

Solving Sparse & High-Dimensional-Output Regression via Compression

Renyan Li¹ Zehui Chen² Guanyi Wang¹

NeurIPS 2024

National University of Singapore¹ Google²

Regression for very high-dim output

- Strong interpretability,
- Computational scalability,
- Provable accuracy.



algorithmic trading



RL for high-dim decision



grounding in LLM

Problem Setting

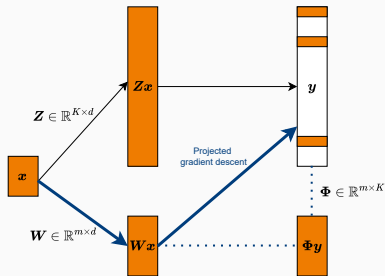
Given input $\mathbf{x} \in \mathbb{R}^d$, predict high-dim output $\mathbf{y} \in \mathbb{R}^K$,

$$\mathbf{y} := \arg \min_{\mathbf{u} \in \mathcal{Y}} \text{dist}(\mathbf{u}, \hat{g}(\mathbf{x}))$$

$$\text{with } \hat{g} := \arg \min_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{y}^i, g(\mathbf{x}^i)),$$

- high-dimensional output, i.e., $K \gg d$.
- sparsity constraint for interpretability, i.e., s -sparse output: $\mathcal{Y} \subseteq \{\mathbf{y} \in \mathbb{R}^K \mid \|\mathbf{y}\|_0 \leq s\}$.
- Using linear models \mathcal{G} , i.e., $g(\mathbf{x}) = \mathbf{Z}\mathbf{x}$, for simplicity.

Framework



- Color of the boxes: orange/white – nonzero/zero components.
- Height of boxes: dimension of corresponding vector.
- Black arrows: the common approach. Blue arrows: our approach.
- Φy : underlying feature for the output.

Propose an **two-stage framework via compression** with

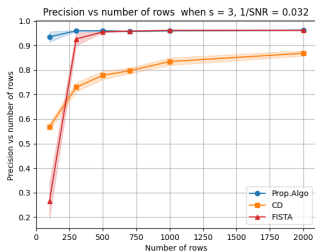
- improved computational efficiency
- the same order of generalization bounds before and after compression

Running Time: Time Complexity Comparison for each prediction.

Dataset	Prop.Algo.	OMP	CD	FISTA
Synthetic Data	≈ 1 second	200-400 seconds	<1 second	<3 seconds
EURLex-4K	<1 second	20-80 seconds	<1 second	≈ 1 second
Wiki10-31K	<5 seconds	500-700 seconds	<5 seconds	5-10 seconds

- Enjoys a better running time for large scale instances

Solution Accuracy



Experiment on synthetic data with $d = 10^4$, $K = 2 \times 10^4$, $s = 3$, and $n = 3 \times 10^4$.

- Outperforms the baselines on precision