in \mathbb{R}^{m \times n}$ be a random matrix: $\alpha \leftarrow R^\top \gamma$

Example: Sketching for scalar Kernel Ridge Regression ($\mathcal{Y} = \mathbb{R}$)

Representer theorem: $\hat{f} = \sum_{i=1}^n \hat{\alpha}_i k_x(\cdot, x_i)$, where

$$\begin{aligned}\hat{\alpha} = (\hat{\alpha}_1, \dots, \hat{\alpha}_n)^\top &= \arg \min_{\alpha \in \mathbb{R}^n} \alpha^\top (K_X^2 + n\lambda K_X) \alpha - 2Y^\top K_X \alpha \\ &= \underbrace{(K_X + n\lambda I_n)^{-1}}_{n \times n} Y\end{aligned}$$

Let $m \ll n$, $R \in \mathbb{R}^{m \times n}$ be a random matrix: $\alpha \leftarrow R^\top \gamma$

$\hat{f} \leftarrow \tilde{f} = \sum_{i=1}^n [R^\top \tilde{\gamma}]_i k_x(\cdot, x_i)$, where

$$\begin{aligned}\tilde{\gamma} = (\tilde{\gamma}_1, \dots, \tilde{\gamma}_m)^\top &= \arg \min_{\gamma \in \mathbb{R}^m} \gamma^\top (RK_X^2 R^\top + n\lambda RK_X R^\top) \gamma - 2Y^\top K_X R^\top \gamma \\ &= \underbrace{(RK_X^2 R^\top + n\lambda RK_X R^\top)^\dagger}_{m \times m} RK_X Y\end{aligned}$$

Low-rank estimator

$$\hat{h}(x) = \sum_{i=1}^n \hat{\alpha}_i(x) \psi_{\mathcal{Y}}(y_i), \quad \text{where} \quad \hat{\alpha}(x) = (K_X + n\lambda I_n)^{-1} k_X^x$$

Low-rank estimator

$$\begin{aligned}\hat{h}(x) &= \sum_{i=1}^n \hat{\alpha}_i(x) \psi_Y(y_i), \quad \text{where} \quad \hat{\alpha}(x) = (K_X + n\lambda I_n)^{-1} k_X^x \\ &= \sqrt{n} S_Y^\# \hat{\alpha}(x)\end{aligned}$$

Low-rank estimator

$$\begin{aligned}\hat{h}(x) &= \sum_{i=1}^n \hat{\alpha}_i(x) \psi_{\mathcal{Y}}(y_i), \quad \text{where} \quad \hat{\alpha}(x) = (K_X + n\lambda I_n)^{-1} k_X^x \\ &= \sqrt{n} S_Y^{\#} \hat{\alpha}(x) \\ &= \sqrt{n} S_Y^{\#} (n S_X S_X^{\#} + n\lambda I_n)^{-1} \sqrt{n} S_X \psi_{\mathcal{X}}(x)\end{aligned}$$

Low-rank estimator

$$\begin{aligned}\hat{h}(x) &= \sum_{i=1}^n \hat{\alpha}_i(x) \psi_{\mathcal{Y}}(y_i), \quad \text{where} \quad \hat{\alpha}(x) = (K_X + n\lambda I_n)^{-1} k_X^x \\ &= \sqrt{n} S_Y^{\#} \hat{\alpha}(x) \\ &= \sqrt{n} S_Y^{\#} (n S_X S_X^{\#} + n\lambda I_n)^{-1} \sqrt{n} S_X \psi_{\mathcal{X}}(x) \\ \hat{h}(x) &= S_Y^{\#} S_X (S_X^{\#} S_X + \lambda I_{\mathcal{H}_X})^{-1} \psi_{\mathcal{X}}(x)\end{aligned}$$

Low-rank estimator

$$\begin{aligned}\hat{h}(x) &= \sum_{i=1}^n \hat{\alpha}_i(x) \psi_{\mathcal{Y}}(y_i), \quad \text{where} \quad \hat{\alpha}(x) = (K_X + n\lambda I_n)^{-1} k_X^x \\ &= \sqrt{n} S_Y^{\#} \hat{\alpha}(x) \\ &= \sqrt{n} S_Y^{\#} (n S_X S_X^{\#} + n\lambda I_n)^{-1} \sqrt{n} S_X \psi_{\mathcal{X}}(x) \\ \hat{h}(x) &= S_Y^{\#} S_X (S_X^{\#} S_X + \lambda I_{\mathcal{H}_X})^{-1} \psi_{\mathcal{X}}(x)\end{aligned}$$

Goal: Given orthogonal projectors \tilde{P}_X and \tilde{P}_Y onto subspaces of \mathcal{H}_X and \mathcal{H}_Y resp.

Low-rank estimator

$$\begin{aligned}\hat{h}(x) &= \sum_{i=1}^n \hat{\alpha}_i(x) \psi_{\mathcal{Y}}(y_i), \quad \text{where} \quad \hat{\alpha}(x) = (K_X + n\lambda I_n)^{-1} k_X^x \\ &= \sqrt{n} S_Y^{\#} \hat{\alpha}(x) \\ &= \sqrt{n} S_Y^{\#} (n S_X S_X^{\#} + n\lambda I_n)^{-1} \sqrt{n} S_X \psi_{\mathcal{X}}(x) \\ \hat{h}(x) &= \textcolor{red}{S}_Y^{\#} \textcolor{red}{S}_X (\textcolor{red}{S}_X^{\#} \textcolor{red}{S}_X + \lambda I_{\mathcal{H}_X})^{-1} \psi_{\mathcal{X}}(x)\end{aligned}$$

Goal: Given orthogonal projectors \tilde{P}_X and \tilde{P}_Y onto subspaces of \mathcal{H}_X and \mathcal{H}_Y resp.

$$\textcolor{red}{S}_X^{\#} \leftarrow \tilde{P}_X S_X^{\#} \quad \text{and} \quad \textcolor{red}{S}_Y^{\#} \leftarrow \tilde{P}_Y S_Y^{\#}$$

Low-rank estimator

$$\begin{aligned}\hat{h}(x) &= \sum_{i=1}^n \hat{\alpha}_i(x) \psi_{\mathcal{Y}}(y_i), \quad \text{where} \quad \hat{\alpha}(x) = (K_X + n\lambda I_n)^{-1} k_X^x \\ &= \sqrt{n} S_Y^\# \hat{\alpha}(x) \\ &= \sqrt{n} S_Y^\# (n S_X S_X^\# + n\lambda I_n)^{-1} \sqrt{n} S_X \psi_{\mathcal{X}}(x) \\ \hat{h}(x) &= \textcolor{red}{S}_Y^\# \textcolor{red}{S}_X (\textcolor{red}{S}_X^\# \textcolor{red}{S}_X + \lambda I_{\mathcal{H}_X})^{-1} \psi_{\mathcal{X}}(x)\end{aligned}$$

Goal: Given orthogonal projectors \tilde{P}_X and \tilde{P}_Y onto subspaces of \mathcal{H}_X and \mathcal{H}_Y resp.

$$\begin{aligned}\textcolor{red}{S}_X^\# &\leftarrow \tilde{P}_X \textcolor{green}{S}_X^\# \quad \text{and} \quad \textcolor{red}{S}_Y^\# \leftarrow \tilde{P}_Y \textcolor{green}{S}_Y^\# \\ \tilde{h}(x) &= \tilde{P}_Y \textcolor{green}{S}_Y^\# \textcolor{red}{S}_X \tilde{P}_X (\tilde{P}_X \textcolor{green}{S}_X^\# \textcolor{red}{S}_X \tilde{P}_X + \lambda I_{\mathcal{H}_X})^{-1} \psi_{\mathcal{X}}(x).\end{aligned}\tag{19}$$

Complexity of IOKR and SISOKR for various types of sketching

Table 3: Time and space complexities at training and inference for the IOKR and SISOKR algorithms with sub-sampling, p -sparsified ($p \in (0, 1]$) or Gaussian sketching, for a test set of size n_{te} and a candidate set of size n_c , such that $n_{te} \leq m_{\mathcal{X}}, m_{\mathcal{Y}} < n \leq n_c$. For the sake of simplicity, we omit the $\mathcal{O}(\cdot)$ in the following.

Method	Training		Inference	
	Time	Space	Time	Space
IOKR	n^3	n^2	$n_{te}nn_c$	nn_c
SISOKR (sub-sampling)	$\max(m_{\mathcal{X}}, m_{\mathcal{Y}})n$	$\max(m_{\mathcal{X}}, m_{\mathcal{Y}})n$	$n_{te}m_{\mathcal{Y}}n_c$	$m_{\mathcal{Y}}n_c$
SISOKR (p -sparsified)	$\max(m_{\mathcal{X}}, m_{\mathcal{Y}})^2pn$	$\max(m_{\mathcal{X}}, m_{\mathcal{Y}})pn$	$\max(n_{te}, nm_{\mathcal{Y}}p)m_{\mathcal{Y}}n_c$	$npm_{\mathcal{Y}}n_c$
SISOKR (Gaussian)	$\max(m_{\mathcal{X}}, m_{\mathcal{Y}})n^2$	n^2	$nm_{\mathcal{Y}}n_c$	nn_c

Related works and differences

Rudi et al. (2015):

1. **scalar** kernel Ridge regression
2. sketching **only** applied in the **input** feature space
3. **Nyström** approximation with **uniform** or **approximate leverage scores** sampling

Related works and differences

Rudi et al. (2015):

1. **scalar** kernel Ridge regression
2. sketching **only** applied in the **input** feature space
3. **Nyström** approximation with **uniform** or **approximate leverage scores** sampling

Ciliberto et al. (2020):

1. **vector-valued** kernel Ridge regression, with possibly infinite-dimensional outputs
2. **no approximation** considered

Related works and differences

Rudi et al. (2015):

1. **scalar** kernel Ridge regression
2. sketching **only** applied in the **input** feature space
3. **Nyström** approximation with **uniform** or **approximate leverage scores** sampling

Ciliberto et al. (2020):

1. **vector-valued** kernel Ridge regression, with possibly infinite-dimensional outputs
2. **no approximation** considered

This work:

1. **vector-valued** kernel Ridge regression, with possibly infinite-dimensional outputs
2. sketching applied in **both** the **input and output** feature space
3. generic **sub-Gaussian** sketches

SISOKR Excess-Risk bound

Theorem (SISOKR excess-risk bound)

Let $\delta \in [0, 1]$, $n \in \mathbb{N}$ sufficiently large such that

$\lambda = n^{-1/(1+\gamma_{\mathcal{X}})} \geq \frac{9\kappa_{\mathcal{X}}^2}{n} \log(\frac{n}{\delta})$. Under **Asm. 1, 2, 3 and 4**, the following holds with probability at least $1 - \delta$

$$\mathbb{E}[\|\tilde{h}(x) - h^*(x)\|_{\mathcal{H}_Y}^2]^{\frac{1}{2}} \leq S(n) + c_2 A_{\rho_X}^{\psi_X}(\tilde{P}_X) + A_{\rho_Y}^{\psi_Y}(\tilde{P}_Y)$$

where

$$S(n) = c_1 \log(4/\delta) n^{-\frac{1}{2(1+\gamma_{\mathcal{X}})}} \quad (\text{regression error})$$

$$A_{\rho_Z}^{\psi_Z}(\tilde{P}_Z) = \mathbb{E}_Z[\|(\tilde{P}_Z - I_{\mathcal{H}_Z})\psi_Z(z)\|_{\mathcal{H}_Z}^2]^{\frac{1}{2}} \quad (\text{sketching reconstruction error})$$

and $c_1, c_2 > 0$ are constants independent of n and δ defined in the proofs.

Sub-Gaussian sketch

Definition

A sub-Gaussian sketch $R_{\mathcal{Z}} \in \mathbb{R}^{m_{\mathcal{Z}} \times n}$ is composed with i.i.d. entries such that

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}] = 0 \quad (20)$$

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}^2] = 1/m \quad (21)$$

$$R_{\mathcal{Z}_{ij}} \sim \frac{\nu_{\mathcal{Z}}^2}{m} - \text{sub-Gaussian}, \quad \text{with} \quad \nu_{\mathcal{Z}} \geq 1 \quad (22)$$

Sub-Gaussian sketch

Definition

A sub-Gaussian sketch $R_{\mathcal{Z}} \in \mathbb{R}^{m_{\mathcal{Z}} \times n}$ is composed with i.i.d. entries such that

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}] = 0 \quad (20)$$

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}^2] = 1/m \quad (21)$$

$$R_{\mathcal{Z}_{ij}} \sim \frac{\nu_{\mathcal{Z}}^2}{m} - \text{sub-Gaussian, with } \nu_{\mathcal{Z}} \geq 1 \quad (22)$$

Examples:

- matrix composed with i.i.d. Gaussian entries

Sub-Gaussian sketch

Definition

A sub-Gaussian sketch $R_{\mathcal{Z}} \in \mathbb{R}^{m_{\mathcal{Z}} \times n}$ is composed with i.i.d. entries such that

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}] = 0 \quad (20)$$

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}^2] = 1/m \quad (21)$$

$$R_{\mathcal{Z}_{ij}} \sim \frac{\nu_{\mathcal{Z}}^2}{m} - \text{sub-Gaussian}, \quad \text{with} \quad \nu_{\mathcal{Z}} \geq 1 \quad (22)$$

Examples:

- matrix composed with i.i.d. Gaussian entries
- matrix composed with i.i.d. bounded random variables

Sub-Gaussian sketch

Definition

A sub-Gaussian sketch $R_{\mathcal{Z}} \in \mathbb{R}^{m_{\mathcal{Z}} \times n}$ is composed with i.i.d. entries such that

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}] = 0 \quad (20)$$

$$\mathbb{E} [R_{\mathcal{Z}_{ij}}^2] = 1/m \quad (21)$$

$$R_{\mathcal{Z}_{ij}} \sim \frac{\nu_{\mathcal{Z}}^2}{m} - \text{sub-Gaussian}, \quad \text{with} \quad \nu_{\mathcal{Z}} \geq 1 \quad (22)$$

Examples:

- matrix composed with i.i.d. Gaussian entries
- matrix composed with i.i.d. bounded random variables
- matrix composed with i.i.d. Gaussian/bounded r.v. multiplied with independent Bernoulli r.v. (El Ahmad et al., 2023)

Sub-Gaussian Sketching Reconstruction Error

Theorem (Sub-Gaussian sketching reconstruction error)

Under **Asm. 1, 2, 3 and 4**, for $\delta \in (0, 1/e]$, $n \in \mathbb{N}$ sufficiently large such that $\frac{9}{n} \log(n/\delta) \leq n^{-\frac{1}{1+\gamma_Z}} \leq \|C_Z\|_{\text{op}}/2$, then if

$$m_Z \geq c_4 \max \left(\nu_Z^2 n^{\frac{\gamma_Z + \mu_Z}{1+\gamma_Z}}, \nu_Z^4 \log(1/\delta) \right), \quad (23)$$

then with probability $1 - \delta$

$$\mathbb{E}_Z[\|(\tilde{P}_Z - I_{\mathcal{H}_Z})\psi_Z(Z)\|_{\mathcal{H}_Z}^2] \leq c_3 n^{-\frac{1-\gamma_Z}{1+\gamma_Z}} \quad (24)$$

where $c_3, c_4 > 0$ are constants independent of n, m_Z, δ defined in the proofs.

Synthetic Least Squares Regression

1) $n = 10,000$, $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$, $d = 300$, $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ linear kernels \implies
 $\mathcal{H}_{\mathcal{X}} = \mathcal{H}_{\mathcal{Y}} = \mathbb{R}^d$

Synthetic Least Squares Regression

1) $n = 10,000$, $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$, $d = 300$, $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ linear kernels \implies
 $\mathcal{H}_{\mathcal{X}} = \mathcal{H}_{\mathcal{Y}} = \mathbb{R}^d$

2) Construct covariance matrices $C_{\mathcal{X}}$ and E such that $\sigma_k(C_{\mathcal{X}}) = k^{-3/2}$
and $\sigma_k(E) = 0.2k^{-1/10}$

Synthetic Least Squares Regression

1) $n = 10,000$, $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$, $d = 300$, $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ linear kernels $\implies \mathcal{H}_{\mathcal{X}} = \mathcal{H}_{\mathcal{Y}} = \mathbb{R}^d$

2) Construct covariance matrices $C_{\mathcal{X}}$ and E such that $\sigma_k(C_{\mathcal{X}}) = k^{-3/2}$ and $\sigma_k(E) = 0.2k^{-1/10}$

3) Draw $H_0 \sim \mathcal{N}(0, I_d)$, and for $i \leq n$, $x_i \sim \mathcal{N}(0, C_{\mathcal{X}})$, $\epsilon_i \sim \mathcal{N}(0, E)$,

$$y_i = C_{\mathcal{X}} H_0 x_i + \epsilon_i \tag{25}$$

Synthetic Least Squares Regression

1) $n = 10,000$, $\mathcal{X} = \mathcal{Y} = \mathbb{R}^d$, $d = 300$, $k_{\mathcal{X}}$ and $k_{\mathcal{Y}}$ linear kernels $\implies \mathcal{H}_{\mathcal{X}} = \mathcal{H}_{\mathcal{Y}} = \mathbb{R}^d$

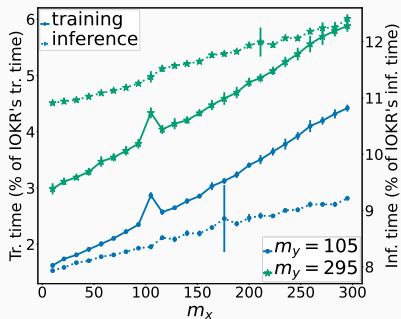
2) Construct covariance matrices $C_{\mathcal{X}}$ and E such that $\sigma_k(C_{\mathcal{X}}) = k^{-3/2}$ and $\sigma_k(E) = 0.2k^{-1/10}$

3) Draw $H_0 \sim \mathcal{N}(0, I_d)$, and for $i \leq n$, $x_i \sim \mathcal{N}(0, C_{\mathcal{X}})$, $\epsilon_i \sim \mathcal{N}(0, E)$,

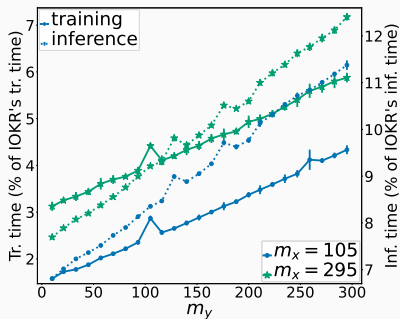
$$y_i = C_{\mathcal{X}} H_0 x_i + \epsilon_i \quad (25)$$

4) $(2 \cdot 10^{-2})$ -SR input and output sketches

Synthetic Least Squares Regression



(a) Training and inference time w.r.t. m_x for $m_y \in \{105, 295\}$



(b) Training and inference time w.r.t. m_y for $m_x \in \{105, 295\}$

Synthetic Least Squares Regression

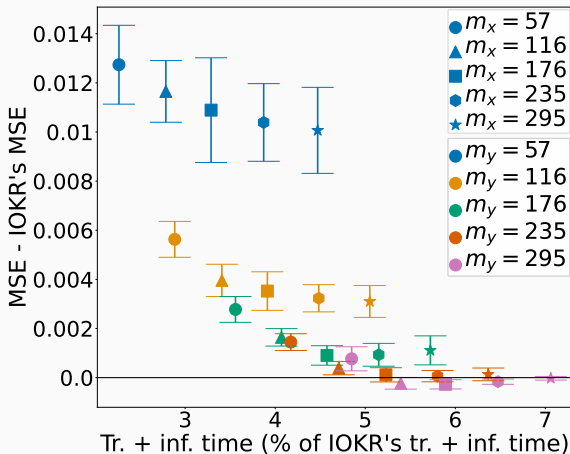
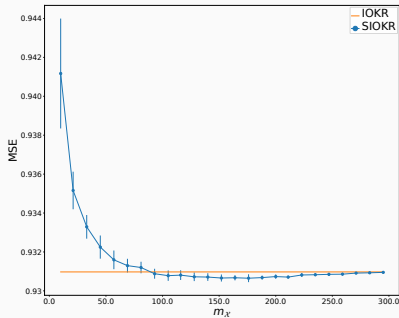
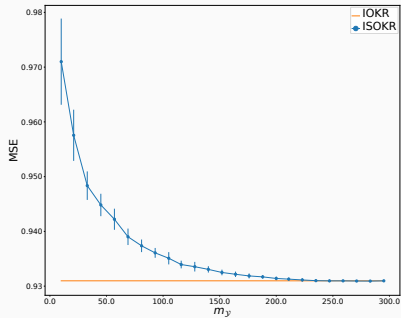


Figure 2: MSE w.r.t. learning time for different values of m_x and m_y

Synthetic Least Squares Regression



(a) SIOKR



(b) ISOKR

Multi-Label Classification

Bibtex and **Bookmarks** (Katakis et al., 2008): tag recommendation problems

Mediamill: detection of semantic concepts in a video

Table 4: Multi-label data sets description.

Data set	n	n_{te}	$n_{features}$	n_{labels}
Bibtex	4880	2515	1836	159
Bookmarks	60000	27856	2150	298
Mediamill	30993	12914	120	101

Multi-Label Classification: Statistical Performance

Table 5: F_1 scores on tag prediction from text data.

Method	Bibtex	Bookmarks	Mediamill
LR	37.2	30.7	NA
SPEN	42.2	34.4	NA
PRLR	44.2	34.9	NA
DVN	44.7	37.1	NA
SISOKR	44.1 ± 0.07	39.3 ± 0.61	57.26 ± 0.04
ISOKR	44.8 ± 0.01	NA	58.02 ± 0.01
SIOKR	44.7 ± 0.09	39.1 ± 0.04	57.33 ± 0.04
IOKR	44.9	NA	58.17

Multi-Label Classification: Computational Performance

Table 6: Training/inference times (in seconds).

Method	Bibtex	Bookmarks	Mediamill
SISOKR	1.41 ± 0.03 / 0.46 ± 0.01	118 ± 1.5 / 20 ± 0.2	66 ± 0.1 / 4 ± 0.01
ISOKR	2.51 ± 0.06 / 0.58 ± 0.01	NA	636 ± 3.7 9 ± 0.2
SIOKR	1.99 ± 0.07 / 1.22 ± 0.03	354 ± 2.1 / 297 ± 2.1	199 ± 0.1 / 121 ± 0.02
IOKR	2.54 / 1.18	NA	621 / 204

Metabolite identification

Inputs: tandem mass spectrum of a metabolite (small molecule

Outputs: molecular structure, i.e. fingerprints, encoded by binary vectors of length $d = 7593$

$n = 6974$ and each molecule is associated to a candidate set: median size = 292 and largest = 36,918 fingerprints

Table 7: MSE and standard errors for the metabolite identification problem. SPEN directly predicts outputs in \mathcal{Y} , then MSE is not defined.

Method	MSE	Tanimoto-Gaussian loss	Top-1 5 10 accuracies
SISOKR	0.813 ± 0.002	0.566 ± 0.007	25.1% 54.2% 64.7%
ISOKR	0.794 ± 0.003	0.509 ± 0.009	28.0% 58.9% 68.9%
SIOKR	0.793 ± 0.002	0.492 ± 0.008	29.5% 61.3% 70.9%
IOKR	0.780 ± 0.002	0.486 ± 0.008	29.6% 61.6% 71.4%
SPEN	NA	0.537 ± 0.008	25.9% 54.1% 64.3%

Table 8: Comparison of training/inference computation times (in seconds).

Method	Metabolite
SISOKR	4.05 ± 0.05 / 1112 ± 29
ISOKR	6.25 ± 50.31 / 1133 ± 32
SIOKR	1.25 ± 0.02 / 1179 ± 37
IOKR	3.54 ± 0.15 / 1191 ± 38

p -Sparsified Sketches: Definition

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

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

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

Computational Property: *Decomposition trick*

Let $m' = \sum_{j=1}^n \mathbb{I}\{R_{:j} \neq 0_s\}$,

$$R = R_{SG} R_{SS},$$

where

- $R_{SG} \in \mathbb{R}^{m \times m'} \leftarrow$ deleting the null columns from R
- $R_{SS} \in \mathbb{R}^{m' \times n} \leftarrow$ sampling the rows of I_n corresponding to the indices of non-zero columns of R .

Example:

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

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

Advantages of sub-sampling sketch

Let $X = \{x_1, \dots, x_5\}$, $k_X^{x_i} = (k(x_i, x_1), \dots, k(x_i, x_5))$ and

$$R_{SS} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

Advantages of sub-sampling sketch

Let $X = \{x_1, \dots, x_5\}$, $k_X^{x_i} = (k(x_i, x_1), \dots, k(x_i, x_5))$ and

$$R_{SS} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$R_{SS}K = \begin{pmatrix} k_X^{x_1} \\ k_X^{x_4} \end{pmatrix} \quad \text{and} \quad R_{SS}KR_{SS}^\top = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_4) \\ k(x_4, x_1) & k(x_4, x_4) \end{pmatrix}$$

\iff

1. Sample $X' = \{x_1, x_4\}$
2. Directly construct sub-Gram matrices $K_{X',X} \in \mathbb{R}^{2 \times 5}$ and $K_{X',X'} \in \mathbb{R}^{2 \times 2}$

Advantages of sub-sampling sketch

Let $X = \{x_1, \dots, x_5\}$, $k_X^{x_i} = (k(x_i, x_1), \dots, k(x_i, x_5))$ and

$$R_{SS} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$

$$R_{SS}K = \begin{pmatrix} k_X^{x_1} \\ k_X^{x_4} \end{pmatrix} \quad \text{and} \quad R_{SS}KR_{SS}^\top = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_4) \\ k(x_4, x_1) & k(x_4, x_4) \end{pmatrix}$$

\iff

1. Sample $X' = \{x_1, x_4\}$
2. Directly construct sub-Gram matrices $K_{X', X} \in \mathbb{R}^{2 \times 5}$ and $K_{X', X'} \in \mathbb{R}^{2 \times 2}$

\implies No need to compute costly matrix multiplications!

\implies No need to compute the whole K and store it in memory!

Time and Space Complexities of $R \cdot K_Z$

Let C_k be the cost of computing $k(x, x')$ for a couple $(x, x') \in \mathcal{X}^2$

- Standard sketch (e.g. Gaussian): $\mathcal{O}(C_k n^2 + n^2 m)$ and $\mathcal{O}(n^2)$,
- p -sparsified sketch: $\mathcal{O}(C_k n^2 m p + n^2 m^2 p)$ and $\mathcal{O}(n^2 m p)$.

\Rightarrow Complexity reduction if $p < 1/m$

Goal of p -sparsified sketches and related works

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

Goal of p -sparsified sketches and related works

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