



Total Least Squares Regression in Input Sparsity Time

Huaian Diao*, Zhao Song‡, David P. Woodruff†, Xin Yang‡

*Northeast Normal University, †Carnegie Mellon University, ‡University of Washington

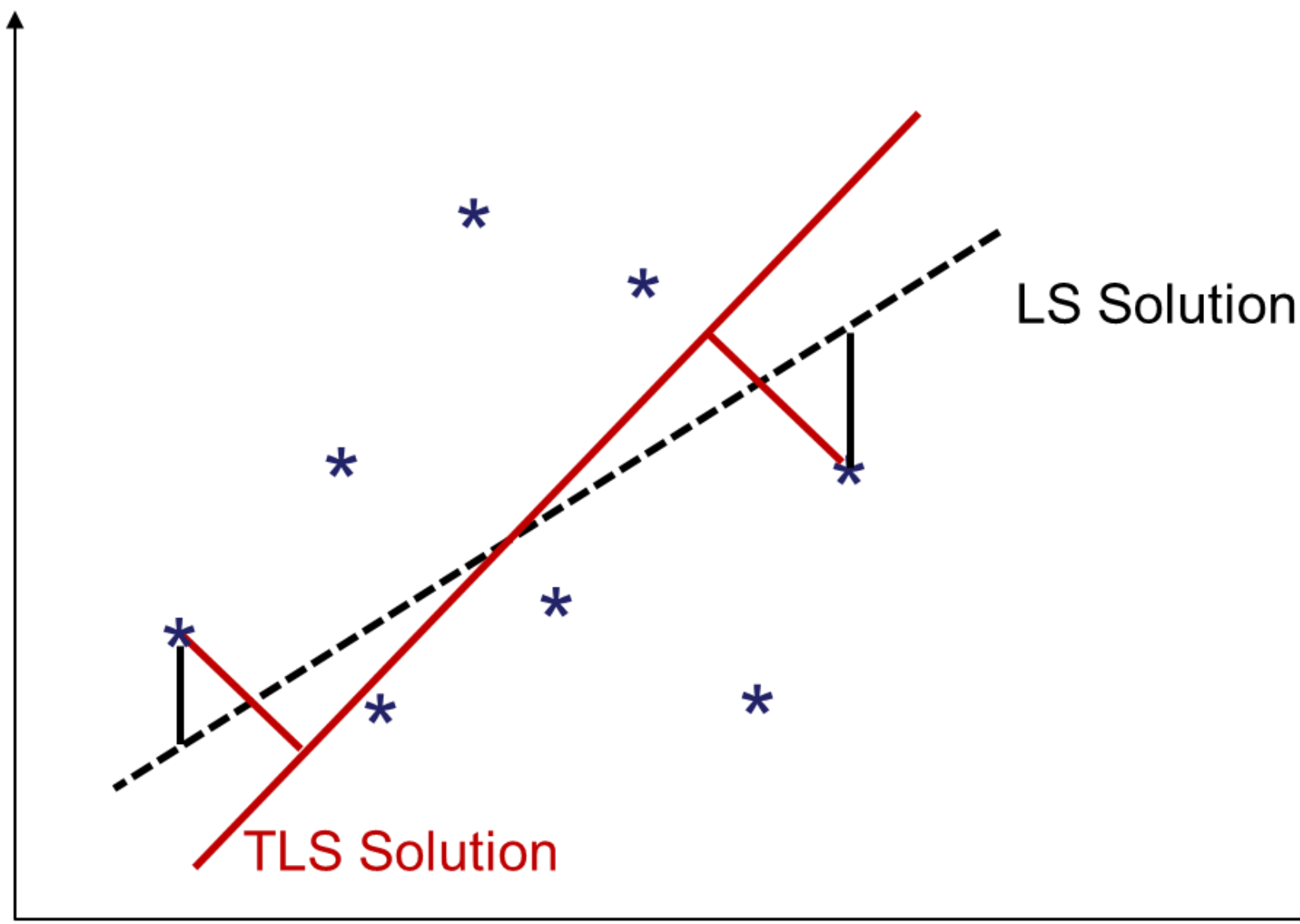
Total Least Squares Problem

Given: $A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{m \times d}$.

Task: Perturb A, B at minimal cost so that the linear system $AX = B$ has a solution.

- Output $X \in \mathbb{R}^{n \times d}$.
- Minimize $\|A - \hat{A}\|_F^2 + \|\hat{A}X - B\|_F^2$ over $\hat{A} \in \mathbb{R}^{m \times d}, X \in \mathbb{R}^{n \times d}$.

Intuition: the ordinary least squares problem $\min_X \|AX - B\|_F^2$ assumes only the labels B are corrupted; but the samples A can also be corrupted.



Compact formulation as low rank approximation: let $C = [A, B] \in \mathbb{R}^{m \times (n+d)}$, then the TLS problem is

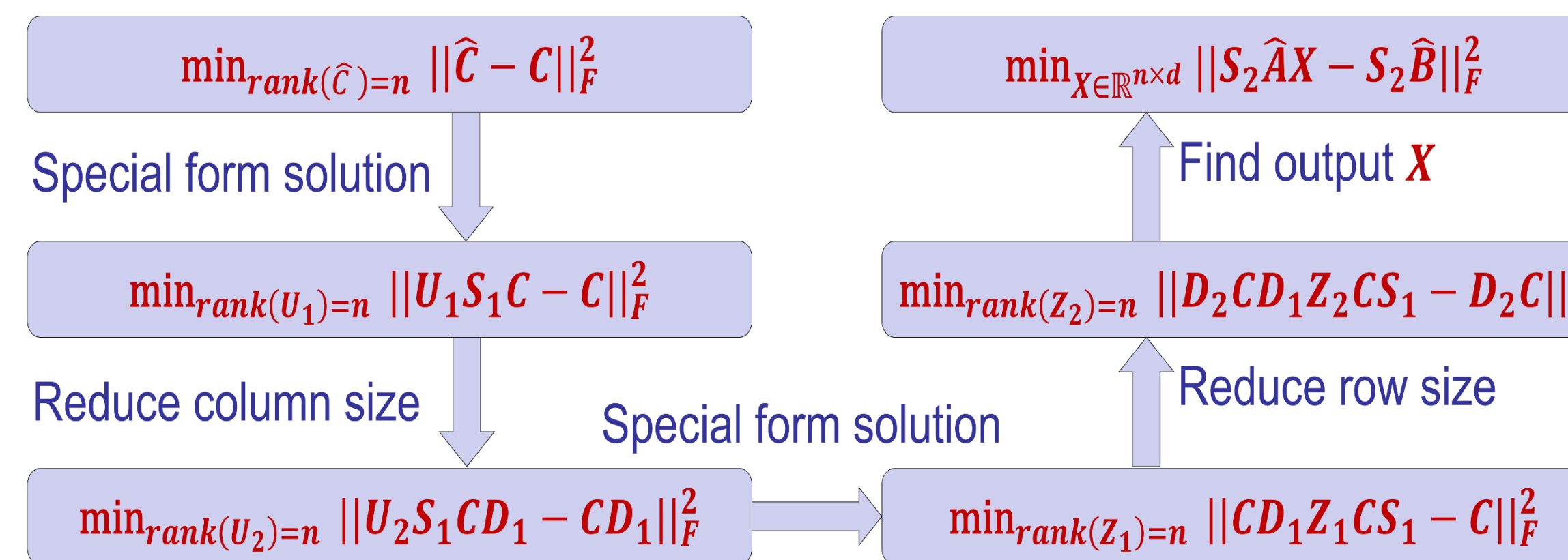
$$\text{OPT} := \min_{\text{rank}(\hat{C})=n} \|C - \hat{C}\|_F^2$$

Main Results

Theorem. For $0 < \epsilon < 1$, with probability **0.9**, we can output $X \in \mathbb{R}^{n \times d}$ so that there exists $\hat{A} \in \mathbb{R}^{m \times n}$ satisfying $\|[\hat{A}, \hat{A}X] - [A, B]\|_F^2 \leq (1 + \epsilon)\text{OPT}$ in time $\tilde{O}(\text{nnz}(A) + \text{nnz}(B)) + d \cdot \text{poly}(n/\epsilon)$. Here nnz is the number of non-zero entries, and \tilde{O} hide some logarithm factors.

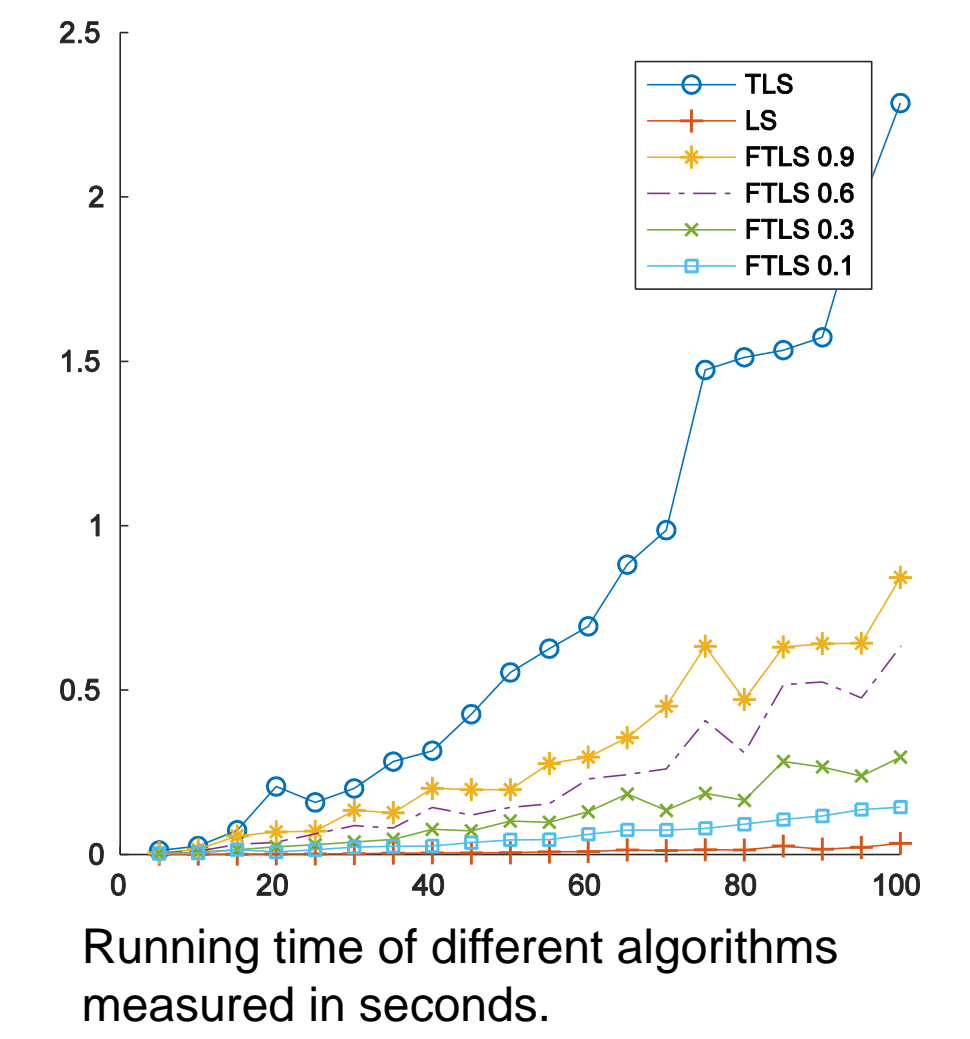
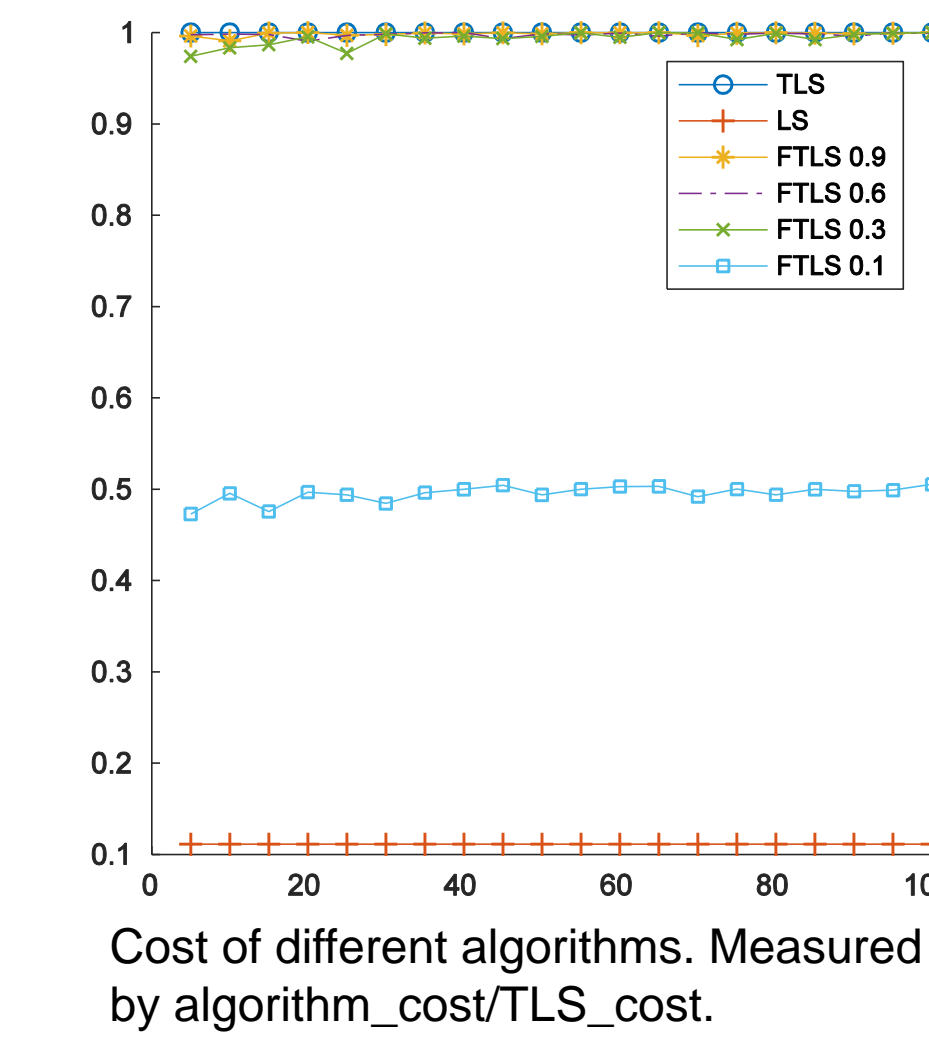
- The output X can be used for future prediction.
- Traditional algorithm for the TLS problem: SVD decomposition, takes $O(m(n+d)^2)$ time.
- Our algorithm runs in input sparsity time. Huge advantage when $m \gg n, d$ and the inputs A, B are sparse.
- Our running time is even smaller than writing down the optimal low rank approximation \hat{C} .
- Key ingredient: chain sketching matrices to reduce dimension.

Algorithm Flowchart



Experimental Results

Synthetic datasets:



Real datasets: regression datasets from UCI datasets.

Method	Cost	C-std	Time	T-std
TLS	0.10	0	1.12	0.05
LS	10^6	0	0.0012	0.0002
FTLS 0.9	0.10	0.0002	0.16	0.0058
FTLS 0.6	0.10	0.0003	0.081	0.0033
FTLS 0.3	0.10	0.0007	0.046	0.0022
FTLS 0.1	0.10	0.0016	0.034	0.0024

Cost and running time of different algorithms for the Airfoil Self-Noise dataset.

Method	Cost	C-std	Time	T-std
TLS	0.93	0	1.36	0.16
LS	666	0	0.0012	0.001
FTLS 0.9	0.93	0.0032	0.30	0.025
FTLS 0.6	0.94	0.0050	0.17	0.01
FTLS 0.3	0.95	0.01	0.095	0.005
FTLS 0.1	0.99	0.03	0.074	0.004

Cost and running time of different algorithms for the Red wine dataset.

- Our algorithm outperforms the ordinary TLS algorithm by running 10~50x faster.
- Our algorithm obtains small cost comparable to the TLS solution.

