

Motivation

Problem: learn a mapping $f : \mathcal{X} \rightarrow \mathcal{Y}$, where \mathcal{Y} is a **structured** space (e.g. graphs, rankings)
Existing works: $\hat{f} = d \circ \hat{h}$: 2-step surrogate method based on input/output kernels [3, 1, 2]
Advantage: **versatility** (i.e., able to handle different output types within a unified framework)
Drawback: **lack of expressiveness** (i.e., not able to handle complex inputs such as texts)
Goal: Build a **versatile** and **expressive** estimator

Some Notations

- $k : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ p.d. kernel, \mathcal{H} its RKHS, $\psi(y) := k(\cdot, y) \in \mathcal{H}$ **relevant representation of the outputs**
- $m \ll n$, $R \in \mathbb{R}^{m \times n}$ **sketching** matrix, i.e. randomly drawn matrix (e.g. **sub-sampling**, **Gaussian**)
- $K = (k(y_i, y_j))_{1 \leq i, j \leq n} \in \mathbb{R}^{n \times n}$, $\tilde{K} = RKR^T \in \mathbb{R}^{m \times m}$ and $\{(\sigma_i(\tilde{K}), \tilde{v}_i), i \in [m]\}$ its **eigenpairs**
- $\hat{\mathcal{H}} = \text{span}((\psi(y_i))_{i=1}^n)$, $\hat{C} = (1/n) \sum_{i=1}^n \psi(y_i) \otimes \psi(y_i) \in \hat{\mathcal{H}}^{\mathcal{H}}$ **empirical covariance operator**
- $\tilde{\mathcal{H}} = \text{span}((\sum_{j=1}^n R_{ij} \psi(y_j))_{i=1}^m)$, $\tilde{C} = \frac{1}{n} \sum_{l=1}^m (\sum_{i=1}^n R_{li} \psi(y_i)) \otimes (\sum_{j=1}^n R_{lj} \psi(y_j)) \in \tilde{\mathcal{H}}^{\mathcal{H}}$

Output Kernel Regression: a surrogate approach

Kernel-induced loss: $\Delta(y, y') := \|\psi(y) - \psi(y')\|_{\mathcal{H}}^2 = k(y, y) - 2k(y, y') + k(y', y')$

Goal: for $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Y}$ a **Deep Neural Network** ($\theta \in \Theta$ denotes its weights), solve

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \|\psi(f_{\theta}(x_i)) - \psi(y_i)\|_{\mathcal{H}}^2 \quad (1)$$

How: let $h_{\theta} : \mathcal{X} \rightarrow \mathcal{H}$ be a **DNN, 2-step surrogate method**:

- $\hat{\theta} = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \|h_{\theta}(x_i) - \psi(y_i)\|_{\mathcal{H}}^2$ (**training step**)
- $f_{\hat{\theta}}(x) = d \circ h_{\hat{\theta}}(x) = \arg \min_{y \in \mathcal{Y}} \|h_{\hat{\theta}}(x) - \psi(y)\|_{\mathcal{H}}^2$ (**inference step**)

Problem: what if $\psi(y)$ is infinite-dimensional or implicit?

Deep Sketched Output Kernel Regression

Solution: consider an orthonormal basis $\tilde{E} = ((\tilde{e}_i)_{i=1}^p)$ of a p -dimensional subspace of \mathcal{H} , where $p \in \mathbb{N}$ is small, and for a DNN $g_W : \mathcal{X} \rightarrow \mathbb{R}^p$ (W its weights),

$$h_{\theta}(x) := g_{\tilde{E}} \circ g_W(x) = \sum_{j=1}^p g_W(x)_j \tilde{e}_j$$

How to build the basis \tilde{E} ?

Let $p = \text{rank}(\tilde{K})$, and $\forall 1 \leq i \leq p$, $\tilde{e}_i = \sqrt{\frac{n}{\sigma_i(\tilde{K})}} \sum_{j=1}^n [R^T \tilde{v}_i]_j \psi(y_j) \in \mathcal{H}$.

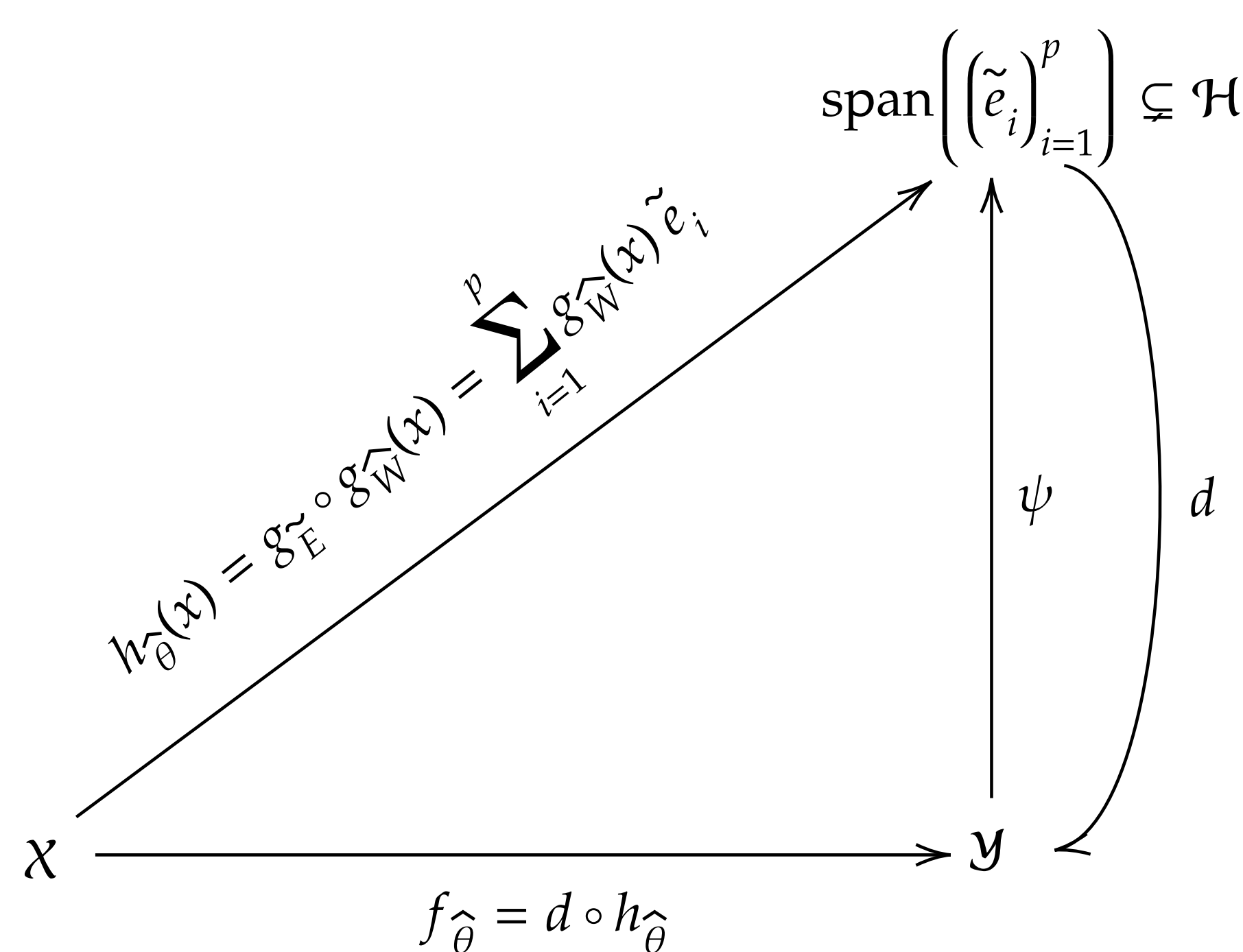
Proposition. The \tilde{e}_i s are the **eigenfunctions**, associated to the eigenvalues $\sigma_i(\tilde{K})/n$, of \tilde{C} whose range is $\tilde{\mathcal{H}}$. Then, $\tilde{E} = ((\tilde{e}_i)_{i=1}^p)$ is an **orthonormal basis** of $\tilde{\mathcal{H}}$.

How to solve the surrogate problem and learn the weights W ?

Proposition. Let $\tilde{E} = ((\tilde{e}_i)_{i=1}^p)$ and $h_{\theta} = g_{\tilde{E}} \circ g_W$. Then

$$\frac{1}{n} \sum_{i=1}^n \|h_{\theta}(x_i) - \psi(y_i)\|_{\mathcal{H}}^2 = \frac{1}{n} \sum_{i=1}^n \|g_W(x_i) - \tilde{\psi}(y_i)\|_2^2,$$

where $\tilde{\psi}(y) = (\tilde{e}_1(y), \dots, \tilde{e}_p(y))^T = \tilde{D}_p^{-1/2} \tilde{V}_p^T R k^y \in \mathbb{R}^p$, $\tilde{D}_p \in \mathbb{R}^p \times \mathbb{R}^p$ and $\tilde{V}_p \in \mathbb{R}^m \times \mathbb{R}^p$ are such that $\tilde{V}_p \tilde{D}_p \tilde{V}_p^T = \tilde{K}$ (SVD of \tilde{K}), and $k^y = (k(y, y_1), \dots, k(y, y_n))$.



Algorithm

1. Training. a. Computations for the basis \tilde{E} .

- Construct $\tilde{D}_p \in \mathbb{R}^p \times \mathbb{R}^p$, $\tilde{V}_p \in \mathbb{R}^m \times \mathbb{R}^p$ such that $\tilde{V}_p \tilde{D}_p \tilde{V}_p^T = \tilde{K}$ (SVD of \tilde{K})
- $\tilde{\Omega} = \tilde{D}_p^{-1/2} \tilde{V}_p^T \in \mathbb{R}^p \times m$

1. Training. b. Solving the surrogate problem.

- $\tilde{\psi}(y_i) = \tilde{\Omega} R k^{y_i} \in \mathbb{R}^p, \forall 1 \leq i \leq n$, $\tilde{\psi}(y_i^{\text{val}}) = \tilde{\Omega} R k^{y_i^{\text{val}}} \in \mathbb{R}^p, \forall 1 \leq i \leq n_{\text{val}}$
- $\hat{W} = \arg \min_{W \in \mathcal{W}} \frac{1}{n} \sum_{i=1}^n \|g_W(x_i) - \tilde{\psi}(y_i)\|_2^2$ (training of g_W with training $\{(x_i, \tilde{\psi}(y_i))\}_{i=1}^n$ and validation $\{(x_i^{\text{val}}, \tilde{\psi}(y_i^{\text{val}}))\}_{i=1}^{n_{\text{val}}}$ pairs and Mean Squared Error loss)

2. Inference.

- $\tilde{\psi}(y_i^c) = \tilde{\Omega} R k^{y_i^c} \in \mathbb{R}^p, \forall 1 \leq i \leq n_c$
- $f_{\hat{\theta}}(x_i^c) = y_j^c$ where $j = \arg \max_{1 \leq j \leq n_c} g_{\hat{W}}(x_i^c)^T \tilde{\psi}(y_j^c), \forall 1 \leq i \leq n_c$

Sketching Size Selection Strategy

Goal: Set the minimal value of m s.t. it captures the information contained in \hat{C}

Solutions:

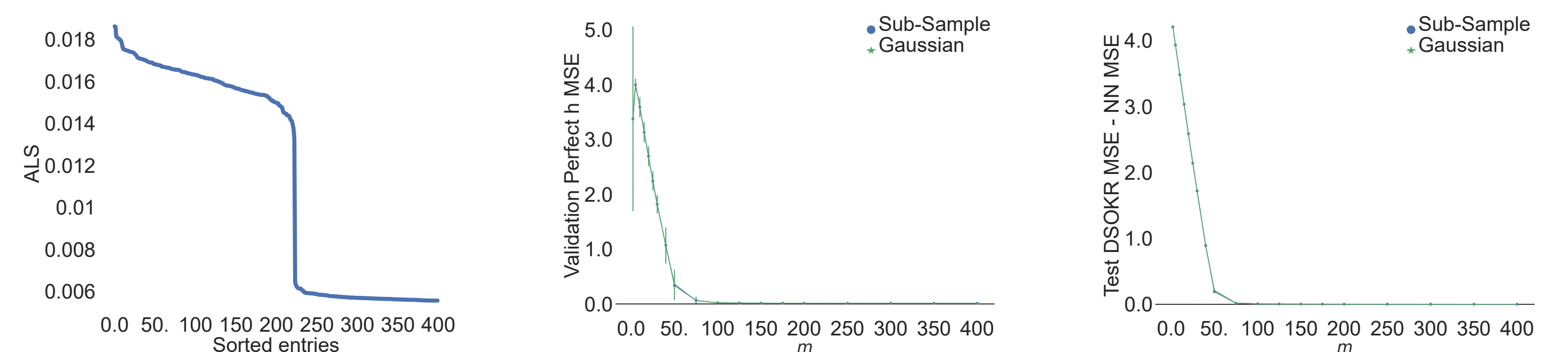
- Approximate leverage scores of \hat{C}
- Set the optimal m according to the performance of the perfect h estimator on the validation set,

$$h : (x, y) \mapsto \sum_{j=1}^p \langle \tilde{e}_j, \psi(y) \rangle_{\mathcal{H}} \tilde{e}_j = \sum_{j=1}^p \tilde{\psi}(y)_j \tilde{e}_j$$

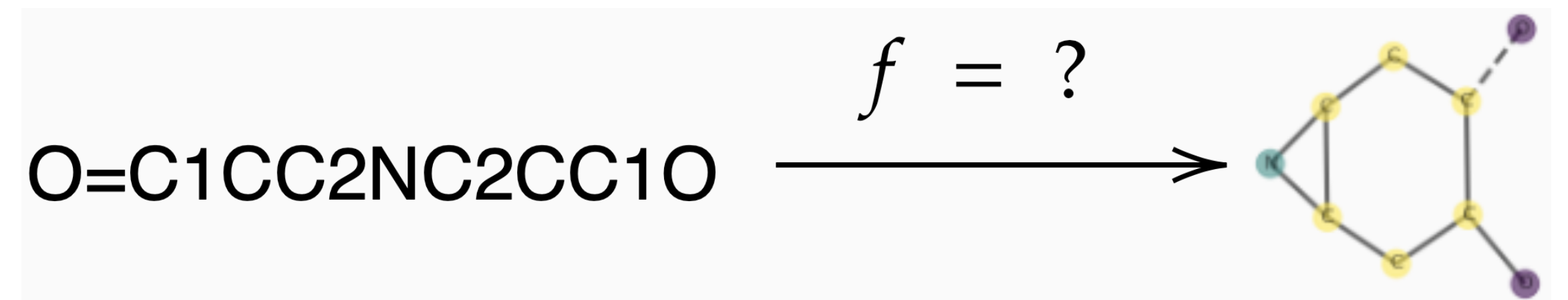
Experiment: Synthetic Least Squares Regression

Setting: $n = 50,000$ training data points, $\mathcal{X} = \mathbb{R}^{2,000}$, $\mathcal{Y} = \mathbb{R}^{1,000}$, k linear kernel so that $\mathcal{H} = \mathcal{Y} = \mathbb{R}^{1,000}$.

Goal: build a dataset such that the outputs lie in a **subspace** of \mathcal{Y} of dimension $d = 50 < 1,000$.



Experiment: SMILES to Molecule on the QMg dataset



	GED w/o edge feature ↓	GED w/ edge feature ↓
SISOKR	3.330 ± 0.080	4.192 ± 0.109
NNBary-FGW	5.115 ± 0.129	-
Sketched ILE-FGW	2.998 ± 0.253	-
DSOKR	1.951 ± 0.074	2.960 ± 0.079

Experiment: Text to Molecule on the ChEBI-20 dataset

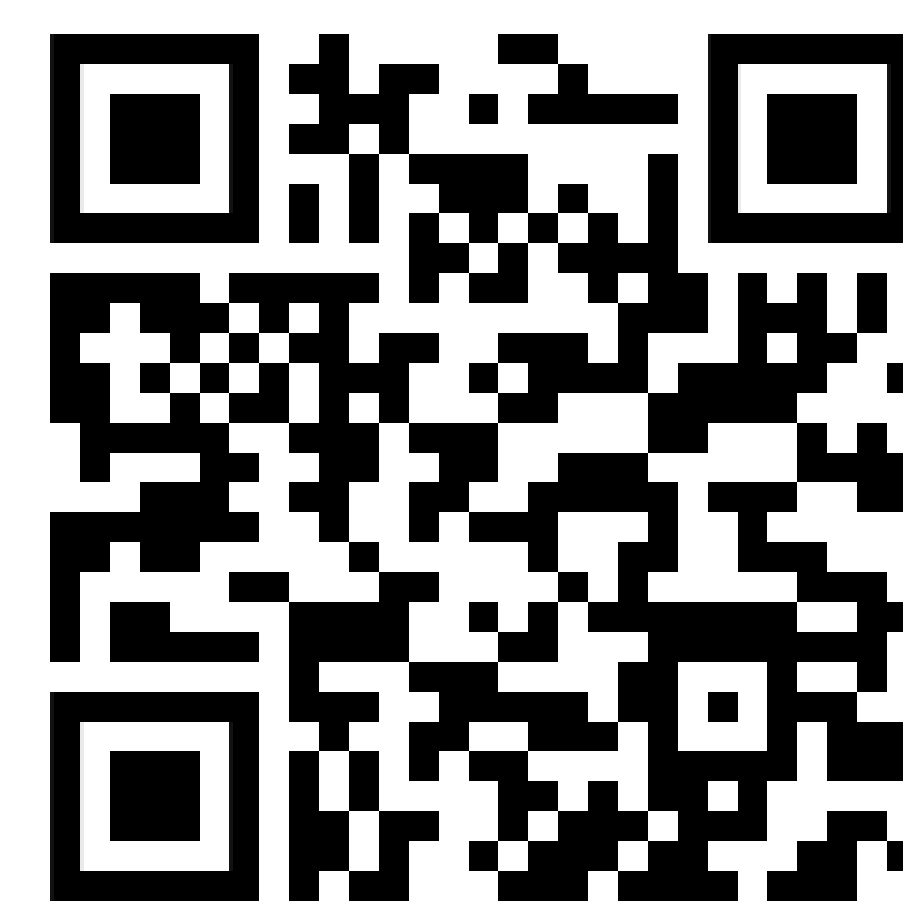
Water is an oxygen hydride consisting of an oxygen atom that is covalently bonded to two hydrogen atoms. $\xrightarrow{f = ?}$

	Hits@1 ↑	Hits@10 ↑	MRR ↑
SISOKR	0.4%	2.8%	0.015
SciBERT Regression	16.8%	56.9%	0.298
CMAM - MLP	34.9%	84.2%	0.513
CMAM - GCN	33.2%	82.5%	0.495
CMAM - Ensemble (MLP×3 + GCN×3)	44.2%	88.7%	0.597
DSOKR - SubSample Sketch	48.2%	87.4%	0.624
DSOKR - Gaussian Sketch	49.0%	87.5%	0.630
DSOKR - Ensemble (SubSample×3)	51.0%	88.2%	0.642
DSOKR - Ensemble (Gaussian×3)	50.5%	87.9%	0.642
DSOKR - Ensemble (SubSample×3 + Gaussian×3)	50.0%	88.3%	0.640

References

- C. Brouard, M. Szafranski, and F. d'Alché Buc. Input output kernel regression: supervised and semi-supervised structured output prediction with operator-valued kernels. *JMLR*, 17(1):6105–6152, 2016.
- C. Ciliberto, L. Rosasco, and A. Rudi. A general framework for consistent structured prediction with implicit loss embeddings. *JMLR*, 21(98):1–67, 2020.
- J. Weston, O. Chapelle, V. Vapnik, A. Elisseeff, and B. Schölkopf. Kernel dependency estimation. In *NeurIPS*, pages 897–904. MIT Press, 2003.

Check-out our code!



Acknowledgements

This project has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement 101120237 (ELIAS), the Télécom Paris research chair on Data Science and Artificial Intelligence for Digitalized Industry and Services (DSADIS) and the PEPR-IA through the project FOUNDRY.