

TODO

Kenny Chiu

February 24, 2022

1 Summary

1.1 Context and background

Lacotte et al. [13] study the performance of iterative Hessian sketch (IHS) [17] for overdetermined least squares problems of the form

$$\mathbf{b}^* = \arg \min_{\mathbf{b} \in \mathbb{R}^d} \left\{ f(\mathbf{b}) = \frac{1}{2} \|\mathbf{X}\mathbf{b} - \mathbf{y}\|^2 \right\}$$

where $\mathbf{X} \in \mathbb{R}^{n \times d}$, $n \geq d$, is a given full rank data matrix and $\mathbf{y} \in \mathbb{R}^n$ is a vector of observations. IHS is an iterative method based on random projections that is effective for large data and ill-conditioned problems. Given step sizes $\{\alpha_t\}$ and momentum parameters $\{\beta_t\}$, the IHS solution is iteratively updated using

$$\mathbf{b}_{t+1} = \mathbf{b}_t - \alpha_t \mathbf{H}_t^{-1} \nabla f(\mathbf{b}_t) + \beta_t (\mathbf{b}_t - \mathbf{b}_{t-1})$$

where $\mathbf{H}_t = \mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X}$ is an approximation of the Hessian $\mathbf{H} = \mathbf{X}^\top \mathbf{X}$ given refreshed (i.i.d.) $m \times n$ sketching (random) matrices $\{\mathbf{S}_t\}$ with $m \ll n$. The performance of IHS with Gaussian sketches (i.e., where $(\mathbf{S}_t)_{ij}$ are i.i.d. $N(0, m^{-1})$) has been studied, but IHS with other sketches have only been empirically studied. In their work, Lacotte et al. [13] draw on results from random matrix and free probability theory and show that the following sketches (asymptotically) converge faster to the optimal solution compared to Gaussian sketches:

1. Truncated Haar sketch, where the rows of \mathbf{S}_t are orthonormal. Orthogonal sketches are preferred over i.i.d. sketches as they do not distort the projection, but orthogonality in general Haar matrices come at the expense of requiring the Gram-Schmidt procedure, which has cost $O(nm^2)$ larger than the $O(nmd)$ cost when using Gaussian sketches.
2. A version of the subsampled randomized Hadamard transform (SRHT), with \mathbf{S}_t constructed from $\mathbf{R}_t = n^{-\frac{1}{2}} \mathbf{B}_t \mathbf{W}_n \mathbf{D}_t \mathbf{P}_t$ where \mathbf{B}_t is a $n \times n$ diagonal matrix of i.i.d. Bernoulli($\frac{m}{n}$) samples, \mathbf{D}_t is a $n \times n$ diagonal matrix of i.i.d. Rademacher samples, \mathbf{P}_t is a $n \times n$ uniformly sampled row permutation matrix, and \mathbf{W}_n is the $n \times n$ Walsh-Hadamard matrix defined recursively as

$$\mathbf{W}_n = \begin{bmatrix} \mathbf{W}_{\frac{n}{2}} & \mathbf{W}_{\frac{n}{2}} \\ \mathbf{W}_{\frac{n}{2}} & -\mathbf{W}_{\frac{n}{2}} \end{bmatrix}$$

where $\mathbf{W}_1 = 1$. \mathbf{S}_t is taken as \mathbf{R}_t with the zeros rows removed (as selected by \mathbf{B}_t). Note that due to this subsampling, \mathbf{S}_t is a $M \times n$ matrix with $\mathbb{E}[M] = m$. By construction, SRHT sketches are orthogonal. Sketching with SRHT only requires $O(nd \log M)$.

1.2 Main contributions

The main contributions of Lacotte et al. [13] include several theoretical results that describe the (asymptotically) optimal value of the parameters for IHS with refreshed Haar or SRHT sketches, the corresponding convergence rates of IHS with these parameters, and closed form

expressions for the inverse moments of SRHT sketches. These results are obtained based on asymptotic results from random matrix theory, in which it is assumed that the matrix dimensions satisfy the aspect ratios $\frac{d}{n} \rightarrow \gamma \in (0, 1)$ and $\frac{m}{n} \rightarrow \xi \in (\gamma, 1)$ as $n, d, m \rightarrow \infty$.

The main results are Theorems 3.1 and 4.1. Theorem 3.1 says that for IHS with refreshed Haar sketches, the optimal convergence rate ρ_H of the relative prediction error is

$$\rho_H = \left(\lim_{n \rightarrow \infty} \frac{\mathbb{E} [\|\mathbf{X}(\mathbf{b}_t - \mathbf{b}^*)\|^2]}{\|\mathbf{X}(\mathbf{b}_0 - \mathbf{b}^*)\|^2} \right)^{\frac{1}{t}} = \rho_G \cdot \frac{\xi(1 - \xi)}{\gamma^2 + \xi - 2\xi\gamma}$$

where ρ_G is the optimal rate of IHS with Gaussian sketches. The aspect ratio scaling factor is less than 1, implying that $\rho_H < \rho_G$ and that IHS with Haar sketches converges faster than with Gaussian sketches. Theorem 4.1 states that the rate ρ_S for IHS with refreshed SRHT sketches is equal to ρ_H under an additional mild assumption on the initialization of the least squares problem (which was not needed for Haar sketches due to existing results from random matrix theory). Theorem 3.1 also states that the optimal convergence rate for IHS with Haar sketches is obtained using momentum values $\beta_t = 0$ (i.e., momentum does not help) and step sizes $\alpha_t = \frac{\theta_{1,H}}{\theta_{2,H}}$ where $\theta_{k,H}$ is the k -th inverse moment of the Haar sketch defined as

$$\theta_{k,H} = \lim_{n \rightarrow \infty} \frac{1}{d} \mathbb{E} [\text{trace} ((\mathbf{U}^\top \mathbf{S}^\top \mathbf{S} \mathbf{U})^{-k})]$$

for $m \times n$ Haar matrix \mathbf{S} and $n \times d$ deterministic matrix \mathbf{U} with orthonormal columns. Closed-form expressions for the first two inverse moments are provided in Lemma 3.2 and are given by

$$\theta_{1,H} = \frac{1 - \gamma}{\xi - \gamma}, \quad \theta_{2,H} = \frac{(1 - \gamma)(\gamma^2 + \xi - 2\gamma\xi)}{(\xi - \gamma)^3}.$$

Theorem 4.1 and Lemma 4.3 together state that the limiting distribution of Haar and SRHT sketches are the same when there is no momentum and therefore so is the optimal step size. However, the optimality of $\beta_t = 0$ for IHS with SRHT sketches is only a conjecture based on numerical simulations.

Other contributions of Lacotte et al. [13] include a complexity analysis of IHS with SRHT sketches and an empirical study of the theoretical results. The complexity analysis concludes that the asymptotic performance of IHS with SRHT sketches is faster than that of the pre-conditioned conjugate gradient method (pCG) [22] by a factor of $\log(d)$. The empirical study verifies that the limiting results can apply in the finite case where the convergence of IHS with Haar and with SRHT sketches are similar and faster than that of Gaussian sketches on ill-conditioned synthetic and real datasets of moderate size ($n \geq 4000$, $d \geq 200$), and that the IHS with SRHT sketches refreshed every iteration has faster convergence than pCG on a similar synthetic dataset.

1.3 Limitations

Limitations of the work by Lacotte et al. [13] include the reliance on asymptotic theory, the empirical evaluation of results on mostly synthetic datasets, the comparison of sketches based

on a single criterion, and the unclear generalizability of results to more complicated problems.

Lacotte et al. [13] obtain the convergence rates of IHS with different sketching matrices by drawing on results from asymptotic random matrix theory. While their simulations show that the theory does apply in moderately-sized datasets, the datasets that they examine are primarily synthetic and designed to satisfy assumptions even if ill-conditioned. However, there is also the counterargument that IHS would only be considered over standard solvers for large data problems, and so these limitations are relatively minor.

Another limitation of their work is that only a single criterion—namely the prediction error between the sketched solution and the optimal solution—is used to compare the performance of the sketching matrices. Other criteria have also been considered in the literature, such as those based on other losses or those based on out-of-sample prediction [8, 17]. While certain criteria are intrinsically related [9], they may still have differing properties and lead to differing results [8].

The main limitation of the work by Lacotte et al. [13] is the simple problem context that the results are derived for. While the theory shows that IHS is promising for large data, overdetermined least squares problems, standard solvers would still be preferred over IHS in large data problems if the appropriate computational resources were available. It is unclear whether the theory could generalize to more complicated problems, such as to undetermined least squares problems or optimization problems with other losses. It would be particularly useful to understand whether there are problems for which IHS would be preferred over conventional solvers in the general case.

1.4 Related literature

Works that analyze the impact of sketch type in sketching methods make up a small portion of the sketching literature. The work by Lacotte et al. [13] is said to be inspired by and therefore most similar to the work by Dobriban and Liu [8], which appears to be the first in the literature to leverage results from asymptotic random matrix theory. Analysis in the asymptotic regime appears to be key in being able to differentiate between the analytical performance of different sketching matrices, which was a challenge in previous works [3, 17, 20]. More recently, Lacotte and Pilanci [12] directly extended their analysis of refreshed sketches in IHS to fixed sketches in a related first-order method that has better guarantees.

Recent related works in the literature also include those that propose extensions of IHS, e.g., IHS with momentum and fixed sketches [15], distributed IHS [6], first-order IHS with adaptive step sizes [25], and Newton sketch [18] (IHS for general convex optimization problems) and its own variants [e.g., 7, 14]. Analyses of the performance in these works generally are intended as a point of comparison against existing methods, are done empirically or make use of conventional analysis techniques rather than asymptotic theory, and do not particularly examine the impact of specific sketching matrices.

2 Mini-proposals

2.1 Proposal 1: A sketched interior point algorithm for quantile regression

Quantile regression [11] offers several advantages over linear regression, such as being able to model different quantiles (as opposed to only a mean), being free from assumptions regarding the parametric form of the response and homoscedasticity, and being transformation equivariant in its response [21]. Given a data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, observations $\mathbf{y} \in \mathbb{R}^n$ and a quantile $\tau \in (0, 1)$, quantile regression fits a linear model on the quantile with the estimated parameters being the solution to the optimization problem

$$\min_{\mathbf{b} \in \mathbb{R}^d} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{b}) (\tau - \mathbb{1}[y_i - \mathbf{x}_i^\top \mathbf{b} < 0]) .$$

The objective is non-differentiable as-is but can be optimized as a linear program. For large data problems, the interior point method transforms the dual program into a constrained optimization problem using log barriers [19], i.e.,

$$\begin{aligned} \arg \max_{\mathbf{a}} \mathbf{y}^\top \mathbf{a} & \implies \arg \max_{\mathbf{a}} \mathbf{y}^\top \mathbf{a} + \mu \sum_{i=1}^n \log a_i \\ \text{s.t. } \begin{cases} \mathbf{X}^\top \mathbf{a} = (1 - \tau) \mathbf{X}^\top \mathbf{1}_n \\ \mathbf{a} \in [0, 1]^n \end{cases} & \text{s.t. } \mathbf{X}^\top \mathbf{a} = (1 - \tau) \mathbf{X}^\top \mathbf{1}_n . \end{aligned}$$

The interior point method solves the dual program by taking a sequence of Newton steps with $\mu \rightarrow 0$. The solution to the Newton step in each iteration satisfies the equation

$$\mathbf{X}^\top \mathbf{W}_t \mathbf{X} \mathbf{b}_t = \mathbf{X}^\top \mathbf{W}_t (\mathbf{y} + \mu \mathbf{A}^{-1} \mathbf{1}_n) .$$

where \mathbf{A} is a $n \times n$ diagonal matrix and \mathbf{W}_t is a $n \times n$ diagonal matrix with positive diagonal entries that changes each iteration. Computing $\mathbf{X}^\top \mathbf{W}_t \mathbf{X}$ is the main computational bottleneck that leads to each iteration having a cost of $O(nd^2)$ [2].

In this proposed project, we consider the case where $d \ll n$ and propose a stochastic interior point algorithm that uses sketching matrices to reduce the cost of the iterative updates. Drawing on methods from the sketching literature [18], the idea is to incorporate partial sketching into the original algorithm where instead of computing $\mathbf{X}^\top \mathbf{W}_t \mathbf{X}$, we compute $\mathbf{X}^\top \mathbf{W}_t^{\frac{1}{2}} \mathbf{S}_t^\top \mathbf{S}_t \mathbf{W}_t^{\frac{1}{2}} \mathbf{X}$. The matrix $\mathbf{S}_t \in \mathbb{R}^{m \times n}$, $m \ll n$, is a dimension-reducing random matrix regenerated every iteration. For example, the subsampled randomized Hadamard transform allows the sketch $\mathbf{S}_t \mathbf{W}_t^{\frac{1}{2}} \mathbf{X}$ to be formed at a cost of $O(nd \log m)$ [13], and so the matrix product above can be computed at a cost of $O(md^2)$. While the sketched solution will only be an approximation of the original, recent work on the convergence of sketching in other optimization problems show promising theoretical and empirical results [e.g., 7, 14, 18]. We note that Yang et al. [24] had previously proposed a stochastic algorithm for quantile

regression. However, their method differs from ours in that they construct a random preconditioning matrix before using standard methods to solve the optimization problem on the conditioned data matrix.

The main contributions of this project would be as follows:

1. A sketched interior point algorithm for optimizing quantile regression problems that is expected to be faster than standard methods currently used in practice.
2. A theoretical analysis of the proposed sketched interior point algorithm that provides convergence guarantees.
3. An empirical comparison of quantile regression models obtained from the proposed sketched algorithm and other existing methods, such as the standard interior point method [19] (implemented in R), the stochastic method by Yang et al. [24], a more modern iteration of the interior point method [26], and a modern quantile regression algorithm based on smoothing [10] (also implemented in R), in large dataset settings.
4. An implementation of the sketch interior point algorithm, e.g., in R, if found to have practical advantages over the existing algorithms.

The main challenge in this project would be the theoretical analysis of the sketched interior point algorithm. The most feasible analysis approach would likely be following that of Pilanci and Wainwright [18] for interior point methods and partial sketches, which would provide a worst-case convergence guarantee for the number of iterations needed to obtain a solution within a desired tolerance. The effect of the sketching matrix may also be of interest, but an analysis approach similar to the asymptotic approach of Lacotte et al. [13] would likely be necessary. However, adapting their least squares approach to that of quantile regression is not straightforward and would likely be more suited for a follow-up project.

Following the completion of this project, there are multiple directions of future work that may be of interest:

1. The proposed sketched interior point algorithm would not be useful for the $n \ll d$ case. A sketch-based method would likely still be possible but would need to use sketches differently, e.g., directly sketching the data matrix as Pham and El Ghaoui [16] did for LASSO, or sketching both the data matrix and the observations as in classical least-squares sketch (although this has been shown to lead to suboptimal performance [17]).
2. Applications of quantile regression or interior point algorithms in general, e.g., composite quantile regression [27] for high-dimensional regression, applications of quadratic programming, may benefit from the faster algorithm.
3. Sketched-based smoothing algorithms for quantile regression. These algorithms approximate the original optimization problem by a differentiable one and therefore sketching should directly follow from the work of Pilanci and Wainwright [18]. Given the more standard setup, it is likely easier to analyze these algorithms than the interior point algorithms.

3 Project report

Abstract

The partial Newton sketch algorithm can be used as an approximate iterative solver for ridge regression. Following Lacotte et al. [13], we attempt to analyze the theoretical properties of partial Newton sketch for ridge regression using results from free probability and random matrix theory. We show that such an approach is not trivial and highlight the aspects of ridge regression that make doing so challenging. We make partial progress in the theory under a hypothetical trace decoupling condition and present some empirical evidence to support this hypothesis.

3.1 Introduction

Ridge regression is a special case of regularized least squares where the penalty function is chosen to be the ℓ_2 -norm of the model parameters. Given data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, observations $\mathbf{y} \in \mathbb{R}^n$ and a regularization parameter $\lambda > 0$, ridge regression obtains estimates of the parameters as the solution to the optimization problem

$$\mathbf{b}^* = \arg \min_{\mathbf{b} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}\mathbf{b} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{b}\|_2^2.$$

While ridge regression can be motivated as a method for reducing overfitting in ordinary least squares (OLS), it also has its computational and analytical benefits over OLS. When \mathbf{X} does not have full column rank (e.g., when $n < d$), then $\mathbf{X}^\top \mathbf{X}$ is singular and the OLS solution is non-unique. When \mathbf{X} is full rank but ill-conditioned, then small changes in \mathbf{X} lead to large changes in $(\mathbf{X}^\top \mathbf{X})^{-1}$ and consequently in the OLS solution. Ridge regression addresses both of these issues by minimizing the variance and mean squared error at the cost of introducing a small bias [4]. The ridge regression solution is unique and is given by

$$\mathbf{b}^* = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} \mathbf{X}^\top \mathbf{y}.$$

In this report, we analyze the theoretical properties of a partial Newton sketch algorithm [18] as an iterative solver for the ridge regression problem. In particular, we attempt to derive an optimal convergence rate and an optimal step size following the approach of Lacotte et al. [13] for iterative Hessian sketch with OLS using asymptotic results from random matrix theory and free probability. We show that while ridge regression can be considered a simple extension to OLS, extending the analysis approach of Lacotte et al. [13] to partial Newton sketch is not trivial.

This report is organized as follows: Section 3.2 provides background about sketching and describes the OLS results by Lacotte et al. [13] that we aim to extend to ridge regression; Section 3.3 highlights relevant work in the literature; Section 3.5 discusses our attempts to analyze Newton sketch for ridge regression and the key differences from OLS that makes the problem challenging; Section 3.6 describes a simulation that empirically supports one of the hypothesized theoretical conditions; and Section 3.7 summarizes our findings and concludes this report.

3.2 Background

This section briefly describes the Newton sketch algorithm and the probability subfields of random matrix theory and free probability that are referred to in this report.

3.2.1 Newton's method and Newton sketch

Given a convex, twice-differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, Newton's method is an efficient iterative method for finding the minimizing solution. The iterative updates are of the form

$$\mathbf{x}_{t+1} = \mathbf{x}_t - (\nabla^2 f(\mathbf{x}_t))^{-1} \nabla f(\mathbf{x}_t)$$

where $\nabla f(\mathbf{x}_t)$ and $\nabla^2 f(\mathbf{x}_t)$ are the gradient and Hessian of f evaluated at \mathbf{x}_t , respectively. Depending on the problem, computing the Hessian $\nabla^2 f(\mathbf{x}_t)$ may be a computational bottleneck. The general Newton sketch algorithm [18] avoids computing $\nabla^2 f(\mathbf{x}_t)$ exactly and instead approximates it with

$$\nabla^2 f(\mathbf{x}_t) \approx \nabla^2 f(\mathbf{x}_t)^{\frac{1}{2}} \mathbf{S}_t^\top \mathbf{S}_t \nabla^2 f(\mathbf{x}_t)^{\frac{1}{2}}$$

where \mathbf{S}_t is a random, rectangular sketching matrix introduced for dimension reduction. In this project, we only consider refreshed sketches where $\mathbf{S}_0, \dots, \mathbf{S}_t$ are i.i.d. realizations (as opposed to a single, fixed realization) to support an analysis based on free probability. We also do not focus on the specific type of sketching matrix used as the challenges we encounter hold generally across sketch types (e.g., i.i.d. Gaussian or orthogonal sketches).

For functions with an additive decomposition of the form $f = f_0 + g$, it may be sufficient for computational improvements to only do a partial Newton sketch [18] where the Hessian is approximated by

$$\nabla^2 f(\mathbf{x}_t) \approx \nabla^2 f_0(\mathbf{x}_t)^{\frac{1}{2}} \mathbf{S}_t^\top \mathbf{S}_t \nabla^2 f_0(\mathbf{x}_t)^{\frac{1}{2}} + \nabla^2 g(\mathbf{x}_t).$$

In particular, the ridge regression loss function has this form where it is analytically more convenient to consider partial Newton sketch updates.

3.2.2 Random matrix theory and free probability

Free probability theory was initially developed by Voiculescu in the 1980's [1] and is concerned with the study of non-commutative random variables. More recently, its connections to random matrix theory were established as a means of studying the limiting spectral distribution of random matrices. In free probability, the notion of freeness is analogous to independence in classical probability theory. A family of random $n \times n$ matrices $\{\mathbf{X}_n^{(1)}, \dots, \mathbf{X}_n^{(I)}\}$ is said to be asymptotically free [5] if

1. $\mathbf{X}_n^{(i)}$ has a limiting spectral distribution for all $i \in \{1, \dots, I\}$, and
2. for all $\{i_1, \dots, i_J\}$ with $i_j \in \{1, \dots, I\}$, $j \in \{1, \dots, J\}$ and $i_1 \neq i_2, \dots, i_{J-1} \neq i_J$, and for all polynomials P_1, \dots, P_J such that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} [\text{trace} (P_j (\mathbf{X}_n^{(i_j)}))] = 0$$

for all $j \in \{1, \dots, J\}$, we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E} \left[\text{trace} \left(\prod_{j=1}^J P_j (\mathbf{X}_n^{(i_j)}) \right) \right] = 0 .$$

The analysis approach of Lacotte et al. [13] uses freeness to calculate the expected normalized trace of matrix products as for asymptotically free random matrices \mathbf{X}_i and \mathbf{X}_j , it holds that as $n \rightarrow \infty$,

$$\frac{1}{n} \mathbb{E} [\text{trace}(\mathbf{X}_i \mathbf{X}_j)] - \frac{1}{n} \mathbb{E} [\text{trace}(\mathbf{X}_i)] \frac{1}{n} \mathbb{E} [\text{trace}(\mathbf{X}_j)] \rightarrow 0 .$$

Other tools from random matrix theory such as transforms of spectral distributions are also used in their analysis, but we do not reach the stage in the analysis where they are useful in this report and so we do not cover them here.

3.3 Related work

Sketched algorithms for ridge regression have been previously considered in the literature. Chowdhury et al. [4] examined ridge regression in the setting where $n \ll d$ and showed that under certain conditions on the sketched approximation, the relative error in the solution of the partial Newton sketch algorithm is bounded above by an exponentially decaying tolerance. Wang et al. [23] compared the classical sketching algorithm (i.e., sketching both the data matrix and the response vector) to the partial Newton sketch algorithm in matrix ridge regression for the case $d \ll n$ and found that both have increased risks relative to the optimal solution. They proposed model averaging as a solution for improving the theoretical properties. It is notable that the analyses by Chowdhury et al. [4] and Wang et al. [23] are both based on conventional statistical learning techniques and that the specific sketching matrix used is not a focus of the analysis.

In this project, we attempt to analyze the partial Newton sketch for ridge regression following the approach that Lacotte et al. [13] used for iterative Hessian sketch (IHS) with OLS. Their approach relies on asymptotic results from random matrix theory and free probability, and this appeared to allow a finer-grain analysis of IHS where the derived convergence rate depends on the specific sketching matrix used. The random matrix theoretic approach to analysis does not seem to have been considered much outside of OLS problems [8, 12]. As we explore in this project, a possible reason for this may be that even simple extensions to standard OLS make it difficult to directly apply existing results from these subfields of probability.

3.4 Newton sketch for ridge regression

Consider the ridge regression loss function

$$f(\mathbf{b}) = \frac{1}{2} \|\mathbf{X}\mathbf{b} - \mathbf{y}\|_2^2 + \frac{\lambda}{2} \|\mathbf{b}\|_2^2$$

for $\mathbf{b} \in \mathbb{R}^d$ given data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ with $d \ll n$, responses $\mathbf{y} \in \mathbb{R}^n$ and a regularization parameter $\lambda > 0$. The gradient and Hessian of the function are respectively given by

$$\begin{aligned}\nabla f(\mathbf{b}) &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) \mathbf{b} - \mathbf{X}^\top \mathbf{y} , \\ \mathbf{H} = \nabla^2 f(\mathbf{b}) &= \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d .\end{aligned}$$

The Newton updates described in Section 3.2.1 therefore have the form

$$\begin{aligned}\mathbf{b}_{t+1} &= \mathbf{b}_t - \alpha_t \mathbf{H}^{-1} \nabla f(\mathbf{b}_t) \\ &= \mathbf{b}_t - \alpha_t (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d)^{-1} ((\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) \mathbf{b}_t - \mathbf{X}^\top \mathbf{y}) .\end{aligned}$$

A partial Newton sketch of the Hessian has the form

$$\mathbf{H}_t = \mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X} + \lambda \mathbf{I}_d$$

where \mathbf{S}_t is a $m \times n$ refreshed sketching matrix with $d < m \ll n$. The partial Newton sketch updates for ridge regression are then given by

$$\mathbf{b}_{t+1} = \mathbf{b}_t - \alpha_t (\mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X} + \lambda \mathbf{I}_d)^{-1} ((\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) \mathbf{b}_t - \mathbf{X}^\top \mathbf{y}) .$$

Note that Chowdhury et al. [4] and Wang et al. [23]) also considered this sketched update for ridge regression. However, our analysis approach differs from theirs in that we adopt the asymptotic random matrix theoretic approach from Lacotte et al. [13]. Also note that we do not consider updates with momentum as Lacotte et al. [13] did as we will show that extending their analysis approach from OLS to ridge regression is already non-trivial.

3.5 Analysis attempt based on random matrix theory

In this section, we show that the proof technique used to obtain Theorems 3.1 and 4.1 of [13] do not easily generalize to the partial Newton sketch updates for ridge regression. We follow the general procedure of the proofs and show how far we can get with the ridge regression setup. We also highlight the key differences between OLS and ridge regression that leads to problems in the proof and discuss possible solutions for rectifying these problems in future work.

The following conjecture formalizes the result analogous to Theorems 3.1 and 4.1 that we would like to prove. Note that additional assumptions will be added to the conjecture as we progress through the proof.

Conjecture 1. *Consider the partial Newton update for ridge regression described in Section 3.4. Let $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^\top$ be the thin singular value decomposition of \mathbf{X} where \mathbf{U} is a $n \times d$ semi-orthogonal matrix, \mathbf{V} is a $d \times d$ orthogonal matrix, and Σ is a $d \times d$ diagonal matrix with the singular values of \mathbf{X} on the diagonal. Define the error vector $\Delta_t = \mathbf{U}^\top \mathbf{X} (\mathbf{b}_t - \mathbf{b}^*)$. For some optimal step size α_t , the sequence of error vectors $\{\Delta_t\}$ satisfies*

$$\rho = \left(\lim_{n \rightarrow \infty} \frac{\mathbb{E} [\|\Delta_t\|_2^2]}{\|\Delta_0\|_2^2} \right)^{\frac{1}{t}}$$

where ρ is the optimal rate of convergence with some closed-form expression.

We begin our attempt to prove Conjecture 1 following the proofs by Lacotte et al. [13]. Using the fact that the ridge regression solution satisfies the equation

$$(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) \mathbf{b}^* = \mathbf{X}^\top \mathbf{Y} ,$$

the update can be rewritten as

$$\begin{aligned} \mathbf{b}_{t+1} &= \mathbf{b}_t - \alpha_t (\mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X} + \lambda \mathbf{I}_d)^{-1} ((\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) \mathbf{b}_t - (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) \mathbf{b}^*) \\ &= \mathbf{b}_t - \alpha_t (\mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X} + \lambda \mathbf{I}_d)^{-1} (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d) (\mathbf{b}_t - \mathbf{b}^*) . \end{aligned}$$

Using the thin SVD of \mathbf{X} , we have the matrix identities

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d &= \mathbf{V} \Sigma^2 \mathbf{V}^\top + \lambda \mathbf{V} \mathbf{V}^\top \\ &= \mathbf{V} (\Sigma^2 + \lambda \mathbf{I}_d) \mathbf{V}^\top , \\ (\mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X} + \lambda \mathbf{I}_d)^{-1} &= (\mathbf{V} \Sigma \mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} \Sigma \mathbf{V}^\top + \lambda \mathbf{V} \mathbf{V}^\top)^{-1} \\ &= \mathbf{V} (\Sigma \mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} \Sigma + \lambda \mathbf{I}_d)^{-1} \mathbf{V}^\top . \end{aligned}$$

However, in order to later on obtain an expression in terms of Δ_t as in the original proof, we require that the data matrix be full column rank. This is a less than ideal assumption to make as one of the advantages of ridge regression is being able to obtain a unique solution with non-full rank data matrices. We return to this point in Section 3.5.1 to discuss how this assumption may be avoided.

Assumption 1. The data matrix \mathbf{X} has full column rank.

Under Assumption 1, the singular values of \mathbf{X} are non-zero and so the above matrices can be rewritten as

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d &= \mathbf{V} \Sigma (\mathbf{I}_d + \lambda \Sigma^{-2}) \Sigma \mathbf{V}^\top , \\ (\mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X} + \lambda \mathbf{I}_d)^{-1} &= \mathbf{V} \Sigma^{-1} (\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} \Sigma^{-1} \mathbf{V}^\top . \end{aligned}$$

Replacing the corresponding matrices in the update with these identities gives

$$\mathbf{b}_{t+1} = \mathbf{b}_t - \alpha_t \mathbf{V} \Sigma^{-1} (\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2}) \Sigma \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) .$$

Multiplying both sides by $\mathbf{U}^\top \mathbf{X}$ gives

$$\begin{aligned} \mathbf{U}^\top \mathbf{X} \mathbf{b}_{t+1} &= \mathbf{U}^\top \mathbf{X} \mathbf{b}_t - \alpha_t \mathbf{U}^\top \mathbf{X} \mathbf{V} \Sigma^{-1} (\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2}) \Sigma \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) \\ &= \mathbf{U}^\top \mathbf{X} \mathbf{b}_t - \alpha_t (\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2}) \Sigma \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) \end{aligned}$$

and then subtracting both sides by $\mathbf{U}^\top \mathbf{X} \mathbf{b}^*$ gives

$$\begin{aligned} \mathbf{U}^\top \mathbf{X} (\mathbf{b}_{t+1} - \mathbf{b}^*) &= \mathbf{U}^\top \mathbf{X} (\mathbf{b}_t - \mathbf{b}^*) - \alpha_t (\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2}) \Sigma \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) \\ &= \left(\mathbf{I}_d - \alpha_t (\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2}) \right) \mathbf{U}^\top \mathbf{X} (\mathbf{b}_t - \mathbf{b}^*) . \end{aligned}$$

Let $\mathbf{Q}_t = \mathbf{I}_d - \alpha_t (\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2})$. Therefore by definition, we have $\Delta_{t+1} = \mathbf{Q}_t \Delta_t$ and

$$\|\Delta_{t+1}\|^2 = \Delta_t^\top \mathbf{Q}_t^\top \mathbf{Q}_t \Delta_t.$$

Taking the expectation with respect to \mathbf{S}_t , we get

$$\mathbb{E} [\|\Delta_{t+1}\|^2] = \Delta_t^\top \mathbb{E} [\mathbf{Q}_t^\top \mathbf{Q}_t] \Delta_t$$

where

$$\begin{aligned} \mathbb{E} [\mathbf{Q}_t^\top \mathbf{Q}_t] &= \mathbf{I}_d - \alpha_t \mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} \right] (\mathbf{I}_d + \lambda \Sigma^{-2}) \\ &\quad - \alpha_t (\mathbf{I}_d + \lambda \Sigma^{-2}) \mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-1} \right] \\ &\quad + \alpha_t^2 (\mathbf{I}_d + \lambda \Sigma^{-2}) \mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} + \lambda \Sigma^{-2})^{-2} \right] (\mathbf{I}_d + \lambda \Sigma^{-2}). \end{aligned}$$

At this point, we run into our first major obstacle that prevents us from applying the key step of the proof of Theorem 3.1. In Theorem 3.1 for OLS, the expression that is obtained from taking the expectation is

$$\mathbb{E} [\|\Delta_{t+1}\|^2] = \Delta_t^\top \mathbb{E} [\mathbf{R}_t^2] \Delta_t$$

where

$$\mathbb{E} [\mathbf{R}_t^2] = \mathbf{I}_d - 2\alpha_t \mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U})^{-1} \right] + \alpha_t^2 \mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U})^{-2} \right].$$

The proof of Theorem 3.1 proceeds to recognize that the matrix $\mathbf{S}_t \mathbf{U}$ can be embedded into a Haar matrix and is therefore rotationally invariant. Using exchangeability arguments, the matrix $\mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U})^{-p} \right]$ has a simple closed-form expression in terms of the inverse moments from which the rest of the proof follows. We do not have rotational invariance in our ridge regression case, and so we follow the proof of Theorem 4.1 from this point onwards. Theorem 4.1 requires an additional assumption on the initialization of the problem in order to avoid computing the expectations directly.

Assumption 2. The initial error vector Δ_0 is random, independent of $\mathbf{S}_0, \dots, \mathbf{S}_t$, and satisfies $\mathbb{E} [\Delta_0 \Delta_0^\top] = \frac{\mathbf{I}_d}{d}$.¹

Under Assumption 2, taking the expectation with respect to \mathbf{S}_t gives

$$\begin{aligned} \mathbb{E} [\|\Delta_{t+1}\|^2] &= \mathbb{E} [\Delta_t^\top \mathbf{Q}_t^\top \mathbf{Q}_t \Delta_t] \\ &= \mathbb{E} [\Delta_0^\top \mathbf{Q}_0^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_0 \Delta_0] \\ &= \mathbb{E} [\text{trace} (\Delta_0^\top \mathbf{Q}_0^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_0 \Delta_0)] \\ &= \text{trace} (\mathbb{E} [\mathbf{Q}_0^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_0 \Delta_0 \Delta_0^\top]) . \end{aligned}$$

¹The convergence rate ρ in the statement of Conjecture 1 is also redefined as $\left(\lim_{n \rightarrow \infty} \frac{\mathbb{E} [\|\Delta_t\|_2^2]}{\mathbb{E} [\|\Delta_0\|_2^2]} \right)^{\frac{1}{t}}$ where $\mathbb{E} [\|\Delta_0\|_2^2] = \mathbb{E} [\Delta_0^\top \Delta_0] = \text{trace} (\mathbb{E} [\Delta_0 \Delta_0^\top]) = 1$.

By the independence of Δ_0 and \mathbf{Q}_i , we then have

$$\begin{aligned}\mathbb{E} [\|\Delta_{t+1}\|^2] &= \text{trace} \left(\mathbb{E} [\mathbf{Q}_0^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_0] \mathbb{E} [\Delta_0 \Delta_0^\top] \right) \\ &= \frac{1}{d} \text{trace} \left(\mathbb{E} [\mathbf{Q}_1^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_0 \mathbf{Q}_0^\top] \right) .\end{aligned}$$

At this point in the proof of Theorem 4.1 for OLS, Lacotte et al. [13] use the asymptotic freeness of $\mathbf{R}_0 \mathbf{R}_0^\top = \mathbf{R}_0^2 = (\mathbf{I}_d - \alpha_0 \mathbf{U}^\top \mathbf{S}_0^\top \mathbf{S}_0 \mathbf{U})^2$ and $\mathbf{R}_t \dots \mathbf{R}_1$ to obtain a closed-form expression for the trace. Asymptotic freeness of these matrices seem to follow from Corollary 4.1 of [5], which says that for $d \times d$ Hermitian random matrices $\{\mathbf{T}_1, \dots, \mathbf{T}_t\}$ with convergent limiting spectral distributions and Haar matrices $\{\mathbf{W}_1, \dots, \mathbf{W}_t\}$ independent of \mathbf{T}_i , the random matrices $\mathbf{W}_1 \mathbf{T}_1 \mathbf{W}_1^\top, \dots, \mathbf{W}_t \mathbf{T}_t \mathbf{W}_t^\top$ are asymptotically free as $n \rightarrow \infty$. Because \mathbf{R}_0^2 and \mathbf{R}_i for $i \neq 0$ are independent and Hermitian, they each have a decomposition of the form $\mathbf{W} \mathbf{T} \mathbf{W}^\top$ where \mathbf{W} is a Haar matrix and \mathbf{T} is diagonal and so the result applies. However, it is less obvious that this result applies in the ridge regression case for

$$\begin{aligned}\mathbf{Q}_0 \mathbf{Q}_0^\top &= \mathbf{I}_d - \alpha_0 (\mathbf{U}^\top \mathbf{S}_0^\top \mathbf{S}_0 \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2}) \\ &\quad - \alpha_0 (\mathbf{I}_d + \lambda \Sigma^{-2}) (\mathbf{U}^\top \mathbf{S}_0^\top \mathbf{S}_0 \mathbf{U} + \lambda \Sigma^{-2})^{-1} \\ &\quad + \alpha_0^2 (\mathbf{U}^\top \mathbf{S}_0^\top \mathbf{S}_0 \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2})^2 (\mathbf{U}^\top \mathbf{S}_0^\top \mathbf{S}_0 \mathbf{U} + \lambda \Sigma^{-2})^{-1}\end{aligned}$$

and $\mathbf{Q}_t \dots \mathbf{Q}_1$ as individual \mathbf{Q}_i and \mathbf{Q}_i^\top are not Hermitian. Due to time constraints on this project, we move forward under the hypothetical assumption that asymptotic freeness is retained in the ridge regression case. *If $\mathbf{Q}_0 \mathbf{Q}_0^\top$ were asymptotically free from $\mathbf{Q}_t \dots \mathbf{Q}_1$* , then from recursive application we obtain

$$\begin{aligned}\lim_{n \rightarrow \infty} \mathbb{E} [\|\Delta_{t+1}\|^2] &= \lim_{n \rightarrow \infty} \frac{1}{d} \text{trace} \left(\mathbb{E} [\mathbf{Q}_2^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{Q}_1^\top] \right) \lim_{n \rightarrow \infty} \frac{1}{d} \text{trace} \left(\mathbb{E} [\mathbf{Q}_0 \mathbf{Q}_0^\top] \right) \\ &= \prod_{i=0}^t \lim_{n \rightarrow \infty} \frac{1}{d} \text{trace} \left(\mathbb{E} [\mathbf{Q}_i \mathbf{Q}_i^\top] \right) \\ &= \prod_{i=0}^t \lim_{n \rightarrow \infty} \frac{1}{d} \text{trace} \left(\mathbb{E} [\mathbf{Q}_i^\top \mathbf{Q}_i] \right) .\end{aligned}$$

Using the expression for $\mathbb{E} [\mathbf{Q}_i^\top \mathbf{Q}_i]$ given earlier, we have

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{1}{d} \text{trace} \left(\mathbb{E} [\mathbf{Q}_i^\top \mathbf{Q}_i] \right) &= 1 - \lim_{n \rightarrow \infty} \frac{2\alpha_i}{d} \text{trace} \left(\mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U} + \lambda \Sigma^{-2})^{-1} \right] (\mathbf{I}_d + \lambda \Sigma^{-2}) \right) \\ &\quad + \lim_{n \rightarrow \infty} \frac{\alpha_i^2}{d} \text{trace} \left(\mathbb{E} \left[(\mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U} + \lambda \Sigma^{-2})^{-2} \right] (\mathbf{I}_d + \lambda \Sigma^{-2})^2 \right) .\end{aligned}$$

Our final obstacle then comes from giving a closed-form expression for the limiting spectral distributions of the above matrices. The proof of Lemma 3.2 (specifically, sub-Lemma A.1) by Lacotte et al. [13] that derives closed-form expressions for the inverse moments of the limiting spectral distribution of $\mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U}$ does not work for the ridge regression case as the proof again relies on embedding $\mathbf{S}_i \mathbf{U}$ into a Haar matrix. A different approach would be

necessary and we leave this for future work due to time constraints on this project.

Assuming that the singular values σ_i of \mathbf{X} along with their inverses σ_i^{-1} are bounded as $n \rightarrow \infty$, the limiting spectral distribution of \mathbf{Q}_i exists when it exists for $\mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U}$. Hypothetically, *if constant closed-form expressions of the limiting traces could be derived*, then the next step would be to minimize $\lim_{n \rightarrow \infty} \mathbb{E} [\|\Delta_{t+1}\|^2] = \rho(\alpha)^{t+1}$ as a function of constant step size α . Note that $\rho(\alpha)$ is quadratic and convex in α . Furthermore, because $\mathbf{Q}_i^\top \mathbf{Q}_i$ is positive semidefinite, its trace is non-negative and therefore $\rho(\alpha) \geq 0$. Thus, minimizing the limiting error norm is equivalent to minimizing $\rho(\alpha)$. Doing so leads to a step size α that is asymptotically optimal (under our assumptions) for a ridge regression solver with partial Newton sketch updates.

3.5.1 Full column rank assumption

We revisit the full rank data matrix assumption (Assumption 1) and consider the implications if it were relaxed. As mentioned, requiring the data matrix to have full column rank is not ideal for ridge regression as one possible reason for using it over OLS is to enforce a unique solution when the data matrix is not full rank. The assumption is made by Lacotte et al. [13] and in our analysis following theirs so that the singular values of \mathbf{X} are positive and therefore Σ is invertible. This allows for factorizing out Σ from the matrix product containing the random matrix, and the remaining product $\mathbf{U}^\top \mathbf{S}_0^\top \mathbf{S}_0 \mathbf{U}$ then has properties that Lacotte et al. [13] exploit. In the ridge regression case, factorizing out Σ does not appear to lead to a particularly convenient form and so we explore the consequences of not doing so.

Using the original identities

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_d &= \mathbf{V} (\Sigma^2 + \lambda \mathbf{I}_d) \mathbf{V}^\top, \\ (\mathbf{X}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{X} + \lambda \mathbf{I}_d)^{-1} &= \mathbf{V} (\Sigma \mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} \Sigma + \lambda \mathbf{I}_d)^{-1} \mathbf{V}^\top, \end{aligned}$$

we obtain an update of the form

$$\mathbf{b}_{t+1} = \mathbf{b}_t - \alpha_t \mathbf{V} (\Sigma \mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} \Sigma + \lambda \mathbf{I}_d)^{-1} (\Sigma^2 + \lambda \mathbf{I}_d) \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) .$$

Notice that multiplying both sides by \mathbf{V}^\top gives

$$\mathbf{V}^\top \mathbf{b}_{t+1} = \mathbf{V}^\top \mathbf{b}_t - \alpha_t (\Sigma \mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} \Sigma + \lambda \mathbf{I}_d)^{-1} (\Sigma^2 + \lambda \mathbf{I}_d) \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*)$$

and then subtracting both sides by $\mathbf{V}^\top \mathbf{b}^*$ gives

$$\begin{aligned} \mathbf{V}^\top (\mathbf{b}_{t+1} - \mathbf{b}^*) &= \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) - \alpha_t (\Sigma \mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} \Sigma + \lambda \mathbf{I}_d)^{-1} (\Sigma^2 + \lambda \mathbf{I}_d) \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) \\ &= \left(\mathbf{I}_d - \alpha_t (\Sigma \mathbf{U}^\top \mathbf{S}_t^\top \mathbf{S}_t \mathbf{U} \Sigma + \lambda \mathbf{I}_d)^{-1} (\Sigma^2 + \lambda \mathbf{I}_d) \right) \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*) . \end{aligned}$$

This form is potentially interesting as we can initiate an analysis similar to the original approach but with the error vector instead defined as $\Delta_t = \mathbf{V}^\top (\mathbf{b}_t - \mathbf{b}^*)$, i.e., as some projection

of the solution error as opposed to some projection of the prediction error. Following the steps in the original approach, under Assumption 2, we again arrive at

$$\mathbb{E} [\|\Delta_{t+1}\|^2] = \frac{1}{d} \text{trace} (\mathbb{E} [\mathbf{Q}_1^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_0 \mathbf{Q}_0^\top])$$

where now $\mathbf{Q}_i = \mathbf{I}_d - \alpha_i (\Sigma \mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U} \Sigma + \lambda \mathbf{I}_d)^{-1} (\Sigma^2 + \lambda \mathbf{I}_d)$. Here, we run into the same challenges with applying results from free probability as in the first approach. While the matrix $\Sigma \mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U} \Sigma$ may appear more complicated to work with than $\mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U}$ in the first approach, its structure has been studied in the free probability literature. For example, Theorem 4.11 of [5] says that for a matrix of the form

$$\mathbf{D}^{\frac{1}{2}} \mathbf{W} \mathbf{T} \mathbf{W}^\top \mathbf{D}^{\frac{1}{2}}$$

where \mathbf{D} and \mathbf{T} are diagonal non-negative matrices and \mathbf{W} is a Haar matrix, the η -transform of the limiting spectral distribution exists and satisfies a certain set of equations. We can obtain the form specified in this result by diagonalizing $\mathbf{S}_i^\top \mathbf{S}_i$. While this particular result does not help with discerning asymptotic freeness, it does show that such matrix structure has been considered in the literature and that it may have further properties that could potentially be useful.

3.6 Empirical experiments

In Section 3.5, we were unable to prove that asymptotic freeness holds for $\mathbf{Q}_1^\top \dots \mathbf{Q}_t^\top \mathbf{Q}_t \dots \mathbf{Q}_1$ where $\mathbf{Q}_i = \mathbf{I}_d - \alpha_i (\mathbf{U}^\top \mathbf{S}_i^\top \mathbf{S}_i \mathbf{U} + \lambda \Sigma^{-2})^{-1} (\mathbf{I}_d + \lambda \Sigma^{-2})$. In this section, we at least empirically investigate the trace decoupling property mentioned in Section 3.2.2 through a simple study where we examine the error

$$\delta(n, \xi) = \frac{1}{d} \text{trace} (\mathbb{E} [\mathbf{Q}_1^\top \mathbf{Q}_2^\top \mathbf{Q}_3^\top \mathbf{Q}_3 \mathbf{Q}_2 \mathbf{Q}_1]) - \frac{1}{d^3} \prod_{i=1}^3 \text{trace} (\mathbb{E} [\mathbf{Q}_i^\top \mathbf{Q}_i])$$

and how it changes as n, d and m grow proportionally according to fixed ratios $\gamma = \frac{d}{n}$ and $\xi = \frac{m}{n}$.

We take $\gamma = 0.05$ and consider two values of $\xi \in \{0.1, 0.2\}$ in our study. We take \mathbf{X} to be a $n \times d$ matrix of i.i.d. standard Gaussians where $n \in \{256, 512, 1024, 2048\}$ (corresponding to $n = 2^p$ for $p \in \{8, \dots, 11\}$) and $d = \lceil \gamma n \rceil$. For $m = \lceil \xi n \rceil$, we take \mathbf{S}_i to be refreshed $m \times n$ sketching matrices of one of the following types:

1. i.i.d. Gaussian: entries are i.i.d. Gaussian random variables with mean 0 and variance m^{-1} .
2. Haar: entries are generated by sampling i.i.d. standard Gaussian random variables and then applying the Gram-Schmidt procedure to obtain a semi-orthogonal matrix.
3. Subsampled randomized Hadamard transform (SRHT) [13]: the matrix \mathbf{S}_i is not formed explicitly. Instead, the left singular vectors \mathbf{U} of \mathbf{X} are transformed through

the map $\mathbf{U} \mapsto \mathbf{B}\mathbf{H}_n\mathbf{D}\mathbf{P}\mathbf{U}$ where \mathbf{P} is a random permutation matrix, \mathbf{D} is a diagonal matrix with Rademacher random variables on the diagonal, \mathbf{H}_n is the $n \times n$ Walsh-Hadamard matrix, and \mathbf{B} is a diagonal matrix with Bernoulli($\frac{m}{n}$) on the diagonal. The rows of the output matrix corresponding to the zero rows of \mathbf{B} are discarded and so the sketch $\mathbf{S}_i\mathbf{U}$ is only approximately a $m \times n$ matrix.

We set $\lambda = 0.1$. As we were also unable to determine the optimal step sizes α_i in Section 3.5, for convenience we use the constant step size

$$\alpha = \frac{(\xi - \gamma)^2}{\gamma^2 + \xi - 2\gamma\xi}$$

which is the optimal (asymptotic) step size for IHS with refreshed orthogonal sketches [13]. This corresponds to step sizes of $\alpha \approx 0.027$ for $\xi = 0.1$ and $\alpha \approx 0.123$ for $\xi = 0.2$.

Figure 1 shows the estimated root mean squared error (RMSE) for each of the three sketch types where the squared error $\delta(n, \xi)^2$ is averaged over 50 independent simulations for each configuration of n and ξ . In all configurations, the RMSE appears to converge to 0 as the dimensions grow which supports the trace decoupling hypothesis. For smaller n , the orthogonal sketches have larger RMSE compared to the i.i.d. Gaussian sketches, though their RMSE appear to decrease more rapidly on average than that of the i.i.d. Gaussian sketches. An interesting note is that the sketch ratio ξ does not seem affect the RMSE for i.i.d. Gaussian sketches, though a larger sketch size does significantly decrease the RMSE for orthogonal sketches.

3.7 Discussion

In this project, we showed that extending the analysis of Lacotte et al. [13] for OLS with IHS to ridge regression with partial Newton sketch is not trivial. The main challenges in attempting to do so come from not being able to directly apply existing results from free probability and random matrix theory to the matrices obtained from the ridge regression updates that are analogous to the ones in the OLS updates. Under certain assumptions and a number of hypothetical conditions, we outlined a procedure to derive the optimal convergence rate and corresponding step size for sketched ridge regression. We examined the full column rank data matrix assumption and showed that relaxing this assumption may lead to a similar analysis but with a different criterion (in terms of the solution error). Through simulations, we empirically showed that the trace decoupling condition hypothesized in the theory appears to hold in at least one simple experimental setting. Though the simulation does not explore all parameters of potential interest (e.g., other values of γ , λ , etc.), the results do instill some optimism that the theory may hold.

One thing to note is that the partial Newton sketch algorithm we considered in this project uses refreshed sketches, i.e., sketching matrices that are resampled every iteration. Results addressing the asymptotic freeness of matrices would not apply to fixed sketches due to the matrices no longer being independent, and thus an analysis based on free probability would likely be unsuitable for the fixed sketch case. Empirical results from Chowdhury et al.

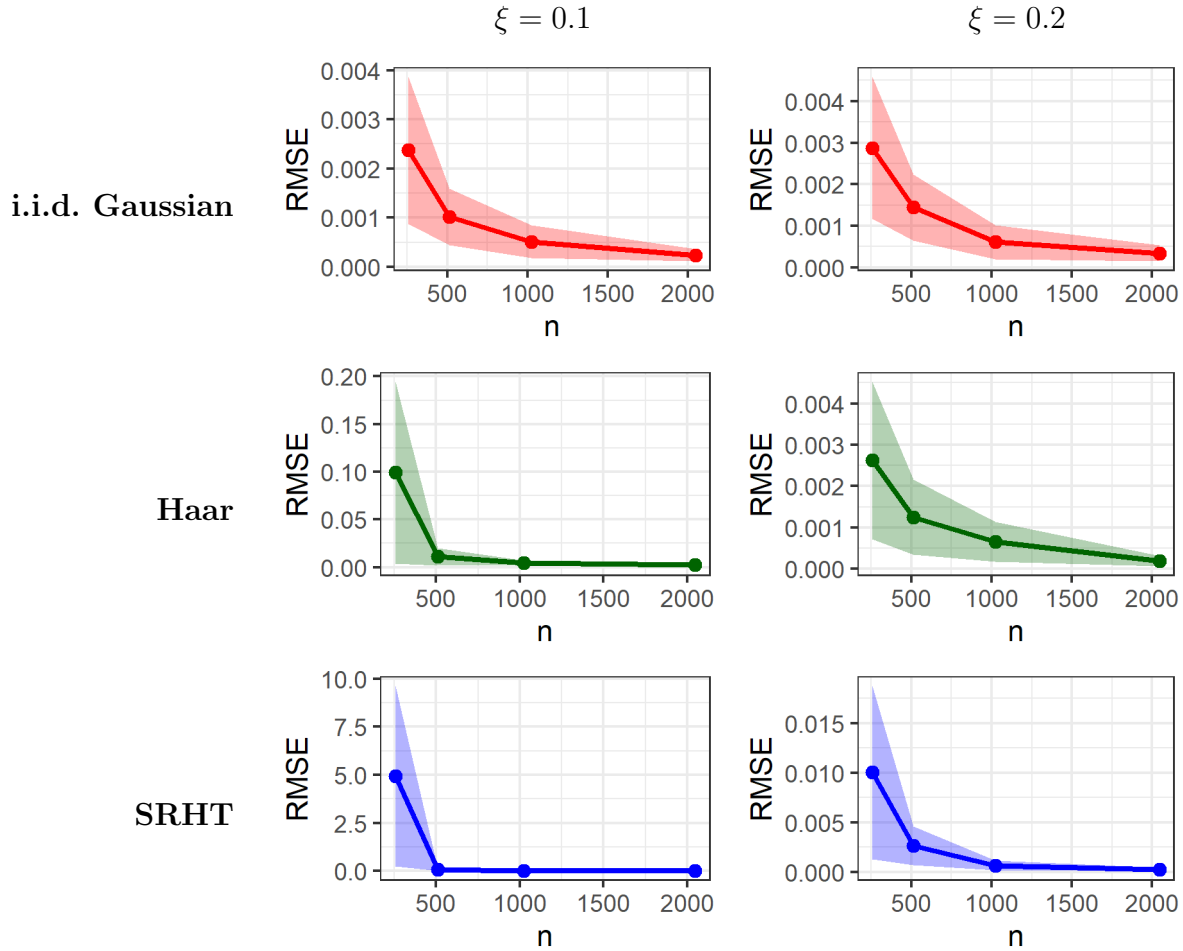


Figure 1: The estimated RMSE and standard error for $\xi \in \{0.1, 0.2\}$ and i.i.d. Gaussian, Haar and SRHT sketches. The mean and standard error at each point is calculated over 50 independent simulations.

[4] do suggest that refreshed sketches enable faster convergence than fixed sketches for ridge regression (in at least one context), and so the fixed sketch case may not be particularly interesting anyways.

Due to the time constraints on this project, there is much work left incomplete and a number of possible directions for future work. One direction is to complete the analysis and verify the asymptotic freeness of the matrices in the proof, as well as to determine whether closed-form expressions of the inverse moments are possible. It may also be of interest to explore whether this analysis approach can be used for other sketch-based methods, though our impression is that given the current state of the free probability and random matrix theory literature, there will only be a narrow set of methods and problems for which such an analysis may be applied. This seems to be the case as to our knowledge, there are only a handful of works [8, 12, 13] in the sketching literature that have used this analysis approach. However, free probability and random matrix theory seems to be an active area of research

and so new developments may change this outlook.

4 Comment

I would like to briefly comment on the iterative Hessian sketch/Newton sketch literature. While the body of works that focus on these algorithms is fairly substantial with many new developments in recent years, there appears to be only one or two research groups driving this portion of the literature and reinforcing their position with many self-citations. Many of these citations are to unpublished papers or concurrent papers. To make matters more confusing, there are instances of cyclical citations where one paper references an old version of another paper, and the new version of the other paper references the paper that references its old version (e.g., [13] and [12]). This made tracking the lineage of the developments in this literature a bit of a nightmare.

References

- [1] Greg W. Anderson, Alice Guionnet, and Ofer Zeitouni. *Free probability*, page 322–413. Cambridge Studies in Advanced Mathematics. Cambridge University Press, 2009. doi: 10.1017/CBO9780511801334.006.
- [2] Colin Chen and Ying Wei. Computational issues for quantile regression. *Sankhyā: The Indian Journal of Statistics*, pages 399–417, 2005.
- [3] Krzysztof Choromanski, Mark Rowland, and Adrian Weller. The unreasonable effectiveness of structured random orthogonal embeddings. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS’17, page 218–227, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.
- [4] Agniva Chowdhury, Jiasen Yang, and Petros Drineas. An iterative, sketching-based framework for ridge regression. In *International Conference on Machine Learning*, pages 989–998. PMLR, 2018.
- [5] Romain Couillet and Mérouane Debbah. *Free probability theory*, page 71–94. Cambridge University Press, 2011. doi: 10.1017/CBO9780511994746.006.
- [6] Michał Dereziński, Burak Bartan, Mert Pilanci, and Michael W Mahoney. Debiasing distributed second order optimization with surrogate sketching and scaled regularization. *arXiv preprint arXiv:2007.01327*, 2020.
- [7] Michał Dereziński, Jonathan Lacotte, Mert Pilanci, and Michael W Mahoney. Newton-LESS: Sparsification without trade-offs for the sketched Newton update. *Advances in Neural Information Processing Systems*, 34, 2021.
- [8] Edgar Dobriban and Sifan Liu. Asymptotics for sketching in least squares. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, Red Hook, NY, USA, 2019. Curran Associates Inc.
- [9] Petros Drineas, Michael W Mahoney, Shan Muthukrishnan, and Tamás Sarlós. Faster least squares approximation. *Numerische mathematik*, 117(2):219–249, 2011.
- [10] Xuming He, Xiaou Pan, Kean Ming Tan, and Wen-Xin Zhou. Smoothed quantile regression with large-scale inference. *Journal of Econometrics*, 2021.
- [11] Roger Koenker and Gilbert Bassett Jr. Regression quantiles. *Econometrica: journal of the Econometric Society*, pages 33–50, 1978.
- [12] Jonathan Lacotte and Mert Pilanci. Optimal randomized first-order methods for least-squares problems. In *International Conference on Machine Learning*, pages 5587–5597. PMLR, 2020.
- [13] Jonathan Lacotte, Sifan Liu, Edgar Dobriban, and Mert Pilanci. Optimal iterative sketching methods with the subsampled randomized Hadamard transform. In

- H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 9725–9735. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/6e69ebbfad976d4637bb4b39de261bf7-Paper.pdf>.
- [14] Jonathan Lacotte, Yifei Wang, and Mert Pilanci. Adaptive Newton sketch: Linear-time optimization with quadratic convergence and effective Hessian dimensionality. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 5926–5936. PMLR, 2021. URL <http://proceedings.mlr.press/v139/lacotte21a.html>.
- [15] Ibrahim Kurban Ozaslan, Mert Pilanci, and Orhan Arikan. Iterative Hessian sketch with momentum. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7470–7474, 2019. doi: 10.1109/ICASSP.2019.8682720.
- [16] Vu Pham and Laurent El Ghaoui. Robust sketching for multiple square-root LASSO problems. In *Artificial Intelligence and Statistics*, pages 753–761. PMLR, 2015.
- [17] Mert Pilanci and Martin J Wainwright. Iterative Hessian sketch: Fast and accurate solution approximation for constrained least-squares. *The Journal of Machine Learning Research*, 17(1):1842–1879, 2016.
- [18] Mert Pilanci and Martin J Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. *SIAM Journal on Optimization*, 27(1): 205–245, 2017.
- [19] Stephen Portnoy and Roger Koenker. The Gaussian hare and the Laplacian tortoise: computability of squared-error versus absolute-error estimators. *Statistical Science*, 12(4):279–300, 1997.
- [20] Garvesh Raskutti and Michael W. Mahoney. A statistical perspective on randomized sketching for ordinary least-squares. *Journal of Machine Learning Research*, 17(213): 1–31, 2016. URL <http://jmlr.org/papers/v17/15-440.html>.
- [21] Robert N Rodriguez and Yonggang Yao. Five things you should know about quantile regression. In *Proceedings of the SAS global forum 2017 conference, Orlando*, pages 2–5, 2017.
- [22] Vladimir Rokhlin and Mark Tygert. A fast randomized algorithm for overdetermined linear least-squares regression. *Proceedings of the National Academy of Sciences*, 105(36):13212–13217, 2008.
- [23] Shusen Wang, Alex Gittens, and Michael W Mahoney. Sketched ridge regression: Optimization perspective, statistical perspective, and model averaging. In *International Conference on Machine Learning*, pages 3608–3616. PMLR, 2017.

-
- [24] Jiyan Yang, Xiangrui Meng, and Michael Mahoney. Quantile regression for large-scale applications. In *International Conference on Machine Learning*, pages 881–887. PMLR, 2013.
 - [25] Aijun Zhang, Hengtao Zhang, and Guosheng Yin. Adaptive iterative Hessian sketch via A -optimal subsampling. *Statistics and Computing*, 30(4):1075–1090, jul 2020. ISSN 0960-3174. doi: 10.1007/s11222-020-09936-8. URL <https://doi.org/10.1007/s11222-020-09936-8>.
 - [26] Pan Zhao and Shenghua Yu. An improved interior point algorithm for quantile regression. *IEEE Access*, 8:139647–139657, 2020.
 - [27] Hui Zou and Ming Yuan. Composite quantile regression and the oracle model selection theory. *The Annals of Statistics*, 36(3):1108–1126, 2008.