

# LOW-RANK APPROXIMATION IN INPUT-SPARSITY TIME, WITH REGULARIZATION ~~AND ROBUSTNESS~~



HAIM AVRON  
TEL AVIV UNIVERSITY

**KEN CLARKSON**

IBM RESEARCH, ALMADEN



DAVID WOODRUFF  
IBM RESEARCH, ALMADEN

# THE PROBLEM

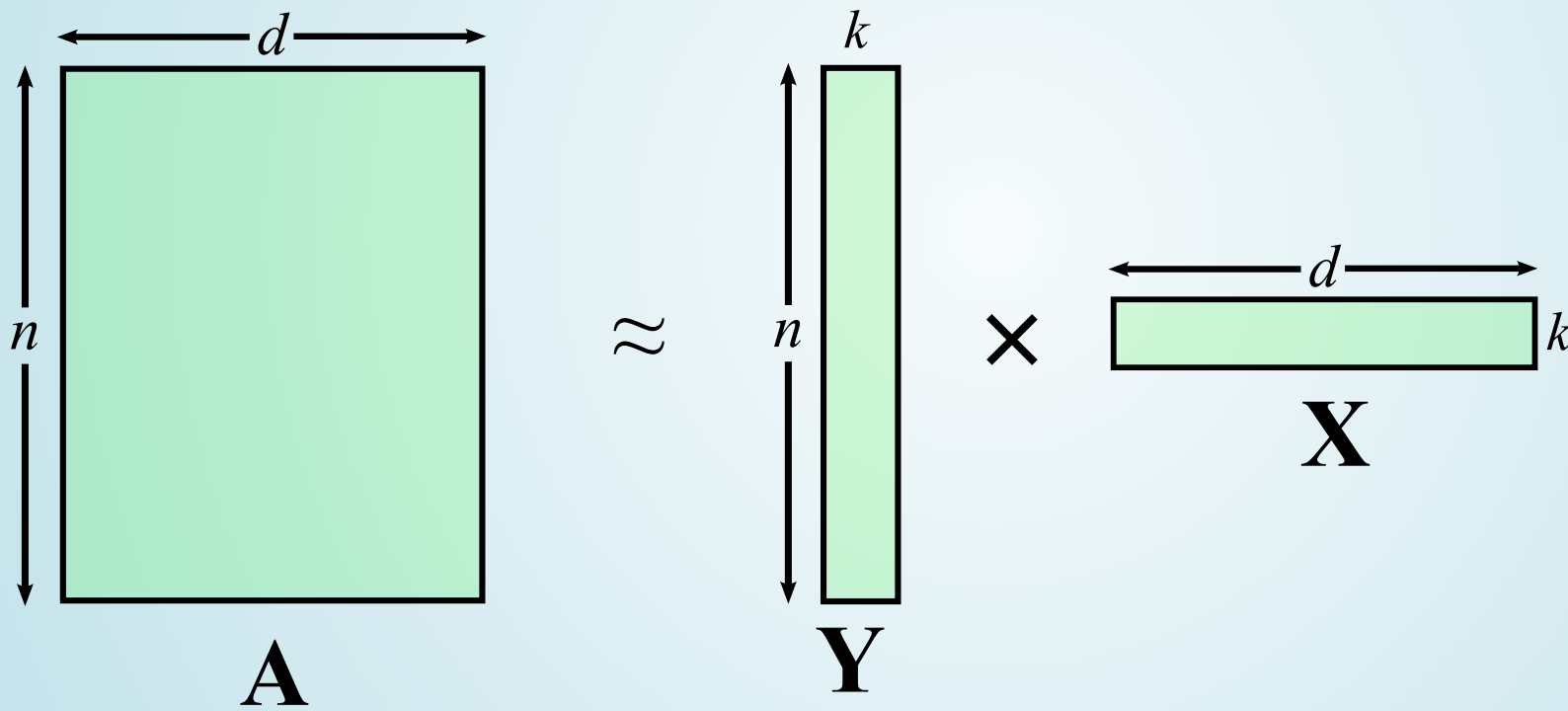
## Regularized low-rank approximation

Given

$\mathbf{A} \in \mathbb{R}^{n \times d}$ , integer  $k$ ,  $\lambda \geq 0$ ,

find

$$\mathbf{Y}_*, \mathbf{X}_* \equiv \operatorname{argmin}_{\substack{\mathbf{Y} \in \mathbb{R}^{n \times k} \\ \mathbf{X} \in \mathbb{R}^{k \times d}}} \|\mathbf{Y}\mathbf{X} - \mathbf{A}\|_F^2 + \lambda(\|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2)$$



## FUN FACTS

Equivalent to nuclear (trace) norm minimization [SS05]:

$$\operatorname{argmin}_{\substack{\mathbf{Y} \in \mathbb{R}^{n \times k} \\ \mathbf{X} \in \mathbb{R}^{k \times d}}} \|\mathbf{YX} - \mathbf{A}\|_F^2 + 2\lambda \|\mathbf{YX}\|_*$$

For best rank- $k$   $\mathbf{A}_k = \mathbf{U}_k \boldsymbol{\Sigma}_k \mathbf{V}_k^\top$ , solution is

$$\mathbf{Y}_* = \mathbf{U}_k \boldsymbol{\Lambda}, \mathbf{X}_* = \boldsymbol{\Lambda} \mathbf{V}_k^\top,$$

where  $\boldsymbol{\Lambda} \equiv (\boldsymbol{\Sigma}_k - \lambda \mathbf{I}_k)_+^{1/2}$ , and  $()_+$  takes negative entries to zero.

The nuclear norm "encourages" low rank;  
constraint insists on  $\operatorname{rank} \leq k$

## WHY THIS PROBLEM?

- Suppose

$$\mathbf{A} = \hat{\mathbf{A}} + \mathbf{N},$$

where  $\text{rank}(\hat{\mathbf{A}}) \leq k$ , and  $\mathbf{N}$  is noise

- Then  $\sigma_i(\mathbf{A}) \approx \sigma_i(\hat{\mathbf{A}}) + \sigma$ ,  
for  $\mathbf{N}$  i.i.d. Gaussian noise appropriate variance
  - Singular values  $\sigma_i()$
  - Singular vectors of  $\mathbf{A}$  lightly perturbed from  $\hat{\mathbf{A}}$
  - So recover  $\hat{\mathbf{A}}$  with the right  $\lambda$
- More generally,  
$$\|\mathbf{YX} - \hat{\mathbf{A}}\|_F^2 \leq \|\mathbf{YX} - \mathbf{A}\|_F^2 - \|\mathbf{N}\|_F^2 + 2\|\mathbf{N}\|_2 \|\mathbf{YX}\|_*,$$
  
so minimizing  $\|\mathbf{YX} - \mathbf{A}\|_F^2 + \|\mathbf{N}\|_2 \|\mathbf{YX}\|_*$  is a plausible proxy  
for minimizing  $\|\mathbf{YX} - \hat{\mathbf{A}}\|_F^2$

# APPROXIMATE SOLUTIONS

Given

$\epsilon > 0, \lambda \geq 0$

a  $(1 + \epsilon)$ -approximation  $\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}$  has

$$\begin{aligned} & \|\tilde{\mathbf{Y}}\tilde{\mathbf{X}} - \mathbf{A}\|_F^2 + \lambda(\|\tilde{\mathbf{X}}\|_F^2 + \|\tilde{\mathbf{Y}}\|_F^2) \\ & \leq (1 + \epsilon) \left[ \|\mathbf{Y}_*\mathbf{X}_* - \mathbf{A}\|_F^2 + \lambda(\|\mathbf{X}_*\|_F^2 + \|\mathbf{Y}_*\|_F^2) \right] \end{aligned}$$

Regularized low-rank approximation, relative error  $\epsilon$

With fixed probability near one,  
 $(1 + \epsilon)$ -approximation can be found in

$$O(\text{nnz}(\mathbf{A})) + \tilde{O}(k^2 \epsilon^{-1} (n + d)) + \text{poly}(\epsilon^{-1} k)$$

time.

input-sparsity-time algorithm

## RIDGE REGRESSION

Given  $\mathbf{Y}_*$ ,

$$\mathbf{X}_* = \operatorname{argmin}_{\mathbf{X} \in \mathbb{R}^{k \times d}} \|\mathbf{Y}_* \mathbf{X} - \mathbf{A}\|_F^2 + \lambda \|\mathbf{X}\|_F^2,$$

that is, *multiple-response ridge regression*.

**Observation.** Suppose  $P_{\mathbf{A}}$  projects to the rowspace of  $\mathbf{A}$ ,  
so for  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{x}^\top P_{\mathbf{A}} \in \operatorname{rowspan}(\mathbf{A})$ . Then

$$\|(\mathbf{Y}_* \mathbf{X} - \mathbf{A})P_{\mathbf{A}}\|_F^2 + \lambda \|\mathbf{X}P_{\mathbf{A}}\|_F^2 \leq \|\mathbf{Y}_* \mathbf{X} - \mathbf{A}\|_F^2 + \lambda \|\mathbf{X}\|_F^2,$$

so without loss of generality, each row of  $\mathbf{X}_*$  is in  $\operatorname{rowspan}(\mathbf{A})$ .

$\mathbf{X}_*$  in  $\operatorname{rowspan}(\mathbf{A})$

To solve these regression problems,  
we reduce to a smaller problem using *sparse embeddings*

# SPARSE EMBEDDINGS

A distribution over matrices  $\mathbf{S} \in \mathbb{R}^{m \times n}$   
so that for any given  $\mathbf{Y} \in \mathbb{R}^{n \times d'}$ ,  
with fixed probability near one,

oblivious subspace embedding

$$\|\mathbf{S}\mathbf{Y}\mathbf{x}\| \approx_{\epsilon} \|\mathbf{Y}\mathbf{x}\| \text{ for all } \mathbf{x} \in \mathbb{R}^{d'}$$

if  $\mathbf{S}$  is a sparse embedding for  $\mathbf{Y}$ ,  
then for all  $\mathbf{X} \in \mathbb{R}^{d' \times d}$ ,  
 $\|\mathbf{S}\mathbf{Y}\mathbf{X}\|_F^2 \approx_{\epsilon} \|\mathbf{Y}\mathbf{X}\|_F^2$

Also called *sketching* matrices

A distribution over matrices  $\mathbf{S} \in \mathbb{R}^{m \times n}$   
so that for any given  $\mathbf{Y} \in \mathbb{R}^{n \times d'}$  and  $\mathbf{A} \in \mathbb{R}^{n \times d}$ ,  
with fixed probability near one,

oblivious affine embeddings

$$\|\mathbf{S}(\mathbf{Y}\mathbf{X} - \mathbf{A})\|_F^2 \approx_{\epsilon} \|\mathbf{Y}\mathbf{X} - \mathbf{A}\|_F^2 \text{ for all } \mathbf{X} \in \mathbb{R}^{d'}$$

There are distributions for which  $\mathbf{S}$  has one non-zero per column,  
so  $\mathbf{S}\mathbf{Y}$  can be computed in  $\text{nnz}(\mathbf{Y})$  time, and  
 $\mathbf{S} \in \mathbb{R}^{m_S \times d'}$  with  $m_S$  small.

sparse oblivious affine embeddings

## MULTIPLE REDUCTIONS

If  $\mathbf{S}$  is an affine embedding, then

$$\tilde{\mathbf{X}} = \operatorname{argmin}_{\mathbf{X} \in \mathbb{R}^{k \times d}} \|\mathbf{S}(\mathbf{Y}_* \mathbf{X} - \mathbf{A})\|_F^2 + \lambda \|\mathbf{X}\|_F^2$$

is a good approximation to  $\mathbf{X}_*$ .

Applying observation above to  $\operatorname{rowspan}(\mathbf{S}\mathbf{A})$ ,

WLOG  $\tilde{\mathbf{X}} = \mathbf{Z}\mathbf{S}\mathbf{A}$  for some  $\mathbf{Z} \in \mathbb{R}^{k \times m_S}$

Applying again "on the other side":

There are sparse affine embeddings  $\mathbf{S}$  and  $\mathbf{R}^\top$  so that

two-sided reduction

$$\mathbf{W}_*, \mathbf{Z}_* \equiv \operatorname{argmin}_{\substack{\mathbf{W} \in \mathbb{R}^{m_R \times k} \\ \mathbf{Z} \in \mathbb{R}^{k \times m_S}}} \|\mathbf{A}\mathbf{R}\mathbf{W}\mathbf{Z}\mathbf{S}\mathbf{A} - \mathbf{A}\|_F^2 + \lambda(\|\mathbf{Z}\mathbf{S}\mathbf{A}\|_F^2 + \|\mathbf{A}\mathbf{R}\mathbf{W}\|_F^2)$$

yield  $(1 + \epsilon)$ -approximation  $\tilde{\mathbf{X}} \equiv \mathbf{Z}_*\mathbf{S}\mathbf{A}$ ,  $\tilde{\mathbf{Y}} \equiv \mathbf{A}\mathbf{R}\mathbf{W}_*$ .

With additional affine *and* sparse embeddings  $\hat{\mathbf{S}}$ ,  $\hat{\mathbf{R}}$ , can reduce to

$$\operatorname{argmin}_{\substack{\mathbf{W} \in \mathbb{R}^{m_R \times k} \\ \mathbf{Z} \in \mathbb{R}^{k \times m_S}}} \|\hat{\mathbf{S}}\mathbf{A}\mathbf{R}\mathbf{W}\mathbf{Z}\mathbf{S}\mathbf{A}\hat{\mathbf{R}} - \hat{\mathbf{S}}\mathbf{A}\hat{\mathbf{R}}\|_F^2 + \lambda(\|\mathbf{Z}\mathbf{S}\mathbf{A}\hat{\mathbf{R}}\|_F^2 + \|\hat{\mathbf{S}}\mathbf{A}\mathbf{R}\mathbf{W}\|_F^2)$$

a subproblem where all matrices have  $\operatorname{poly}(k/\epsilon)$  rows and columns.



## RIDGE REGRESSION, REPRISÉ

Suppose we want to apply sparse embedding machinery to

$$\mathbf{x}_* = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \lambda \|\mathbf{x}\|^2,$$

that is, solve instead

$$\min_{\mathbf{x} \in \mathbb{R}^d} \|\mathbf{S}(\mathbf{A}\mathbf{x} - \mathbf{b})\|^2 + \lambda \|\mathbf{x}\|^2.$$

When  $\lambda$  is very large,  $\lambda \|\mathbf{x}\|^2$  dominates, and  $\|\mathbf{S}(\mathbf{A}\mathbf{x} - \mathbf{b})\|^2 \approx \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$  can be loose.

So ridge regression should get easier as  $\lambda$  increases.

## THE STATISTICAL DIMENSION OF RIDGE REGRESSION

The statistical dimension (regression degrees of freedom, effective dimension) of ridge regression is

$$\text{sd}_\lambda(\mathbf{A}) \equiv \sum_i \frac{1}{1 + \lambda/\sigma_i^2(\mathbf{A})},$$

where  $\sigma_i(\mathbf{A})$  is the  $i$ 'th singular value of  $\mathbf{A}$ .

Note that

$\text{sd}_\lambda(\mathbf{A}) \leq \text{rank}(\mathbf{A}) \leq \min\{n, d\}$ , and  
 $\text{sd}_\lambda(\mathbf{A})$  is decreasing in  $\lambda$ .

## SMALLER EMBEDDINGS FOR RIDGE REGRESSION

smaller ridge

A ridge regression problem with  $n$  rows can be reduced to one with  $\tilde{O}(\epsilon^{-1} \text{sd}_\lambda(\mathbf{A}))$  rows, in  $O(\text{nnz}(\mathbf{A})) + \tilde{O}(d(\epsilon^{-1} \text{sd}_\lambda(\mathbf{A}) + \text{sd}_\lambda(\mathbf{A})^2))$  time. The resulting problem yields a  $(1 + \epsilon)$ -approximation.

- This extends prior work from sampling to sketching
- For low-rank approximation, some  $k \rightarrow O(\epsilon^{-1} \text{sd}_\lambda(\mathbf{Y}_*))$
- There is another reduction that reduces  $d$  to  $\text{poly}(\text{sd}_\lambda(\mathbf{A})\epsilon^{-1} (\sigma_1(A)^2/\lambda))$

## CONCLUSIONS

- Better dependence on  $k$  and  $\epsilon$  in "lower order terms" for this class
- Algorithms get faster as  $\lambda$  gets large
- We don't actually know  $\text{sd}_\lambda(\mathbf{A})$
- Also: can get pre-conditioners for kernel ridge regression
- Results generalize to regularizers  $f(\mathbf{X}, \mathbf{Y})$  satisfying some orthogonal invariance conditions

Thank you for your attention!