

NEW SUBSAMPLING ALGORITHMS FOR FAST LEAST SQUARES REGRESSION

PARAMVEER S. DHILLON¹, YICHAO LU², DEAN FOSTER² AND LYLE UNGAR¹
 {¹CIS,²STATISTICS (WHARTON)} UNIVERSITY OF PENNSYLVANIA, PHILADELPHIA, PA, U.S.A

BACKGROUND

- Problem: Estimation of ordinary least squares (OLS) regression when $n \gg p$ (n observations, p features).
- OLS Regression: $Y = \mathbf{X}w_0 + \epsilon$, $\epsilon \sim \mathcal{N}(0, 1)$.
- MLE solution $\rightarrow \hat{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top Y$.
 - + Running Time (FLOPS): $O(np^2)$.
 - + Error bound: $\|w_0 - \hat{w}\| \rightarrow O(\sqrt{\frac{p}{n}})$
- Current state of the art: Preconditioning based approaches [(Drineas)⁺ 07, (Rokhlin)⁺ 08].
 1. Transform the data with randomized Hadamard (SRHT) or Fast Fourier Transform (FFT).
 2. Uniformly subsample the resulting matrix ($n_{subs} = O(p)$).
 3. Estimate the OLS on this smaller matrix.
 - + Running Time (FLOPS): $O(\max(np \log p, n_{subs}p^2))$.
 - + Error bound: $\|w_0 - \hat{w}\| \rightarrow O(\sqrt{\frac{p}{n_{subs}}})$.

THE ALGORITHMS

- Either precondition the data matrix \mathbf{X} and then subsample (*fixed design*) or directly subsample (*sub-gaussian random design*) (If you believe the data is i.i.d.).
- 1. **Full Subsampling Algorithm (FS)** Subsample \mathbf{X} and Y , then $\hat{w}_{FS} = (\mathbf{X}_{subs}^\top \mathbf{X}_{subs})^{-1} \mathbf{X}_{subs}^\top Y_{subs}$.
 - + Similar to [(Drineas)⁺ 07], but novel error analysis.
- 2. **Covariance Subsampling Algorithm (CovS)** $\rightarrow \hat{w}_{CovS} = (\mathbf{X}_{subs}^\top \mathbf{X}_{subs})^{-1} \mathbf{X}^\top Y$.
- 3. **Uluru** \rightarrow Two stage algorithm
 - (a) Stage 1: Use **FS** to estimate \hat{w}_{FS} .
 - (b) Stage 2: Use **CovS** to estimate $\hat{w}_{correct}$ on the remaining observations ($n_{rem} = n \setminus n_{subs}$).
 - (c) Perform Sampling Correction: $\hat{w}_{Uluru} = \hat{w}_{FS} + \hat{w}_{correct}$.

Uluru



Methods	Running Time O(FLOPS)	Error bound
OLS	$O(np^2)$	$O(\sqrt{p/n})$
FS	$O(nr p^2)$	$O(\sqrt{p/nr})$
CovS	$O(nr p^2 + np)$	*
Uluru	$O(nr p^2 + np)$	$O(\sqrt{p/n})$

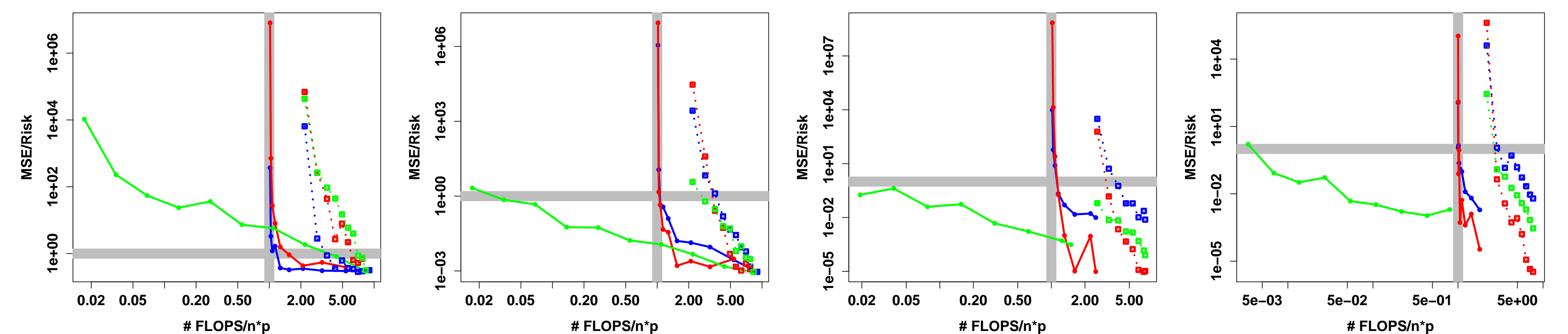
- We do not increase the error of **Uluru** by using less data in estimating the covariance matrix. So our estimate of the *quadratic* term is as solid as the rock formation **Uluru**!

THEORY (SUMMARY)

- When $n_{subs} \ll n_{rem}$, keeping only the dominating terms, the results can be summarized as: With failure probability less than some fixed number, the algorithms have the following error bounds.
 1. **FS** $\rightarrow O(\sigma \sqrt{\frac{p}{n_{subs}}})$.
 2. **CovS** $\rightarrow O(\sqrt{\frac{p}{n_{subs}}} \|w\| + \sigma \sqrt{\frac{p}{n}})$.
 3. **Uluru** $\rightarrow O(\sigma \frac{p}{n_{subs}} + \sigma \sqrt{\frac{p}{n}})$.
 - If the second term for the error of the **Uluru** algorithm dominates, i.e. if $r (= \frac{n_{subs}}{n}) > O(\sqrt{p/n})$ then the error bound of **Uluru** $\approx O(\sigma \sqrt{\frac{p}{n}})$ (completely independent of r !).
 - The threshold for r only depends on the properties of design matrix (n, p) and not on the noise level σ .
 - **FS** and **CovS** do not have this property.

EXPERIMENTS

- Results for synthetic datasets (Plots 1-2, low signal and high signal) and for real world datasets (Plots 3-4, CPUSMALL, CADATA). Color scheme: + (Green)-FS, + (Blue)-CovS, + (Red)-Uluru. The solid lines indicate no preconditioning (i.e. random design) and dashed lines indicate fixed design with Randomized Hadamard preconditioning. The FLOPS reported are the theoretical values.



CONCLUSION

Uluru has a runtime of $O(np)$ and obtains error bound of $O(\sqrt{\frac{p}{n}})$ which is the same as full OLS and is independent of amount of subsampling.