

A Appendix

We now provide additional details for our results in Section 2 of the paper.

Lemma 1. *Let A, B be full column rank matrices of size $n \times d$, with $\log(n) = d^{o(1)}$. Let S be an SRHT with¹ $m = \tilde{O}((d + \log(1/\delta))/\epsilon^2)$ rows. For any matrix B of size $n \times d$ we have*

$$\|X\|_2 = \|(SA)^\dagger SB\|_2 \lesssim \|A^\dagger B\|_2 + \epsilon \|\Sigma^{-1}\|_2 \sqrt{(1 + d/k)(\|B\|_2^2 + \|B\|_F^2/k)},$$

with probability $1 - 1/\text{poly}(d)$.

Proof. From Equation 5 in the body and the triangle inequality, we have

$$(SA)^\dagger SB = V\Sigma^{-1} \left(\sum_{k=0}^{\infty} T^k \right) U^T S^T SB \quad (1)$$

$$\|(SA)^\dagger SB\|_2 \leq \|A^\dagger B\|_2 + \|(SA)^\dagger SB - A^\dagger B\|_2 \quad (2)$$

$$\leq \|A^\dagger B\|_2 + \|V\Sigma^{-1} \left(\sum_{k=0}^{\infty} T^k \right) U^T S^T SB - V\Sigma^{-1} U^T B\|_2 \quad (3)$$

$$\leq \|A^\dagger B\|_2 + \|\Sigma^{-1}\|_2 \sum_{k=0}^{\infty} \epsilon^k \|U^T S^T SB - U^T B\|_2 \quad (4)$$

$$\leq \|A^\dagger B\|_2 + \|\Sigma^{-1}\|_2 \frac{\epsilon}{1 - \epsilon} \sqrt{(1 + d/k)(\|B\|_2^2 + \|B\|_F^2/k)} \quad (5)$$

where we used equation 8 in the body in the last step. \square

A.1 AD

Let us manually derive the AD for the least squares regression problem.

$$\begin{aligned} \text{LLS}(A, b) &= \text{LS}(A^T A, A^T b) \\ &= \text{LS}(M, m) \\ (M, m) &\equiv (A^T A, A^T b) \\ B &= A^T \\ (B_1, B_2) &= (B, B) \\ C &= B_1 A \\ d &= B_2 b \\ \bar{A}_M &= A\bar{M} + A\bar{M}^T \\ \bar{A}_m &= b\bar{m}^T \\ \bar{A} &= \bar{A}_M + \bar{A}_m \\ &= A\bar{M} + A\bar{M}^T + b\bar{m}^T \\ \bar{b} &= A\bar{m} \\ (\bar{M}, \bar{m}) &= \mathcal{J}^T(\text{LS})(M, m)(\bar{y}) \\ &= (-\bar{m}y^T, \text{LS}(M^T, \bar{y})) \\ &= (-\text{LS}(A^T A, \bar{y})y^T, \text{LS}(A^T A, \bar{y})) \end{aligned}$$

This gives us the final reverse mode AD:

$$\bar{A} = -A(A^T A)^{-1} \bar{y} y^T - A y \bar{y}^T (A^T A)^{-1} + b y^T (A^T A)^{-1} \quad (6)$$

$$\bar{b} = A(A^T A)^{-1} \bar{y} \quad (7)$$

¹For a function f , we use the notation $\tilde{O}(f)$ to denote $f \cdot \text{polylog}(f)$.

A.2 Approximation bounds

Let us derive some additional bounds which were missing in the main paper:

$$\|Ay\bar{y}^T M^{-1} - Ay_D\bar{y}^T M_S^{-1}\|_F \leq \|Ay\bar{y}^T M^{-1} - Ay_D\bar{y}^T M^{-1}\|_F + \|Ay_D\bar{y}^T (M - M_S^{-1})\|_F \quad (8)$$

$$\leq \|U(I - (U^T S^T S U)^{-1})U^T b\|_F + \epsilon \|Ay_D\| \|\bar{y}\| \|\Sigma^{-1}\|_2 \|\Sigma^{-1}\|_F \quad (9)$$

$$\lesssim \epsilon (\|b\|_2 + \|b\|_2 \|\bar{y}\|_2 \|\Sigma^{-1}\|_2 \|\Sigma^{-1}\|_F) \quad (10)$$

Table 1: Cheat sheet to derive AD.

Original	Forward Transform	Reverse Transform
$z = a + b$	$\dot{z} = \dot{a} + \dot{b}$	$(\bar{a}, \bar{b}) = (\bar{z}, \bar{z})$
$z = ab$	$\dot{z} = \dot{a}b + a\dot{b}$	$(\bar{a}, \bar{b}) = (\bar{z}b, a\bar{z})$
$(z_1, z_2) = (a, a)$	$(\dot{z}_1, \dot{z}_2) = (\dot{a}, \dot{a})$	$\bar{a} = \bar{z}_1 + \bar{z}_2$
$Y = AXB$	$\dot{Y} = \dot{A}YB$	$\bar{X} = A^T \bar{Y} B^T$
$y = \text{LS}(M, m)$	$\dot{y} = \mathcal{J} \text{LS}(M, m)(y, \dot{M}, \dot{m})$ $= \text{LS}(M, \dot{m} - \dot{M}y)$	$(\bar{M}, \bar{m}) = \mathcal{J}^T \text{LS}(M, m)(y, \bar{y})$ $= (-\bar{m}y^T, \text{LS}(M^T, \bar{y}))$

Table 2: Forward mode AD Transformations.

Type	Primal	Forward Transform
Regular	$y = \text{LLS}(A, b)$	$\dot{y} = \mathcal{J} \text{LLS}(A, b)(y, \dot{A}, \dot{b})$ $= \text{LS}(A^T \dot{A}, \dot{A}^T b + A^T \dot{b} - (\dot{A}^T A + A^T \dot{A})y)$
“Diff + Sketch”	$y_D = \text{LLS}(A, b, S)$	$\dot{y}_D = \mathcal{J} \text{LLS}(A, b)(y_D, \dot{A}, \dot{b}, S)$ $= \text{LS}(A^T S^T S \dot{A}, \dot{A}^T b + A^T \dot{b} - (\dot{A}^T A + A^T \dot{A})y_D)$
“Sketch + Diff”	$y_S = \text{LLS}(A, b, S)$	$\dot{y}_S = \mathcal{J} \text{LLS}(A, b)(y_S, \dot{A}, \dot{b}, S)$ $= \text{LS}(A^T S^T S \dot{A}, \dot{A}^T S^T S b + A^T S^T \dot{b}$ $- (\dot{A}^T S^T S \dot{A} + A^T S^T S \dot{A})y_S)$

Table 3: Reverse mode AD Transformations.

Type	Primal	Reverse Transform
Regular	$y = \text{LLS}(A, b)$	$(\bar{A}, \bar{b}) = \mathcal{J}^T \text{LLS}(A, b)(y, \bar{y})$ $= (-A^T \bar{y} \bar{y}^T - Ay\bar{y}^T M^{-1} + b\bar{y}^T M^{-1}, A^T \bar{y})$
“Diff + Sketch”	$y_D = \text{LLS}(A, b, S)$	$(\bar{A}, \bar{b}) = \mathcal{J}^T \text{LLS}(A, b)(y, \bar{y})$ $= (-AM_S^{-1} \bar{y} y_D^T - Ay_D \bar{y}^T M_S^{-1} + b\bar{y}^T M_S^{-1}, A^T \bar{y})$
“Sketch + Diff”	$y_S = \text{LLS}_S(A, b, S)$	$(\bar{A}_S, \bar{b}_S) = \mathcal{J}^T \text{LLS}_S(A, b, S)(y_S, \bar{y})$ $= (-S^T A_S^T \bar{y} y_S^T - S^T S A y_S \bar{y}^T M_S^{-1} + S^T S b \bar{y}^T M_S^{-1}, S^T A_S^T \bar{y})$

A.3 “Sketch and Differentiate”

Lemma 2. *The reverse mode approximation error for the term \bar{b} when we approximate it by sketching matrix S can be bounded with probability $1 - \delta$ as follows: $\|\bar{b} - \bar{b}_S\|_2 \leq \|\Sigma^{-1}\|_2 \|\bar{y}\|_2 (\epsilon + (1 + \epsilon)\|I - S^T S\|_2)$.*

Proof. Let us use Lemma 1 and sub-multiplicativity to obtain the following:

$$\begin{aligned}
\|\bar{b} - \bar{b}_S\|_2 &= \|AM^{-T}\bar{y} - S^T SAM_S^{-T}\bar{y}_S\|_2 \\
&= \|AM^{-1}\bar{y} - AM_S^{-1}\bar{y} + AM_S^{-1}\bar{y} - S^T SAM_S^{-1}\bar{y}\|_2 \\
&\leq \|AM^{-1} - AM_S^{-1}\|_2 \|\bar{y}\|_2 + \|I - S^T S\|_2 \|AM_S^{-1}\|_2 \|\bar{y}\|_2 \\
&\leq \epsilon \|\Sigma^{-1}\|_2 \|\bar{y}\|_2 + \|I - S^T S\|_2 \|AM_S^{-1}\|_2 \|\bar{y}\|_2 \\
&\leq \|\Sigma^{-1}\|_2 \|\bar{y}\|_2 (\epsilon + (1 + \epsilon)\|I - S^T S\|_2)
\end{aligned} \tag{11}$$

□

where we used a lemma from the main paper. So, the error can be large ($\|I - S^T S\|_2$).

Lemma 3. *The reverse mode approximation error for the term \bar{A} when we approximate it using the sketching matrix S can be bounded with probability $1 - \delta$.*

Proof.

$$\begin{aligned}
\|\bar{A} - \bar{A}_S\|_F &= \|-2AM^{-T}\bar{y}y^T + b\bar{y}^T M^{-1} - (-2S^T SAM_S^{-T}\bar{y}_S y_S^T + S^T S b\bar{y}_S^T M_S^{-1})\|_F \\
&\leq \|2AM^{-1}\bar{y}y^T - 2S^T SAM_S^{-1}\bar{y}_S y_S^T\|_F + \|b\bar{y}^T M^{-1} - S^T S b\bar{y}_S^T M_S^{-1}\|_F
\end{aligned} \tag{12}$$

$$\begin{aligned}
\|AM^{-1}\bar{y}y^T - S^T SAM_S^{-1}\bar{y}_S y_S^T\|_F &\leq \|AM^{-1}\bar{y}y^T - AM_S^{-1}\bar{y}_S y_S^T\|_F + \|AM_S^{-1}\bar{y}_S y_S^T - S^T SAM_S^{-1}\bar{y}_S y_S^T\|_F \\
&\leq \epsilon \|\Sigma^{-1}\|_F \|\bar{y}\| \|y\| + \|AM_S^{-1}\bar{y}_S y_S^T - S^T SAM_S^{-1}\bar{y}_S y_S^T\|_F
\end{aligned} \tag{13}$$

□

A.4 “Differentiate and Sketch”

Lemma 4. *The reverse mode approximation error for the term \bar{b} when we sketch only the computationally expensive terms by S , with probability at least $1 - \delta$, satisfies: $\|\bar{b} - \bar{b}_S\|_2 \lesssim \epsilon \|\Sigma^{-1}\|_2 \|\bar{y}\|_2$.*

Proof. Let us use the sketching properties and sub-multiplicativity to obtain the following:

$$\begin{aligned}
\|\bar{b} - \bar{b}_S\|_2 &= \|AM^{-T}\bar{y} - AM_S^{-T}\bar{y}_S\|_2 \\
&\approx \|U(I - U^T S^T S U)\Sigma^{-1}V^T\bar{y}\|_2 \quad \bar{y} \approx \bar{y}_S \\
&\lesssim \epsilon \|U\|_2 \|\Sigma^{-1}\|_2 \|\bar{y}\|_2 \\
&\lesssim \epsilon \|\Sigma^{-1}\|_2 \|\bar{y}\|_2
\end{aligned} \tag{14}$$

□

Lemma 5. *The reverse mode approximation error for the term \bar{A} when we sketch only the computationally expensive terms by S , with probability $1 - 1/\text{poly}(d)$, satisfies: $\|\bar{A} - \bar{A}_S\|_2 \lesssim \epsilon \|\bar{y}\|_2 (\|\Sigma^{-1}\|_2 \|y\|_2 + \frac{1}{1-\epsilon} \|\Sigma^{-1}\|_2 \|Ay - b\|_2 \|A^\dagger\|_2)$.*

Proof. The approximation error can be split into 3 terms such that $\|\bar{A} - \bar{A}_S\| \leq Q_1 + Q_2 + Q_3$ where:

$$\begin{aligned}
Q_1 &= \|b\bar{y}^T M^{-1} - b\bar{y}_S^T M_S^{-1}\|_F \\
&\leq \epsilon \|b\bar{y}^T\|_F \|\Sigma^{-1}\|_2 \|\Sigma^{-1}\|_F
\end{aligned}$$

Let us bound Q_2 as follows:

$$\begin{aligned}
Q_2 &= \|AM^{-1}\bar{y}y^T - AM_{S'}^{-1}\bar{y}y_S^T\|_F = \|A(M^{-1} - M_{S'}^{-1})\bar{y}y^T + AM_{S'}^{-1}\bar{y}y^T - AM_{S'}^{-1}\bar{y}y_S^T\|_F \\
&\leq \|A(M^{-1} - M_{S'}^{-1})\bar{y}y^T\|_F + \|AM_{S'}^{-1}\bar{y}(y - y_S)^T\|_F \\
&\leq \epsilon\|\Sigma^{-1}\|_2\|\bar{y}y^T\|_F + \|AM_{S'}^{-1}\|_2\|\bar{y}(y - y_S)^T\|_F \\
&\leq \epsilon\|\Sigma^{-1}\|_2\|\bar{y}\|_2\|y\|_2 + \|AM_{S'}^{-1}\|_2\|\bar{y}\|_2\|(y - y_S)\|_2 \\
&\leq \epsilon\|\bar{y}\|_2(\|\Sigma^{-1}\|_2\|y\|_2 + \|AM_{S'}^{-1}\|_2\|Ay - b\|_2\|A^\dagger\|_2) \\
&\leq \epsilon\|\bar{y}\|_2(\|\Sigma^{-1}\|_2\|y\|_2 + (1 + \epsilon)\|\Sigma^{-1}\|_2\|Ay - b\|_2\|A^\dagger\|_2)
\end{aligned} \tag{15}$$

where we used the following result [Price et al. \(2017\)](#):

$$\|y - y_S\|_2 \leq \epsilon\|Ay - b\|_2\|A^\dagger\|_2 \tag{16}$$

and the last term Q_3 can be bounded as:

$$Q_3 = \|Ay\bar{y}^T M^{-1} - Ay_S\bar{y}^T M_S^{-1}\| = \|Ay\bar{y}^T(M^{-1} - M_S^{-1}) + Ay\bar{y}^T M_S^{-1} - Ay_S\bar{y}^T M_S^{-1}\| \tag{17}$$

$$\begin{aligned}
&\leq \epsilon\|Ay\bar{y}^T\| + \|A(y - y_S)\bar{y}^T M_S^{-1}\| \\
&\leq \epsilon\|Ay\bar{y}^T\| + \epsilon\|A\|\|Ay - b\|\|\bar{y}^T M_S^{-1}\|
\end{aligned} \tag{18}$$

Note that all three terms Q_1, Q_2, Q_3 are $O(\epsilon)$. \square

B Experiments

We plot the performance of the two proposed approaches for obtaining forward and reverse mode AD in the case of linear regression. We generate a linear regression problem by choosing the entries of matrix A and vector b from i.i.d. $N(0, 1)$ (Normal distribution with mean 0 and variance 1). The differences from the two approaches, “sketch+differentiate” and “differentiate+sketch” are shown in Figure 1.

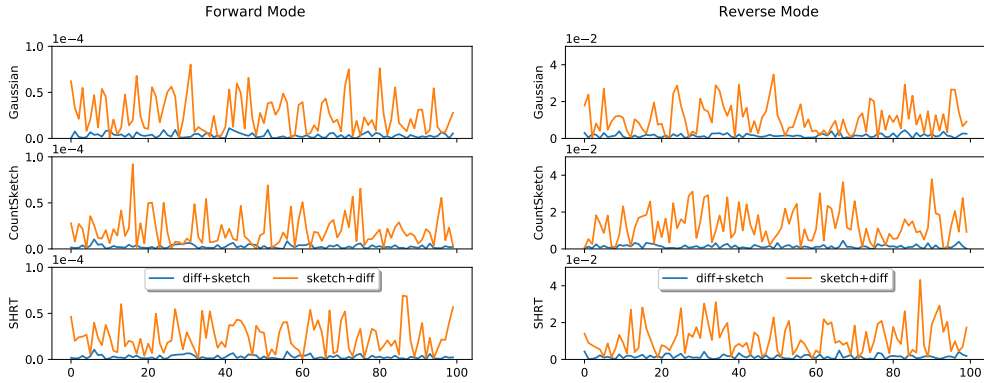


Figure 1: Numerical observation that differentiation and sketching do not commute, and that differentiation-then-sketch is more accurate. We show the forward mode along with its approximation corresponding to the three sketching matrices of Gaussian, Count-sketch and Subsampled Randomized Hadamard Transform (SRHT), on a randomly generated least squares problem of size 100000×100 , along with a random perturbation. Reverse mode is shown for a subsample of 100 randomly chosen values for the variable b , where we used sign as the cost function.

References

Eric Price, Zhao Song, and David P. Woodruff. Fast regression with an l_∞ guarantee. In *44th International Colloquium on Automata, Languages, and Programming, ICALP 2017, July 10-14, 2017, Warsaw, Poland*, pages 59:1–59:14, 2017. doi: 10.4230/LIPIcs.ICALP.2017.59.