

Least Squares Estimation Using Sketched Data with Heteroskedastic Errors

ICML 2022

Sokbae (Simon) Lee, Columbia University

Serena Ng, Columbia University and NBER

Sketching

- Want a **sketch** $\tilde{A} = \Pi A \in \mathbb{R}^{m \times d}$ that preserves features of $A \in \mathbb{R}^{n \times d}$
- Aim to have a much smaller m than n
- Random Sampling:
 - rows in \tilde{A} are rows of A .
 - e.g., Bernoulli sampling; uniform sampling w or w/o replacement; leverage score sampling
- Random Projections:
 - rows in \tilde{A} are linear combinations of rows of A .
 - e.g., Gaussian; SRHT; Countsketch

Literature: Algorithmic Perspective

- The early works of Sarlos (2006), Drineas, Mahoney, and Muthukrishnan (2006) and Drineas, Mahoney, Muthukrishnan, and Sarlos (2011) consider sketching of the least squares estimator from an [algorithmic perspective](#) (worst case analysis with the fixed data).
- See, e.g., Woodruff (2014), Drineas and Mahoney (2018) and Martinsson and Tropp (2020) for a review.

Literature: Statistical Perspective

- However, recent works due to Ma, Mahoney, and Yu (2015), Raskutti and Mahoney (2016), and Dobriban and Liu (2019) show that an optimal worst-case error may not yield an optimal mean-squared error.
- This led to interest in better understanding the statistical implications of sketching. For example, Geppert, Ickstadt, Munteanu, Quedenfeld, and Sohler (2017) considers Bayesian estimation while Ahfock, Astle, and Richardson (2020) and Ma, Zhang, Xing, Ma, and Mahoney (2020) provide **asymptotic distribution theory for the sketched least squares estimators under homoskedasticity**.

Regression Model

- Given i.i.d. observations $\{(y_i, X_i) : i = 1, \dots, n\}$, we consider a linear regression model:

$$y = X\beta_0 + e, \quad \mathbb{E}(Xe) = 0$$

- $\sqrt{n}(\hat{\beta}_{OLS} - \beta_0) \rightarrow_d N(0, V_1)$ as $n \rightarrow \infty$, where V_1 is the sandwich variance defined as

$$V_1 := [\mathbb{E}(X_i X_i^T)]^{-1} \mathbb{E}(e_i^2 X_i X_i^T) [\mathbb{E}(X_i X_i^T)]^{-1}.$$

- Under homoskedasticity, V_1 becomes

$$V_0 := \mathbb{E}(e_i^2) [\mathbb{E}(X_i X_i^T)]^{-1}.$$

Sketched OLS

- A sketch of the data (y, X) is (\tilde{y}, \tilde{X}) , where $\tilde{y} = \Pi y$, $\tilde{X} = \Pi X$, and Π is usually an $m \times n$ random matrix.
- The sketched least squares estimator is
$$\tilde{\beta}_{OLS} := (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{y}.$$

Assumption

- (i) *The data $\mathcal{D}_n := \{(y_i, X_i) \in \mathbb{R}^{1+p} : i = 1, \dots, n\}$ are independent and identically distributed (i.i.d.).
Furthermore, X has singular value decomposition
$$X = U_X \Sigma_X V_X^T.$$*
- (ii) *$\mathbb{E}(y_i^4) < \infty$, $\mathbb{E}(\|X_i\|^4) < \infty$, and $\mathbb{E}(X_i X_i^T)$ has full rank p .*
- (iii) *The random matrix Π is independent of \mathcal{D}_n .*
- (iv) *$m = m_n \rightarrow \infty$ but $m/n \rightarrow 0$ as $n \rightarrow \infty$, while p is fixed.*

Two Leading Examples

For simplicity, we focus on Bernoulli sampling (BS) from Random Sampling and Countsketch (CS) from Random Projections.

- Bernoulli sampling (BS): $\Pi = \sqrt{\frac{n}{m}}B$, where B is a diagonal sampling matrix of i.i.d. Bernoulli random variables with success probability m/n .
- Countsketch (CS) : only one non-zero entry in each column of Π . The non-zero entry takes on value $\{+1, -1\}$ randomly drawn with equal probability, and is located uniformly at random for each column.

Theorem (OLS)

Let Assumption 1 hold and $\mathbb{E}(e_i|X_i) = 0$.

- (i) Under BS, $m^{1/2}(\tilde{\beta}_{OLS} - \hat{\beta}_{OLS}) \rightarrow_d N(0, V_1)$.
- (ii) Under CS, $m^{1/2}(\tilde{\beta}_{OLS} - \hat{\beta}_{OLS}) \rightarrow_d N(0, V_0)$.

Theorem 1 indicates that both sampling schemes yield asymptotically normal estimates, but for different reasons have different asymptotic variances.

Practical Inference

- In applications, researchers would like to test a hypothesis about β_0 using a sketched estimate, and our results provide all the quantities required for inference.
- Since $m/n \rightarrow 0$,

$$\begin{aligned}m^{1/2}(\tilde{\beta} - \hat{\beta}) &= m^{1/2}(\tilde{\beta} - \beta_0) - (m/n)^{1/2} n^{1/2}(\hat{\beta} - \beta_0) \\ &= m^{1/2}(\tilde{\beta} - \beta_0) + o_p(1).\end{aligned}$$

- Then, asymptotic normality of $m^{1/2}(\tilde{\beta} - \hat{\beta})$ provides a guide to conduct inference for β_0 :

$$\tilde{V}_m^{-1/2}(\tilde{\beta}_{OLS} - \beta_0) \approx N(0, I_p),$$

where the form of \tilde{V}_m will be given below.

Monte Carlo Experiments

	(1)	(2)	(3)	(4)
	SIZE		POWER	
	S.E.0	S.E.1	S.E.0	S.E.1
(I) HOMOSKEDASTIC DESIGN				
BERNOULLI	0.046	0.050	0.490	0.496
UNIFORM	0.047	0.052	0.489	0.490
LEVERAGE	0.045	0.053	0.483	0.513
COUNTSKETCH	0.049	0.051	0.479	0.489
SRHT	0.056	0.061	0.492	0.498
SRFT	0.055	0.057	0.484	0.489
(II) HETEROSKEDASTIC DESIGN				
BERNOULLI	0.310	0.047	0.713	0.436
UNIFORM	0.301	0.053	0.719	0.435
LEVERAGE	0.183	0.051	0.727	0.529
COUNTSKETCH	0.054	0.057	0.813	0.812
SRHT	0.054	0.056	0.804	0.809
SRFT	0.050	0.052	0.799	0.806

- AHFOCK, D. C., W. J. ASTLE, AND S. RICHARDSON (2020):
“Statistical properties of sketching algorithms,” *Biometrika*, 108(2),
283–297.
- DOBRIBAN, E., AND S. LIU (2019): “Asymptotics for sketching in
least squares,” in *Proceedings of the 33rd International Conference on
Neural Information Processing Systems*, pp. 3675–3685.
- DRINEAS, P., M. MAHONEY, AND S. MUTHUKRISHNAN (2006):
“Sampling Algorithms for L2 Regression and Applications,”
*Proceedings of the 17th Annual ACM-SIAM Symposium on Discrete
Algorithms*, pp. 1127–1136.
- DRINEAS, P., M. MAHONEY, S. MUTHUKRISHNAN, AND T. SARLOS
(2011): “Faster Least Squares Approximation,” *Numerical
Mathematics*, 117, 219–249.

- DRINEAS, P., AND M. W. MAHONEY (2018): “Lectures on randomized numerical linear algebra,” in *The Mathematics of Data*, ed. by M. W. Mahoney, J. C. Duchi, and A. C. Gilbert, pp. 1–48. AMS/IAS/SIAM.
- GEPPERT, L., K. ICKSTADT, A. MUNTEANU, J. QUDEDENFELD, AND C. SOHLER (2017): “Random Projections for Bayesian Regressions,” *Statistical Computing*, 27:, 79–101.
- MA, P., M. W. MAHONEY, AND B. YU (2015): “A statistical perspective on algorithmic leveraging,” *Journal of Machine Learning Research*, 16(1), 861–911.

- MA, P., X. ZHANG, X. XING, J. MA, AND M. MAHONEY (2020): “Asymptotic Analysis of Sampling Estimators for Randomized Numerical Linear Algebra Algorithms,” in *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, ed. by S. Chiappa, and R. Calandra, vol. 108 of *Proceedings of Machine Learning Research*, pp. 1026–1035. PMLR.
- MARTINSSON, P.-G., AND J. A. TROPP (2020): “Randomized numerical linear algebra: Foundations and algorithms,” *Acta Numerica*, 29, 403–572.
- RASKUTTI, G., AND M. W. MAHONEY (2016): “A statistical perspective on randomized sketching for ordinary least-squares,” *Journal of Machine Learning Research*, 17(1), 7508–7538.

- SARLOS, T. (2006): “Improved Approximation Algorithms for Large Matrices via Random Projections,” *Proceedings of the 47 IEEE Symposium on Foundations of Computer Science*.
- WOODRUFF, D. P. (2014): “Sketching as a tool for numerical linear algebra,” *Foundations and Trends in Theoretical Computer Science*, 10(1–2), 1–157.