



Berkeley

SIAM Conference on  
Computational Science  
and Engineering



February 27-March 3, 2017  
Hilton Atlanta, Atlanta, Georgia, USA

# Sketched Ridge Regression: Optimization and Statistical Perspectives

**Shusen Wang**

UC Berkeley

Alex Gittens

RPI

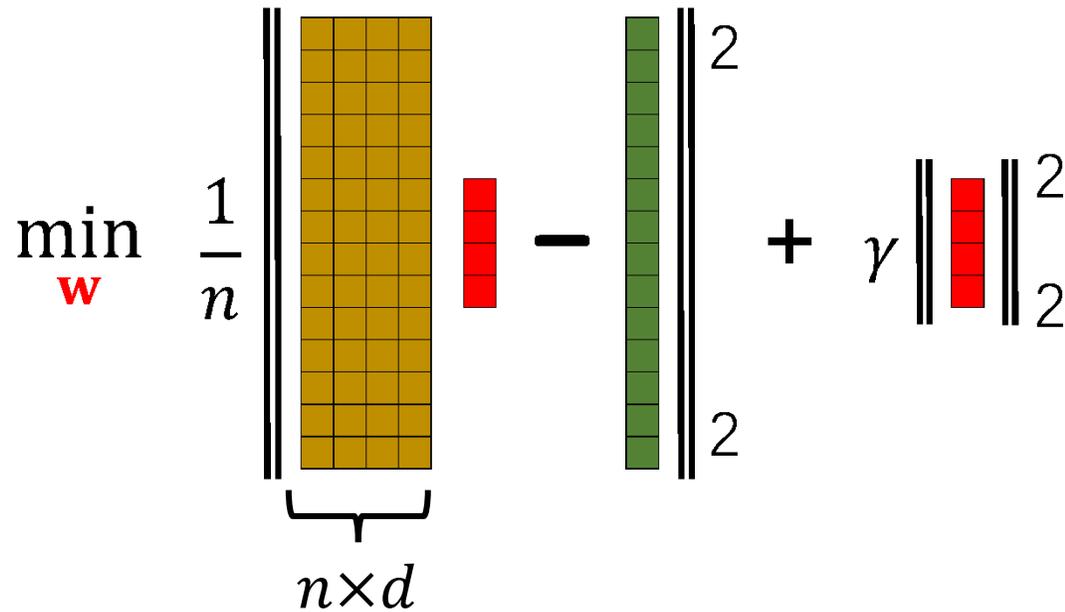
Michael Mahoney

UC Berkeley

# Overview

# Ridge Regression

$$\min_{\mathbf{w}} \left\{ f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\}$$



Over-determined:  
 $n \gg d$







# Approximate Ridge Regression

$$\min_{\mathbf{w}} \left\{ f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\}$$

$$\min_{\mathbf{w}} \frac{1}{n} \left\| \begin{array}{c} \text{4x4 grid} \\ \text{4x1 red bar} \end{array} - \begin{array}{c} \text{4x1 green bar} \end{array} \right\|_2^2 + \gamma \left\| \begin{array}{c} \text{4x1 red bar} \end{array} \right\|_2^2$$

sketch size

- Sketched solution:  $\mathbf{w}^S$
- Sketch size  $\tilde{O}\left(\frac{d}{\epsilon}\right)$

# Approximate Ridge Regression

$$\min_{\mathbf{w}} \left\{ f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\}$$

$$\min_{\mathbf{w}} \frac{1}{n} \left\| \begin{array}{c} \text{sketch size} \\ \text{sketch} \end{array} \right\|_2^2 + \gamma \|\mathbf{w}\|_2^2$$

- Sketched solution:  $\mathbf{w}^S$
- Sketch size  $\tilde{O}\left(\frac{d}{\epsilon}\right)$
- $f(\mathbf{w}^S) \leq (1 + \epsilon) \min_{\mathbf{w}} f(\mathbf{w})$

Optimization Perspective

# Approximate Ridge Regression

$$\min_{\mathbf{w}} \left\{ f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\}$$

$$\min_{\mathbf{w}} \frac{1}{n} \left\| \begin{array}{c} \text{4x4 grid} \\ \text{red column} \end{array} \right\|_2^2 - \left\| \begin{array}{c} \text{green column} \end{array} \right\|_2^2 + \gamma \left\| \begin{array}{c} \text{red column} \end{array} \right\|_2^2$$

## Statistical Perspective

- Bias
- Variance

# Related Work

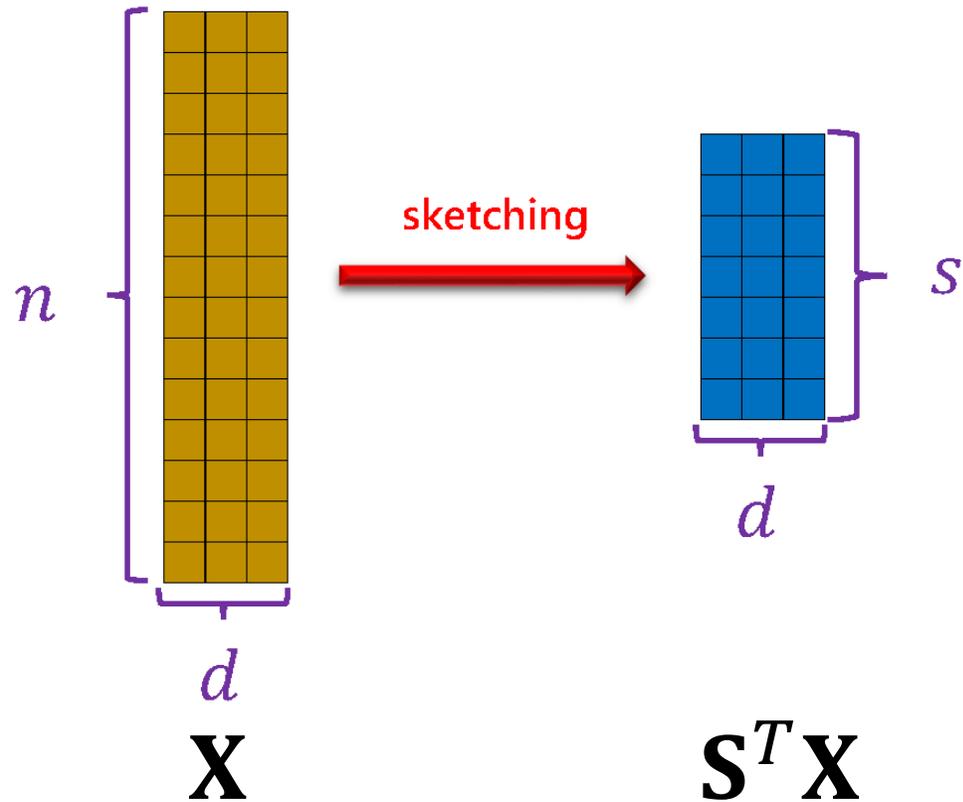
- Least squares regression:  $\min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$

## Reference

- Drineas, Mahoney, and Muthukrishnan: Sampling algorithms for  $l_2$  regression and applications. In *SODA*, 2006.
- Drineas, Mahoney, Muthukrishnan, and Sarlos: Faster least squares approximation. *Numerische Mathematik*, 2011.
- Clarkson and Woodruff: Low rank approximation and regression in input sparsity time. In *STOC*, 2013.
- Ma, Mahoney, and Yu: A statistical perspective on algorithmic leveraging. *Journal of Machine Learning Research*, 2015.
- Pilanci and Wainwright: Iterative Hessian sketch: fast and accurate solution approximation for constrained least squares. *Journal of Machine Learning Research*, 2015.
- Raskutti and Mahoney: A statistical perspective on randomized sketching for ordinary least-squares. *Journal of Machine Learning Research*, 2016.
- Etc ...

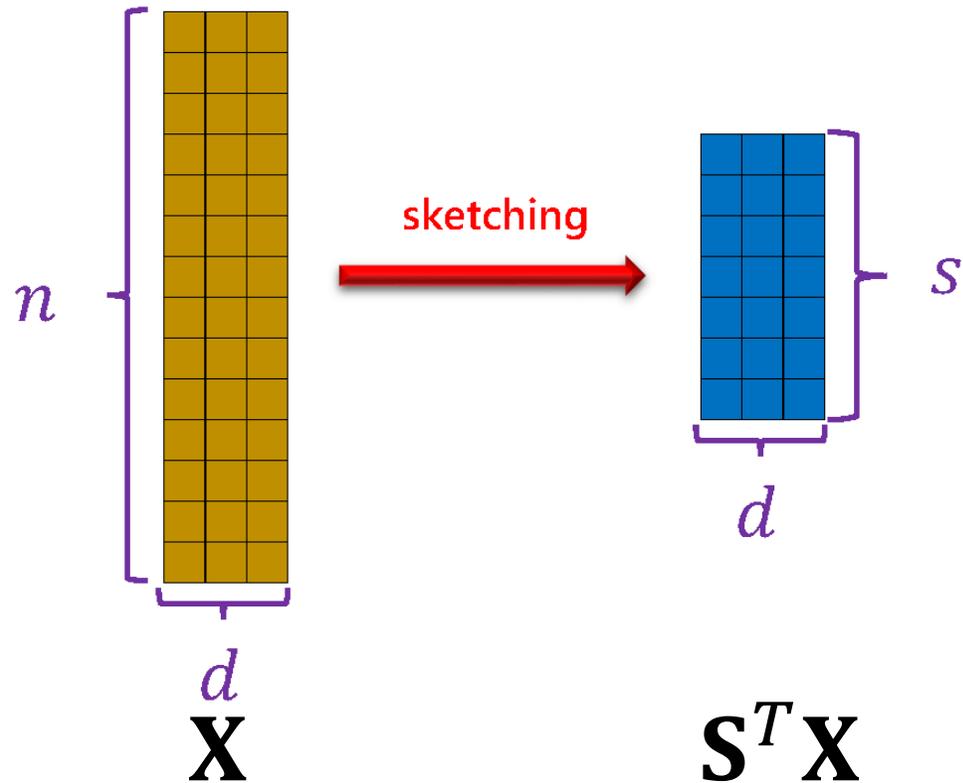
# Sketched Ridge Regression

# Matrix Sketching



- Turn big matrix into smaller one.
- $\mathbf{X} \in \mathbb{R}^{n \times d} \Rightarrow \mathbf{S}^T \mathbf{X} \in \mathbb{R}^{s \times d}$ .
- $\mathbf{S} \in \mathbb{R}^{n \times s}$  is called *sketching matrix*, e.g.,
  - Uniform sampling
  - Leverage score sampling
  - Gaussian projection
  - Subsampled randomized Hadamard transform (SRHT)
  - Count sketch (sparse embedding)
  - Etc.

# Matrix Sketching



- Some matrix sketching methods are efficient.
  - Time cost is  $o(nds)$  — lower than multiplication.
- Examples:
  - Leverage score sampling:  $O(nd \log n)$  time
  - SRHT:  $O(nd \log s)$  time

# Ridge Regression

- Objective function:

$$f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2$$

- Optimal solution:

$$\begin{aligned} \mathbf{w}^* &= \underset{\mathbf{w}}{\operatorname{argmin}} f(\mathbf{w}) \\ &= (\mathbf{X}^T \mathbf{X} + n\gamma \mathbf{I}_d)^\dagger (\mathbf{X}^T \mathbf{y}) \end{aligned}$$

- Time cost:  $O(nd^2)$  or  $O(ndt)$

# Sketched Ridge Regression

- Goal: *efficiently* and *approximately* solve

$$\operatorname{argmin}_{\mathbf{w}} \left\{ f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\}.$$

# Sketched Ridge Regression

- Goal: *efficiently* and *approximately* solve

$$\operatorname{argmin}_{\mathbf{w}} \left\{ f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\}.$$

- Approach: reduce the size of  $\mathbf{X}$  and  $\mathbf{y}$  by matrix sketching.

$$\min_{\mathbf{w}} \frac{1}{n} \left\| \begin{array}{c} \text{sketch of } \mathbf{X} \\ \text{sketch of } \mathbf{y} \end{array} \right\|_2^2 + \gamma \|\mathbf{w}\|_2^2$$

# Sketched Ridge Regression

- Sketched solution:

$$\begin{aligned}\mathbf{w}^s &= \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{1}{n} \|\mathbf{S}^T \mathbf{X} \mathbf{w} - \mathbf{S}^T \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\} \\ &= (\mathbf{X}^T \mathbf{S} \mathbf{S}^T \mathbf{X} + n\gamma \mathbf{I}_d)^\dagger (\mathbf{X}^T \mathbf{S} \mathbf{S}^T \mathbf{y})\end{aligned}$$

$$\min_{\mathbf{w}} \frac{1}{n} \left\| \begin{array}{c} \text{4x4 grid} \\ \text{red bar} \end{array} - \begin{array}{c} \text{green bar} \end{array} \right\|_2^2 + \gamma \left\| \begin{array}{c} \text{red bar} \end{array} \right\|_2^2$$

# Sketched Ridge Regression

- Sketched solution:

$$\begin{aligned}\mathbf{w}^s &= \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \frac{1}{n} \|\mathbf{S}^T \mathbf{X} \mathbf{w} - \mathbf{S}^T \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \right\} \\ &= (\mathbf{X}^T \mathbf{S} \mathbf{S}^T \mathbf{X} + n\gamma \mathbf{I}_d)^\dagger (\mathbf{X}^T \mathbf{S} \mathbf{S}^T \mathbf{y})\end{aligned}$$

- Time:  $O(sd^2) + T_s$ 
  - $T_s$  is the cost of sketching  $\mathbf{S}^T \mathbf{X}$
  - E.g.  $T_s = O(nd \log s)$  for SRHT.
  - E.g.  $T_s = O(nd \log n)$  for leverage score sampling.

# **Theory: Optimization Perspective**

# Optimization Perspective

- Recall the objective function  $f(\mathbf{w}) = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \gamma \|\mathbf{w}\|_2^2$ .
- Bound  $f(\mathbf{w}^S) - f(\mathbf{w}^*)$ .
- $\frac{1}{n} \|\mathbf{X}\mathbf{w}^S - \mathbf{X}\mathbf{w}^*\|_2^2 \leq f(\mathbf{w}^S) - f(\mathbf{w}^*)$ .

# Optimization Perspective

For the sketching methods

- SRHT or leverage sampling with  $s = \tilde{O}\left(\frac{\beta d}{\epsilon}\right)$ ,
- uniform sampling with  $s = O\left(\frac{\mu \beta d \log d}{\epsilon}\right)$ ,

$f(\mathbf{w}^s) - f(\mathbf{w}^*) \leq \epsilon f(\mathbf{w}^*)$  holds w.p. 0.9.

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : the design matrix
- $\gamma$ : the regularization parameter
- $\beta = \frac{\|\mathbf{X}\|_2^2}{n\gamma + \|\mathbf{X}\|_2^2} \in (0, 1]$
- $\mu \in \left[1, \frac{n}{d}\right]$ : the row coherence of  $\mathbf{X}$

# Optimization Perspective

For the sketching methods

- SRHT or leverage sampling with  $s = \tilde{O}\left(\frac{\beta d}{\epsilon}\right)$ ,
- uniform sampling with  $s = O\left(\frac{\mu \beta d \log d}{\epsilon}\right)$ ,

$f(\mathbf{w}^s) - f(\mathbf{w}^*) \leq \epsilon f(\mathbf{w}^*)$  holds w.p. 0.9.

$$\implies \frac{1}{n} \|\mathbf{X}\mathbf{w}^s - \mathbf{X}\mathbf{w}^*\|_2^2 \leq \epsilon f(\mathbf{w}^*).$$

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : the design matrix
- $\gamma$ : the regularization parameter
- $\beta = \frac{\|\mathbf{X}\|_2^2}{n\gamma + \|\mathbf{X}\|_2^2} \in (0, 1]$
- $\mu \in \left[1, \frac{n}{d}\right]$ : the row coherence of  $\mathbf{X}$

# **Theory: Statistical Perspective**

# Statistical Model

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : fixed design matrix
- $\mathbf{w}_0 \in \mathbb{R}^d$ : the *true* and *unknown* model
- $\mathbf{y} = \mathbf{X}\mathbf{w}_0 + \boldsymbol{\delta}$ : observed response vector
  - $\delta_1, \dots, \delta_n$  are random noise
  - $\mathbb{E}[\boldsymbol{\delta}] = \mathbf{0}$  and  $\mathbb{E}[\boldsymbol{\delta}\boldsymbol{\delta}^T] = \xi^2 \mathbf{I}_n$

# Bias-Variance Decomposition

- Risk:  $R(\mathbf{w}) = \frac{1}{n} \mathbb{E} \|\mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w}_0\|_2^2$ 
  - $\mathbb{E}$  is taken w.r.t. the random noise  $\delta$ .

# Bias-Variance Decomposition

- Risk:  $R(\mathbf{w}) = \frac{1}{n} \mathbb{E} \|\mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w}_0\|_2^2$ 
  - $\mathbb{E}$  is taken w.r.t. the random noise  $\delta$ .
  - Risk measures prediction error.

# Bias-Variance Decomposition

- Risk:  $R(\mathbf{w}) = \frac{1}{n} \mathbb{E} \|\mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w}_0\|_2^2$
- $R(\mathbf{w}) = \text{bias}^2(\mathbf{w}) + \text{var}(\mathbf{w})$

# Bias-Variance Decomposition

- Risk:  $R(\mathbf{w}) = \frac{1}{n} \mathbb{E} \|\mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w}_0\|_2^2$
- $R(\mathbf{w}) = \text{bias}^2(\mathbf{w}) + \text{var}(\mathbf{w})$

Optimal  
Solution

- $\text{bias}(\mathbf{w}^*) = \gamma\sqrt{n} \|(\boldsymbol{\Sigma}^2 + n\gamma\mathbf{I}_d)^{-1}\boldsymbol{\Sigma}\mathbf{V}^T\mathbf{w}_0\|_2,$
- $\text{var}(\mathbf{w}^*) = \frac{\xi^2}{n} \|(\mathbf{I}_d + n\gamma\boldsymbol{\Sigma}^{-2})^{-1}\|_2^2,$

Sketched  
Solution

- $\text{bias}(\mathbf{w}^s) = \gamma\sqrt{n} \|(\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{S}\mathbf{S}^T\mathbf{U}\boldsymbol{\Sigma} + n\gamma\mathbf{I}_d)^\dagger\boldsymbol{\Sigma}\mathbf{V}^T\mathbf{w}_0\|_2,$
- $\text{var}(\mathbf{w}^s) = \frac{\xi^2}{n} \|(\mathbf{U}^T\mathbf{S}\mathbf{S}^T\mathbf{U} + n\gamma\boldsymbol{\Sigma}^{-2})^\dagger\mathbf{U}^T\mathbf{S}\mathbf{S}^T\|_2^2,$

- Here  $\mathbf{X} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$  is the SVD.

# Statistical Perspective

For the sketching methods

- SRHT or leverage sampling with  $s = \tilde{O}\left(\frac{d}{\epsilon^2}\right)$ ,
- uniform sampling with  $s = O\left(\frac{\mu d \log d}{\epsilon^2}\right)$ ,

the followings hold w.p. 0.9:

$$1 - \epsilon \leq \frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon,$$

$$(1 - \epsilon) \frac{n}{s} \leq \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \leq (1 + \epsilon) \frac{n}{s}.$$

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : the design matrix
- $\mu \in \left[1, \frac{n}{d}\right]$ : the row coherence of  $\mathbf{X}$

Good!

Bad! Because  $n \gg s$ .

# Statistical Perspective

For the sketching methods

- SRHT or leverage sampling with  $s = \tilde{O}\left(\frac{d}{\epsilon^2}\right)$ ,
- uniform sampling with  $s = O\left(\frac{\mu d \log d}{\epsilon^2}\right)$ ,

the followings hold w.p. 0.9:

$$1 - \epsilon \leq \frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon,$$

$$(1 - \epsilon) \frac{n}{s} \leq \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \leq (1 + \epsilon) \frac{n}{s}.$$

- $\mathbf{X} \in \mathbb{R}^{n \times d}$ : the design matrix
- $\mu \in \left[1, \frac{n}{d}\right]$ : the row coherence of  $\mathbf{X}$

If  $\mathbf{y}$  is noisy

⇒ variance dominates bias

⇒  $R(\mathbf{w}^s) \gg R(\mathbf{w}^*)$ .

# Conclusions

- Use sketched solution to initialize numerical optimization.
  - $Xw^S$  is close to  $Xw^*$ .

Optimization Perspective

# Conclusions

- Use sketched solution to initialize numerical optimization.
  - $\mathbf{X}\mathbf{w}^S$  is close to  $\mathbf{X}\mathbf{w}^*$ .

Optimization Perspective

- $\mathbf{w}^{(t)}$ : output of the  $t$ -th iteration of CG algorithm.
- $\frac{\|\mathbf{X}\mathbf{w}^{(t)} - \mathbf{X}\mathbf{w}^*\|_2^2}{\|\mathbf{X}\mathbf{w}^{(0)} - \mathbf{X}\mathbf{w}^*\|_2^2} \leq 2 \left( \frac{\sqrt{\kappa(\mathbf{X}^T\mathbf{X})} - 1}{\sqrt{\kappa(\mathbf{X}^T\mathbf{X})} + 1} \right)^t$ .
- Initialization is important.

# Conclusions

- Use sketched solution to initialize numerical optimization.
  - $\mathbf{Xw}^S$  is close to  $\mathbf{Xw}^*$ .
- Never use sketched solution to replace the optimal solution.
  - Much higher variance  $\rightarrow$  bad generalization.

Optimization Perspective

Statistical Perspective

# Model Averaging

# Model Averaging

- Independently draw  $\mathbf{S}_1, \dots, \mathbf{S}_g$ .
- Compute the sketched solutions  $\mathbf{w}_1^{\mathbf{S}}, \dots, \mathbf{w}_g^{\mathbf{S}}$ .
- Model averaging:  $\mathbf{w}^{\mathbf{S}} = \frac{1}{g} \sum_{i=1}^g \mathbf{w}_i^{\mathbf{S}}$ .

# Optimization Perspective

- For sufficiently large  $s$ ,

$$\frac{f(\mathbf{w}_1^s) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \leq \epsilon \quad \text{holds w.h.p.}$$

**Without model averaging**

# Optimization Perspective

- For sufficiently large  $s$ ,

$$\frac{f(\mathbf{w}_1^s) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \leq \epsilon \quad \text{holds w.h.p.}$$

Without model averaging

- Using the **same** matrix sketching and **same**  $s$ ,

$$\frac{f(\mathbf{w}^s) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \leq \frac{\epsilon}{g} + \epsilon^2 \quad \text{holds w.h.p.}$$

With model averaging

# Optimization Perspective

- For sufficiently large  $s$ ,

$$\frac{f(\mathbf{w}_1^s) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \leq \epsilon \text{ holds w.h.p.}$$

Without model averaging

- Using the **same** matrix sketching and **same**  $s$ ,

$$\frac{f(\mathbf{w}^s) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \leq \frac{\epsilon}{g} + \epsilon^2 \text{ holds w.h.p.}$$

With model averaging

# Optimization Perspective

- For sufficiently large  $s$ ,

$$\frac{f(\mathbf{w}_1^s) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \leq \epsilon \text{ holds w.h.p.}$$

Without model averaging

- Using the **same** matrix sketching and **same**  $s$ ,

$$\frac{f(\mathbf{w}^s) - f(\mathbf{w}^*)}{f(\mathbf{w}^*)} \leq \frac{\epsilon}{g} + \epsilon^2 \text{ holds w.h.p.}$$

With model averaging

If  $s \gg d \implies \epsilon^2$  is very small  $\implies$  error bound  $\propto \frac{\epsilon}{g}$ .

# Statistical Perspective

- Risk:  $R(\mathbf{w}) = \frac{1}{n} \mathbb{E} \left\| \mathbf{X}\mathbf{w} - \mathbf{X}\mathbf{w}_0 \right\|_2^2 = \text{bias}^2(\mathbf{w}) + \text{var}(\mathbf{w})$
- Model averaging :
  - $\text{bias}(\mathbf{w}^s) = \gamma \sqrt{n} \left\| \frac{1}{g} \sum_{i=1}^g (\mathbf{\Sigma} \mathbf{U}^T \mathbf{S}_i \mathbf{S}_i^T \mathbf{U} \mathbf{\Sigma} + n\gamma \mathbf{I}_d)^\dagger \mathbf{\Sigma} \mathbf{V}^T \mathbf{w}_0 \right\|_2$ ,
  - $\text{var}(\mathbf{w}^s) = \frac{\xi^2}{n} \left\| \frac{1}{g} \sum_{i=1}^g (\mathbf{U}^T \mathbf{S}_i \mathbf{S}_i^T \mathbf{U} + n\gamma \mathbf{\Sigma}^{-2})^\dagger \mathbf{U}^T \mathbf{S}_i \mathbf{S}_i^T \right\|_2^2$ .
  - Here  $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  is the SVD.

# Statistical Perspective

- For sufficiently large  $s$ , the followings hold w.h.p.:

$$\frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon \quad \text{and} \quad \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \leq \frac{n}{s} (1 + \epsilon).$$

**Without model averaging**

# Statistical Perspective

- For sufficiently large  $s$ , the followings hold w.h.p.:

$$\frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon \quad \text{and} \quad \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \leq \frac{n}{s} (1 + \epsilon).$$

Without model averaging

- Using the **same** sketching methods and **same**  $s$ , the followings hold w.h.p.:

$$\frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon \quad \text{and} \quad \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \approx \frac{n}{s} \left( \frac{1}{\sqrt{g}} + \epsilon \right)^2$$

With model averaging

# Statistical Perspective

- For sufficiently large  $s$ , the followings hold w.h.p.:

$$\frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon \quad \text{and} \quad \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \leq \frac{n}{s} (1 + \epsilon).$$

Without model averaging

- Using the **same** sketching methods and **same**  $s$ , the followings hold w.h.p.:

$$\frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon \quad \text{and} \quad \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \approx \frac{n}{s} \left( \frac{1}{\sqrt{g}} + \epsilon \right)^2$$

With model averaging

# Statistical Perspective

- For sufficiently large  $s$ , the followings hold w.h.p.:

$$\frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon \quad \text{and} \quad \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \leq \frac{n}{s} (1 + \epsilon).$$

Without model averaging

- Using the **same** sketching methods and **same**  $s$ , the followings hold w.h.p.:

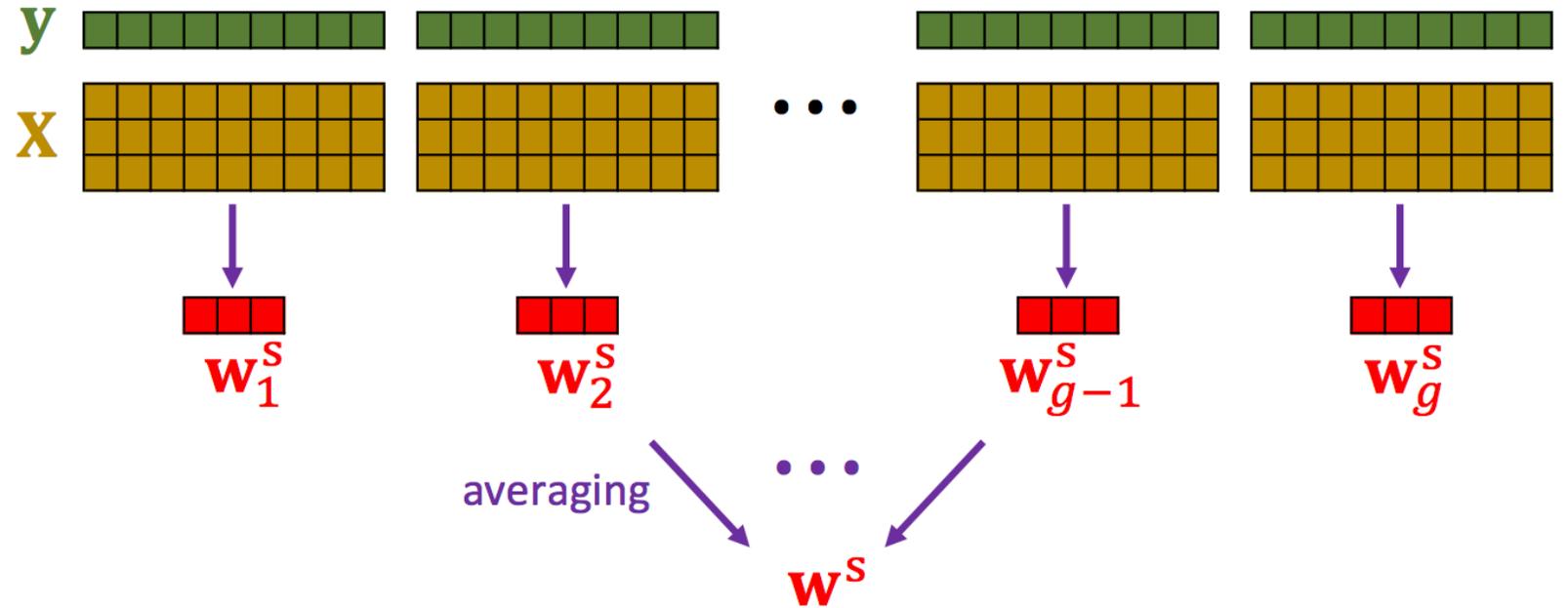
$$\frac{\text{bias}(\mathbf{w}^s)}{\text{bias}(\mathbf{w}^*)} \leq 1 + \epsilon \quad \text{and} \quad \frac{\text{var}(\mathbf{w}^s)}{\text{var}(\mathbf{w}^*)} \approx \frac{n}{s} \left( \frac{1}{\sqrt{g}} + \epsilon \right)^2$$

With model averaging

If  $\epsilon$  is small, then  $\text{var}(\mathbf{w}^s) \propto \frac{1}{g}$ .

# Applications to Distributed Optimization

- $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$  are (randomly) split among  $g$  machines.
- Equivalent to uniform sampling with  $s = \frac{n}{g}$ .



# Optimization Perspective

- **Application to distributed optimization:**
  - If  $s = \frac{n}{g} \gg d$ ,  $\mathbf{w}^s$  is very close to  $\mathbf{w}^*$  (provably).
  - $\mathbf{w}^s$  is good initialization of distributed optimization algorithms.

# Statistical Perspective

- **Application to distributed machine learning:**
  - If  $s = \frac{n}{g} \gg d$ , then  $R(\mathbf{w}^S)$  is comparable to  $R(\mathbf{w}^*)$ .
  - If low-precision solution suffices, then  $\mathbf{w}^S$  is a good substitute of  $\mathbf{w}^*$ .
  - One-shot solution.

# Thank You!

The paper is at [arXiv:1702.04837](https://arxiv.org/abs/1702.04837)