

On Random Embeddings and Their Application to Optimisation



Zhen Shao
St Anne's College
University of Oxford

A DPhil thesis submitted for the degree of
Doctor of Philosophy

Michaelmas 2021

Acknowledgements

I would like to thank my supervisor, Prof Coralia Cartis, for her patience, support and teaching over the last four years, without which this thesis would not be possible.

I would also like to thank Chris Breward and Colin Please for making Industrially Focused Mathematical Modelling CDT possible, and the collaborations I had with Numerical Algorithm Group Ltd., in particular Dr Jan Fiala.

Throughout my DPhil, I have been generously supported by my office mates, my friends in the Mathematical Institute and St Anne's College. Thank you all for being with me in this journey.

I am also grateful for all the teachings and support I received during my undergraduate years at Oxford. In particular, my tutors at Pembroke College.

Finally, I would like to thank my parents, for raising me up and giving me the best environment for my education.

Abstract

Random embeddings project high-dimensional spaces to low-dimensional ones; they are careful constructions which allow the approximate preservation of key properties, such as the pair-wise distances between points. Often in the field of optimisation, one needs to explore high-dimensional spaces representing the problem data or its parameters and thus the computational cost of solving an optimisation problem is connected to the size of the data/variables. This thesis studies the theoretical properties of norm-preserving random embeddings, and their application to several classes of optimisation problems.

Our investigations into random projections present subspace embedding properties for s -hashing ensembles — sparse random matrices with s non-zero entries per column — that are optimal in the projection dimension m of the sketch, namely, $m = \mathcal{O}(d)$ where d is the dimension of the subspace. A diverse set of results are presented that address the case when the input matrix has sufficiently low coherence; how the acceptable coherence changes with the number s of non-zeros per column in the s -hashing matrices, or is reduced through suitable transformations. In particular, we propose a new random embedding, the Hashed Randomised Hadamard Transform, that improves upon the Subsampled Randomised Hadamard Transform by replacing sub-sampling with hashing.

We apply these sketching techniques to linear least squares problems, as part of a Blendenpik-type algorithm, that uses a sketch of the data matrix to build a high quality preconditioner and then solves a preconditioned formulation of the original problem. We also include suitable linear algebra tools for rank-deficient and for sparse problems that lead to our implementation, Ski-LLS, outperforming not only sketching-based routines on randomly-generated input, but also state of the art direct solver SPQR and iterative code HSL on certain subsets of the sparse Florida matrix collection; namely, on least squares problems that are significantly over-determined, or moderately sparse, or difficult.

Instead of sketching in the data/observational space as in the linear least squares case above, we then consider sketching in the variable/parameter domain for a more generic problem and algorithm. We propose a general random-subspace first-order framework for unconstrained non-convex optimisation that requires a weak probabilistic assumption on the subspace gradient, which we show to be satisfied by various random matrix

ensembles, such as Gaussian and hashing sketching. We show that, when safeguarded with trust region or quadratic regularisation techniques, this random subspace approach satisfies, with high probability, a complexity bound of order $\mathcal{O}(\epsilon^{-2})$ to drive the (full) gradient norm below ϵ ; matching in the accuracy order, deterministic counterparts of these methods and securing almost sure convergence. We then particularise this framework to random subspace Gauss-Newton methods for nonlinear least squares problems, that only require the calculation of the Jacobian matrix in a subspace, with similar complexity guarantees.

We further investigate second-order methods for non-convex optimisation, and propose a Random Subspace Adaptive Regularised Cubic (R-ARC) method, which we analyse under various assumptions on the objective function and the sketching matrices. We show that, when the sketching matrix achieves a subspace embedding of the augmented matrix of the gradient and the Hessian with sufficiently high probability, then the R-ARC method satisfies, with high probability, a complexity bound of order $\mathcal{O}(\epsilon^{-3/2})$ to drive the (full) gradient norm below ϵ ; matching in the accuracy order the deterministic counterpart (ARC). We also show that the same complexity bound is obtained when the Hessian matrix has sparse rows and appropriate sketching matrices are chosen. We also investigate R-ARC's convergence to second order critical points. We show that the R-ARC method also drives the Hessian in the subspace to be approximately positive semi-definite with high probability, for a variety of sketching matrices; and furthermore if the Hessian matrix has low rank and scaled Gaussian sketching matrices are used, the R-ARC drives the (full) Hessian to be approximately positive semi-definite, with high probability, at the rate $\mathcal{O}(\epsilon^{-3})$, again matching in the accuracy order its deterministic counterpart.

Contents

1	Introduction	1
1.1	Background	1
1.2	Random embedding	2
1.3	Linear least squares	7
1.3.1	Dense linear least squares	7
1.3.2	Sparse linear least squares	9
1.3.3	Sketching for linear least squares	9
1.3.4	Alternative approaches for preconditioning and solving large-scale LLS problems	10
1.4	Minimising a general unconstrained objective function	11
1.4.1	First order methods	12
1.4.2	Second order methods	14
1.4.3	Applications to non-linear least squares problems	16
1.5	Contributions of this thesis	17
1.5.1	New theoretical analysis of hashing matrices and development of new random embeddings	17
1.5.2	Analysis, state-of-the-art implementation and benchmarking of new large-scale linear least squares solver using random embeddings	18
1.5.3	First-order subspace methods and their application to non-linear least squares	19
1.5.4	Random subspace cubic regularisation algorithm, R-ARC	20
1.6	Structure of thesis	20
2	Random embedding	22
2.1	Introduction and relevant literature	22
2.2	Technical Background	24
2.2.1	Random embeddings	24
2.2.2	Generic properties of subspace embeddings	25
2.2.3	Sparse matrix distributions and their embeddings properties	27
2.2.3.1	Coherence-dependent embedding properties of sparse random matrices	28

2.2.3.2	Non-uniformity of vectors and their relation to embedding properties of sparse random matrices	30
2.3	Hashing sketching with $m = \mathcal{O}(r)$	31
2.4	Relaxing the coherence requirement using s -hashing matrices	37
2.4.1	The embedding properties of s -hashing matrices	37
2.4.2	A general embedding property for an s -hashing variant	38
2.4.3	The Hashed-Randomised-Hadamard-Transform sketching	41
3	Sketching for linear least squares	43
3.1	Introduction and relevant literature	43
3.2	Algorithmic framework and analysis	44
3.2.1	A generic algorithmic framework for solving linear least squares with sketching	45
3.2.2	Analysis of Algorithm 1	46
3.3	Implementation details	50
3.3.1	Ski-LLS, an implementation of Algorithm 1	50
3.3.2	Discussion of our implementation	52
3.4	Numerical study	53
3.4.1	Test Set	53
3.4.2	Numerical illustrations	55
3.4.3	Compilation and running environment for timed experiments	56
3.4.4	Tuning to set the default parameters	56
3.4.5	Residual accuracy of Ski-LLS	60
3.5	Numerical performance	61
3.5.1	Solvers compared and their parameters	61
3.5.2	Running time performance on randomly generated full rank dense A	62
3.5.3	Running time performance on randomly generated full rank sparse A	64
3.5.4	Large scale benchmark of Ski-LLS-sparse on the Florida Matrix Collection . .	64
4	First order subspace method for general objectives	68
4.1	Introduction	68
4.2	General algorithmic framework and its convergence result	69
4.2.1	Generic algorithmic framework and assumptions	69
4.2.2	A probabilistic convergence result	72
4.2.3	Corollaries of Theorem 4.2.1	73
4.3	Proof of Theorem 4.2.1	75
4.4	An algorithmic framework based on sketching	83
4.4.1	A generic random subspace method based on sketching	83

4.4.2	The random matrix distribution \mathcal{S} in Algorithm 3	86
4.4.2.1	Gaussian sketching matrices	86
4.4.2.2	s -hashing matrices	88
4.4.2.3	(Stable) 1-hashing matrices	88
4.4.2.4	Sampling matrices	90
4.5	Random subspace quadratic regularisation and subspace trust region methods	91
4.5.1	Random subspace quadratic regularisation with sketching	93
4.5.2	Iteration complexity of random subspace quadratic regularisation methods . .	95
4.5.2.1	Using scaled Gaussian matrices	95
4.5.2.2	Using stable 1-hashing matrices	96
4.5.2.3	Using sampling matrices	97
4.5.3	Random subspace trust region methods with sketching	98
4.5.4	Iteration complexity of random subspace trust region methods	99
4.5.4.1	Using scaled Gaussian matrices	100
4.5.4.2	Using stable 1-hashing matrices	101
4.5.4.3	Using sampling matrices	102
4.6	Randomised Subspace Gauss-Newton (R-SGN) for non-linear least squares	103
5	Second order subspace methods for general objectives	109
5.1	Introduction	109
5.2	R-ARC: random subspace adaptive cubic regularisation method	109
5.3	Fast convergence rate assuming subspace embedding of the Hessian matrix	111
5.3.1	Auxiliary results	112
5.3.2	Satisfying the assumptions of Theorem 4.2.1	113
5.3.3	Iteration complexity of Algorithm 6 to decrease $\nabla f(x_k)$ below ϵ	115
5.3.3.1	Discussion	116
5.4	Fast convergence rate assuming the sparsity of the Hessian matrix	117
5.5	Convergence to second order (subspace) critical points	119
5.6	Convergence to second order (full space) critical points	121
6	Conclusion and future directions	126
	Bibliography	127

List of Figures

1.1	Randomly sampling and then re-scaling gives a good estimate of the norm when the vector components have similar magnitude.	3
1.2	(Randomised) Fourier transform makes the magnitude of the entries of a vector more similar	4
3.1	Hashing combined with a randomised Discrete Hartley Transform produces more accurate sketched matrix SA for a given m/d ratio comparing to sampling combined with a randomised Discrete Hartley Transform; the accuracy of the sketch is reflected in the quality of the preconditioner R constructed from the matrix SA , see (3.2.15).	55
3.2	When the data matrix A is an ill-conditioned sparse Gaussian matrix, using 1, 2, 3–hashing produces similarly good preconditioners.	56
3.3	When the data matrix A has higher coherence, using s –hashing with $s > 1$ is crucial to produce an acceptable preconditioner.	56
3.4	Runtime of Ski-LLS-dense on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, Ski-LLS-dense is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. We choose $\gamma = 1.7$ to approximately minimize the runtime across the above plots.	57
3.5	Runtime of Ski-LLS-dense without R-CPQR on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, Ski-LLS-dense without R-CPQR is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. Note that using LAPACK QR instead of R-CPQR results in slightly shorter running time (c.f. Figure 3.4). We choose $\gamma = 1.7$ to approximately minimize the runtime across the above plots.	57

3.6	Runtime of Blendenpik on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, Blendenpik is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. Note that Blendenpik handles coherent dense A significantly less well than our dense solvers. We choose $\gamma = 2.2$ to approximately minimize the runtime across the above plots. .	58
3.7	Runtime of LSRN on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, LSRN is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. Note that LSRN runs more than 5 times slower comparing to Blendenpik or Ski-LLS in the serial testing environment, due to the use of SVD and Gaussian sketching. We choose $\gamma = 1.1$ to approximately minimize the runtime across the above plots.	58
3.8	Running time of Ski-LLS-sparse on sparse matrices $A \in \mathbb{R}^{m \times d}$ from Test Set 2 with $n = 80000, d = 4000$ and $n = 120000, d = 3000$ using different values of s and $\gamma = m/d$. We choose $m = 1.4d, s = 2$ in consideration of the above plot and the residual accuracy in Figure 3.9 but also taking into account some experiments of Ski-LLS-sparse we have done on the Florida matrix collection.	59
3.9	Corresponding residual on sparse matrices $A \in \mathbb{R}^{m \times d}$ from Test Set 2 with $n = 80000, d = 4000$ and $n = 120000, d = 3000$ using different values of s and $\gamma = m/d$. Note that using 1-hashing ($s = 1$) results in inaccurate solutions.	59
3.10	Runtime for LSRN on sparse matrices $A \in \mathbb{R}^{m \times d}$ from Test Set 2 with $n = 80000, d = 4000$ and $n = 120000, d = 3000$ using different values of $\gamma = m/d$. We choose $m = 1.1d$ in consideration of the above plot and the residual accuracy in Figure 3.11 but also taking into account some experiments of Ski-LLS-sparse we have done on the Florida matrix collection.	60
3.11	Residual values obtained by LSRN on the same sparse problems as in Figure 3.10. .	60
3.12	Time taken by solvers to compute the solution of problem (1.3.1) for A being coherent dense matrices of various sizes (x-axis)	63
3.13	Time taken by solvers to compute the solution of problem (1.3.1) for A being semi-coherent dense matrices of various sizes (x-axis)	63
3.14	Time taken by solvers to compute the solution of problem (1.3.1) for A being incoherent dense matrices of various sizes (x-axis)	63
3.15	Running time comparison of Ski-LLS with LS_HSL and LS_SPQR for randomly generated incoherent sparse matrices of different sizes.	65

3.16	Running time comparison of Ski-LLS with LS_HSL and LS_SPQR for randomly generated semi-coherent sparse matrices of different sizes.	65
3.17	Running time comparison of Ski-LLS with LS_HSL and LS_SPQR for randomly generated coherent sparse matrices of different sizes.	65
3.18	Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with $n \geq 30d$	66
3.19	Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with $n \geq 10d$	66
3.20	Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with with $n \geq 10d$ and the unpreconditioned LSQR takes more than 5 seconds to solve.	67
3.21	Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with $n \geq 10d$ and $\text{nnz}(A) \geq 0.01nd$	67
4.1	ARTIF (left), BRATU2D (middle) and OSCIGRNE (right) objective value against cumulative Jacobian action size for R-SGN with coordinate sampling.	106
4.2	ARTIF (left), BRATU2D (middle) and OSCIGRNE (right) objective value against cumulative Jacobian action size for R-SGN with Gaussian sketching.	106
4.3	ARTIF (left), BRATU2D (middle) and OSCIGRNE (right) objective value against cumulative Jacobian action size for R-SGN with 3-hashing sketching.	106
4.4	R-SGN on the CHEMOTHERAPY dataset	107
4.5	R-SGN on the GISETTE dataset	107

Chapter 1

Introduction

1.1 Background

This thesis is about random embeddings and their application to improving the efficiency of optimisation algorithms for different problem classes. In particular, regarding random embeddings, the novelty of our work is in the analysis of sparse projections and in proposing a new general random embedding with attractive theoretical properties and numerical performance. Then we transfer these results and existing random embeddings to improve optimisation algorithms for linear and non-linear least squares problems, as well as for general objectives using first and second order information.

Numerical optimisation designs algorithms that find an extreme value of a given function. Such computational routines find numerous applications in data science, finance and machine learning. The computational cost of an optimisation algorithm typically grows with the dimension of the function being optimised, which in many applications increases as the data set becomes larger or the model for the data becomes more complex. For example, the classical computation of the solution of fitting a linear model to a data set (linear least squares) grows linearly with the size of the data set and quadratically with the number of variables in the model. Given the ever-increasing amount of data and complexity of models, recent research trends attempt to make classical optimisation algorithms faster and more scalable, [32, 78, 5, 17, 73, 74]. This work explores two topics in this context: algorithms for linear least squares that compute an accurate solution (up to the machine precision), and with computational complexities lower than classical methods; and algorithms for general unconstrained objective functions that compute an approximate solution with first and second order guarantees of optimality, with probability arbitrarily close to one and matching, in the order of desired accuracy of the solution, the complexity of classical optimization methods.

We begin with a simple example of how random embeddings can help solve linear least squares

$$\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2 \tag{1.1.1}$$

faster. Consider the problem $\min_x f(x) = (x-6)^2 + (2x-5)^2 + (3x-7)^2 + (4x-10)^2$, corresponding to $A = \begin{pmatrix} 1 & 2 & 3 & 4 \end{pmatrix}^T$ and $b = \begin{pmatrix} 6 & 5 & 7 & 10 \end{pmatrix}^T$. Solving $f'(x) = 0$ or equivalently $A^T A x = A^T b$, we obtain $x = \frac{77}{30} \approx 2.567$. Sketching with random embedding S transforms the problem (1.1.1) to

$$\min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2^2, \quad (1.1.2)$$

where $S \in \mathbb{R}^{m \times n}$ is some matrix we choose. We give two examples for S .

Example 1 (Sampling). Let¹ $S = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}$. Then $SA = \begin{pmatrix} 1 & 3 \end{pmatrix}^T$ gives the 1st and 3rd row of the matrix A . $Sb = \begin{pmatrix} 6 & 7 \end{pmatrix}^T$ gives the 1st and 3rd entry of the vector b . Solving (1.1.2) gives us $x = \frac{81}{30} = 2.700$.

Example 2 (Hashing). Let² $S = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix}$. Then $SA = \begin{pmatrix} 3 & 7 \end{pmatrix}^T$ where the 1st row of SA is the sum of the 1st and 2nd rows of A ; the 2nd row of SA is the sum of the 3rd and 4th rows of A . $Sb = \begin{pmatrix} 11 & 17 \end{pmatrix}^T$ where the 1st entry of Sb is the sum of the 1st and 2nd entries of b ; the 2nd entry of Sb is the sum of the 3rd and 4th entries of b . Solving (1.1.2) gives us $x = \frac{152}{58} \approx 2.621$.

In both examples we reduce the number of rows of A, b from four to two, but using hashing gives a more accurate result because it uses each row of A and entry in b . In later sections of this thesis we show that computing the solution of problem (1.3.1) with hashing sketching leads to improved performance. We also show how to use random embeddings to compute an accurate instead of just an approximate solution of linear least squares.

For the remainder of this chapter, we first review key concepts about random embeddings, and compare and contrast some well-known random embeddings. Then we introduce the problem of linear least squares, classical techniques for its solution, and random embedding-based approaches known as sketching. We then introduce general non-convex optimisation problems; classical first and second order methods to solve them; and existing theoretical results on these ‘full space’ methods. We also introduce non-linear least squares as a particular type of non-convex optimisation problems. This chapter ends with a summary of the structure and contributions of this thesis to random embeddings, their applications to linear least squares, and to general non-convex optimisations. Our detailed contributions and relevant literature reviews can be found in individual chapters.

1.2 Random embedding

JL Lemma Dimensionality reduction techniques using random embeddings rely crucially on the Johnson-Lindenstrauss lemma which first appeared in 1984 [56]. It states that to calculate the approximate 2-norm³ of a set of vectors, it suffices to randomly project them to lower dimensions,

¹Namely, S has one non-zero per row.

²Namely, S has one non-zero per column.

³By 2-norm of a vector, we mean its usual Euclidean norm.

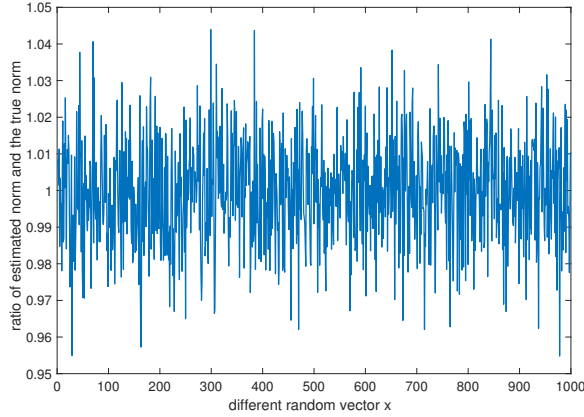


Figure 1.1: Randomly sampling and then re-scaling gives a good estimate of the norm when the vector components have similar magnitude.

calculate their length in the projected space. This is equivalent to multiplying the vectors representing the high-dimensional points (on the left) by an under-determined matrix with entries following some probabilistic distributions. We call such a matrix, a random embedding or projection. In particular, random embeddings for a set of points that approximately preserve their 2-norms are called Johnson-Lindenstrauss (JL)-embeddings (formally defined in Definition 2.2.2). More specifically, we may choose scaled Gaussian matrices as the embedding ⁴.

Lemma 1.2.1 (JL Lemma [56, 25]). *Given a fixed, finite set $Y \subseteq \mathbb{R}^n$, $\epsilon, \delta > 0$, let $S \in \mathbb{R}^{m \times n}$ have entries independently distributed as the normal $N(0, n^{-1})$, with $m = \mathcal{O}\left(\epsilon^{-2} \log\left(\frac{|Y|}{\delta}\right)\right)$ and where $|Y|$ refers to the cardinality of the set Y . Then we have, with probability at least $1 - \delta$, that*

$$(1 - \epsilon)\|y\|_2^2 \leq \|Sy\|_2^2 \leq (1 + \epsilon)\|y\|_2^2 \quad \text{for all } y \in Y. \quad (1.2.1)$$

Intuitively, we are able to use the above dimensionality reduction technique because we are only concerned with Euclidean distances, expressed as sum of squares. If a vector x has n entries with similar magnitudes, to calculate its 2-norm, we only need to sample some of its entries, say m entries, then calculate the sum of squares of those entries, and rescale by n/m to obtain the approximate 2-norm of x . This is illustrated in Figure 1.1, where we set x to be a random Gaussian vector with independent identically distributed entries. We see that the error in the norm estimation is within 5%.

In general, the magnitudes of the entries are dissimilar. However, we can preprocess x by applying a random, norm-preserving transformation, before sampling and re-scaling. In Figure 1.2, we apply a (randomised) Fourier transform to a vector x with a non-uniform distribution of the magnitude of entries. We observe that the square of the entries of x are more uniformly distributed after the transform.

⁴In the original paper [56], the lemma appears as an existence result concerning Lipschitz mappings. Here we state the ‘modern’ form that is proved in [25] and that is more relevant to this thesis.

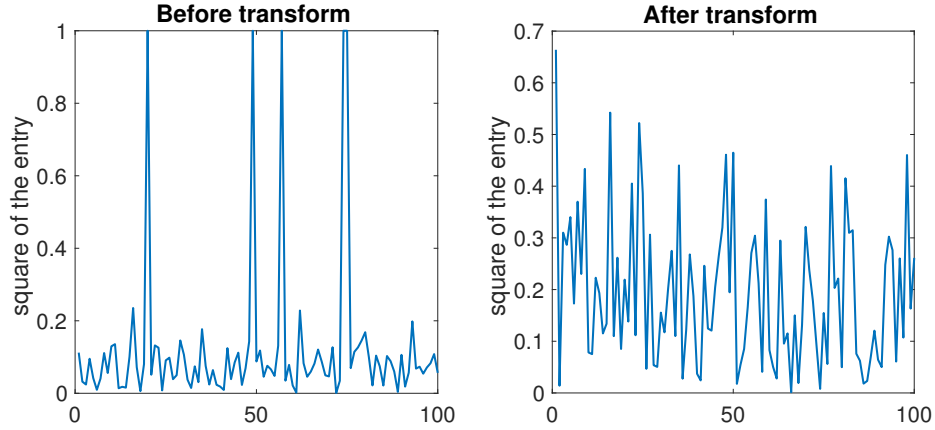


Figure 1.2: (Randomised) Fourier transform makes the magnitude of the entries of a vector more similar

Multiplying by a square Gaussian matrix has a similar effect in making the magnitude of the entries of a vector more similar. While multiplying by an under-determined Gaussian matrix is a composition of multiplying by a square Gaussian matrix and then an under-determined sampling matrix (with one non-zero entry per row in a random column, whose value is one). This is the intuition behind the JL Lemma. For more details on random matrices, see [100].

Subspace embedding Instead of a discrete group of points, Subspace embeddings (formally defined in Definition 2.2.3) also aim to approximately preserve the 2-norm of each point in a column subspace of some given matrix A . Subspace embeddings are useful when the point whose 2-norm is to be approximately preserved is unknown but lies in a given subspace; such as in the application of using random embeddings to solve linear least squares faster, where the optimal $Ax^* - b$ is unknown but lies in the subspace generated by the columns of A and the vector b . Subspace embeddings also find applications in computing a high quality preconditioner of a linear system, and solving the low-rank matrix approximation problem [74]. Often, a random matrix distribution can be both an (oblivious)⁵ JL-embedding and an (oblivious) subspace embedding, see [101] where the author derives the oblivious subspace embedding property of the scaled Gaussian matrices from its oblivious JL-embedding property.

Oblivious subspace embedding A crucial advantage of random embeddings comparing to deterministic ones is that their embedding properties are data independent. For example, it is well known that the singular value decomposition (SVD) gives the most efficient low-dimensional embedding of a column subspace (Eckart–Young theorem). However, for each given matrix, its SVD needs to be computed before the embedding can be applied; which is computationally expensive and the cost scales with the data size. By contrast, random embeddings are independent of the data matrix and hence no data-specific computation is required (aside from constructing the random embedding

⁵Data independent, see Definition 2.2.4.

by sampling from the given distribution and applying the random embedding to the data matrix). Therefore, due to this property, random embeddings are oblivious embeddings (formally defined in Definition 2.2.4). A consequence of the embedding being oblivious to the data is that there is, in general, a positive probability that the randomly drawn embedding fails to embed the data (in the sense of providing a JL-embedding or a subspace-embedding). However the failure probability is exponentially small and can be bounded above by appropriately setting the dimension of the embedded space. Moreover, the iterative nature of our subspace algorithms in later chapters takes into account that in some iterations, the embedding may fail. But the probability that those algorithms fail to converge at the expected theoretical rate approaches zero exponentially fast with the total number of iterations.

Next, we briefly review a list of commonly used random embeddings, which are represented by random matrices.

Popular random matrices and their properties Sampling matrices have one non-zero entry per row in a random column.

Definition 1.2.1. We define $S \in \mathbb{R}^{m \times n}$ to be a scaled sampling matrix if, independently for each $i \in [m]$, we sample $j \in [n]$ uniformly at random and let $S_{ij} = \sqrt{\frac{n}{m}}$.

The scaling factor is included so that given $x \in \mathbb{R}^n$, we have $\mathbb{E}[\|Sx\|_2] = \|x\|_2$ for any scaled sampling matrix S . Sampling matrices are computationally inexpensive to apply to vectors/matrices so that embeddings based on them can be computed efficiently. However, the success of sampling matrices is highly dependent on the data. Even if we have $\mathbb{E}[\|Sx\|_2] = \|x\|_2$, $\|Sx\|_2$ may have high variance, such as when x has a single non-zero entry in its first row.

Non-uniformity of a vector, formally defined in Definition 2.2.6, provides a measure of how different the magnitudes of the entries are; and the success of sampling matrices as an oblivious JL embedding depends on this. Similarly, the success of the sampling matrices as an oblivious subspace embedding depends on the coherence of a matrix (formally defined in Definition 2.2.5), which provides a measure of the non-uniformity of vectors in the matrix column subspace (Lemma 2.2.4).

There are broadly two types of approaches to tackle the high variance challenge of using sampling matrices. The first type is based on transforming the vector/matrix to one with the same norm/column subspace but with higher uniformity. For example, it is well known that for any fixed vector $x \in \mathbb{R}^n$, pre-multiplication by a square Gaussian matrix (with each entry following $N(0, n^{-1})$) transforms the vector into one with independent normally distributed entries while preserving $\|x\|_2$ in expectation. In high dimensions, the resulting vector has high uniformity (due to entries having the same distribution and the scaling factor) and is thus suitable for applying sampling. A scaled Gaussian matrix can be thought as the product of a scaled sampling matrix (with the scaling being $\sqrt{\frac{n}{m}}$) and a square Gaussian matrix (with each entry following $N(0, n^{-1})$).

Definition 1.2.2. We say $S \in \mathbb{R}^{m \times n}$ is a scaled Gaussian matrix if S_{ij} are independently distributed as $N(0, m^{-1})$.

(Scaled) Gaussian matrices with appropriate dimensions have been shown to be an oblivious JL/subspace embeddings [101]. However, Gaussian matrices are computationally expensive to apply, especially when embedding a linear subspace represented by a dense basis due to the cost of dense matrix-matrix multiplication.

Subsampled-Randomised-Hadamard-Transform (SRHT) [2, 97] uses an alternative non-uniformity reduction technique based on the Hadamard transform, which is similar to the Fourier transform. A Fast-Fourier-type algorithm exists that allows applying SRHT in $O(n \log(n))$ time for $x \in \mathbb{R}^n$ [2], thus having a better complexity than the naive matrix-matrix multiplication, while still achieving comparable theoretical properties as the scaled Gaussian matrix. For subspace embedding of matrices in $\mathbb{R}^{n \times d}$, the embedding dimension of SRHT has a $\log(d)$ multiplicative factor compared to that of scaled Gaussian matrices [97]. We have SRHT formally defined as below.

Definition 1.2.3. A Subsampled-Randomised-Hadamard-Transform (SRHT) [2, 97] is an $m \times n$ matrix of the form $S = S_s H D$ with $m \leq n$, where

- D is a random $n \times n$ diagonal matrix with ± 1 independent entries.
- H is an $n \times n$ Walsh-Hadamard matrix defined by

$$H_{ij} = n^{-1/2} (-1)^{\langle (i-1)_2, (j-1)_2 \rangle}, \quad (1.2.2)$$

where $(i-1)_2, (j-1)_2$ are binary representation vectors of the numbers $(i-1), (j-1)$ respectively⁶.

- S_s is a random scaled $m \times n$ sampling matrix (defined in Definition 1.2.1), independent of D .

A crucial drawback of SRHT is that if the column space is represented by a sparse matrix, the embedded matrix, although of smaller dimensions, is generally dense. Though sampling matrices preserve sparsity, we have mentioned above their downsides concerning high variance.

The second way to tackle the disadvantages of sampling is to use another sparse embedding ensemble instead. The 1-hashing matrices have one non-zero per column instead of one non-zero per row as in the sampling matrix; moreover, the value of the non-zero is ± 1 with equal probability so that $\mathbb{E} [\|Sx\|_2]^2 = \|x\|_2^2$. We have the following formal definition.

Definition 1.2.4 (1-hashing [21, 58]). $S \in \mathbb{R}^{m \times n}$ is a 1-hashing matrix if independently for each $j \in [n]$, we sample i uniformly at random and let $S_{ij} = \pm 1$ with equal probability.

Conceptually, unlike sampling, which discards a number of rows of the vector/matrix to be embedded, hashing uses every single row. The dimensionality reduction is achieved by hashing

⁶For example, $(3)_2 = (1, 1)$.

those rows into m slots, and adding them with sign-randomisation if multiple rows are hashed into a single slot. Therefore intuitively, hashing is more robust than sampling because it uses all the rows, and theoretical results have been established to show 1-hashing matrices with appropriate dimensions are oblivious JL/subspace embeddings without any requirement on the non-uniformity of the input [21, 79].

Finally, 1-hashing can be generalised to s -hashing, that is, matrices with s non-zeros per column, defined below.

Definition 1.2.5. [21] $S \in \mathbb{R}^{m \times n}$ is a s -hashing matrix if independently for each $j \in [n]$, we sample without replacement $i_1, i_2, \dots, i_s \in [m]$ uniformly at random and let $S_{i_k j} = \pm 1/\sqrt{s}$ for $k = 1, 2, \dots, s$.

Conceptually, s -hashing sketches each row of the input s times (into s different rows of the output), with each row being scaled by $\frac{1}{\sqrt{s}}$, and has better theoretical properties when used as a JL/subspace embedding than 1-hashing [22]. However, we note that while 1-hashing does not increase the number of non-zeros in the vector/matrix to be embedded, s -hashing may increase it by up to s times.

1.3 Linear least squares

Linear Least Squares (LLS) problems arising from fitting observational data to a linear model are mathematically formulated as the optimisation problem,

$$\min_{x \in \mathbb{R}^d} \|Ax - b\|_2^2, \quad (1.3.1)$$

where $A \in \mathbb{R}^{n \times d}$ is the data matrix that has (potentially unknown) rank r , $b \in \mathbb{R}^n$ is the vector of observations, and $n \geq d \geq r$. Thus (1.3.1) represents an optimisation problem where we have n data points and a model of d variables.

Problem (1.3.1) is equivalent to solving the normal equations

$$A^T A x = A^T b. \quad (1.3.2)$$

Numerous techniques have been proposed for the solution of (1.3.2), and they traditionally involve the factorization of $A^T A$, just A , or iterative methods. The ensuing cost in the worst case is $\mathcal{O}(nd^2)$, which is prohibitive for large n and d [37]. We briefly survey here the main classical techniques for solving LLS (1.3.1)/(1.3.2) following [85]. For iterative methods and preconditioning, see [9], while for sparse input matrices, see also [27].

1.3.1 Dense linear least squares

We say problem (1.3.1) is a dense Linear Least Squares (dense LLS) if the matrix A is a dense matrix. Namely, the matrix A does not have sufficiently many zero entries for specialised sparse

linear algebra routines to be advantageous. To solve dense LLS, we may employ direct methods based on factorizations, or iterative methods typically based on conjugate gradient techniques.

Direct methods for dense LLS Cholesky factorization computes $A^T A = LL^T$, where $L \in \mathbb{R}^{d \times d}$ is a lower triangular matrix. Then the normal equations (1.3.2) are solved by forward-backward substitutions involving the matrix L . The main costs are computing $A^T A$ and factorizing it, both taking $\mathcal{O}(nd^2)$ though many practical algorithms compute the factorization without forming $A^T A$ explicitly.⁷ This method is affected by the potential ill-conditioning of $A^T A$ (since the condition number of $A^T A$ is the square of the condition number of A) and so may not solve (1.3.2) accurately.

Employing the QR factorization aims to solve (1.3.1) directly without using (1.3.2). Computing $A = QR$, where $Q \in \mathbb{R}^{n \times n}$ is orthogonal and $R \in \mathbb{R}^{n \times d}$ is upper triangular, we have that $\|Ax - b\|_2^2 = \|Rx - Q^T b\|_2^2$. As R is both over-determined and upper triangular, its last $n - d$ rows are zeros. Therefore, $\|Rx - Q^T b\|_2^2$ is minimised by making the first d rows of Rx equal to the first d rows of $Q^T b$ which involves solving a linear system of equations involving the upper triangular matrix R . Hence the dominant cost is the QR factorization, which is $\mathcal{O}(nd^2)$.

When A is rank-deficient or approximately rank-deficient, Cholesky factorization may break down and (un-pivoted) QR factorization may give a rank-deficient R , introducing numerical instabilities in solving systems involving R . Singular Value Decomposition(SVD)-based methods are the most robust in dealing with rank-deficient problems, as an SVD reveals the spectrum, and therefore the extent of rank-deficiency, of the matrix A .

A (full) SVD of A is computed as $A = U_f \Sigma_f V_f^T$, where $U_f \in \mathbb{R}^{n \times d}$, $V_f \in \mathbb{R}^{d \times d}$ have orthonormal columns, and Σ_f is diagonal with non-negative and decreasing diagonal entries. The rank deficiency in A is then controlled by carefully selecting a cut-off point in the diagonal of Σ_f . After which the method proceeds similarly to QR-based approach by replacing A in (1.3.1) with its factorization and using the fact that left multiplying matrices with orthonormal columns/right multiplying orthogonal matrices does not change the 2-norm. However SVD-based methods are more computationally expensive than QR-based ones [37].

Iterative methods for dense LLS LSQR [86] and related algorithms such as LSMR [34] apply conjugate gradient method to solve the normal equations (1.3.2), only requiring matrix-vector multiplications involving A or A^T . In the worst case, $\mathcal{O}(d)$ iterations with $\mathcal{O}(nd)$ floating-point arithmetic operations per iteration are required. Therefore the worst case cost of iterative methods for dense LLS is $\mathcal{O}(nd^2)$ ⁸. But if the spectrum of A (the distribution of the singular values of A) is favorable these methods may take less iterations [37].

⁷Factorising $A^T A$ takes $\mathcal{O}(d^3)$ only but given that $n \geq d$, it is still $\mathcal{O}(nd^2)$.

⁸We note that iterative methods may not converge in $\mathcal{O}(d)$ iterations if A has a large condition number due to the effect of floating-point arithmetic.

Preconditioning techniques lower the condition number of the system by transforming the problem (1.3.1) into an equivalent form before applying iterative methods. For example, a sophisticated preconditioner for (1.3.1) is the incomplete Cholesky preconditioner [94], that uses an incomplete Cholesky factor of A to transform the problem. In general, some preconditioning should be used together with iterative methods [39].

1.3.2 Sparse linear least squares

When the matrix A is sparse (that is, there is a significant number of zeros in A such that specialised sparse linear algebra algorithms could be advantageous), we refer to the problem as sparse Linear Least Squares (sparse LLS).

Direct methods for sparse LLS We refer the reader to [27], where sparse Cholesky and QR factorizations are described. The main difference compared to the dense counterparts is that when the positions of a large number of zero entries of A are given, it is possible to use that information alone to predict positions of some zero entries in the resulting factors so that no computation is required to compute their values. Therefore, sparse factorizations could be faster than their dense counterparts on sparse A .

Iterative methods for sparse LLS The LSQR and LSMR algorithms mentioned in solving dense LLS automatically take advantage of the sparsity of A , as the matrix-vector multiplications involving A or A^T are faster when A is sparse.

1.3.3 Sketching for linear least squares

Over the past fifteen years, sketching techniques have been investigated for improving the computational efficiency and scalability of methods for the solution of (1.3.1); see, for example, the survey papers [73, 101]. Sketching uses a carefully-chosen random matrix $S \in \mathbb{R}^{m \times n}$, $m \ll n$, to sample/project the measurement matrix A to lower dimensions, while approximately preserving the geometry of the entire column space of A ; this quality of S (and of its associated distribution) is captured by the (oblivious) subspace embedding property [101] in Definition 2.2.3. The sketched matrix SA is then used to either directly compute an approximate solution to problem (1.3.1) or to generate a high-quality preconditioner for the iterative solution of (1.3.1); the latter has been the basis of state-of-the-art randomized linear algebra codes such as Blendenpik [5] and LSRN [78], where the latter improves on the former by exploiting sparsity and allowing rank-deficiency of the input matrix⁹.

Using sketching to compute an approximate solution of (1.3.1) in a computationally efficient way is proposed by Sarlos [93]. Using sketching to compute a preconditioner for the iterative solution

⁹LSRN also allows and exploits parallelism but this is beyond our scope here.

of (1.3.1) via a QR factorization of the sketched matrix is proposed by Rokhlin [92]. The choice of the sketching matrix is very important as it needs to approximately preserve the geometry of the column space of a given matrix (a subspace embedding, see Definition 2.2.3) with high-probability, while allowing efficient computation of the matrix-matrix product SA . Sarlos [93] proposed using the fast Johnson-Lindenstrauss transform (FJLT) discovered by Ailon and Chazelle [2]. The FJLT (and similar embeddings such as the SRHT) is a structured matrix that takes $O(nd \log(d))$ flops to apply to A , while requiring the matrix S to have about $m = O(d \log(d))$ rows to be a subspace embedding (also see Tropp [97], Ailon and Liberty [3]). More recently, Clarkson and Woodruff [21] proposed the hashing matrix that has one non-zero per column in random rows, with the value of the non-zero being ± 1 . This matrix takes $O(nnz(A))$ flops to apply to A , but needs $m = \Theta(d^2)$ rows to be a subspace embedding (also see Meng et al [77], Nelson et al [81, 80]). Recent works have also shown that increasing number of non-zeros per column of the hashing matrix leads to reduced requirement of number of rows (Cohen [22], Nelson [79]). Further work on hashing by Bourgain [10] showed that if the coherence¹⁰ is low, hashing requires fewer rows. These sketching matrices have found applications in practical implementations of sketching algorithms, namely, Blendenpik [5] used a variant of FJLT; LSRN [78] used Gaussian sketching; Iyer [53] and Dahiya [24] experimented with 1-hashing¹¹. The state-of-the-art sketching solvers Blendenpik [5] and LSRN [78] demonstrated several times speed-ups comparing to solvers based on QR/SVD factorizations in LAPACK [4], and LSRN additionally showed significant speed-up comparing to the solver based on sparse QR in SuiteSparse [28] when the measurement matrix A is sparse. However, Blendenpik and LSRN have not fully exploited the power of sketching, namely, Blendenpik only solves problem (1.3.1) when the measurement matrix A has full column rank, and LSRN uses Gaussian sketching matrices with dense SVD even for a sparse measurement matrix A . We propose a new solver in Chapter 3.

1.3.4 Alternative approaches for preconditioning and solving large-scale LLS problems

On the preconditioning side for linear least squares, [18] considered alternative regularization strategies to compute an Incomplete Cholesky preconditioner for rank-deficient least squares. LU factorization may alternatively be used for preconditioning. After a factorization $PAQ = LU$ where P, Q are permutation matrices, the normal equation $A^T A x = A^T b$ is transformed as $L^T L y = L^T c$ with $y = U Q^T x$ and $c = P b$. In [49], L is further preconditioned with L_1^{-1} where L_1 is the upper square part of L . On the other hand, [36] proposed and implemented a new sparse QR factorization method, with C++ code and encouraging performance on Inverse Poisson Least Squares problems. For a survey, see [39, 38].

¹⁰Maximum row norm of the left singular matrix U from the compact SVD of the matrix $A = U \Sigma V^T$. Formally defined in Definition 2.2.5.

¹¹hashing matrices with 1 non-zero per column

In addition to the sketch-precondition-solve methodology we use, large-scale linear least squares may alternatively be solved by first-order methods, zeroth-order methods (including Stochastic Gradient Descent (SGD)) and classical sketching (namely, solve the randomly embedded linear least square problem directly, as in see [93]). First order methods construct iterates $x_{t+1} = x_t - \mu_t H_t^{-1} g(x_t) + \beta_t(x_t - x_{t-1})$, where $H_t = A^T S_t^T S_t A$, $g(x_t) = A^T A x_t - A^T b$ and the last term represents the momentum. This is proposed in [64, 62], deriving optimal sequences μ_t, β_t for Gaussian and subsampled randomised Hadamard transforms, for S_t fixed or refreshed at each iteration. See also [43] for a randomised method for consistent linear least squares (namely, the residual at the optimal solution is zero). On the other hand, because linear least squares is a convex problem, SGD can be used, with [69] investigating using SGD with heavy ball momentum and [57] investigating using SGD with sketched Hessian. Using gradient-based sampling instead of uniform or leverage-scores-based sampling is explored in [110]. Finally, [70] provides techniques for a posteriori error estimates for classical/explicit sketching.

1.4 Minimising a general unconstrained objective function

We consider the following problem

$$\min_{x \in \mathbb{R}^d} f(x), \tag{1.4.1}$$

where f is smooth and non-convex. We will be satisfied if our algorithm returns an (approximate) local minimum of the objective function f – a point at which, if f is continuously differentiable, its gradient $\nabla f(x)$ is (approximately) zero; if f is twice continuously differentiable, then in addition to its gradient being zero, its Hessian $\nabla^2 f(x)$ is (approximately) positive semi-definite. This may not be the global minimum – namely, the smallest value of f over the entire \mathbb{R}^d . Finding the latter is much more computationally challenging and the remit of the field of global optimization [68]. Though global optimization is beyond our scope here, local optimization algorithms are often key ingredients in the development of techniques for the former.

Starting from a(ny) initial guess $x_0 \in \mathbb{R}^d$, classical (local) algorithms for solving (1.4.1) are iterative approaches that generate iterates x_k , $k \geq 0$, recursively, based on an update of the form

$$x_{k+1} = x_k - s_k,$$

where the step s_k is chosen so that the objective f typically decreases at each iteration. Depending, for example, on the problem information used in the computation of s_k , algorithms can be classified into those that use only up to first order information (i.e. $f(x), \nabla f(x)$); and those that also use second order information (i.e. $\nabla^2 f(x)$).

1.4.1 First order methods

For a(ny) user-provided tolerance $\epsilon > 0$, first order methods find an iterate x_k such that $\|\nabla f(x_k)\|_2 < \epsilon$. This ensures the approximate achievement of the necessary optimality condition $\nabla f(x) = 0$ that holds at any local minimiser of problem (1.4.1). However, note that this condition is not sufficient, e.g. x with $\nabla f(x) = 0$ could be a saddle point or even a local maximiser. However, as the iterates also progressively decrease f , we increase the chance that we find an approximate local minimiser. Indeed, several recent works have shown that for a diverse set of landscapes, first order methods such as the gradient descent methods escape/do not get trapped in saddle points and approach local minimisers; see for example, [55, 66]. We briefly review the three main classical first order methods: steepest descent with line search, the trust region method and the adaptive regularisation method.

Steepest descent with linesearch Steepest descent method seek the update

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k), \quad (1.4.2)$$

where the step α_k is determined by a line-search detailed below.

Given a constant $0 < \beta < 1$, the linesearch algorithm starts with some initial guess of $\alpha_k > 0$, and repeatedly decreases α_k until the following Armijo condition is satisfied:

$$f(x_k) - f(x_k - \alpha_k \nabla f(x_k)) \geq \beta \alpha_k \|\nabla f(x_k)\|_2^2. \quad (1.4.3)$$

It can be shown that assuming $f(x)$ is continuously differentiable, one can always find $\alpha_k > 0$ satisfying (1.4.3); moreover, provided the gradient of f is Lipschitz continuous¹² and f is bounded below, the steepest descent algorithm with linesearch requires at most $\mathcal{O}(\epsilon^{-2})$ evaluations of the objective function and its gradient to converge to an iterate x_k such that $\|\nabla f(x_k)\|_2 \leq \epsilon$ [16, 83]. This complexity bound is sharp, as shown, for example, in Theorem 2.2.3 of [16].

The trust region method In the trust region method, the step s_k is calculated by minimizing a local model $m_k(s) = f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T B_k s$, where B_k is a symmetric matrix (that is required to be uniformly bounded above with k). The matrix B_k could be some approximation of the Hessian/curvature information, if possible.

We then compute s_k by approximately minimising $m_k(s)$ within a trust region $\|s\|_2 \leq \Delta_k$ so that the decrease in m_k achieved by taking the step s_k is at least as much as that can be achieved by considering the steepest descent direction in the trust region. The trust region radius Δ_k is initialised at some value Δ_0 and subsequently dynamically adjusted: for the computed step s_k ,

¹²For some $L > 0$, we have that $\|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x - y\|_2$ for all $x, y \in \mathbb{R}^d$.

if we have sufficient decrease in f , $f(x_k) - f(x_k + s_k) \geq \eta[m_k(0) - m_k(s_k)]$, for some (iteration-independent) constant $0 < \eta < 1$ and $0 < \gamma < 1$, then $\Delta_{k+1} = \gamma^{-1}\Delta_k$ and $x_{k+1} = x_k + s_k$. Otherwise, $\Delta_{k+1} = \gamma\Delta_k$ and we do not take the step s_k ($x_{k+1} = x_k$).

It has been shown that the first order trust region method also has a global complexity of $\mathcal{O}(\epsilon^{-2})$ in terms of gradient and objective function evaluations [46]. This complexity bound is also sharp [16].

The adaptive regularisation method Like the trust region method, adaptive regularisation uses a local model around the current iterate x_k to compute a suitable step s_k . Unlike the trust region method, which explicitly restricts the size of the potential step, a regularisation term is imposed to a first-order Taylor model to implicitly restrict the size of the step:

$$m_k(s) = f(x_k) + \nabla f(x_k)^T s + \frac{1}{2}\sigma_k \|s\|_2^2 = T_{f,1}(x_k, s) + \frac{1}{2}\sigma_k \|s\|_2^2, \quad (1.4.4)$$

where $T_{f,1}(x_k, s)$ is the first-order Taylor series expansion of f around x_k . We minimise the local regularised model to compute a trial step s_k . Then, as in the trust region method, we evaluate the objective at the trial step $x_k + s_k$ and dynamically adjust the regularisation parameter by computing $\rho_k = \frac{f(x_k) - f(x_k + s_k)}{T_{f,1}(x_k, 0) - T_{f,1}(x_k, s_k)}$; and, if $\rho_k \geq \eta$, we set $x_{k+1} = x_k + s_k$ and set $\sigma_{k+1} = \max(\gamma\sigma_k, \sigma_{min})$, otherwise we do not take the step ($x_{k+1} = x_k$) and increase the regularisation by setting $\sigma_{k+1} = \frac{1}{\gamma}\sigma_k$, where $\gamma, \sigma_{min} \in (0, 1)$ are constants.

The first order adaptive regularisation method also has a (sharp) complexity of $\mathcal{O}(\epsilon^{-2})$ for both gradient and objective function evaluations, under the same assumptions on f as for the trust region methods [16].

Subspace first order methods The coordinate descent method [102] iteratively computes the next iterate by fixing most components of the variable x at their values from the current iterate, while approximately minimising the objective with respect to the remaining components; thus effectively searches a potential step in a restricted subspace of the full variable space. The coordinate descent method is convergent for some convex problems [102], but fails to converge for nonconvex problems [89]. The coordinate descent has found many applications in solving large-scale problems [7, 90]. Randomised coordinate descent methods have been an intense topic of recent investigations due to the demands of large scale problems; see [82, 91, 65, 48, 107, 109, 33, 105, 106, 71, 87]. For a survey see [102]. In Chapter 4 of this thesis, we will study a probabilistic subspace first order algorithm for general non-convex problems, that only needs directional gradient evaluations $S_k \nabla f(x_k)$ (so that the algorithm only searches the step in a subspace), where $S_k \in \mathbb{R}^{l \times d}$ is some random matrix to be specified with l being a user-chosen constant. Our work builds on the framework in [17]. However, here we establish more general results and use subspaces explicitly to save gradient evaluation/computation cost. We show that the appropriate replacement of the full gradient with

subspace gradients $S_k \nabla f(x_k)$ does not harm the worst-case complexity; although in our specific algorithm, since S_k is a random matrix, there is a probabilistic component in our convergence result. That is, we have $\|\nabla f(x_k)\|_2 < \epsilon$ with $k = \mathcal{O}(\epsilon^{-2})$ with probability proportional to $1 - e^{-\mathcal{O}(k)}$. The failure probability is very small for any reasonable value of k , and we show some numerical examples illustrating the effectiveness of our algorithm compared to the classical first order based methods when applied to non-linear least squares.

1.4.2 Second order methods

In order to improve both the performance and the optimality guarantees of first order methods, we add curvature information both in the construction and in the termination of algorithms, when this is available. Given accuracy tolerances $\epsilon_S, \epsilon_2 > 0$, we may strengthen our goal to try to find a point where simultaneously,

$$\|\nabla f(x_k)\|_2 < \epsilon_S, \lambda_{\min}(\nabla^2 f(x_k)) > -\epsilon_2 \quad (1.4.5)$$

where $\lambda_{\min}(\cdot)$ denotes the left-most eigenvalue of a matrix. These conditions secure approximate second order criticality conditions and strengthen the guarantee that we are close to a local minimiser, where $\nabla^2 f(x)$ is positive semidefinite. Clearly, in order to achieve this aim, the algorithm needs to be provided with second order information, the Hessian $\nabla^2 f(x_k)$, at each iteration. Let us review the classical optimisation methods for smooth, non-convex objective where both first and second order information is available.

Newton's method In Newton's method, the iterates are constructed according to

$$x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} \nabla f(x_k), \quad (1.4.6)$$

that is, the step s_k satisfies the linear system $\nabla^2 f(x_k) s_k = -\nabla f(x_k)$. Note that here one assumes that the matrix $\nabla^2 f(x_k)$ is positive definite for each k . For general functions, one may add regularisation terms, or use a linesearch or trust region, which we will discuss later.

Newton's method is attractive because it has a quadratic convergence property once x_k gets into a neighbourhood of a nondegenerate solution. However such a neighbourhood is typically not known a priori, before the run of the algorithm. It turns out that starting from an arbitrary starting point, the complexity of Newton's method can be the same as the steepest descent method, $\mathcal{O}(\epsilon^{-2})$, even if we assume the Newton direction is always well-defined [16].

The second order trust region method Trust-region methods could also use additional second order information, and the local model at each iteration becomes

$$m_k(s) = f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T \nabla^2 f(x_k) s. \quad (1.4.7)$$

We then compute an approximate minimiser of the model, subject to s being smaller than the trust-region radius. As before, it suffices to compute an approximate solution for convergence. In addition to requiring the model decrease to be at least as much as that of in the steepest descent direction, for second order criticality, we also require that the model decrease is at least as much as that obtained in the direction of the eigenvector of $\nabla^2 f(x_k)$ corresponding to the smallest eigenvalue (if such an eigenvalue is negative, otherwise this second condition is not required). Then, the objective at the trial step is again evaluated, and the ratio of the function decrease with the model decrease is compared to a pre-defined constant, and steps are taken/not taken; trust region radius is increased/decreased accordingly.

The second order trust-region algorithm has been shown to converge to a first order critical point $\|\nabla f(x_k)\|_2 < \epsilon_S$ in $\mathcal{O}(\epsilon_S^{-2})$ iterations; moreover, this bound is sharp. Thus the first order complexity of the (first order) trust region method is not improved by upgrading the trust region model to include accurate second order information. However, one can further show that the second order trust-region algorithm converges to a second order critical point $\lambda_{\min}(\nabla^2 f(x_k)) > -\epsilon_2$ and $\|\nabla f(x_k)\|_2 \leq \epsilon_S$ in $\mathcal{O}(\max(\epsilon_S^{-2}, \epsilon_2^{-3}))$ iterations. We see that the main advantage of the second order trust region over the first order one is that it allows one to quantifiably compute an approximate second order critical point [16].

The second order adaptive regularisation method The second order adaptive regularisation method, however, is able to not only allow convergence to a second order critical point, but also allow an improved speed of convergence to a first order critical point. The algorithm is the following: at each iteration, we build the model as

$$m_k(s) = f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T \nabla^2 f(x_k) s + \frac{1}{3} \sigma_k \|s\|_2^3, \quad (1.4.8)$$

where σ_k is an adaptive parameter whose value increases or decreases depending on the amount of objective function decrease achieved by a step calculated from (approximately) minimising such a model. Compared to the first order regularisation method, the approximate minimisation here requires that $\|\nabla_s m_k(s_k)\|_2 \leq \frac{1}{2} \theta_1 \|s_k\|_2^2$ and $\lambda_{\min}(\nabla_s^2 m_k(s_k)) \geq \theta_2 \|s_k\|_2$ (for two iteration independent constants $\theta_1, \theta_2 > 0$). Assuming that f is bounded below and that its Hessian is Lipschitz continuous¹³, it can be shown that this algorithm converges to a point where $\|\nabla f(x_k)\|_2 < \epsilon_S$ in $\mathcal{O}(\epsilon_S^{-3/2})$ evaluations of the objective function, gradient and Hessian; and to a point where $\|\nabla f(x_k)\|_2 < \epsilon_S$ and $\lambda_{\min}(\nabla^2 f(x_k)) > -\epsilon_2$ in $\mathcal{O}(\max(\epsilon_S^{-3/2}, \epsilon_2^{-3}))$ evaluations of the objective function, gradient and Hessian. Moreover, both of these bounds are sharp and the first-order one is provably optimal for second order methods [16, 11]. Minimising (a slight variant of) (1.4.8) to compute a step was first suggested in [47]. Independently, [84] considered using (1.4.8) from a different perspective. The above mentioned method was first proposed in [15] that improves on previous

¹³For some $L > 0$, we have that $\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L \|x - y\|_2$ for all $x, y \in \mathbb{R}^d$.

attempts, namely, allowing inexact model minimisation and without requiring the knowledge of problem-specific constants.

Subspace second order adaptive regularisation methods Methods that only use the gradient/Hessian in a subspace have been studied in the recent years. The sketched Newton algorithm [88] requires a sketching matrix that is proportional to the rank of the Hessian. Sketched online Newton [72] uses streaming sketches to scale up a second-order method, comparable to Gauss–Newton, for solving online learning problems. The randomised subspace Newton [42] efficiently sketches the full Newton direction for a family of generalised linear models, such as logistic regression. Other randomised versions of Newton’s method include [44, 41, 8]. The global convergence of the above methods, however, require the objective function f to be convex (or even strongly convex). For general non-convex optimisation, [17, 45] give generic frameworks that apply to the first order/second order general optimisation methods. Our main focus is non-convex functions and so we build on these works. Chapter 5 of this thesis proposes a second order adaptive regularisation method when operating in a random subspace. Specifically, both the gradient and the Hessian will be replaced by their subspace equivalent. We are able to show that under suitable assumptions on the subspace sketching matrix $S_k \in \mathbb{R}^{l \times d}$ (that could be a scaled Gaussian matrix with $l = \mathcal{O}(r)$, where r is the rank of the Hessian $\nabla^2 f(x_k)$), both the fast convergence rate $\mathcal{O}(\epsilon_S^{-3/2})$ to the first order critical point, and the convergence to the second order critical point with a rate $\mathcal{O}(\epsilon_2^{-3})$ can be retained.

1.4.3 Applications to non-linear least squares problems

Non-linear least squares are a subclass of general unconstrained optimisation problems. We aim to solve

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{2} \sum_{i=1}^n (r_i(x))^2 = \frac{1}{2} \|r(x)\|_2^2, \quad (1.4.9)$$

where $r(x) = (r_1, r_2, \dots, r_n)(x)$ is a vector-valued smooth residual function $r : \mathbb{R}^d \rightarrow \mathbb{R}^n$. This formulation has a wide range of applications in weather forecasting, finance and machine learning problems. We briefly overview some classical solution methods here, following [85].

The Gauss–Newton method The Gauss–Newton method is a simplification of Newton’s method that exploits the structure of non-linear least squares problems. In particular, we can approximate the Hessian as $\nabla^2 f(x_k) = J(x_k)^T J(x_k) + \sum_{i=1}^n r_i(x_k) \nabla^2 r_i(x_k) \approx J(x_k)^T J(x_k)$ where

$$J(x) = \left(\frac{\partial r_i(x)}{\partial x_j} \right)_{ij} \in \mathbb{R}^{n \times d}.$$

This approximation is justified in the case when $r(x) \approx 0$ at a solution x or when r is approximately linear in the variables. Despite using only first-derivative information about r , the Gauss–Newton

method has been shown to enjoy a local super-linear convergence rate to a first order critical point. When the Gauss-Newton direction is safeguarded with a trust region or regularization technique (which is often referred to as a Levenberg-Marquardt method), it can be shown to have global convergence provided for example, that J is Lipschitz continuous. To ensure global convergence of linesearch variants of Gauss-Newton, we additionally need to require that the Jacobian’s singular values are uniformly bounded away from zero – a very strong assumption.

Subspace Gauss Newton method Expanding on our work in [14], in Chapter 4, we present such an algorithmic variant and its numerical performance when compared to the full Gauss Newton method. Subspace Gauss-Newton variants can also be found in [42].

1.5 Contributions of this thesis

1.5.1 New theoretical analysis of hashing matrices and development of new random embeddings

In Chapter 2, we study theoretical properties of hashing matrices and propose a new oblivious subspace embedding based on hashing. We show that hashing matrices — with one nonzero entry per column and of size proportional to the rank of the data matrix — generate a subspace embedding with high probability, provided the given data matrix has low coherence. We then show that s -hashing matrices, with $s > 1$ nonzero entries per column, satisfy similarly good sketching properties for a larger class of low coherence data matrices.

More specifically, we show that a hashing matrix $S \in \mathbb{R}^{m \times n}$ with $m = \mathcal{O}(r)$ is an oblivious subspace embedding for matrices A with low coherence. Hashing matrices have been shown empirically to be almost as good as Gaussian matrices [24] in terms of their embedding properties, but the theoretical results show that they need at least $m = \mathcal{O}(r^2)$ rows [81]. Our result explains the phenomenon observed in [24]. In addition, it was observed empirically that one needs at least 2 non-zeros per hashing column for the projection to be accurate. We show that using s non-zeros per column instead of 1 non-zero per column relaxes the coherence requirement by \sqrt{s} . Thus we expect more performance improvement if we increase $s = 1$ to $s = 2$ than from $s = 2$ to $s = 3$. Previous work on s -hashing has independently discovered the \sqrt{s} factor [10], but our result relies on a single, intuitive proof, and is not tied to any particular proof for the case of 1–hashing. So if the coherence requirement bound for 1-hashing is subsequently improved, our work allows the result on s -hashing to improve automatically.

Cascading this result, we also introduce and analyse a new random embedding: Hashed-Randomised-Hadamard-Transform (HRHT), that combines the coherence reduction properties of randomised Hadamard Transform with the coherence-based theory of hashing embeddings. Compared to Subsampled-Randomised-Hadamard-Transform, HRHT is able to achieve subspace embedding with the embed-

ding dimension $\mathcal{O}(r)$ where r is the rank of the matrix to be embedded, matching the optimal bound known for the Gaussian embedding. Experiments using random embeddings for preconditioning linear least squares show the improved performance of HRHT over SRHT.

1.5.2 Analysis, state-of-the-art implementation and benchmarking of new large-scale linear least squares solver using random embeddings

Chapter 3 concerns the solution of large-scale Linear Least Squares (LLS) problems, by applying random embeddings to reduce the dimensionality of the observation/sample space. The sketched matrix SA is used to generate a high-quality preconditioner for the iterative solution of (1.3.1) and has been the basis of state-of-the-art randomized linear algebra codes such as Blendenpik [5] and LSRN [78], where the latter improves on the former by exploiting input sparsity, parallelization and allowing rank-deficiency of the input matrix.

We propose and analyse a sketching framework for potentially rank-deficient LLS. Our framework includes the algorithm used by Blendenpik [5] and LSRN [78]; and additionally it allows one to use a wide range of rank-revealing factorizations to build a preconditioner using the sketched matrix SA . Our analysis shows that one can recover a minimal residual solution with a rank-revealing QR factorization with sketching, or the minimal norm solution with a total orthogonal factorization with sketching. This framework allows us to use (randomised) column pivoted QR factorization for dense LLS so that our solver solves rank-deficient LLS satisfactorily without using the expensive SVD. We are also able to use a sparse rank-revealing factorization for sparse LLS, obtaining a significant speed-up over LSRN [78], state-of-the-art sparse solvers LS_SPQR [28] and incomplete Cholesky factorization preconditioned Krylov subspace method LS_HSL [94].

Numerically, we developed a solver SKi-LLS (SKetchIng-Linear-Least-Square) combining our theoretical and algorithmic ideas and state-of-the-art C++ implementations. For dense inputs, the solver is more robust than Blendenpik (as it solves rank-deficient or approximately rank-deficient problems); while being quicker than Blendenpik for matrices with high coherence and comparable in speed with Blendenpik for other matrices. In order to overcome the speed deterioration of the column-pivoted QR comparing to the un-pivoted QR, we used a recent development of randomised column pivoted QR that exploits randomisation and the importance of memory usage and cache in modern computing architecture [76]. For sparse inputs, by using a sparse QR factorization code developed by Davis [28], our solver is more than 10 times faster than LSRN, LS_SPQR and LS_HSL for sparse Gaussian inputs. We extensively compared our solver with LSRN, LS_SPQR and LS_HSL on the Florida Matrix Collection [29], and our solver is extremely competitive on strongly-over-determined inputs or ill-conditioned inputs.

1.5.3 First-order subspace methods and their application to non-linear least squares

In Chapter 4, we analyse a general randomised algorithmic framework for minimizing a general objective function (1.4.1), that improves upon the one introduced in [17], so that an arbitrarily high probability convergence/complexity result can be derived. We formulate more specific conditions on the reduced local models that are based on random embeddings of the variable space (in contrast to embedding the observational space in the linear least squares case). Compared to [45], our algorithm applies more generally¹⁴, also to quadratic regularisation (see later sections).

Compared to [17, 45], we use a weaker/different definition of a ‘true’ iteration, when the approximate problem information is sufficiently accurate; this definition is based on the random embedding satisfying a (one-sided) JL-embedding property (see (1.5.1) below), which is novel. Using the latter property and typical smoothness assumptions on the problem, we show that our framework of random subspace methods has complexity $\mathcal{O}(\epsilon^{-2})$ to generate an approximate first-order stationary point, with exponentially high probability. To ensure this, the random subspace needs to sketch the gradient, replacing $\nabla f(x_k)$ with $S_k \nabla f(x_k)$ in the algorithm, where $S_k \in \mathbb{R}^{l \times d}$ satisfies, with positive probability,

$$\|S_k \nabla f(x_k)\|_2 \geq (1 - \epsilon_S) \|\nabla f(x_k)\|_2. \quad (1.5.1)$$

We show that the above is achieved when S_k is a scaled Gaussian matrix, a hashing (sparse-embedding) matrix, a sampling matrix, and many more.

We note again that this framework marks a significant departure from probabilistic model assumptions [17, 45], since our model gradient, $S_k \nabla f(x_k)$, does not even have the same dimension as the true gradient¹⁵. The intuition behind this requirement is that in classical trust-region or adaptive regularization, the norm of the gradient is a key ingredient in the recipe for ensuring a convergence result, and hence by preserving the norm of the gradient with sketching, we are able to obtain a similar worst case complexity. Another interesting observation is that in the case of S_k being the sampling matrix, for which our method reduces to a randomised block-coordinate approach, we show how the success of the algorithm on non-convex smooth problems is connected with the ‘non-uniformity’ of the gradient; thus leading to almost sure convergence under some strong assumptions related to the objective’s gradient. We then particularize this framework to Gauss-Newton techniques for nonlinear least squares problems, where the Jacobian is computed in a subspace.

¹⁴Also, in the case of trust region methods, our framework does not need to compare the norm of the model gradient with the trust region radius at each iteration in order to decide if the trust region radius should be increased (see [45], Algorithm 2.5).

¹⁵Hence the probabilistic model condition which bounds $\|\nabla m(x_k) - \nabla f(x_k)\|_2$ is not applicable here.

1.5.4 Random subspace cubic regularisation algorithm, R-ARC

In Chapter 5, we further analyse the random subspace framework when second order information is added and applied to general non-convex optimisation. We propose and analyse an algorithm that is a random subspace version of the second order adaptive cubic regularisation method. We show that the subspace variant matches the optimal convergence rate $\mathcal{O}(\epsilon_1^{-3/2})$ of the full-dimensional variant to generate $\|\nabla f(x_k)\|_2 \leq \epsilon_1$ under suitable assumptions: either the embedding matrix S_k provides a subspace embedding of the Hessian $\nabla^2 f(x_k)$, or the Hessian is sparse in the sense that only few rows are non-zero.

We further analyse convergence to second order critical points of the second order adaptive regularisation method. We first show that in general, the algorithm converges to a point where the subspace Hessian $S_k \nabla^2 f(x_k) S_k^T$ is approximately positive semi-definite. Then we prove that if scaled Gaussian matrices are used as random embeddings, the algorithm converges to a point where the full Hessian is approximately positive semi-definite, at a rate $\mathcal{O}(\epsilon_2^{-3})$ that matches the full-dimensional second order cubic regularisation method.

1.6 Structure of thesis

In Chapter 2, we first give the necessary technical background on random embeddings, which will be used throughout this thesis. We then state and prove our theorem on the coherence requirement needed to use 1-hashing with $m = \mathcal{O}(d)$ as an oblivious subspace embedding (defined in Definition 2.2.4). Cascading this result, we show how increasing the number of non-zeros per column from 1 to s relaxes the coherence requirement for s -hashing matrices by a factor of \sqrt{s} , and propose a new random matrix distribution for subspace embeddings that has at most s non-zeros per column. Then, we propose a carefully constructed random matrix distribution that uses hashing matrices, achieving $m = \mathcal{O}(d)$ as a subspace embedding with high probability for any sufficiently over-determined matrix A .

In Chapter 3, we propose and analyse an algorithmic framework that uses random embedding (sketching) for potentially rank-deficient linear least squares. Then we introduce our linear least squares solver Ski-LLS which implements the framework and discuss its key features and implementation details. We test and benchmark Ski-LLS against state of the art algorithms and test problems in the remainder of Chapter 3.

In Chapter 4, we move onto problem (1.4.1). We first propose and analyse an algorithmic framework that relies on stochastic reduced local models. We then show how sketching-based subspace methods fit into this framework, and derive results for quadratic-regularisation and trust-region algorithms for general unconstrained objective optimisation. We then apply this framework to

Gauss-Newton method and nonlinear least squares, obtaining a subspace Gauss-Newton method and illustrating its performance numerically.

In Chapter 5, we propose and analyse the subspace cubic-regularisation based approach for solving (1.4.1). We first show how subspace embedding of the Hessian of the objective function allows the same convergence rate as the (full-space) cubic-regularisation methods. We then show how the sparsity of the Hessian allows a similar convergence result to be derived. We then go on to analyse the convergence to second-order critical points using the subspace cubic-regularisation based approach. We show that using scaled Gaussian embeddings allows convergence of R-ARC to a second order critical point with a rate essentially the same as the full-space method.

Finally in Chapter 6, we summarise the main results in this thesis and set some future directions.

Notation. Throughout the thesis, we let $\langle \cdot, \cdot \rangle$ and $\|\cdot\|_2$ denote the usual Euclidean inner product and norm, respectively, and $\|\cdot\|_\infty$, the l_∞ norm. Also, for some $n \in \mathbb{N}$, $[n] = \{1, 2, \dots, n\}$. For a(ny) symmetric positive definite matrix \overline{W} , we define the norm $\|x\|_{\overline{W}} := x^T \overline{W} x$, for all x , as the norm induced by \overline{W} . The notation $\Theta(\cdot)$ denotes both lower and upper bounds of the respective order. $\Omega(\cdot)$ denotes a lower bound of the respective order. $\mathcal{O}(\cdot)$ denotes an upper bound of the respective order.

Chapter 2

Random embedding

2.1 Introduction and relevant literature

This chapter is based and expands materials in [95, 12].

Main contributions This chapter aims to explore the theoretical properties of sparse sketching matrices S for improved efficiency and scalability of methods for solving the LLS problem (1.3.1), when A is sparse or dense. After introducing the necessary technical background in Section 2, we firstly investigate projection and computational properties of 1-hashing matrices, random sparse matrices with 1 non-zero per column, that were first proposed in the randomised linear algebra context by Clarkson and Woodruff [21]. Sparse matrices allow faster computation of the matrix-matrix product SA , leading to faster embeddings than their dense counterparts. Moreover, sparse matrices preserve the sparsity of the data matrix A , allowing the sparsity of the embedded matrix SA to be exploited by specialized numerical linear algebra routines for faster computation.

It has been observed numerically in [24] that 1-hashing matrices, with the same projection dimensions as Gaussian matrices (matrices with i.i.d. Gaussian entries) are as efficient in projecting and solving the LLS problem (1.3.1). However, it was shown in [81, 77] that 1-hashing matrices require at least an $\mathcal{O}(r^2)$ projection dimension to work effectively (as an oblivious subspace embedding, defined in Definition 2.2.4) comparing to an $\mathcal{O}(r)$ projection dimension required by Gaussian matrices [101], where r is the rank of A in (1.3.1). Thus a gap exists between the theory and the empirical observation. Our main result on 1-hashing matrices shows that 1-hashing matrices can have the same theoretical properties as the Gaussian matrices, namely, being an oblivious subspace embedding for matrices of rank r with the projection dimension being $\mathcal{O}(r)$, given that the matrix A has low coherence (defined in Definition 2.2.5).

Cascading on this result, we then show in Section 2.4, firstly that s -hashing matrices, which are sparse matrices with (fixed) s non-zeros per column first proposed as a candidate matrix distribution for oblivious subspace embeddings in [79], achieves being an oblivious subspace embedding with the projection dimension of $\mathcal{O}(r)$ for matrices A of rank r , but with the coherence requirement on A

being relaxed by \sqrt{s} comparing to 1-hashing matrices. Our numerical illustration (Figure 3.3) shows that 2-hashing is more effective as a sketching matrix in solving (1.3.1).

Secondly in Section 2.4, we propose a new matrix distribution called s -hashing variant matrices that has at most s non-zeros per column for oblivious subspace embeddings. Using a novel result that allows us to connect any coherence-restricted embedding result for 1-hashing matrices to s -hashing matrices, we show that s -hashing variant matrices have a similar coherence restricted embedding property as s -hashing matrices.

At the end of Section 2.4, we combine s -hashing (variant) matrices with coherence reduction transformations that were first proposed in [2] to derive a new random matrix distribution that will be an oblivious subspace embedding with the projection dimension $\mathcal{O}(r)$ for any matrix A sufficiently over-determined; and will take $\mathcal{O}(nd \log(n))$ flops to apply. This so-called Hashed-Randomised-Hadamard-Transform (HRHT) improves upon the previously proposed Subsampled-Randomized-Hadamard-Transform (SRHT) by lowering the projection dimension from $\mathcal{O}(r \log(r))$ to $\mathcal{O}(r)$ while maintaining the complexity of embedding time up to a $\log(n)/\log(d)$ multiplicative factor.

Related literature After the Johnson-Lindenstrauss Lemma appeared in 1984, Indyk and Motawani proved scaled Gaussian matrices is a JL-embedding [50] for which an elementary proof was provided by Dasgupta and Gupta [25]. Achlioptas [1] showed that matrices with all entries being ± 1 with equal probability, or indeed matrices with all entries being $+1, -1, 0$ with equal probability are JL-embeddings. However these random ensembles take $\mathcal{O}(nd^2)$ to apply to an $n \times d$ matrix in general. The Fast-Johnson-Lindenstrauss-Transform (FJLT) as a JL-embedding was proposed in [2] and as a subspace-embedding was proposed in [93]. The FJLT is based on Fast Fourier Transform-like algorithms and is faster to apply to matrices and vectors. The construction and analysis of FJLT are subsequently improved in [97, 3, 5, 92], ending with [97] analysing a variant of FJLT called Subsampled Randomised Hadamard Transform (SRHT) using matrix concentration inequalities. SRHT takes $\mathcal{O}(nd \log(d))$ flops to apply to A , while requiring S to have about $m = \mathcal{O}(r \log(r))$ rows to be an oblivious subspace embedding, where r is the rank of A . Clarkson and Woodruff [21] proposed and analysed using the 1-hashing matrices as a candidate distribution for subspace embeddings, and Nelson and Nguyen [79] proposed and analysed using the s -hashing matrices. These sparse matrices are subsequently analysed in [77, 81, 80, 22], showing that for 1-hashing matrices, $m = \Omega(r^2)$ is required for being an oblivious subspace embedding (see also Example 3); and $m = \mathcal{O}(r^2)$ is sufficient. And for s -hashing matrices, $m = \mathcal{O}(r \log(r))$ is sufficient for being an oblivious subspace embedding with $s = \mathcal{O}(\log(r))$. Recently, Bourgain, Dirksen and Nelson [10] showed a coherence dependent result of s -hashing matrices (see Theorem 2.2.1).

Comparing to the existing results, our results on 1 and s -hashing matrices are the first oblivious subspace embedding results on 1 and s -hashing matrices with $m = \mathcal{O}(r)$; though our result does have a strict coherence requirement. Our result on Hashed-Randomised-Hadamard-Transform has a lower embedding dimension than the SRHT. (Note that it has also been shown in [97] that the embedding dimension of SHRT could not be further lowered due to the Coupon Collector’s Problem.)

Finally, we mention some recent works on random embeddings. Recent results concerning oblivious (tensor) subspace embeddings [51] could be particularized to oblivious (vector) subspace embeddings, leading to a matrix distribution $S \in \mathbb{R}^{m \times n}$ with $m = \mathcal{O}(r \log^4 r \log n)$ (where r is the rank of the matrix to be embedded) that requires $\mathcal{O}(n \log n)$ operations to apply to any vector. This has slightly worse space and time complexity than sub-sampled randomized Hadamard transform. Regarding sparse embeddings, [19] proposed a ‘stable’ 1-hashing matrix that has the (ϵ, δ) -oblivious JL embedding property with the optimal $m = \mathcal{O}(\epsilon^{-2} \log(1/\delta))$ (same as scaled Gaussian matrices) while each row also has approximately the same number of non-zeros. The algorithm samples n non-zero row indices for n columns of S by sampling without replacement from the set $\{[m], [m], \dots, [m]\}$ where $[m] = \{1, 2, \dots, m\}$ is repeated $\lceil \frac{n}{m} \rceil$ times. [67] proposed learning the positions and values of non-zero entries in 1-hashing matrices by assuming the data comes from a fixed distribution.

2.2 Technical Background

In this section, we review some important concepts and their properties that we then use throughout the thesis. We employ several variants of the notion of random embeddings for finite or infinite sets, as we define next.

2.2.1 Random embeddings

We start with a very general concept of embedding a (finite or infinite) number of points; throughout, we let $\epsilon \in (0, 1)$ be the user-chosen/arbitrary error tolerance in the embeddings and $n, k \in \mathbb{N}$.¹

Definition 2.2.1 (Generalised JL² embedding [101]). *A generalised ϵ -JL embedding for a set $Y \subseteq \mathbb{R}^n$ is a matrix $S \in \mathbb{R}^{m \times n}$ such that*

$$-\epsilon \|y_i\|_2 \cdot \|y_j\|_2 \leq \langle Sy_i, Sy_j \rangle - \langle y_i, y_j \rangle \leq \epsilon \|y_i\|_2 \cdot \|y_j\|_2, \quad \text{for all } y_i, y_j \in Y. \quad (2.2.1)$$

If we let $y_i = y_j$ in (2.2.1), we recover the common notion of an ϵ -JL embedding, that approximately preserves the length of vectors in a given set.

¹Note that here ϵ is not the error tolerance of the algorithms that we will discuss later in this thesis, e.g. linear least squares, general non-convex optimisations. While the error tolerance in the embedding influences the performance of embedding-based algorithms, it is not necessary to have a small error in the embedding in order to achieve a small error tolerance in the actual algorithm. Because the inaccuracy of the embedding may be mitigated by repeated iterations of the algorithm, or an indirect use of the embedding. In particular, although $\epsilon \in (0, 1)$, we do not require the embedding accuracy ϵ to be close to zero in this thesis.

²Note that ‘JL’ stands for Johnson-Lindenstrauss, recalling their pioneering lemma [56].

Definition 2.2.2 (JL embedding [101]). *An ϵ -JL embedding for a set $Y \subseteq \mathbb{R}^n$ is a matrix $S \in \mathbb{R}^{m \times n}$ such that*

$$(1 - \epsilon)\|y\|_2^2 \leq \|Sy\|_2^2 \leq (1 + \epsilon)\|y\|_2^2 \quad \text{for all } y \in Y. \quad (2.2.2)$$

Often, in the above definitions, the set $Y = \{y_1, \dots, y_k\}$ is a finite collection of vectors in \mathbb{R}^n . But an infinite number of points may also be embedded, such as in the case when Y is an entire subspace. Then, an embedding approximately preserves pairwise distances between any points in the column space of a matrix $B \in \mathbb{R}^{n \times k}$.

Definition 2.2.3 (ϵ -subspace embedding [101]). *An ϵ -subspace embedding for a matrix $B \in \mathbb{R}^{n \times k}$ is a matrix $S \in \mathbb{R}^{m \times n}$ such that*

$$(1 - \epsilon)\|y\|_2^2 \leq \|Sy\|_2^2 \leq (1 + \epsilon)\|y\|_2^2 \quad \text{for all } y \in Y = \{y : y = Bz, z \in \mathbb{R}^k\}. \quad (2.2.3)$$

In other words, S is an ϵ -subspace embedding for B if and only if S is an ϵ -JL embedding for the column subspace Y of B .

Oblivious embeddings are matrix distributions such that given a(ny) subset/column subspace of vectors in \mathbb{R}^n , a random matrix drawn from such a distribution is an embedding for these vectors with high probability. We let $1 - \delta \in [0, 1]$ denote a(ny) success probability of an embedding.

Definition 2.2.4 (Oblivious embedding [101, 93]). *A distribution \mathcal{S} on $S \in \mathbb{R}^{m \times n}$ is an (ϵ, δ) -oblivious embedding if given a fixed/arbitrary set of vectors, we have that, with probability at least $1 - \delta$, a matrix S from the distribution is an ϵ -embedding for these vectors.*

Using the above definitions of embeddings, we have distributions that are *oblivious JL-embeddings* for a(ny) given/fixed set Y of some vectors $y \in \mathbb{R}^n$, and distributions that are *oblivious subspace embeddings* for a(ny) given/fixed matrix $B \in \mathbb{R}^{n \times k}$ (and for the corresponding subspace Y of its columns). We note that depending on the quantities being embedded, in addition to ϵ and δ dependencies, the size m of S may depend on n and the ‘dimension’ of the embedded sets; for example, in the case of a finite set Y of k vectors in \mathbb{R}^n , m additionally may depend on k while in the subspace embedding case, m may depend on the rank r of B .

2.2.2 Generic properties of subspace embeddings

A necessary condition for a matrix S to be an ϵ -subspace embedding for a given matrix is that the sketched matrix has the same rank.

Lemma 2.2.1. *If the matrix S is an ϵ -subspace embedding for a given matrix B for some $\epsilon \in (0, 1)$, then $\text{rank}(SB) = \text{rank}(B)$, where $\text{rank}(\cdot)$ denotes the rank of the argument matrix.*

Proof. Let $B \in \mathbb{R}^{n \times k}$. By rank-nullity theorem, $\text{rank}(B) + \dim \ker(B) = \text{rank}(SB) + \dim \ker(SB) = k$. Clearly, $\dim \ker(SB) \geq \dim \ker(B)$. If the previous inequality is strict, then there exists $z \in \mathbb{R}^k$ such that $\|SBz\|_2 = 0$ and $\|Bz\|_2 > 0$, contradicting the assumption that S is an ϵ -subspace embedding for B according to (2.2.3). \square

Given any matrix $A \in \mathbb{R}^{n \times d}$ of rank r , the compact singular value decomposition (SVD) of A provides a perfect subspace embedding. In particular, let

$$A = U\Sigma V^T, \quad (2.2.4)$$

where $U \in \mathbb{R}^{n \times r}$ with orthonormal columns, $\Sigma \in \mathbb{R}^{r \times r}$ is diagonal matrix with strictly positive diagonal entries, and $V \in \mathbb{R}^{d \times r}$ with orthonormal columns [37]. Then the matrix U^T is a ϵ -subspace embedding for A for any $\epsilon \in (0, 1)$.

Next, we connect the embedding properties of S for A with those for U in (2.2.4), using a proof technique in Woodruff [101].

Lemma 2.2.2. *Let $A \in \mathbb{R}^{n \times d}$ with rank r and SVD-decomposition factor $U \in \mathbb{R}^{n \times r}$ defined in (2.2.4), and let $\epsilon \in (0, 1)$. Then the following equivalences hold:*

- (i) *a matrix S is an ϵ -subspace embedding for A if and only if S is an ϵ -subspace embedding for U , namely,*

$$(1 - \epsilon)\|Uz\|_2^2 \leq \|SUz\|_2^2 \leq (1 + \epsilon)\|Uz\|_2^2, \quad \text{for all } z \in \mathbb{R}^r. \quad (2.2.5)$$

- (ii) *A matrix S is an ϵ -subspace embedding for A if and only if for all $z \in \mathbb{R}^r$ with $\|z\|_2 = 1$, we have³*

$$(1 - \epsilon)\|Uz\|_2^2 \leq \|SUz\|_2^2 \leq (1 + \epsilon)\|Uz\|_2^2. \quad (2.2.6)$$

Proof. (i) Let $A = U\Sigma V^T$ be defined as in (2.2.4). If S is an ϵ -subspace embedding for A , let $z \in \mathbb{R}^r$ and define $x = V\Sigma^{-1}z \in \mathbb{R}^d$. Then we have $Uz = Ax$ and

$$\|SUz\|_2^2 = \|SAx\|_2^2 \leq (1 + \epsilon)\|Ax\|_2^2 = (1 + \epsilon)\|Uz\|_2^2, \quad (2.2.7)$$

where we have used $Uz = Ax$ and (2.2.3). Similarly, we have $\|SUz\|_2^2 \geq (1 - \epsilon)\|Uz\|_2^2$. Hence S is an ϵ -subspace embedding for U .

Conversely, given S is an ϵ -subspace embedding for U , let $x \in \mathbb{R}^d$ and $z = \Sigma V^T x \in \mathbb{R}^r$. Then we have $Ax = Uz$, and $\|SAx\|_2^2 = \|SUz\|_2^2 \leq (1 + \epsilon)\|Uz\|_2^2 = (1 + \epsilon)\|Ax\|_2^2$. Similarly $\|SAx\|_2^2 \geq (1 - \epsilon)\|Ax\|_2^2$. Hence S is an ϵ -subspace embedding for A .

³We note that since $\|z\|_2 = 1$ and U has orthonormal columns, $\|Uz\|_2 = \|z\|_2 = 1$ in (2.2.6).

- (ii) Since the equivalence in (i) holds, note that (2.2.5) clearly implies (2.2.6). The latter also implies the former if (2.2.6) is applied to $z/\|z\|_2$ for any nonzero $z \in \mathbb{R}^r$.

□

Remark 1. Lemma 2.2.2 shows that to obtain a subspace embedding for an $n \times d$ matrix A it is sufficient (and necessary) to embed correctly its left-singular matrix that has rank r . Thus, the dependence on d in subspace embedding results can be replaced by dependence on r , the rank of the input matrix A . As rank deficient matrices A are important in this thesis, we opt to state our results in terms of their r dependency (instead of d).

The matrix U in (2.2.4) can be seen as the ideal ‘sketching’ matrix for A ; however, there is not much computational gain in doing this as computing the compact SVD has similar complexity as computing a minimal residual solution to (1.3.1) directly.

2.2.3 Sparse matrix distributions and their embeddings properties

In terms of optimal embedding properties, it is well known that (dense) scaled Gaussian matrices S with $m = \mathcal{O}(\epsilon^{-2}(r + \log(1/\delta)))$ provide an (ϵ, δ) -oblivious subspace embedding for $n \times d$ matrices A of rank r [101]. However, the computational cost of the matrix-matrix product SA is $\mathcal{O}(nd^2)$, which is similar to the complexity of solving the original LLS problem (1.3.1); thus it seems difficult to achieve computational gains by calculating a sketched solution of (1.3.1) in this case. In order to improve the computational cost of using sketching for solving LLS problems, and to help preserve input sparsity (when A is sparse), sparse random matrices have been proposed, namely, such as random matrices with one non-zero per row. However, uniformly sampling rows of A (and entries of b in (1.3.1)) may miss choosing some (possibly important) row/entry. A more robust proposal, both theoretically and numerically, is to use hashing matrices, with one (or more) nonzero entries per column, which when applied to A (and b), captures all rows of A (and entries of b) by adding two (or more) rows/entries with randomised signs. The definition of s -hashing matrices and was given in Definition 1.2.5 and when $s = 1$, we have 1-hashing matrices defined in Definition 1.2.4.

Still, in general, the optimal dimension dependence present in Gaussian subspace embeddings cannot be replicated even for hashing distributions, as our next example illustrates.

Example 3 (The 1-hashing matrix distributions fails to yield an oblivious subspace embedding with $m = \mathcal{O}(r)$). Consider the matrix

$$A = \begin{pmatrix} I_{r \times r} & 0 \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{n \times d}. \quad (2.2.8)$$

If S is a 1-hashing matrix with $m = \mathcal{O}(r)$, then $SA = (S_1 \ 0)$, where the S_1 block contains the first r columns of S . To ensure that the rank of A is preserved (cf. Lemma 2.2.1), a necessary condition for S_1 to have rank r is that the r non-zeros of S_1 are in different rows. Since, by definition, the

respective row is chosen independently and uniformly at random for each column of S , the probability of S_1 having rank r is no greater than

$$\left(1 - \frac{1}{m}\right) \cdot \left(1 - \frac{2}{m}\right) \cdot \dots \cdot \left(1 - \frac{r-1}{m}\right) \leq e^{-\frac{1}{m} - \frac{2}{m} - \dots - \frac{r-1}{m}} = e^{-\frac{r(r-1)}{2m}}, \quad (2.2.9)$$

For the probability⁴ (2.2.9) to be at least $1/2$, we must have $m \geq \frac{r(r-1)}{2 \log(2)}$.

The above example improves upon the lower bound in Nelson et al. [80] by slightly relaxing the requirements on m and n ⁵. We note that in the order of r (or equivalently⁶, d), the lower bound $m = \mathcal{O}(r^2)$ for 1-hashing matches the upper bound given in Nelson and Nguyen [79], Meng and Mahoney [77].

When S is an s -hashing matrix, with $s > 1$, the tight bound $m = \Theta(r^2)$ can be improved to $m = \Theta(r \log r)$ for s sufficiently large. In particular, Cohen [22] derived a general upper bound that implies, for example, subspace embedding properties of s -hashing matrices provided $m = \mathcal{O}(r \log r)$ and $s = \mathcal{O}(\log r)$; the value of s may be further reduced to a constant (that is not equal to 1) at the expense of increasing m and worsening its dependence of d . A lower bound for guaranteeing oblivious embedding properties of s -hashing matrices is given in [81]. Thus we can see that for s -hashing (and especially for 1-hashing) matrices, their general subspace embedding properties are suboptimal in terms of the dependence of m on d when compared to the Gaussian sketching results. To improve the embedding properties of hashing matrices, we must focus on special structure input matrices.

2.2.3.1 Coherence-dependent embedding properties of sparse random matrices

A feature of the problematic matrix (2.2.8) is that its rows are separated into two groups, with the first r rows containing all the information. If the rows of A were more ‘uniform’ in the sense of equally important in terms of relevant information content, hashing may perform better as a sketching matrix. Interestingly, it is not the uniformity of the rows of A but the uniformity of the rows of U , the left singular matrix from the compact SVD of A , that plays an important role. The concept of coherence is a useful proxy for the uniformity of the rows of U and A ⁷.

Definition 2.2.5. (*Matrix coherence [73]*) The coherence of a matrix $A \in \mathbb{R}^{n \times d}$, denoted $\mu(A)$, is the largest Euclidean norm of the rows of U defined in (2.2.4). Namely,

$$\mu(A) = \max_{i \in [n]} \|U_i\|_2, \quad (2.2.10)$$

⁴The argument in the example relating to 1-hashing sketching is related to the birthday paradox, as mentioned (but not proved) in Nelson and Nguyen [80].

⁵We note that in fact, [79] considers a more general set up, namely, any matrix distribution with column sparsity one.

⁶See Remark 1.

⁷We note that sampling matrices were shown to have good subspace embedding properties for input matrices with low coherence [2, 97]. Even if the coherence is minimal, the size of the sampling matrix has a $d \log d$ dependence where the $\log d$ term cannot be removed due to the coupon collector problem [97].

where U_i denotes the i th row of U .⁸

Some useful properties follow.

Lemma 2.2.3. *Let $A \in \mathbb{R}^{n \times d}$ have rank $r \leq d \leq n$. Then*

$$\sqrt{\frac{r}{n}} \leq \mu(A) \leq 1. \quad (2.2.11)$$

Furthermore, if $\mu(A) = \sqrt{\frac{r}{n}}$, then $\|U_i\|_2 = \sqrt{\frac{r}{n}}$ for all $i \in [n]$ where U is defined in (2.2.4).

Proof. Since the matrix $U \in \mathbb{R}^{n \times r}$ has orthonormal columns, we have that

$$\sum_{i=1}^n \|U_i\|_2^2 = r. \quad (2.2.12)$$

Therefore the maximum 2-norm of U must not be less than $\sqrt{\frac{r}{n}}$, and thus $\mu(A) \geq \sqrt{\frac{r}{n}}$. Furthermore, if $\mu(A) = \sqrt{\frac{r}{n}}$, then (2.2.12) implies $\|U_i\|_2 = \sqrt{\frac{r}{n}}$ for all $i \in [n]$.

Next, by expanding the set of columns of U to a basis of \mathbb{R}^n , there exists $U_f \in \mathbb{R}^{n \times n}$ such that $U_f = (U \hat{U})$ orthogonal where $\hat{U} \in \mathbb{R}^{n \times (n-d)}$ has orthonormal columns. The 2-norm of i th row of U is bounded above by the 2-norm of i th row of U_f , which is one. Hence $\mu(A) \leq 1$. \square

We note that for A in (2.2.8), we have $\mu(A) = 1$. The maximal coherence of this matrix sheds some light on the ensuing poor embedding properties we noticed in Example 1.

Bourgain et al [10] gives a general result that captures the coherence-restricted subspace embedding properties of s -hashing matrices.

Theorem 2.2.1 (Bourgain et al [10]). *Let $A \in \mathbb{R}^{n \times d}$ with coherence $\mu(A)$ and rank r ; and let $0 < \epsilon, \delta < 1$. Assume also that*

$$m \geq c_1 \max \left\{ \delta^{-1}, \left[(r + \log m) \min \{ \log^2(r/\epsilon), \log^2(m) \} + r \log(1/\delta) \right] \epsilon^{-2} \right\} \quad (2.2.13)$$

$$\text{and } s \geq c_2 \left[\log(m) \log(1/\delta) \min \{ \log^2(r/\epsilon), \log^2(m) \} + \log^2(1/\delta) \right] \mu(A)^2 \epsilon^{-2}, \quad (2.2.14)$$

where c_1 and c_2 are positive constants. Then a(ny) s -hashing matrix $S \in \mathbb{R}^{m \times n}$ is an ϵ -subspace embedding for A with probability at least $1 - \delta$.

Substituting $s = 1$ in (2.2.14), we can use the above Theorem to deduce an upper bound μ of acceptable coherence values of the input matrix A , namely,

$$\mu(A) \leq c_2^{-1/2} \epsilon \left[\log(m) \log(1/\delta) \min \{ \log^2(r/\epsilon), \log^2(m) \} + \log^2(1/\delta) \right]^{-1/2} := \mu. \quad (2.2.15)$$

Thus Theorem 2.2.1 implies that the distribution of 1-hashing matrices with m satisfying (2.2.13) is an oblivious subspace embedding for any input matrix A with $\mu(A) \leq \mu$, where μ is defined in (2.2.15).

⁸Note that the concept of coherence is different to the (in)coherence used in compressed sensing literature [31], in particular our notion of coherence is not invariant under a different coordinate representation.

2.2.3.2 Non-uniformity of vectors and their relation to embedding properties of sparse random matrices

In order to prove some of our main results, we need a corresponding notion of coherence of vectors, to be able to measure the ‘importance’ of their respective entries; this is captured by the so-called non-uniformity of a vector.

Definition 2.2.6 (Non-uniformity of a vector). *Given $x \in \mathbb{R}^n$, the non-uniformity of x , $\nu(x)$, is defined as*

$$\nu(x) = \frac{\|x\|_\infty}{\|x\|_2}. \quad (2.2.16)$$

We note that for any vector $x \in \mathbb{R}^n$, we have $\frac{1}{\sqrt{n}} \leq \nu(x) \leq 1$.

Lemma 2.2.4. *Given $A \in \mathbb{R}^{n \times d}$, let $y = Ax$ for some $x \in \mathbb{R}^d$. Then*

$$\nu(y) \leq \mu(A). \quad (2.2.17)$$

Proof. Let $A = U\Sigma V^T$ be defined as in (2.2.4), and let $z = \Sigma V^T x \in \mathbb{R}^r$. Then $y = Ax = Uz$. Therefore

$$\|y\|_\infty = \|Uz\|_\infty = \max_{1 \leq i \leq n} |\langle U_i, z \rangle| \leq \max_{1 \leq i \leq n} \|U_i\|_2 \|z\|_2 \leq \mu(A) \|z\|_2, \quad (2.2.18)$$

where U_i denotes the i^{th} row of U . Furthermore, $\|y\|_2 = \|Uz\|_2 = \|z\|_2$ which then implies $\nu(y) = \|y\|_\infty / \|y\|_2 \leq \mu(A)$. \square

The next lemmas are crucial to our results in the next section; the proof of the first lemma can be found in the paper [35].

We also note, in subsequent results, the presence of *problem-independent constants*, also called absolute constants that will be implicitly or explicitly defined, depending on the context. Our convention here is as expected, that the same notation denotes the same constant across all results in this chapter.

The following expression will be needed in our results,

$$\bar{\nu}(\epsilon, \delta) := C_1 \sqrt{\epsilon} \min \left\{ \frac{\log(E/\epsilon)}{\log(1/\delta)}, \sqrt{\frac{\log(E)}{\log(1/\delta)}} \right\}, \quad (2.2.19)$$

where $\epsilon, \delta \in (0, 1)$ and $E, C_1 > 0$.

Lemma 2.2.5 ([35], Theorem 2). *Suppose that $\epsilon, \delta \in (0, 1)$, and E satisfies $C \leq E < \frac{2}{\delta \log(1/\delta)}$, where $C > 0$ and C_1 are problem-independent constants. Let $m \leq n \in \mathbb{N}$ with $m \geq E\epsilon^{-2} \log(1/\delta)$.*

Then, for any $x \in \mathbb{R}^n$ with

$$\nu(x) \leq \bar{\nu}(\epsilon, \delta), \quad (2.2.20)$$

where $\bar{\nu}(\epsilon, \delta)$ is defined in (2.2.19), a randomly generated 1-hashing matrix $S \in \mathbb{R}^{m \times n}$ is an ϵ -JL embedding for $\{x\}$ with probability at least $1 - \delta$.

Lemma 2.2.6. Let $\epsilon, \delta \in (0, 1)$, $\nu \in (0, 1]$ and $m, n \in \mathbb{N}$. Let \mathcal{S} be a distribution of $m \times n$ random matrices. Suppose that for any given y with $\nu(y) \leq \nu$, a matrix $S \in \mathbb{R}^{m \times n}$ randomly drawn from \mathcal{S} is an ϵ -JL embedding for $\{y\}$ with probability at least $1 - \delta$. Then for any given set $Y \subseteq \mathbb{R}^n$ with $\max_{y \in Y} \nu(y) \leq \nu$ and cardinality $|Y| \leq 1/\delta$, a matrix S randomly drawn from \mathcal{S} is an ϵ -JL embedding for Y with probability at least $1 - |Y|\delta$.

Proof. Let $Y = \{y_1, y_2, \dots, y_{|Y|}\}$. Let B_i be the event that S is an ϵ -JL embedding for $y_i \in Y$. Then $\mathbb{P}(B_i) \geq 1 - \delta$ by assumption. We have

$$\mathbb{P}(S \text{ is an } \epsilon\text{-JL embedding for } Y) = \mathbb{P}(\cap_i B_i) = 1 - \mathbb{P}((\cap_i B_i)^c) \quad (2.2.21)$$

$$= 1 - \mathbb{P}(\cup_i B_i^c) \quad (2.2.22)$$

$$\geq 1 - \sum_i \mathbb{P}(B_i^c) \quad (2.2.23)$$

$$= 1 - \sum_i [1 - \mathbb{P}(B_i)] \quad (2.2.24)$$

$$\geq 1 - \sum_i [1 - (1 - \delta)] = 1 - |Y|\delta. \quad (2.2.25)$$

□

Lemma 2.2.7. Let $\epsilon \in (0, 1)$, and $Y \subseteq \mathbb{R}^n$ be a finite set such that $\|y\|_2 = 1$ for each $y \in Y$. Define

$$Y_+ = \{y_1 + y_2 : y_1, y_2 \in Y\} \quad (2.2.26)$$

$$Y_- = \{y_1 - y_2 : y_1, y_2 \in Y\}. \quad (2.2.27)$$

If $S \in \mathbb{R}^{m \times n}$ is an ϵ -JL embedding for $\{Y_+ \cup Y_-\}$, then S is a generalised ϵ -JL embedding for Y .

Proof. Let $y_1, y_2 \in Y$. We have that

$$\begin{aligned} |\langle Sy_1, Sy_2 \rangle - \langle y_1, y_2 \rangle| &= |(\|S(y_1 + y_2)\|^2 - \|S(y_1 - y_2)\|^2) \\ &\quad - (\|(y_1 + y_2)\|^2 - \|(y_1 - y_2)\|^2)|/4 \\ &\leq \epsilon(\|(y_1 + y_2)\|^2 + \|(y_1 - y_2)\|^2)/4 \\ &= \epsilon(\|y_1\|^2 + \|y_2\|^2)/2 \\ &= \epsilon, \end{aligned} \quad (2.2.28)$$

where to obtain the inequality, we use that S is an ϵ -JL embedding for $\{Y_+ \cup Y_-\}$; the last equality follows from $\|y_1\|_2 = \|y_2\|_2 = 1$. □

2.3 Hashing sketching with $m = \mathcal{O}(r)$

Our first result shows that if the coherence of the input matrix is sufficiently low, the distribution of 1-hashing matrices with $m = \mathcal{O}(r)$ is an (ϵ, δ) -oblivious subspace embedding.

The following expression will be useful later,

$$\bar{\mu}(\epsilon, \delta) := C_1 \sqrt{C_2} \sqrt{\epsilon} \min \left\{ \frac{\log(E/(C_2\epsilon))}{4r + \log(1/\delta)}, \sqrt{\frac{\log(E)}{4r + \log(1/\delta)}} \right\}, \quad (2.3.1)$$

where $\epsilon, \delta \in (0, 1)$, $r, E > 0$ are to be chosen/defined depending on the context, and $C_1, C_2 > 0$ are problem-independent constants.

Theorem 2.3.1. *Suppose that $\epsilon, \delta \in (0, 1)$, $r \leq d \leq n, m \leq n \in \mathbb{N}$, $E > 0$ satisfy*

$$C \leq E \leq \frac{2e^{4r}}{[4r + \log(1/\delta)] \delta}, \quad (2.3.2)$$

$$m \geq EC_2^{-2} \epsilon^{-2} [4r + \log(1/\delta)], \quad (2.3.3)$$

where $C > 0$ and $C_1, C_2 > 0$ are problem-independent constants. Then for any matrix $A \in \mathbb{R}^{n \times d}$ with rank r and

$$\mu(A) \leq \bar{\mu}(\epsilon, \delta), \quad (2.3.4)$$

where $\bar{\mu}(\epsilon, \delta)$ is defined in (2.3.1), a randomly generated 1-hashing matrix $S \in \mathbb{R}^{m \times n}$ is an ϵ -subspace embedding for A with probability at least $1 - \delta$.

The proof of Theorem 2.3.1 relies on the fact that the coherence of the input matrix gives a bound on the non-uniformity of the entries for all vectors in its column space (Lemma 2.2.4), adapting standard arguments in [101] involving set covers.

Definition 2.3.1. *A γ -cover of a set M is a subset $N \subseteq M$ with the property that given any point $y \in M$, there exists $w \in N$ such that $\|y - w\|_2 \leq \gamma$.*

Consider a given real matrix $U \in \mathbb{R}^{n \times r}$ with orthonormal columns, and let

$$M := \{Uz \in \mathbb{R}^n : z \in \mathbb{R}^r, \|z\|_2 = 1\}. \quad (2.3.5)$$

The next two lemmas show the existence of a γ -net N for M , and connect generalised JL embeddings for N with JL embeddings for M .

Lemma 2.3.1. *Let $0 < \gamma < 1$, $U \in \mathbb{R}^{n \times r}$ have orthonormal columns and M be defined in (2.3.5). Then there exists a γ -cover N of M such that $|N| \leq (1 + \frac{2}{\gamma})^r$.*

Proof. Let $\tilde{M} = \{z \in \mathbb{R}^r : \|z\|_2 = 1\}$. Let $\tilde{N} \subseteq \tilde{M}$ be the maximal set such that no two points in \tilde{N} are within distance γ from each other. Then it follows that the r -dimensional balls centred at points in \tilde{N} with radius $\gamma/2$ are all disjoint and contained in the r -dimensional ball centred at the origin with radius $(1 + \gamma/2)$. Hence

$$\frac{\text{Volume of the } r\text{-dimensional ball centred at the origin with radius } (1 + \gamma/2)}{\text{Total volume of the } r\text{-dimensional balls centred at points in } \tilde{N} \text{ with radius } \gamma/2} \quad (2.3.6)$$

$$= \frac{1}{|\tilde{N}|} \frac{(1 + \frac{\gamma}{2})^r}{(\frac{\gamma}{2})^r} \geq 1, \quad (2.3.7)$$

which implies $|\tilde{N}| \leq (1 + \frac{2}{\gamma})^r$.

Let $N = \{Uz \in \mathbb{R}^n : z \in \tilde{N}\}$. Then $|N| \leq |\tilde{N}| \leq (1 + \frac{2}{\gamma})^r$ and we show N is a γ -cover for M . Given $y_M \in M$, there exists $z_M \in \tilde{M}$ such that $y_M = Uz_M$. By definition of \tilde{N} , there must be $z_N \in \tilde{N}$ such that $\|z_M - z_N\|_2 \leq \gamma$ as otherwise \tilde{N} would not be maximal. Let $y_N = Uz_N \in N$. Since U has orthonormal columns, we have $\|y_M - y_N\|_2 = \|z_M - z_N\|_2 \leq \gamma$.

□

Lemma 2.3.2. *Let $\epsilon, \gamma \in (0, 1)$, $U \in \mathbb{R}^{n \times d}$, $M \subseteq \mathbb{R}^n$ associated with U be defined in (2.3.5). Suppose $N \subseteq M$ is a γ -cover of M and $S \in \mathbb{R}^{m \times n}$ is a generalised ϵ_1 -JL embedding for N , where $\epsilon_1 = \frac{(1-\gamma)(1-\gamma^2)}{1+2\gamma-\gamma^2}\epsilon$. Then S is an ϵ -JL embedding for M .*

To prove Lemma 2.3.2, we need the following Lemma.

Lemma 2.3.3. *Let $\gamma \in (0, 1)$, $U \in \mathbb{R}^{n \times r}$ having orthonormal columns and $M \subseteq \mathbb{R}^n$ associated with U be defined in (2.3.5). Let N be a γ -cover of M , $y \in M$. Then for any $k \in \mathbb{N}$, there exists $\alpha_0, \alpha_1, \dots, \alpha_k \in \mathbb{R}$, $y_0, y_1, y_2, \dots, y_k \in N$ such that*

$$\|y - \sum_{i=0}^k \alpha_i y_i\|_2 \leq \gamma^{k+1}, \quad (2.3.8)$$

$$|\alpha_i| \leq \gamma^i, i = 0, 1, \dots, k. \quad (2.3.9)$$

Proof. We use induction. Let $k = 0$. Then by definition of a γ -cover, there exists $y_0 \in N$ such that $\|y - y_0\| < \gamma$. Letting $\alpha_0 = 1$, we have covered the $k = 0$ case.

Now assume (2.3.8) and (2.3.9) are true for $k = K \in \mathbb{N}$. Namely there exists $\alpha_0, \alpha_1, \dots, \alpha_K \in \mathbb{R}$, $y_0, y_1, y_2, \dots, y_K \in N$ such that

$$\|y - \sum_{i=0}^K \alpha_i y_i\|_2 \leq \gamma^{K+1} \quad (2.3.10)$$

$$|\alpha_i| \leq \gamma^i, i = 0, 1, \dots, K. \quad (2.3.11)$$

Because $y, y_0, y_1, \dots, y_K \in N \subseteq M$, there exists $z, z_0, z_1, \dots, z_K \in \mathbb{R}^r$ such that $y = Uz, y_0 = Uz_0, y_1 = Uz_1, \dots, y_K = Uz_K$ with $\|z\| = \|z_0\| = \dots = \|z_K\| = 1$. Therefore

$$\frac{y - \sum_{i=0}^K \alpha_i y_i}{\|y - \sum_{i=0}^K \alpha_i y_i\|_2} = \frac{U(z - \sum_{i=0}^K \alpha_i z_i)}{\|U(z - \sum_{i=0}^K \alpha_i z_i)\|_2} = U \frac{z - \sum_{i=0}^K \alpha_i z_i}{\|z - \sum_{i=0}^K \alpha_i z_i\|_2} \in M, \quad (2.3.12)$$

where we have used that the columns of U are orthonormal.

Since N is a γ -cover for M , there exists $y_{K+1} \in N$ such that

$$\left\| \frac{y - \sum_{i=0}^K \alpha_i y_i}{\|y - \sum_{i=0}^K \alpha_i y_i\|_2} - y_{K+1} \right\|_2 \leq \gamma. \quad (2.3.13)$$

Multiplying both sides by $\alpha_{K+1} := \|y - \sum_{i=0}^K \alpha_i y_i\|_2 \leq \gamma^{K+1}$, we have

$$\|y - \sum_{i=0}^{K+1} \alpha_i y_i\|_2 \leq \gamma^{K+2}, \quad (2.3.14)$$

$$|\alpha_i| \leq \gamma^i, i = 0, 1, \dots, K+1. \quad (2.3.15)$$

□

Proof of Lemma 2.3.2. Let $y \in M$ and $k \in \mathbb{N}$, and consider the approximate representation of y provided in Lemma 2.3.3, namely, assume that (2.3.8) and (2.3.9) hold. Then we have

$$\begin{aligned} \|S \sum_{i=0}^k \alpha_i y_i\|_2^2 &= \sum_{i=0}^k \|S \alpha_i y_i\|_2^2 + \sum_{0 \leq i < j \leq k} 2 \langle S \alpha_i y_i, S \alpha_j y_j \rangle \\ &= \sum_{i=0}^k \|S \alpha_i y_i\|_2^2 + \sum_{0 \leq i < j \leq k} 2 \langle \alpha_i y_i, \alpha_j y_j \rangle + \\ &\quad + \left[\sum_{0 \leq i < j \leq k} 2 \langle S \alpha_i y_i, S \alpha_j y_j \rangle - \sum_{0 \leq i < j \leq k} 2 \langle \alpha_i y_i, \alpha_j y_j \rangle \right] \\ &\leq (1 + \epsilon_1) \sum_{i=0}^k \|\alpha_i y_i\|_2^2 + \sum_{0 \leq i < j \leq k} 2 \langle \alpha_i y_i, \alpha_j y_j \rangle + 2 \sum_{0 \leq i < j \leq k} \epsilon_1 \|\alpha_i y_i\|_2 \|\alpha_j y_j\|_2 \\ &= \left\| \sum_{i=0}^k \alpha_i y_i \right\|_2^2 + \epsilon_1 \left[\sum_{i=0}^k \|\alpha_i y_i\|_2^2 + 2 \sum_{0 \leq i < j \leq k} \|\alpha_i y_i\|_2 \|\alpha_j y_j\|_2 \right], \end{aligned}$$

where to deduce the inequality, we use that S is a generalised ϵ_1 -JL embedding for N . Using $\|y_i\|_2 = 1$ and $|\alpha_i| \leq \gamma^i$, we have

$$\begin{aligned} \frac{1}{\epsilon_1} \left\{ \left\| S \sum_{i=0}^k \alpha_i y_i \right\|_2^2 - \left\| \sum_{i=0}^k \alpha_i y_i \right\|_2^2 \right\} &= \sum_{i=0}^k \|\alpha_i y_i\|_2^2 + 2 \sum_{0 \leq i < j \leq k} \|\alpha_i y_i\|_2 \|\alpha_j y_j\|_2 \\ &\leq \sum_{i=0}^k \gamma^{2i} + 2 \sum_{0 \leq i < j \leq k} \gamma^i \gamma^j \\ &\leq \frac{1 - \gamma^{k+1}}{1 - \gamma} + \frac{2\gamma(1 - \gamma^{k-i})(1 - \gamma^{2k})}{(1 - \gamma)(1 - \gamma^2)}, \end{aligned} \quad (2.3.16)$$

where we have used

$$\begin{aligned} \sum_{0 \leq i < j \leq k} \gamma^i \gamma^j &= \sum_{i=0}^{k-1} \gamma^i \sum_{j=i+1}^k \gamma^j = \sum_{i=0}^{k-1} \gamma^{2i+1} \sum_{j=0}^{k-i-1} \gamma^j \\ &= \frac{\gamma(1 - \gamma^{k-i})}{1 - \gamma} \sum_{i=0}^{k-1} \gamma^{2i} = \frac{\gamma(1 - \gamma^{k-i})(1 - \gamma^{2k})}{(1 - \gamma)(1 - \gamma^2)}. \end{aligned}$$

Letting $k \rightarrow \infty$ in (2.3.16), we deduce

$$\frac{1}{\epsilon_1} \left\{ \left\| S \sum_{i=0}^{\infty} \alpha_i y_i \right\|_2^2 - \left\| \sum_{i=0}^{\infty} \alpha_i y_i \right\|_2^2 \right\} \leq \frac{1}{1 - \gamma} + \frac{2\gamma}{(1 - \gamma)(1 - \gamma^2)} = \frac{1 + 2\gamma - \gamma^2}{(1 - \gamma)(1 - \gamma^2)},$$

Letting $k \rightarrow \infty$ in (2.3.8) implies $y = \sum_{i=0}^{\infty} \alpha_i y_i$, and so the above gives

$$\|Sy\|_2^2 - \|y\|_2^2 \leq \epsilon_1 \frac{1 + 2\gamma - \gamma^2}{(1 - \gamma)(1 - \gamma^2)} = \epsilon \|y\|_2^2,$$

where to get the equality, we used $\|y\|_2 = 1$ and the definition of ϵ_1 . The lower bound in the ϵ_1 -JL embedding follows similarly. \square

We are ready to prove Theorem 2.3.1.

Proof of Theorem 2.3.1. Let $A \in \mathbb{R}^{n \times d}$ with rank r and satisfying (2.3.4). Let $U \in \mathbb{R}^{n \times r}$ be an SVD factor of A as defined in (2.2.4), which by definition of coherence, implies

$$\mu(U) = \mu(A) \leq \bar{\mu}(\epsilon, \delta), \quad (2.3.17)$$

where $\bar{\mu}(\epsilon, \delta)$ is defined in (2.3.1). We let $\gamma, \epsilon_1, \delta_1 \in (0, 1)$ be defined as

$$\gamma = \frac{2}{e^2 - 1}, \quad C_2 = \frac{(1 - \gamma)(1 - \gamma^2)}{1 + 2\gamma - \gamma^2}, \quad \epsilon_1 = C_2 \epsilon \quad \text{and} \quad \delta_1 = e^{-4r} \delta, \quad (2.3.18)$$

and note that $C_2 \in (0, 1)$ and

$$\bar{\nu}(\epsilon_1, \delta_1) = \bar{\mu}(\epsilon, \delta), \quad (2.3.19)$$

where $\bar{\nu}(\cdot, \cdot)$ is defined in (2.2.19). Let $M \in \mathbb{R}^n$ be associated to U as in (2.3.5) and let $N \subseteq M$ be the γ -cover of M as guaranteed by Lemma 2.3.1, with γ defined in (2.3.18) which implies that $|N| \leq e^{2r}$.

Let $S \in \mathbb{R}^{m \times n}$ be a randomly generated 1-hashing matrix with $m \geq E\epsilon_1^{-2} \log(1/\delta_1) = EC_2^{-2} \epsilon^{-2} [4r + \log(1/\delta)]$, where to obtain the last equality, we used (2.3.18).

To show that the sketching matrix S is an ϵ -subspace embedding for A (with probability at least $1 - \delta$), it is sufficient to show that S is an ϵ_1 -generalised JL embedding for $N \subseteq M$ (with probability at least $1 - \delta$). To see this, recall (2.3.5) and Lemma 2.2.2(ii) which show that S is an ϵ -subspace embedding for A if and only if S is an ϵ -JL embedding for M . Our sufficiency claim now follows by invoking Lemma 2.3.2 for S , N and M .

We are left with considering in detail the cover set $N = \{y_1, y_2, \dots, y_{|N|}\}$ and the following useful ensuing sets

$$\begin{aligned} N_+ &= \{y_i + y_j : i, j \in [1, |N|]\} \quad \text{and} \quad N_- = \{y_i - y_j : i, j \in [1, |N|]\}, \\ N_-^{(1)} &= \{y_i - y_j : 1 \leq i < j \leq |N|\} \quad \text{and} \quad N_-^{(2)} = \{y_i - y_j : 1 \leq j < i \leq |N|\}. \end{aligned}$$

Now let $Y := N_+ \cup N_-^{(1)}$ and show that

$$\nu(y) \leq \bar{\nu}(\epsilon_1, \delta_1) \quad \text{for all} \quad y \in Y. \quad (2.3.20)$$

To see this, assume first that $y = y_i + y_j \in N_+$, with $y_i, y_j \in N \subseteq M$. Thus there exist $z_i, z_j \in R^r$ such that $y_i = Uz_i$ and $y_j = Uz_j$, and so $y = U(z_i + z_j)$. Using Lemma 2.2.4, $\nu(y) \leq \mu(U) = \mu(A)$,

which together with (2.3.17) and (2.3.19), gives (2.3.20) for points $y \in N_+$; the proof for $y \in N_-^{(1)}$ follows similarly.

Lemma 2.2.5 with $(\epsilon, \delta) := (\epsilon_1, \delta_1)$ provides that for any $x \in \mathbb{R}^n$ with $\nu(x) \leq \bar{\nu}(\epsilon_1, \delta_1)$, S is an ϵ_1 -JL embedding for $\{x\}$ with probability at least $1 - \delta_1$. This and (2.3.20) imply that the conditions of Lemma 2.2.6 are satisfied for $Y = N_+ \cup N_-^{(1)}$, from which we conclude that S is an ϵ_1 -JL embedding for Y with probability at least $1 - |Y|\delta_1$. Note that

$$|Y| \leq |N_+| + |N_-^{(1)}| \leq \frac{1}{2}|N|(|N|+1) + \frac{1}{2}|N|(|N|-1) = |N|^2.$$

This, the definition of δ_1 in (2.3.18) and $|N| \leq e^{2r}$ imply that $1 - |Y|\delta_1 \geq 1 - \delta$. Therefore S is an ϵ_1 -JL embedding for $N_+ \cup N_-^{(1)}$ with probability at least $1 - \delta$.

Finally, Definition 2.2.2 of JL-embeddings implies that the sign of the embedded vector is irrelevant and that $\{0\}$ is always embedded, and so if S is an ϵ_1 -JL embedding for $N_+ \cup N_-^{(1)}$, it is also an ϵ_1 -JL embedding for $N_+ \cup N_-$. Lemma 2.2.7 now provides us with the desired result that then, S is a generalised ϵ_1 -JL embedding for N . \square

Next we discuss the results in Theorem 2.3.1.

Conditions for a well-defined coherence requirement While our result guarantees optimal dimensionality reduction for the sketched matrix, using a very sparse 1-hashing matrix for the sketch, it imposes implicit restrictions on the number n of rows of A . Recalling (2.2.11), we note that condition (2.3.4) is well-defined when

$$\sqrt{\frac{r}{n}} \leq \bar{\mu}(\epsilon, \delta). \quad (2.3.21)$$

Using the definition of $\bar{\mu}(\epsilon, \delta)$ in (2.3.1) and assuming reasonably that $\log(1/\delta) = \mathcal{O}(r)$, we have the lower bound

$$\bar{\mu}(\epsilon, \delta) \geq C_1 \sqrt{C_2} \sqrt{\epsilon} \frac{\min \{ \log(E/(C_2\epsilon)), \sqrt{\log E} \}}{4r + \log(1/\delta)},$$

and so (2.3.21) is satisfied if

$$n \geq \frac{r(4r + \log(1/\delta))^2}{C_1^2 C_2 \epsilon \min \{ \log^2(E/(C_2\epsilon)), \log E \}} = \mathcal{O} \left(\frac{r^3}{\epsilon \log^2(\epsilon)} \right). \quad (2.3.22)$$

Comparison with data-independent bounds Existing results show that $m = \Theta(r^2)$ is both necessary and sufficient in order to secure an oblivious subspace embedding property for 1-hashing matrices with no restriction on the coherence of the input matrix [80, 79, 77]. Aside from requiring more projected rows than in Theorem 2.3.1, these results implicitly impose $n \geq \mathcal{O}(r^2)$ for the size/rank of data matrix in order to secure meaningful dimensionality reduction.

Table 2.1: Summary of results for 1-hashing

Result	μ (coherence of A)	m (size of sketching S)
[77]	–	$\Theta(r^2)$
[10]	$\mathcal{O}\left(\log^{-3/2}(r)\right)$	$\mathcal{O}\left(r \log^2(r)\right)$
Theorem 2.3.1	$\mathcal{O}\left(r^{-1}\right)$	$\mathcal{O}(r)$

Comparison with data-dependent bounds To the best of our knowledge, the only data-dependent result for hashing matrices is [10] (see Theorem 2.2.1). From (2.2.13), we have that $m \geq c_1 r \min\{\log^2(r/\epsilon), \log^2(m)\} \epsilon^{-2}$ and hence $m = \Omega(r \log^2 r)$; while Theorem 2.3.1 only needs $m = \mathcal{O}(r)$. However, the coherence requirement on A in Theorem 2.2.1 is weaker than (2.3.4) and so [10] applies to a wider range of inputs at the expense of a larger value of m required for the sketching matrix.

Summary and look ahead Table 2.1 summarises existing results and we see stricter coherence assumptions lead to improved dimensionality reduction properties. In the next section, we investigate relaxing coherence requirements by using hashing matrices with increased column sparsity (s -hashing) and coherence reduction transformations.

2.4 Relaxing the coherence requirement using s -hashing matrices

This section investigates the embedding properties of s -hashing matrices when $s \geq 1$. Indeed, [10] shows that s -hashing relaxes their particular coherence requirement by \sqrt{s} . Theorem 2.4.1 presents a similar result for our particular coherence requirement (2.3.4) that again guarantees embedding properties for $m = \mathcal{O}(r)$. Then we present a new s -hashing variant that allows us to give a general result showing that (any) subspace embedding properties of 1-hashing matrices immediately translate into similar properties for these s -hashing matrices when applied to a larger class of data matrices, with larger coherence. A simplified embedding result with $m = \mathcal{O}(r)$ is then deduced for this s -hashing variant. Finally, s -hashing or s -hashing variant is combined with the randomised Hadamard transform with Theorem 2.4.4 and 2.4.5 guaranteeing embedding properties for $m = \mathcal{O}(r)$ given that the data matrix A has $n = \mathcal{O}(r^3)$.

Numerical benefits of s -hashing (for improved preconditioning) are investigated in later sections; see Figures 3.2 and 3.3 for example.

2.4.1 The embedding properties of s -hashing matrices

Our next result shows that using s -hashing (Definition 1.2.5) relaxes the particular coherence requirement in Theorem 2.3.1 by \sqrt{s} .

Theorem 2.4.1. *Let $r \leq d \leq n \in \mathbb{N}^+$. Let $C_1, C_2, C_3, C_M, C_\nu, C_s > 0$ be problem-independent constants. Suppose that $\epsilon, \delta \in (0, C_3)$, $m, s \in \mathbb{N}^+$ and $E > 0$ satisfy⁹*

$$1 \leq s \leq C_s C_2^{-1} \epsilon^{-1} [4r + \log(1/\delta)], \quad (2.4.1)$$

$$C_M \leq E \leq C_2^2 \epsilon^2 s [4r + \log(1/\delta)]^{-1} e^{C_s (C_2 \epsilon s)^{-1} [4r + \log(1/\delta)]}, \quad (2.4.2)$$

$$m \geq \{EC_2^{-2} \epsilon^{-2} [4r + \log(1/\delta)], se\}. \quad (2.4.3)$$

Then for any matrix $A \in \mathbb{R}^{n \times d}$ with rank r and $\mu(A) \leq \sqrt{s} C_\nu C_1^{-1} \bar{\mu}(\epsilon, \delta)$, where $\bar{\mu}(\epsilon, \delta)$ is defined in (2.3.1), a randomly generated s -hashing matrix $S \in \mathbb{R}^{m \times n}$ is an ϵ -subspace embedding for A with probability at least $1 - \delta$.

Theorem 2.4.1 parallels Theorem 2.3.1; and its proof relies on the following lemma which parallels Lemma 2.2.5.

Lemma 2.4.1 ([54], Theorem 1.5). *Let $C_1, C_3, C_M, C_\nu, C_s > 0$ be problem-independent constants. Suppose that $\epsilon, \delta \in (0, C_3)$, $m, s \in \mathbb{N}^+$, $E \in \mathbb{R}$ satisfy*

$$1 \leq s \leq C_s \epsilon^{-1} \log(1/\delta),$$

$$C_M \leq E < \epsilon^2 s \log^{-1}(1/\delta) e^{C_s (\epsilon s)^{-1} \log(1/\delta)},$$

$$m \geq \max \{E \epsilon^{-2} \log(1/\delta), se\}.$$

Then for any $x \in \mathbb{R}^n$ with $\nu(x) \leq \sqrt{s} C_\nu C_1^{-1} \bar{\nu}(\epsilon, \delta)$, where $\bar{\nu}(\epsilon, \delta)$ is defined in (2.2.19), a randomly generated s -hashing matrix $S \in \mathbb{R}^{m \times n}$ is an ϵ -JL embedding for $\{x\}$ with probability at least $1 - \delta$.

The proof of Theorem 2.4.1 follows the same argument as Theorem 2.3.1, replacing 1-hashing with s -hashing and using Lemma 2.4.1 instead of Lemma 2.2.5. We omit the details.

2.4.2 A general embedding property for an s -hashing variant

Note that in both Theorem 2.2.1 and Theorem 2.4.1, allowing column sparsity of hashing matrices to increase from 1 to s results in coherence requirements being relaxed by \sqrt{s} . We introduce an s -hashing variant that allows us to generalise this result.

Definition 2.4.1. *We say $T \in \mathbb{R}^{m \times n}$ is an s -hashing variant matrix if independently for each $j \in [n]$, we sample with replacement $i_1, i_2, \dots, i_s \in [m]$ uniformly at random and add $\pm 1/\sqrt{s}$ to $T_{i_k j}$, where $k = 1, 2, \dots, s$.¹⁰*

Both s -hashing and s -hashing variant matrices reduce to 1-hashing matrices when $s = 1$. For $s \geq 1$, the s -hashing variant has at most s non-zeros per column, while the usual s -hashing matrix has precisely s nonzero entries per same column.

The next lemma connects s -hashing variant matrices to 1-hashing matrices.

⁹Note that the expressions of the lower bounds in (2.3.3) and (2.4.3) are identical apart from the choice of E and the condition $m \geq se$.

¹⁰We add $\pm 1/\sqrt{s}$ to $T_{i_k j}$ because we may have $i_k = i_l$ for some $l < k$, as we have sampled with replacement.

Lemma 2.4.2. *An s -hashing variant matrix $T \in \mathbb{R}^{m \times n}$ (as in Definition 2.4.1) could alternatively be generated by calculating $T = \frac{1}{\sqrt{s}} [S^{(1)} + S^{(2)} + \dots + S^{(s)}]$, where $S^{(k)} \in \mathbb{R}^{m \times n}$ are independent 1-hashing matrices for $k = 1, 2, \dots, s$.*

Proof. In Definition 2.4.1, an s -hashing variant matrix T is generated by the following procedure:

```

for  $j = 1, 2, \dots, n$  do
  for  $k = 1, 2, \dots, s$  do
    Sample  $i_k \in [m]$  uniformly at random and add  $\pm 1/\sqrt{s}$  to  $T_{i_k, j}$ .

```

Due to the independence of the entries, the 'for' loops in the above routine can be swapped, leading to the equivalent formulation,

```

for  $k = 1, 2, \dots, s$  do
  for  $j = 1, 2, \dots, n$  do
    Sample  $i_k \in [m]$  uniformly at random and add  $\pm 1/\sqrt{s}$  to  $T_{i_k, j}$ .

```

For each $k \leq s$, the 'for' loop over j in the above routine generates an independent random 1-hashing matrix $S^{(k)}$ and adds $(1/\sqrt{s}) S^{(k)}$ to T .

□

We are ready to state and prove the main result in this section.

Theorem 2.4.2. *Let $s, r \leq d \leq n \in \mathbb{N}^+$, $\epsilon, \delta \in (0, 1)$. Suppose that $m \in \mathbb{N}^+$ is chosen such that the distribution of 1-hashing matrices $S \in \mathbb{R}^{m \times ns}$ is an (ϵ, δ) -oblivious subspace embedding for any matrix $B \in \mathbb{R}^{ns \times d}$ with rank r and $\mu(B) \leq \mu$ for some $\mu > 0$. Then the distribution of s -hashing variant matrices $T \in \mathbb{R}^{m \times n}$ is an (ϵ, δ) -oblivious subspace embedding for any matrix $A \in \mathbb{R}^{n \times d}$ with rank r and $\mu(A) \leq \mu\sqrt{s}$.*

Proof. Applying Lemma 2.4.2, we let

$$T = \frac{1}{\sqrt{s}} [S^{(1)} + S^{(2)} + \dots + S^{(s)}] \quad (2.4.4)$$

be a randomly generated s -hashing variant matrix where $S^{(k)} \in \mathbb{R}^{m \times n}$ are independent 1-hashing matrices, $k \in \{1, \dots, s\}$. Let $A \in \mathbb{R}^{n \times d}$ with rank r and with $\mu(A) \leq \mu\sqrt{s}$; let $U \in \mathbb{R}^{n \times r}$ be an SVD-factor of A as defined in (2.2.4). Let

$$W = \frac{1}{\sqrt{s}} \begin{pmatrix} U \\ \vdots \\ U \end{pmatrix} \in \mathbb{R}^{ns \times r}. \quad (2.4.5)$$

As U has orthonormal columns, the matrix W also has orthonormal columns and hence the coherence of W coincides with the largest Euclidean norm of its rows

$$\mu(W) = \frac{1}{\sqrt{s}} \mu(U) = \frac{1}{\sqrt{s}} \mu(A) \leq \mu. \quad (2.4.6)$$

Let $S = (S^{(1)} \dots S^{(s)}) \in \mathbb{R}^{m \times ns}$. We note that the j -th column of S is generated by sampling $i \in [m]$ and setting $S_{ij} = \pm 1$. Moreover, as $S^{(k)}$, $k \in \{1, \dots, s\}$, are independent, the sampled entries are independent. Therefore, S is distributed as a 1-hashing matrix. Furthermore, due to our assumption on the distribution of 1-hashing matrices, m is chosen such that $S \in \mathbb{R}^{m \times ns}$ is an (ϵ, δ) -oblivious subspace embedding for $(ns) \times r$ matrices of coherence at most μ . Applying this to input matrix W , we have that with probability at least $1 - \delta$,

$$(1 - \epsilon)\|z\|_2^2 = (1 - \epsilon)\|Wz\|_2^2 \leq \|SWz\|_2^2 \leq (1 + \epsilon)\|Wz\|_2^2 = (1 + \epsilon)\|z\|_2^2, \quad (2.4.7)$$

for all $z \in \mathbb{R}^r$, where in the equality signs, we used that W has orthonormal columns. On the other hand, we have that

$$SW = \frac{1}{\sqrt{s}} (S^{(1)} \dots S^{(s)}) \begin{pmatrix} U \\ \vdots \\ U \end{pmatrix} = \frac{1}{\sqrt{s}} [S^{(1)}U + S^{(2)}U + \dots + S^{(s)}U] = TU.$$

This and (2.4.7) provide that, with probability at least $1 - \delta$,

$$(1 - \epsilon)\|z\|_2^2 \leq \|TUz\|_2^2 \leq (1 + \epsilon)\|z\|_2^2, \quad (2.4.8)$$

which implies that T is an ϵ -subspace embedding for A by Lemma 2.2.2. \square

Theorem 2.3.1 and Theorem 2.4.2 imply an s -hashing variant version of Theorem 2.4.1.

Theorem 2.4.3. *Suppose that $\epsilon, \delta \in (0, 1)$, $s, r \leq d \leq n, m \leq n \in \mathbb{N}^+$, $E > 0$ satisfy (2.3.2) and (2.3.3). Then for any matrix $A \in \mathbb{R}^{n \times d}$ with rank r and $\mu(A) \leq \bar{\mu}(\epsilon, \delta)\sqrt{s}$, where $\bar{\mu}(\epsilon, \delta)$ is defined in (2.3.1), a randomly generated s -hashing variant matrix $S \in \mathbb{R}^{m \times n}$ is an ϵ -subspace embedding for A with probability at least $1 - \delta$.*

Proof. Theorem 2.3.1 implies that the distribution of 1-hashing matrices $S \in \mathbb{R}^{m \times ns}$ is an (ϵ, δ) -oblivious subspace embedding for any matrix $B \in \mathbb{R}^{ns \times d}$ with rank r and $\mu(B) \leq \bar{\mu}(\epsilon, \delta)$. We also note that this result is invariant to the number of rows in B (as long as the column size of S matches the row count of B), and so the expressions for m , $\bar{\mu}(\epsilon, \delta)$ and the constants therein remain unchanged.

Theorem 2.4.2 then provides that the distribution of s -hashing variant matrices $S \in \mathbb{R}^{m \times n}$ is an (ϵ, δ) -oblivious subspace embedding for any matrix $A \in \mathbb{R}^{n \times d}$ with rank r and $\mu(A) \leq \bar{\mu}(\epsilon, \delta)\sqrt{s}$; the desired result follows. \square

Theorem 2.4.1 and Theorem 2.4.3 provide similar results, and we find that the latter provides simpler constant expressions (such as for E).

2.4.3 The Hashed-Randomised-Hadamard-Transform sketching

Here we consider the Randomised-Hadamard-Transform [2], to be applied to the input matrix A before sketching, as another approach that allows reducing the coherence requirements under which good subspace embedding properties can be guaranteed. It is common to use the Subsampled-RHT (SHRT) [2], but the size of the sketch needs to be at least $\mathcal{O}(r \log r)$; this prompts us to consider using hashing instead of subsampling in this context (as well), and obtain an optimal order sketching bound. Figure 3.2 illustrates numerically the benefit of HRHT sketching for preconditioning compared to SRHT.

Definition 2.4.2. A Hashed-Randomised-Hadamard-Transform (HRHT) is an $m \times n$ matrix of the form $S = S_h H D$ with $m \leq n$, where

- D is a random $n \times n$ diagonal matrix with ± 1 independent entries.
- H is an $n \times n$ Walsh-Hadamard matrix defined by

$$H_{ij} = n^{-1/2} (-1)^{\langle (i-1)_2, (j-1)_2 \rangle}, \quad (2.4.9)$$

where $(i-1)_2, (j-1)_2$ are binary representation vectors of the numbers $(i-1), (j-1)$ respectively¹¹.

- S_h is a random $m \times n$ s -hashing or s -hashing variant matrix, independent of D .

Our next results show that if the input matrix is sufficiently over-determined, the distribution of HRHT matrices with optimal sketching size and either choice of S_h , is an (ϵ, δ) -oblivious subspace embedding.

Theorem 2.4.4 (s -hashing version). $r \leq d \leq n \in \mathbb{N}^+$. Let $C_1, C_2, C_3, C_M, C_\nu, C_s > 0$ be problem-independent constants. Suppose that $\epsilon, \delta \in (0, C_3)$, $m, s \in \mathbb{N}^+$ and $E > 0$ satisfy (2.4.1), (2.4.2) and (2.4.3). Let $\delta_1 \in (0, 1)$ and suppose further that

$$n \geq \frac{\left(\sqrt{r} + \sqrt{8 \log(n/\delta_1)} \right)^2}{s C_\nu^2 C_1^{-2} \bar{\mu}(\epsilon, \delta)^2}, \quad (2.4.10)$$

where $\bar{\mu}(\epsilon, \delta)$ is defined in (2.3.1). Then for any matrix $A \in \mathbb{R}^{n \times d}$ with rank r , an HRHT matrix $S \in \mathbb{R}^{m \times n}$ with an s -hashing matrix S_h , is an ϵ -subspace embedding for A with probability at least $(1 - \delta)(1 - \delta_1)$.

Theorem 2.4.5 (s -hashing variant distribution). Suppose that $\epsilon, \delta \in (0, 1)$, $r \leq d \leq n, m \leq n \in \mathbb{N}$, $E > 0$ satisfy (2.3.2) and (2.3.3). Let $\delta_1 \in (0, 1)$ and suppose further that

$$n \geq \frac{\left(\sqrt{r} + \sqrt{8 \log(n/\delta_1)} \right)^2}{s \bar{\mu}(\epsilon, \delta)^2}, \quad (2.4.11)$$

¹¹For example, $(3)_2 = (1, 1)$.

where $\bar{\mu}(\epsilon, \delta)$ is defined in (2.3.1). Then for any matrix $A \in \mathbb{R}^{n \times d}$ with rank r , an HRHT matrix $S \in \mathbb{R}^{m \times n}$ with an s -hashing variant matrix S_h , is an ϵ -subspace embedding for A with probability at least $(1 - \delta)(1 - \delta_1)$.

The proof of Theorem 2.4.4 and Theorem 2.4.5 relies on the analysis in [97] of Randomised-Hadamard-Transforms, which are shown to reduce the coherence of any given matrix with high probability.

Lemma 2.4.3. [97] Let $r \leq n \in \mathbb{N}^+$ and $U \in \mathbb{R}^{n \times r}$ have orthonormal columns. Suppose that H, D are defined in Definition 2.4.2 and $\delta_1 \in (0, 1)$. Then

$$\mu(HDU) \leq \sqrt{\frac{r}{n}} + \sqrt{\frac{8 \log(n/\delta_1)}{n}}$$

with probability at least $1 - \delta_1$.

We are ready to prove Theorem 2.4.4.

Proof of Theorem 2.4.4 and Theorem 2.4.5. Let $A = U\Sigma V$ be defined in (2.2.4), $S = S_h HD$ be an HRHT matrix. Define the following events:

- $B_1 = \left\{ \mu(HDU) \leq \sqrt{r/n} + \sqrt{8 \log(n/\delta_1)/n} \right\},$
- $B_2 = \left\{ \mu(HDA) \leq \sqrt{r/n} + \sqrt{8 \log(n/\delta_1)/n} \right\},$
- $B_3 = \{ \mu(HDA) \leq \hat{\mu}(s, \epsilon, \delta) \},$
- $B_4 = \{ S_h \text{ is an } \epsilon\text{-subspace embedding for } HDA \},$
- $B_5 = \{ S_h HD \text{ is an } \epsilon\text{-subspace embedding for } A \},$

where $\hat{\mu}(s, \epsilon, \delta) = \sqrt{s} C_\nu C_1^{-1} \bar{\mu}(\epsilon, \delta)$ if S_h is an s -hashing matrix and $\hat{\mu}(s, \epsilon, \delta) = \sqrt{s} \bar{\mu}(\epsilon, \delta)$ if S_h is an s -hashing variant matrix, and where $\bar{\mu}(\epsilon, \delta)$ is defined in (2.3.1).

Observe that B_4 implies B_5 because B_4 gives

$$(1 - \epsilon) \|Ax\|^2 \leq (1 - \epsilon) \|HDAx\|^2 \leq \|S_h HDAx\|^2 \leq (1 + \epsilon) \|HDAx\|^2 \leq (1 + \epsilon) \|Ax\|^2, \quad (2.4.12)$$

where the first and the last equality follows from HD being orthogonal. Moreover, observe that $B_1 = B_2$ because $\mu(HDA) = \max_i \|(HDU)_i\|_2 = \mu(HDU)$, where the first equality follows from $HDA = (HDU)\Sigma V^T$ being an SVD of HDA . Furthermore, B_2 implies B_3 due to (2.4.10) in the s -hashing case; and (2.4.11) in the s -hashing variant case.

Thus $\mathbb{P}(B_5) \geq \mathbb{P}(B_4) = \mathbb{P}(B_4|B_3)\mathbb{P}(B_3) \geq P(B_4|B_3)\mathbb{P}(B_2) = \mathbb{P}(B_4|B_3)\mathbb{P}(B_1)$. If S_h is an s -hashing matrix, Theorem 2.4.1 gives $\mathbb{P}(B_4|B_3) \geq 1 - \delta$. If S_h is an s -hashing variant matrix, Theorem 2.4.3 gives $\mathbb{P}(B_4|B_3) \geq 1 - \delta$. Therefore in both cases, we have

$$P(B_5) \geq \mathbb{P}(B_4|B_3)\mathbb{P}(B_1) \geq (1 - \delta)\mathbb{P}(B_1) \geq (1 - \delta)(1 - \delta_1), \quad (2.4.13)$$

where the third inequality uses Lemma 2.4.3. □

Chapter 3

Sketching for linear least squares

3.1 Introduction and relevant literature

This chapter is based and expands materials in [95, 12].

Main contribution This chapter builds on the insight from the theoretical results in the last chapter to propose, analyse and benchmark a sketching based solver of (1.3.1). We first propose and analyse a rank-deficient generic sketching framework for (1.3.1), which includes the algorithm used by two previous sketching-based solvers, Blendenpik [5] and LSRN [78] but additionally allows more flexibility of the choice of factorizations of the sketched matrix SA for building a preconditioner for (1.3.1). Our analysis shows that under this framework, one can compute a minimal residual solution of (1.3.1) with sketching if a rank-revealing factorization is used; or the minimal norm solution of (1.3.1) if a total orthogonal factorization is used. Next, based on this algorithmic framework, we propose Ski-LLS, a software package for solving (1.3.1) where we carefully distinguish whether A is dense or sparse. If A is dense, Ski-LLS combines our novel hashed coherence reduction transformation¹ analysed in Theorem 2.4.4 with a recently proposed randomized column pivoted QR factorization [76], achieving better robustness and faster speed than Blendenpik and LSRN. If A is sparse, Ski-LLS combines s -hashing analysed in Theorem 2.4.1 with the state-of-the-art sparse QR factorization in [28], achieving 10 times faster speed on random sparse ill-conditioned problems and competitive performance on a test set of 181 matrices from the Florida Matrix Collection [29] comparing to the state-of-the-art direct and iterative solvers for sparse (1.3.1), which are based on sparse QR factorization [28] and incomplete Cholesky factorization preconditioned LSQR [94] respectively.

Relevant literature Classically, dense (1.3.1) is solved by LAPACK [4], and sparse (1.3.1) is either solved by a sparse direct method implemented, say in [28] or a preconditioned LSQR (see a comparison of different preconditioners in [39]). Sarlo [93] first proposed using sketching matrices

¹For better numerical performance we use DHT as in Blendenpik instead of the Hadamard Transform analysed in the theory for Ski-LLS.

that are oblivious subspace embeddings to solve (1.3.1) by solving $\min_{x \in \mathbb{R}^d} \|SAx - Sb\|_2$. This approach requires the row of S to grow proportionally to the inverse square of the residual accuracy; hence is impractical for obtaining high accuracy solutions. Instead, Rokhlin [92] proposed using the sketch SA to compute a preconditioner of (1.3.1); and then solve (1.3.1) using preconditioned LSQR. This approach allows machine precision solutions to be computed in a small number of LSQR iterations if the matrix S is a subspace embedding of A . This algorithmic idea was carefully implemented in Blendenpik [5], achieving four times speed-up against LAPACK. Noting that Blendenpik only solves full rank (1.3.1) and does not take advantage of sparse A , Meng, Saunders and Mahoney [78] proposed LSRN, which takes advantage of sparse A and computes an accurate solution even when A is rank-deficient by using Gaussian matrices to sketch and the SVD to compute a preconditioner. However, the run-time comparisons are conducted in a multi-core parallel environment, unlike Blendenpik, which uses the serial environment.

Recently, the numerical performance of using 1-hashing matrices as the sketching matrix to solve (1.3.1) was explored in [24]. [53, 52] further explored using Blendenpik-like solvers in a distributed computing environment. [63] explores using random embedding to solve $L2$ -regularised least squares.

To the best of our knowledge, Ski-LLS is the first solver that uses s -hashing (with $s > 1$); uses a sparse factorization when solving sparse (1.3.1); and uses the hashing combined with coherence reduction transformations for dense problems. This work also presents the first large scale comparison of sketching-based LLS solvers with the state-of-the-art classical sparse solvers on the Florida Matrix Collection.

3.2 Algorithmic framework and analysis

We now turn our attention to the LLS problem (1.3.1) we are interested in solving. Building on the Blendenpik [5] and LSRN [78] techniques, we introduce a generic algorithmic framework for (1.3.1) that can employ any rank-revealing factorization of SA , where S is a(ny) sketching matrix; we then analyse its convergence.

3.2.1 A generic algorithmic framework for solving linear least squares with sketching

Algorithm 1 Generic Sketching Algorithm for Linear Least Squares

Initialization

Given $A \in \mathbb{R}^{n \times d}$ and $b \in \mathbb{R}^n$, set positive integers m and it_{max} , and accuracy tolerances τ_a and τ_r , and an $m \times n$ random matrix distribution \mathcal{S} .

1. Randomly draw a sketching matrix $S \in \mathbb{R}^{m \times n}$ from \mathcal{S} , compute the matrix-matrix product $SA \in \mathbb{R}^{m \times d}$ and the matrix-vector product $Sb \in \mathbb{R}^m$.

2. Compute a factorization of SA of the form,

$$SA = QR\hat{V}^T, \quad (3.2.1)$$

where

- $R = \begin{pmatrix} R_{11} & R_{12} \\ 0 & 0 \end{pmatrix} \in \mathbb{R}^{d \times d}$, where $R_{11} \in \mathbb{R}^{p \times p}$ is nonsingular.
- $Q = (Q_1 \ Q_2) \in \mathbb{R}^{m \times d}$, where $Q_1 \in \mathbb{R}^{m \times p}$ and $Q_2 \in \mathbb{R}^{m \times (d-p)}$ have orthonormal columns.
- $\hat{V} = (V_1 \ V_2) \in \mathbb{R}^{d \times d}$ is an orthogonal matrix with $V_1 \in \mathbb{R}^{d \times p}$.

3. Compute $x_s = V_1 R_{11}^{-1} Q_1^T Sb$. If $\|Ax_s - b\|_2 \leq \tau_a$, terminate with solution x_s .

4. Else, iteratively, compute

$$y_\tau \approx \arg \min_{y \in \mathbb{R}^p} \|Wy - b\|_2, \quad (3.2.2)$$

where

$$W = AV_1 R_{11}^{-1}, \quad (3.2.3)$$

using LSQR [86] with (relative) tolerance τ_r and maximum iteration count it_{max} . Return $x_\tau = V_1 R_{11}^{-1} y_\tau$.

Remark 2. (i) The factorization $SA = QR\hat{V}^T$ allows column-pivoted QR, or other rank-revealing factorization, complete orthogonal decomposition ($R_{12} = 0$) and the SVD ($R_{12} = 0$, R_{11} diagonal). It also includes the usual QR factorisation if SA is full rank; then the R_{12} block is absent.

(ii) Often in implementations, the factorization (3.2.1) has $R = \begin{pmatrix} R_{11} & R_{12} \\ 0 & R_{22} \end{pmatrix}$, where $R_{22} \approx 0$ and is treated as the zero matrix.

(iii) For computing x_s in Step 3, we note that in practical implementations, R_{11} in (3.2.1) is upper triangular, enabling efficient calculation of matrix-vector products involving R_{11}^{-1} ; then, there is no need to form/calculate R_{11}^{-1} explicitly.

(iv) For the solution of (3.2.2), we use the termination criterion $\|y_\tau - y_*\|_{W^T W} \leq \tau_r \|y_\tau - y_*\|_{W^T W}$ in the theoretical analysis, where y_* is defined in (3.2.5). In practical implementations different termination criteria need to be employed (see Section 3.3.2).

3.2.2 Analysis of Algorithm 1

Given problem (1.3.1), we denote its minimal Euclidean norm solution as follows

$$x_{*,2} = \arg \min_{x^* \in \mathbb{R}^d} \|x_*\|_2 \quad \text{subject to} \quad \|Ax_* - b\|_2 = \min_x \|Ax - b\|_2. \quad (3.2.4)$$

and let

$$y_* = \arg \min_{y \in \mathbb{R}^p} \|Wy - b\|_2, \quad \text{where } W \text{ is defined in (3.2.3).} \quad (3.2.5)$$

The following two lemmas provide basic properties of Algorithm 1.

Lemma 3.2.1. *$W \in \mathbb{R}^{n \times p}$ defined in (3.2.3) has full rank p .*

Proof. Note SW has rank p because $SW = Q_1$, where Q_1 is defined in (3.2.1). By rank-nullity theorem in \mathbb{R}^p , $\text{rank}(W) + \dim \ker(W) = \text{rank}(SW) + \dim \ker(SW)$ where $\ker(W)$ denotes the null space of W ; and since $\dim \ker(SW) \geq \dim \ker(W)$, we have that $\text{rank}(SW) \leq \text{rank}(W)$. So $\text{rank}(W) \geq p$. It follows that $\text{rank}(W) = p$ because $W \in \mathbb{R}^{n \times p}$ can have at most rank p . \square

Lemma 3.2.2. *In Algorithm 1, if S is an ϵ -subspace embedding for A for some $\epsilon \in (0, 1)$, then $p = r$ where r is the rank of A .*

Proof. Lemma 2.2.1 gives $r = \text{rank}(A) = \text{rank}(SA) = p$. \square

If the LLS problem (1.3.1) has a sufficiently small optimal residual, then Algorithm 1 terminates early in Step 3 with the solution x_s of the sketched problem $\min \|SAx - Sb\|_2$; then, no LSQR iterations are required.

Lemma 3.2.3 (Explicit Sketching Guarantee). *Given problem (1.3.1), suppose that the matrix $S \in \mathbb{R}^{m \times n}$ in Algorithm 1 is an ϵ -subspace embedding for the augmented matrix $(A \ b)$ for some $0 < \epsilon < 1$. Then*

$$\|Ax_s - b\|_2 \leq \frac{1 + \epsilon}{1 - \epsilon} \|Ax_* - b\|_2, \quad (3.2.6)$$

where x_s is defined in Step 3 of Algorithm 1 and x_* is a(ny) solution of (1.3.1).

The proof is similar to the result in [101] that shows that any solution of the sketched problem $\min_x \|SAx - Sb\|_2$ satisfies (3.2.6). For completeness, the proof is included here.

Proof. We have that $x_s \in \arg \min \|SAx - Sb\|_2$ by checking the optimality condition $(SA)^T SAx_s = (SA)^T Sb$. Hence we have that

$$\|Ax_s - b\|_2^2 \leq \frac{1}{1 - \epsilon} \|SAx_s - Sb\|_2^2 \leq \frac{1}{1 - \epsilon} \|SAx_* - Sb\|_2^2 \leq \frac{1 + \epsilon}{1 - \epsilon} \|Ax_* - b\|_2^2, \quad (3.2.7)$$

where the first and the last inequality follow from S being a subspace embedding for $(A \ b)$, while the second inequality is due to x_s minimizing $\|SAx - Sb\|_2$. \square

The following technical lemma is needed in the proof of our next theorem.

Lemma 3.2.4. *Let $A \in \mathbb{R}^{n \times d}$ and $V_1 \in \mathbb{R}^{d \times p}$ be defined in Algorithm 1. Then $\ker(V_1^T) \cap \text{range}(A^T) = \{0\}$, where $\ker(V_1^T)$ and $\text{range}(A^T)$ denote the null space of V_1^T and range subspace generated by the rows of A , respectively.*

Proof. Let $z \in \ker(V_1^T) \cap \text{range}(A^T)$. Then $V_1^T z = 0$ and $z = A^T w$ for some $w \in \mathbb{R}^n$. Let U, Σ, V be the SVD factors of A as defined in (2.2.4). Since S is an ϵ -subspace embedding for A , $\text{rank}((SU)^T) = \text{rank}(SU) = \text{rank}(SA) = r$, where r is the rank of A and hence there exists $\hat{w} \in \mathbb{R}^m$ such that $(SU)^T \hat{w} = U^T w$. Note that

$$0 = V_1^T z = V_1^T A^T w = V_1^T V \Sigma U^T w = V_1^T V \Sigma U^T S^T \hat{w} = V_1^T A^T S^T \hat{w} = R_{11}^T Q_1^T \hat{w}, \quad (3.2.8)$$

which implies $Q_1^T \hat{w} = 0$ because R_{11}^T is nonsingular. It follows that

$$z = A^T w = V \Sigma U^T w = V \Sigma U^T S^T \hat{w} = (SA)^T \hat{w} = V \begin{pmatrix} R_{11}^T Q_1^T \\ R_{12}^T Q_1^T \end{pmatrix} \hat{w} = 0, \quad (3.2.9)$$

where we have used $Q_1^T \hat{w} = 0$ for the last equality. \square

Theorem 3.2.1 shows that when the LSQR algorithm in Step 4 converges, Algorithm 1 returns a minimal residual solution of (1.3.1).

Theorem 3.2.1 (Implicit Sketching Guarantee). *Given problem (1.3.1), suppose that the matrix $S \in \mathbb{R}^{m \times n}$ in Algorithm 1 is an ϵ -subspace embedding for the augmented matrix $(A \ b)$ for some $0 < \epsilon < 1$. If $y_\tau = y_*$ in Step 4 of Algorithm 1 (by setting $\tau_r := 0$), where y_* is defined in (3.2.5), then x_τ in Step 5 satisfies $x_\tau = x_*$, where x_* is a solution of (1.3.1).*

Proof. Using the optimality conditions (normal equations) for the LLS in (3.2.5), and $y_\tau = y_*$, we deduce $W^T W y_\tau = W^T b$, where W is defined in (3.2.5). Substituting the definition of x_τ from Step 5 of Algorithm 1, we deduce

$$(R_{11}^{-1})^T V_1^T A^T A x_\tau = (R_{11}^{-1})^T V_1^T A^T b.$$

Multiplying the last displayed equation by R_{11}^T , we obtain

$$V_1^T (A^T A x_\tau - A^T b) = 0. \quad (3.2.10)$$

It follows from (3.2.10) that $A^T A x_\tau - A^T b \in \ker(V_1^T) \cap \text{range}(A^T)$. But Lemma 3.2.4 implies that the latter set intersection only contains the origin, and so $A^T A x_\tau - A^T b = 0$; this and the normal equations for (1.3.1) imply that x_τ is an optimal solution of (1.3.1). \square

The following technical lemma is needed for our next result; it re-states Theorem 3.2 from [78] in the context of Algorithm 1.

Lemma 3.2.5. [78] *Given problem (1.3.1), let $x_{*,2}$ be its minimal Euclidean norm solution defined in (3.2.4) and $P \in \mathbb{R}^{d \times p}$, a nonsingular matrix. Let $x_\tau := Py_\tau$, where y_τ is assumed to be the minimal Euclidean norm solution of $\min_{y \in \mathbb{R}^p} \|APy - b\|_2$. Then $x_\tau = x_{*,2}$ if $\text{range}(P) = \text{range}(A^T)$.*

Theorem 3.2.2 further guarantees that if $R_{12} = 0$ in (3.2.1) such as when a complete orthogonal factorization is used, then the minimal Euclidean norm solution of (1.3.1) is obtained.

Theorem 3.2.2 (Minimal-Euclidean Norm Solution Guarantee). *Given problem (1.3.1), suppose that the matrix $S \in \mathbb{R}^{m \times n}$ in Algorithm 1 is an ϵ -subspace embedding for the augmented matrix $(A \ b)$ for some $0 < \epsilon < 1$. If $R_{12} = 0$ in (3.2.1) and $y_\tau = y_*$ in Step 4 of Algorithm 1 (by setting $\tau_r := 0$), where y_* is defined in (3.2.5), then x_τ in Step 5 satisfies $x_\tau = x_{*,2}$, where $x_{*,2}$ is the minimal Euclidean norm solution (3.2.4) of (1.3.1).*

Proof. The result follows from Lemma 3.2.5 with $P := V_1 R_{11}^{-1}$, provided $\text{range}(V_1 R_{11}^{-1}) = \text{range}(A^T)$. To see this, note that

$$\text{range}(V_1 R_{11}^{-1}) = \text{range}(V_1) = \text{range}((SA)^T),$$

where the last equality follows from $(SA)^T = V_1 R_{11}^T Q_1^T + V_2 R_{12} Q_1^T$ and $R_{12} = 0$. Using the SVD decomposition (2.2.4) of A , we further have

$$\text{range}(V_1 R_{11}^{-1}) = \text{range}(A^T S^T) = \text{range}(V \Sigma U^T S^T) = \text{range}(V \Sigma (SU)^T).$$

Since S is an ϵ -subspace embedding for A , it is also an ϵ -subspace embedding for U by Lemma 2.2.2 and therefore by Lemma 2.2.1, $\text{rank}(SU) = \text{rank}(U) = r$. Since $SU \in \mathbb{R}^{m \times r}$ has full column rank, we have that $\text{range}(V \Sigma (SU)^T) = \text{range}(V) = \text{range}(A^T)$. \square

Theorem 3.2.3 gives an iteration complexity bound for the inner solver in Step 4 of Algorithm 1, as well as particularising this result for a special starting point for which an optimality guarantee can be given. It relies crucially on the quality of the preconditioner provided by the sketched factorization in (3.2.1), and its proof uses standard LSQR results.

Theorem 3.2.3 (Rate of convergence). *Given problem (1.3.1), suppose that the matrix $S \in \mathbb{R}^{m \times n}$ in Algorithm 1 is an ϵ -subspace embedding for the augmented matrix $(A \ b)$ for some $0 < \epsilon < 1$. Then:*

(i) *Step 4 of Algorithm 1 takes at most*

$$\tau \leq O\left(\frac{|\log \tau_r|}{|\log \epsilon|}\right) \quad (3.2.11)$$

LSQR iterations to return a solution y_τ such that

$$\|y_\tau - y_*\|_{W^T W} \leq \tau_r \|y_0 - y_*\|_{W^T W}, \quad (3.2.12)$$

where y_ and W are defined in (3.2.5).*

(ii) If we initialize $y_0 := Q^T S b$ for the LSQR method in Step 4, then at termination of Algorithm 1, we can further guarantee that

$$\|Ax_\tau - b\|_2 \leq \left(1 + \frac{2\epsilon\tau_r}{1-\epsilon}\right) \|Ax_* - b\|_2, \quad (3.2.13)$$

where x_τ is computed in Step 4 of Algorithm 1 and x_* is a solution of (1.3.1).

Proof. (i) Using results in [37], LSQR applied to (3.2.2) converges as follows

$$\frac{\|y_j - y_*\|_{W^T W}}{\|y_0 - y_*\|_{W^T W}} \leq 2 \left(\frac{\sqrt{\kappa[W^T W]} - 1}{\sqrt{\kappa[W^T W]} + 1} \right)^j, \quad (3.2.14)$$

where y_j denotes the j th iterate of LSQR and $\kappa(W^T W)$ refers to the condition number of $W^T W$. Since S is an ϵ -subspace embedding for A , we have that the largest singular value of W satisfies

$$\sigma_{\max}(W) = \max_{\|y\|=1} \|AV_1 R_{11}^{-1} y\| \leq (1-\epsilon)^{-1/2} \max_{\|y\|=1} \|SAV_1 R_{11}^{-1} y\| = (1-\epsilon)^{-1/2} \max_{\|y\|=1} \|Q_1 y\| = (1-\epsilon)^{-1/2},$$

where we have used that $SAV_1 R_{11}^{-1} = Q_1$ from (3.2.1). Similarly, it can be shown that the smallest singular value of W satisfies $\sigma_{\min}(W) \geq (1+\epsilon)^{-1/2}$. Hence

$$\kappa(W^T W) \leq \frac{1+\epsilon}{1-\epsilon}. \quad (3.2.15)$$

Hence we have

$$\frac{\sqrt{\kappa[W^T W]} - 1}{\sqrt{\kappa[W^T W]} + 1} \leq \frac{\sqrt{1+\epsilon} - \sqrt{1-\epsilon}}{\sqrt{1+\epsilon} + \sqrt{1-\epsilon}} = \frac{(\sqrt{1+\epsilon} - \sqrt{1-\epsilon})(\sqrt{1+\epsilon} + \sqrt{1-\epsilon})}{(\sqrt{1+\epsilon} + \sqrt{1-\epsilon})^2} \leq \epsilon.$$

Thus (3.2.14) implies $\|y_\tau - y_*\|_{W^T W} \leq \tau_r \|y_0 - y_*\|_{W^T W}$ whenever $\tau \geq \frac{\log(2) + |\log \tau_r|}{|\log \epsilon|}$.

(ii) If we initialize $y_0 := Q^T S b$ for the LSQR method in Step 4, then we have

$$\begin{aligned} \|y_0 - y_*\|_{W^T W} &= \|Ax_s - Ax_*\|_2 = \|Ax_s - b - (Ax_* - b)\|_2 \leq \|Ax_s - b\|_2 + \|Ax_* - b\|_2 \\ &\leq \left(\sqrt{\frac{1+\epsilon}{1-\epsilon}} - 1 \right) \|Ax_* - b\|_2 \leq \frac{2\epsilon}{1-\epsilon} \|Ax_* - b\|_2, \end{aligned}$$

where we have used $\sqrt{\frac{1+\epsilon}{1-\epsilon}} \leq \frac{1+\epsilon}{1-\epsilon}$ to get the last inequality. Using part (i), after at most $\frac{\log(2) + |\log \tau_r|}{|\log \epsilon|}$ LSQR iterations, we have that

$$\|y_\tau - y_*\|_{W^T W} \leq \frac{2\epsilon\tau_r}{1-\epsilon} \|Ax_* - b\|_2. \quad (3.2.16)$$

Note that $\|Ax_\tau - Ax_*\|_2 = \|y_\tau - y^*\|_{W^T W}$. Using the triangle inequality, we deduce

$$\|Ax_\tau - b\|_2 = \|Ax_\tau - Ax_* + Ax_* - b\|_2 \leq \left(1 + \frac{2\epsilon\tau_r}{1-\epsilon}\right) \|Ax_* - b\|_2. \quad (3.2.17)$$

□

3.3 Implementation details

3.3.1 Ski-LLS, an implementation of Algorithm 1

Sketching-for-Linear-Least-Squares (Ski-LLS) implements Algorithm 1 for solving (1.3.1). We distinguish two cases based on whether the data matrix A is stored as a dense matrix or a sparse matrix.

Dense A When A is stored as a dense matrix ², we employ the following implementation of Algorithm 1. The resulting solver is called Ski-LLS-dense.

1. In Step 1 of Algorithm 1, we let

$$S = S_h F D, \quad (3.3.1)$$

where

- (a) D is a random $n \times n$ diagonal matrix with ± 1 independent entries, as in Definition 2.4.2.
 - (b) F is a matrix representing the normalized Discrete Hartley Transform (DHT), defined as $F_{ij} = \sqrt{1/n} [\cos(2\pi(i-1)(j-1)/n) + \sin(2\pi(i-1)(j-1)/n)]$ ³. We use the (DHT) implementation in FFTW 3.3.8 ⁴.
 - (c) S_h is an s -hashing matrix, defined in Definition 1.2.5. We use the sparse matrix-matrix multiplication routine in SuiteSparse 5.3.0 ⁵ to compute $S_h \times (FDA)$.
2. In Step 2 of Algorithm 1, we use the randomized column pivoted QR (R-CPQR) proposed in [76, 75] ⁶.
 3. In Step 3 of Algorithm 1, since R_{11} from R-CPQR is upper triangular, we do not explicitly compute its inverse, but instead, use back-solve from the LAPACK provided by Intel MKL 2019 ⁷.
 4. In Step 4 of Algorithm 1, we use the LSQR routine implemented in LSRN [78] ⁸.

The user can choose the value of the following parameters: m (default is $1.7d$), s (default is 1), τ_a (default is 10^{-8}), it_{max} (default value is 10^4). $rcond$ (default value is 10^{-12}), which is a parameter used in Step 2 of Algorithm 1. The R-CPQR we use computes $SA = Q\tilde{R}\hat{V}^T$, which is then used to compute R_{11} by letting $p = \max \left\{ q : \tilde{R}_{qq} \geq rcond \right\}$, R_{11} be the upper left $p \times p$ block of \tilde{R} .

²This does not necessarily imply that every entry of A is non-zero, however, we presume a large number of the entries are zero such that specialized sparse numerical linear algebras are ineffective.

³Here we use the same transform (DHT) as that in Blendenpik for comparison of other components of Ski-LLS-dense, instead of the Walsh-Hadamard transform defined in Definition 2.4.2.

⁴Available at <http://www.fftw.org>.

⁵Available at <https://people.engr.tamu.edu/davis/suitesparse.html>.

⁶The implementation can be found at <https://github.com/flame/hqrrp/>. The original code only has a 32-bit integer interface. We wrote a 64-bit integer wrapper as our code has 64-bit integers.

⁷See <https://software.intel.com/content/www/us/en/develop/tools/oneapi/components/onemkl.html>.

⁸Available at <https://web.stanford.edu/group/SOL/software/lsrcn/>. We fixed some very minor bugs in the code.

wisdom (default value is 1). The DHT we use is faster with pre-tuning, see Blendenpik [6] for a detailed discussion. If the DHT has been pre-tuned, the user needs to set *wisdom* = 1, otherwise set *wisdom* = 0. In all our experiment, the default is to tune the DHT using the crudest tuning mechanism offered by FFTW, which typically takes less than one minute.

We also offer an implementation without using R-CPQR for dense full-rank problems. The only difference is that in Step 2 of Algorithm 1, we assume that the matrix A has full-rank $r = d$. Hence we use DGEQRF from LAPACK to compute a QR factorization of SA (the same routine is used in Blendenpik) instead of R-CPQR. It has the same list of parameters with the same default values, except the parameter *rcond* is absent because it does not use R-CPQR.

Sparse A When A is stored as a sparse matrix ⁹, we employ the following implementation of Algorithm 1. The resulting solver is called Ski-LLS-sparse.

1. In Step 1 of Algorithm 1, we let S be an s -hashing matrix, defined in Definition 1.2.5.
2. In Step 2 of Algorithm 1, we use the sparse QR factorization (SPQR) proposed in [28] and implemented in SuiteSparse.
3. In Step 3 of Algorithm 1, since R_{11} from SPQR is upper triangular, we do not explicitly compute its inverse, but instead, use the sparse back-substitution routine from SuiteSparse.
4. In Step 4 of Algorithm 1, we use the LSQR routine implemented in LSRN, extended to include the use of sparse preconditioner and sparse numerical linear algebras from SuiteSparse.

The user can choose the value of the following parameters: m (default value is $1.4d$), s (default value is 2), τ_a (default value is 10^{-8}), τ_r (default value is 10^{-6}), it_{max} (default value is 10^4). And $rcond_{thres}$ (default value 10^{-10}), which checks the conditioning of R_{11} computed by SPQR. If $\kappa(R_{11}) \geq 1/rcond_{thres}$, we use the perturbed back-solve for upper triangular linear systems involving R_{11} (see the next point). *perturb* (default value 10^{-10}). When $\kappa(R_{11}) \geq 1/rcond_{thres}$, any back-solve involving R_{11} or its transpose will be modified in the following way: When divisions by a diagonal entry r_{ii} of R_{11} is required where $1 \leq i \leq p$, we divide by $r_{ii} + perturb$ instead. ¹⁰ *ordering* (default value 2) which is a parameter to the SPQR routine that influences the permutation matrix \hat{V} and the sparsity of R . ¹¹.

⁹Here we assume that the user stored the matrix in a sparse matrix format because a large number of entries are zero. Throughout computations, we maintain the sparse matrix format for effective numerical linear algebras.

¹⁰This is a safe-guard when SPQR fails to detect the rank of A . This happens infrequently [28].

¹¹Note that this is slightly different from the SPQR default, which is to use to use COLAMD if $m2 \leq 2*n2$; otherwise try AMD. Let f be the flops for $\text{chol}((S^*P)^*(S^*P))$ with the ordering P found by AMD. Then if $f/\text{nnz}(R) \geq 500$ and $\text{nnz}(R)/\text{nnz}(S) \geq 5$ then try METIS, and take the best ordering found (AMD or METIS), where typically $m2 = m$, $n2 = n$ for $SA \in \mathbb{R}^{m \times n}$. In contrast, Ski-LLS by default always use the AMD ordering.

3.3.2 Discussion of our implementation

Subspace embedding properties achieved via numerical calibration Our analysis of Algorithm 1 in Section 3.2 relies crucially on S being an ϵ -subspace embedding of A . For dense matrices, Blendenpik previously used SR-DHT, defined in (3.4.7) with theoretical guarantees of the oblivious ϵ -subspace embedding property for full rank A if $m = \mathcal{O}(d \log(d))$. Theorem 2.4.4 shows when using hashing instead of sampling with randomised Walsh-Hadamard transform, hashing achieves being an oblivious ϵ -subspace embedding with $m = \mathcal{O}(d)$ (note that $r = d$ for full rank A) under the additional dimensional assumption of A . In Ski-LLS-dense, HR-DHT is used instead of HRHT analyzed in Theorem 2.4.2 because as mentioned in Blendenpik paper [5], DHT is more flexible (Walsh-Hadamard transform only allows n to be an integer power of 2 so that padding is needed); and SR-DHT based sketching solver has stabler and shorter running time comparing to when DHT is replaced by Walsh-Hadamard transform. Moreover, we aim to show in addition to the theoretical advantage of hashing (Theorem 2.4.4), numerically using hashing instead of sampling combined with coherence-reduction transformations yields a more effective solver for (1.3.1) in terms of running time.

Therefore to compare to Blendenpik, we chose to use HR-DHT instead of HRHT. We then use numerical calibration as used in Blendenpik to determine the default value of m for Ski-LLS-dense such that ϵ -subspace embedding of A is achieved with sufficiently high (all the matrices in the calibration set) probability. (See the next section and Appendices). Note that the U-shaped curve appears in Figure 3.4, because as $\gamma := m/d$ grows, we have better subspace embeddings so that ϵ decreases, resulting in fewer LSQR iterations according to (3.2.11). However the factorization cost in Step 2 and the sketching cost in Step 1 will grow as m grows. Thus a trade-off is achieved when m is neither too big nor too small.

For sparse matrices, Theorem 2.4.1 guarantees the oblivious ϵ -subspace embedding property s -hashing matrices for matrices A with low coherence. However as Figure 3.8, 3.9 suggest, s -hashing with $s > 1$ and $m = \mathcal{O}(d)$ tends to embed higher coherence A as well. The specific default values of m, s are again chosen using numerical calibration; and the characteristic U-shape is because of a similar trade-off as in the dense case.

What if S is not an ϵ -subspace embedding of A Note that even \mathcal{S} is an oblivious subspace embedding for matrices $A \in \mathbb{R}^{n \times d}$, for a given $A \in \mathbb{R}^{n \times d}$, there is a chance that a randomly generated matrix S from \mathcal{S} fails to embed A . However, in this case, Ski-LLS will still compute an accurate solution of (1.3.1) given that A has full rank and SA has the same rank as A . Because then the preconditioner $V_1 R_{11}^{-1}$ is an invertible square matrix. The situation is less clear when A is rank-deficient and S fails to embed A . However, with the default parameters chosen from numerical

calibrations, the accuracy of Ski-LLS is excellent for A being both random dense/sparse matrices and for A in the Florida matrix collection.

Approximate numerical factorization in Step 2 In both our dense and sparse solvers, Step 2 $SA = QR\hat{V}^T$ is not guaranteed to be accurate when A is rank-deficient. This is because R-CPQR, like CPQR, does not guarantee detection of rank although in almost all cases the numerical rank is correctly determined (in the sense that if one follows the procedure described in the definition of the parameter $rcond$, the factorization $SA = QR\hat{V}^T$ will be accurate up to approximately $rcond$ error). Similarly, SPQR performs heuristic rank-detection for speed efficiency and therefore rank-detection and the resulting accuracy is not guaranteed. Also, we have not analysed the implication of floating-point arithmetic for Ski-LLS. The accuracy of Ski-LLS, however, is demonstrated in a range of dense and sparse test problems, see later sections.

Practical LSQR termination criterion The termination criterion proposed in Step 4 of the algorithm is not practical as we do not know y^* . In practice, we terminate Step 4 of Algorithm 1 if $\frac{\|W^T(Wy_k - b)\|}{\|W\|\|Wy_k - b\|} \leq \tau_r$ where W is defined in (3.2.3) similarly to what is used in LSRN [78]. See the original LSQR paper [86], Section 6 for a justification.

3.4 Numerical study

3.4.1 Test Set

The matrix A

1. The following are three different types of random dense matrices that are the same type of test matrices used by Avron et.al. [5] for comparing Blendenpik with LAPACK least square solvers. They have different ‘non-uniformity’ of rows.

- (a) Incoherent dense, defined by

$$A = U\Sigma V^T \in \mathbb{R}^{n \times d}, \quad (3.4.1)$$

where $U \in \mathbb{R}^{n \times d}$, $V \in \mathbb{R}^{d \times d}$ are matrices generated by orthogonalising columns of two independent matrices with i.i.d. $N(0,1)$ entries. $\Sigma \in \mathbb{R}^{d \times d}$ is a diagonal matrix with diagonal entries equally spaced from 1 to 10^6 (inclusive).

- (b) Semi-coherent dense, defined by

$$A = \begin{pmatrix} B \\ 0 \end{pmatrix} I_{d/2}^0 + 10^{-8} J_{n,d} \in \mathbb{R}^{n \times d}, \quad (3.4.2)$$

where B is an incoherent dense matrix defined in (3.4.1), $J_{n,d} \in \mathbb{R}^{n \times d}$ is a matrix of all ones.

(c) Coherent dense, defined by

$$A = \begin{pmatrix} I_{d \times d} \\ 0 \end{pmatrix} + 10^{-8} J_{n,d} \in \mathbb{R}^{n \times d}, \quad (3.4.3)$$

where $J_{n,d}$ is a matrix of all ones.

2. The following are three different types of random sparse matrices with different ‘non-uniformity’ of rows.

(a) Incoherent sparse, defined by

$$A = \text{sprandn}(n, d, 0.01, 1e-6) \in \mathbb{R}^{n \times d}, \quad (3.4.4)$$

where ‘sprandn’ is a command in MATLAB that generates a matrix with approximately $0.01nd$ normally distributed non-zero entries and a condition number approximately equals to 10^6 .

(b) Semi-coherent sparse, defined by

$$A = \hat{D}^5 B, \quad (3.4.5)$$

where $B \in \mathbb{R}^{n \times d}$ is an incoherent sparse matrix defined in (3.4.4) and \hat{D} is a diagonal matrix with independent $N(0, 1)$ entries on the diagonal.

(c) Coherent sparse, defined by

$$A = \hat{D}^{20} B, \quad (3.4.6)$$

where $B \in \mathbb{R}^{n \times d}$, \hat{D} are the same as in (3.4.5).

3. (Florida matrix collection) A total of 181 matrices in the Florida (SuiteSparse) matrix collection [29] satisfying:

(a) If the matrix is under-determined, we transpose it to make it over-determined.

(b) We only take a matrix $A \in \mathbb{R}^{n \times d}$ if $n \geq 30000$ and $n \geq 2d$.

Remark 3. Note that here ‘coherence’ is a more general concept indicating the non-uniformity of the rows of A . Although ‘coherent dense/sparse’ A tends to have higher values of $\mu(A)$ than that of ‘incoherent dense/sparse’ A , the value of $\mu(A)$ may be similar for semi-coherent and coherent test matrices. The difference is that for ‘coherent’ test matrices, the row norms of A (and U from the SVD of A) tend to be more non-uniform. In the dense matrix cases, the rows of A are progressively less uniform due to the presence of the identity blocks. In the sparse matrix cases, the rows of A are progressively less uniform due to the scalings from the increasing powers of a diagonal Gaussian matrix.

The vector b In all the test, the vector $b \in \mathbb{R}^n$ in (1.3.1) is chosen to be a vector of all ones.

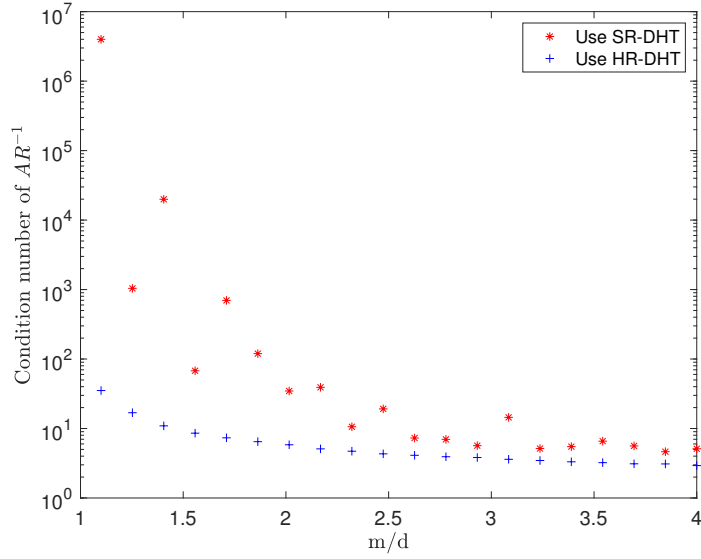


Figure 3.1: Hashing combined with a randomised Discrete Hartley Transform produces more accurate sketched matrix SA for a given m/d ratio comparing to sampling combined with a randomised Discrete Hartley Transform; the accuracy of the sketch is reflected in the quality of the preconditioner R constructed from the matrix SA , see (3.2.15).

3.4.2 Numerical illustrations

The case for using hashing instead of sampling In Figure 3.1, we generate a random coherent dense matrix $A \in \mathbb{R}^{4000 \times 400}$ defined in (3.4.3), and for each m/d (where $d = 400$), we sketch the matrix using $S \in \mathbb{R}^{m \times n}$ being a HR-DHT defined in (3.3.1) and using SR-DHT defined by

$$S = S_s F D, \quad (3.4.7)$$

where $S_s \in \mathbb{R}^{m \times n}$ is a scaled sampling matrix, whose individual rows contain a single non-zero entry at a random column with value $\sqrt{\frac{n}{m}}$; F, D are defined the same as in (3.3.1). This is the sketching used in Blendenpik. We then compute an (non-pivoted) QR factorization of each sketch $SA = QR$, and the condition number of AR^{-1} .

We see that using hashing instead of sampling allows the use of smaller m to reach a given preconditioning quality.

The case for using s -hashing with $s > 1$ In Figure 3.2, we let $A \in \mathbb{R}^{4000 \times 400}$ be a random incoherent sparse matrix defined in (3.4.4), while in Figure 3.3, A be defined as

$$A = \begin{pmatrix} B & 0 \\ 0 & I_{d/2} \end{pmatrix} + 10^{-8} J_{n,d} \in \mathbb{R}^{n \times d}, \quad (3.4.8)$$

where $B \in \mathbb{R}^{n \times d}$ is a random incoherent sparse matrix, and $J_{n \times d}$ is a matrix of all ones¹². Comparing Figure 3.2 with Figure 3.3, we see that using s -hashing matrices with $s > 1$ is essential to produce a good preconditioner.

¹²We use this type of random sparse matrix instead of one of the types defined in (3.4.5) because this matrix better showcases the failure of 1-hashing.

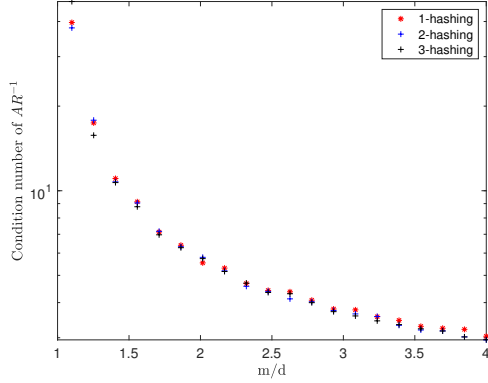


Figure 3.2: When the data matrix A is an ill-conditioned sparse Gaussian matrix, using 1, 2, 3-hashing produces similarly good preconditioners.

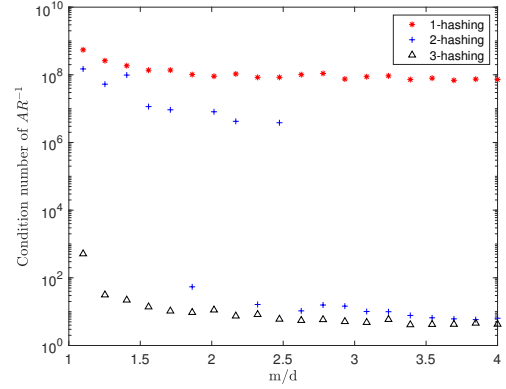


Figure 3.3: When the data matrix A has higher coherence, using s -hashing with $s > 1$ is crucial to produce an acceptable preconditioner.

3.4.3 Compilation and running environment for timed experiments

The above numerical illustrations are done in MATLAB as it does not involve running time. For all the other studies, unless otherwise mentioned, we use Intel C compiler `icc` with optimisation flag `-O3` to compile all the C code, and Intel Fortran compiler `ifort` with `-O3` to compile Fortran-based code. All code has been compiled in sequential mode and linked with sequential dense/sparse linear algebra libraries provided by Intel MKL, 2019 and Suitesparse 5.3.0. The machine used has Intel(R) Xeon(R) CPU E5-2667 v2 @ 3.30GHz with 8GB RAM.

3.4.4 Tuning to set the default parameters

The default parameter values m for Ski-LLS-dense (both with and without R-CPQR) solvers and m, s for Ski-LLS-sparse are chosen using a calibrating random matrix set. See the below graphs.

Calibration for Dense Solvers In Figure 3.4, Figure 3.5, Figure 3.6, Figure 3.7 we tested Ski-LLS-dense, Ski-LLS-dense without R-CPQR, Blendenpik and LSRN on the same calibration set and chose the optimal parameters for them for fair comparison. The default parameters chosen are $m = 1.7d$, $m = 1.7d$, $m = 2.2d$ and $m = 1.1d$ for Ski-LLS-dense, Ski-LLS-dense without R-CPQR, Blendenpik and LSRN respectively.

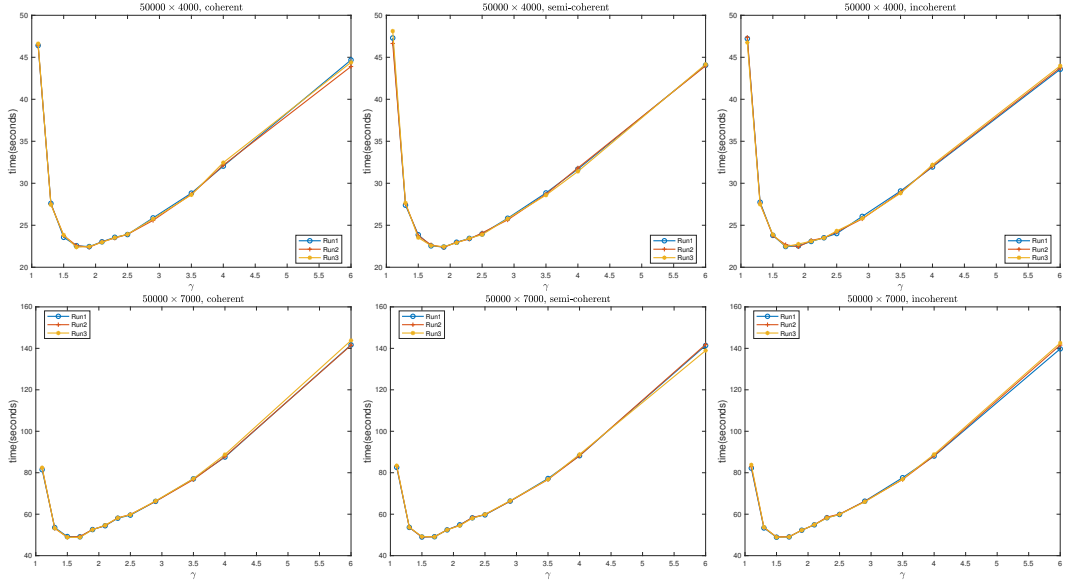


Figure 3.4: Runtime of Ski-LLS-dense on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, Ski-LLS-dense is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. We choose $\gamma = 1.7$ to approximately minimize the runtime across the above plots.

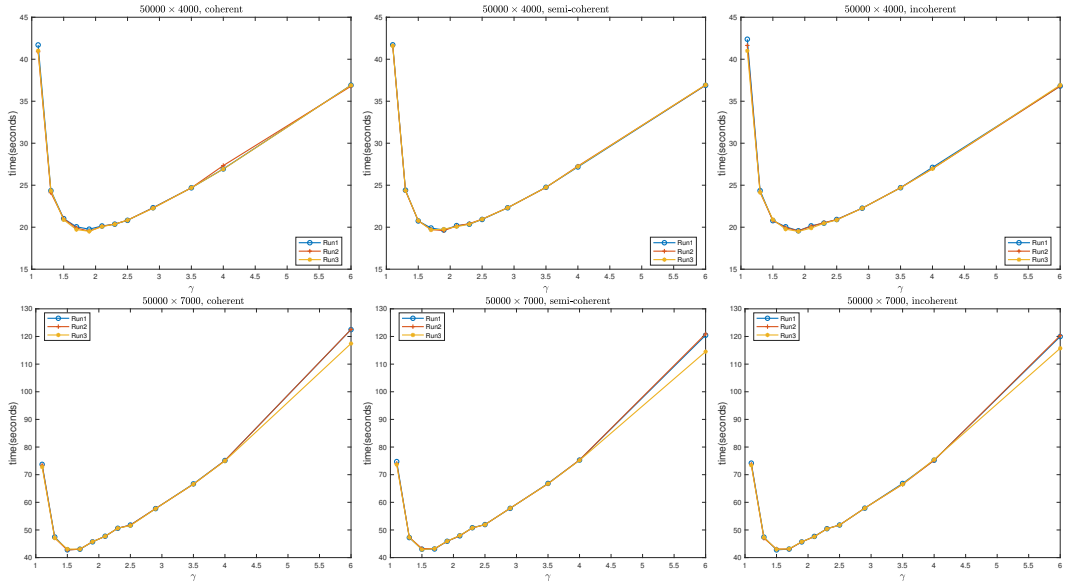


Figure 3.5: Runtime of Ski-LLS-dense without R-CPQR on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, Ski-LLS-dense without R-CPQR is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. Note that using LAPACK QR instead of R-CPQR results in slightly shorter running time (c.f. Figure 3.4). We choose $\gamma = 1.7$ to approximately minimize the runtime across the above plots.

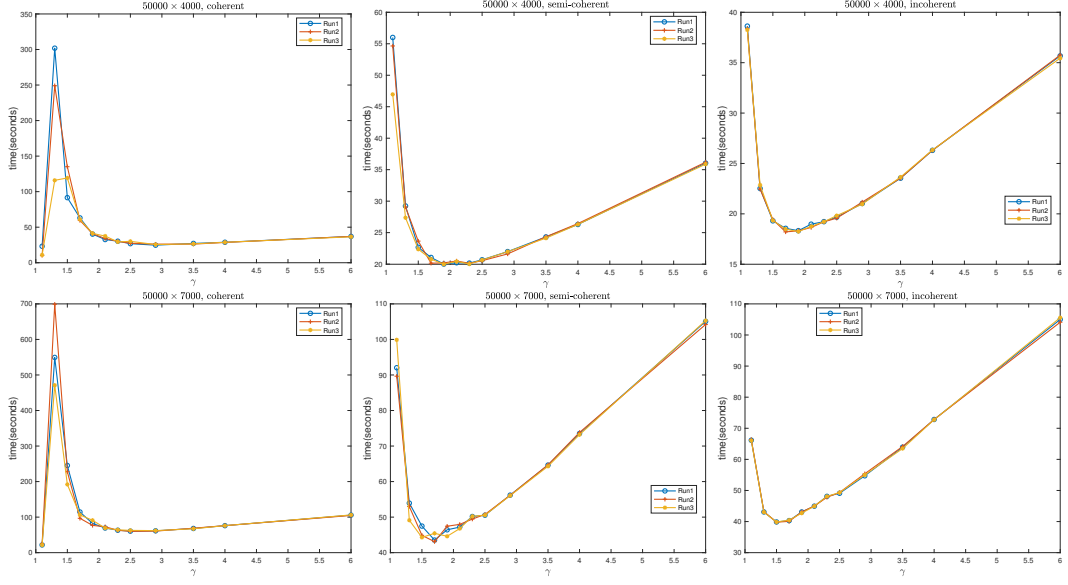


Figure 3.6: Runtime of Blendenpik on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, Blendenpik is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. Note that Blendenpik handles coherent dense A significantly less well than our dense solvers. We choose $\gamma = 2.2$ to approximately minimize the runtime across the above plots.

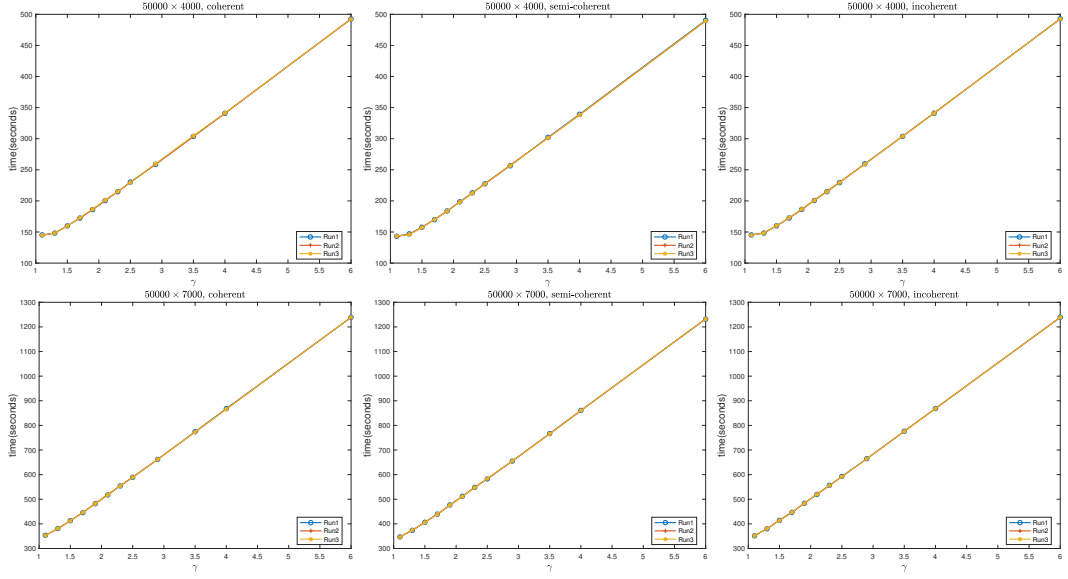


Figure 3.7: Runtime of LSRN on dense matrices $A \in \mathbb{R}^{n \times d}$ from Test Set 1 with $n = 50000, d = 4000$ and $n = 50000, d = 7000$ and different values of $\gamma = m/d$. For each plot, LSRN is run three times on (the same) randomly generated A . We see that the runtime has low variance despite the randomness in the solver. Note that LSRN runs more than 5 times slower comparing to Blendenpik or Ski-LLS in the serial testing environment, due to the use of SVD and Gaussian sketching. We choose $\gamma = 1.1$ to approximately minimize the runtime across the above plots.

Calibration for Ski-LLS-sparse In In Figures 3.8 and Figure 3.9 Figure 3.10 and Figure 3.11, we tested Ski-LLS-sparse and LSRN on the same calibration set and choose the optimal parameters

from them for fair comparison. Note that in the below calibration, $\tau_r = 10^{-4}$ is used instead of the default value of Ski-LLS-sparse. There is no τ_a because the solver at that time has not implemented Step 3 of Algorithm 1. The SPQR ordering used is the SuiteSparse default instead of Ski-LLS default (AMD). The other parameters, $it_{max}, rcond_{thres}, perturb$ are the same as the default.

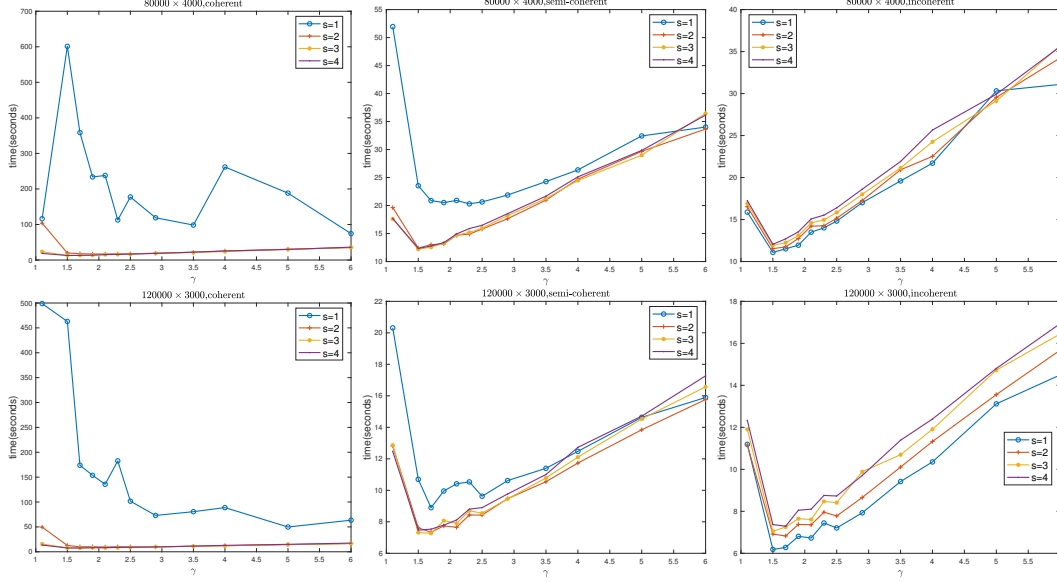


Figure 3.8: Running time of Ski-LLS-sparse on sparse matrices $A \in \mathbb{R}^{m \times d}$ from Test Set 2 with $n = 80000, d = 4000$ and $n = 120000, d = 3000$ using different values of s and $\gamma = m/d$. We choose $m = 1.4d, s = 2$ in consideration of the above plot and the residual accuracy in Figure 3.9 but also taking into account some experiments of Ski-LLS-sparse we have done on the Florida matrix collection.

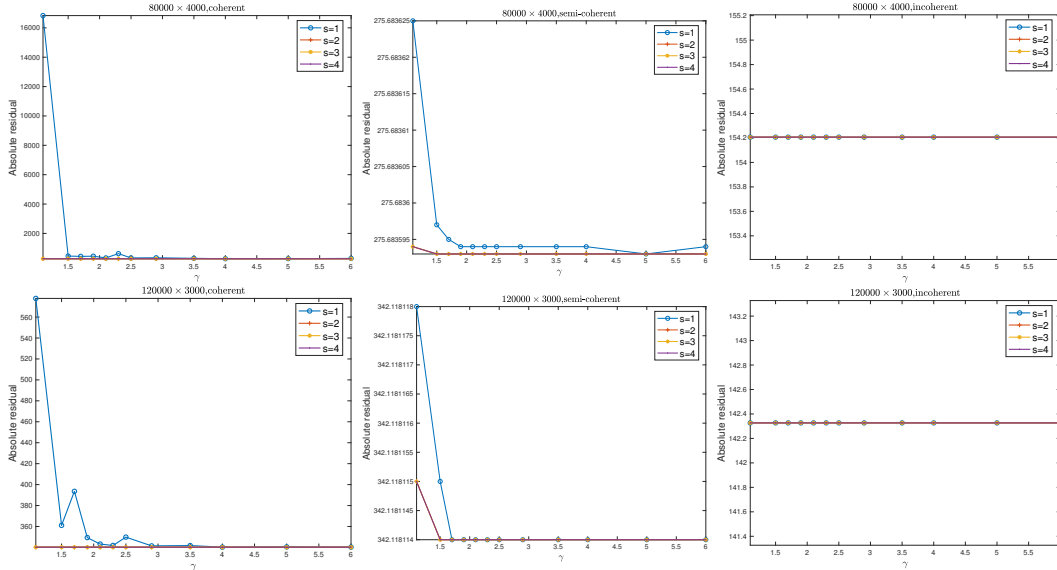


Figure 3.9: Corresponding residual on sparse matrices $A \in \mathbb{R}^{m \times d}$ from Test Set 2 with $n = 80000, d = 4000$ and $n = 120000, d = 3000$ using different values of s and $\gamma = m/d$. Note that using 1-hashing ($s = 1$) results in inaccurate solutions.

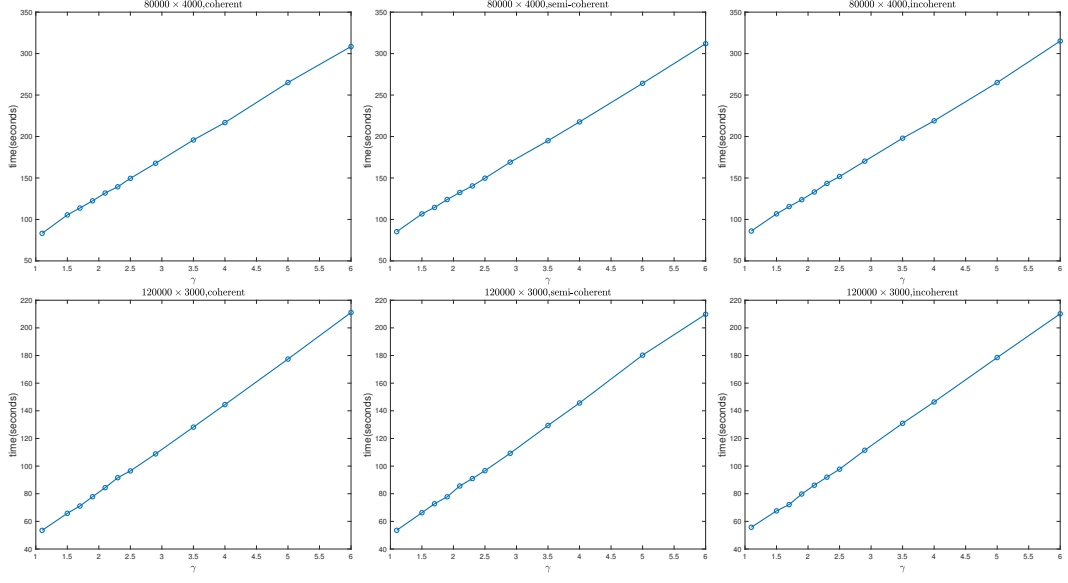


Figure 3.10: Runtime for LSRN on sparse matrices $A \in \mathbb{R}^{m \times d}$ from Test Set 2 with $n = 80000, d = 4000$ and $n = 120000, d = 3000$ using different values of $\gamma = m/d$. We choose $m = 1.1d$ in consideration of the above plot and the residual accuracy in Figure 3.11 but also taking into account some experiments of Ski-LLS-sparse we have done on the Florida matrix collection.

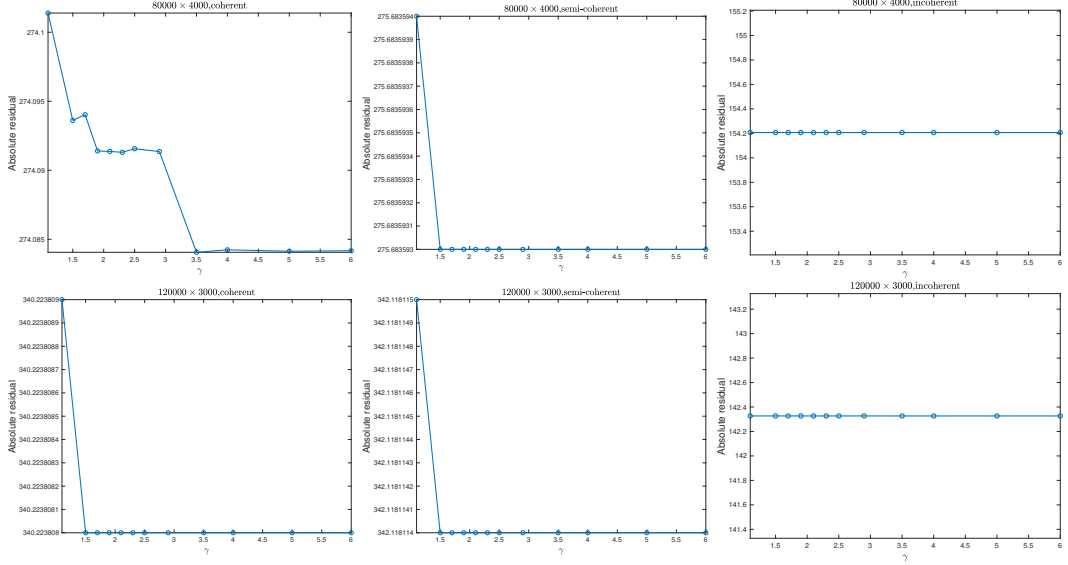


Figure 3.11: Residual values obtained by LSRN on the same sparse problems as in Figure 3.10.

3.4.5 Residual accuracy of Ski-LLS

Since our theory in Section 3.2 is only an approximation of the practical implementation as discussed in Section 3.3.2, we numerically test Ski-LLS's accuracy of solving (1.3.1). We choose 14 matrices A in the Florida matrix collection with different dimensions and rank-deficiencies (Table 3.2). We use LAPACK's SVD-based linear least squares solver (SVD), LSRN, Blendenpik, Ski-LLS-dense and

	lp_ship12l	Franz1	GL7d26	cis-n4c6-b2	lp_modszk1	rel5	ch5-5-b1
SVD	18.336	26.503	50.875	6.1E-14	33.236	14.020	7.3194
LSRN	18.336	26.503	50.875	3.2E-14	33.236	14.020	7.3194
Blendenpk	NaN	9730.700	NaN	3.0E+02	NaN	NaN	340.9200
Ski-LLS-dense	18.336	26.503	50.875	5.3E-14	33.236	14.020	7.3194
Ski-LLS-sparse	18.336	26.503	50.875	6.8E-14	33.236	14.020	7.3194
	n3c5-b2	ch4-4-b1	n3c5-b1	n3c4-b1	connectus	landmark	cis-n4c6-b3
SVD	9.0E-15	4.2328	3.4641	1.8257	282.67	1.1E-05	30.996
LSRN	6.7E-15	4.2328	3.4641	1.8257	282.67	1.1E-05	30.996
Blendenpk	1.3E+02	66.9330	409.8000	8.9443	NaN	NaN	3756.200
Ski-LLS-dense	5.2E-15	4.2328	3.4641	1.8257	282.67	1.1E-05	30.996
Ski-LLS-sparse	6.9E-15	4.2328	3.4641	1.8257	282.67	1.1E-05	30.996

Table 3.1: Residuals of solvers for a range of rank-deficient problems taken from the Florida matrix collection [29]. The matrices are all sparse but we convert them into dense format before applying a dense solver such as Blendenpk. We see both Ski-LLS-dense and Ski-LLS-sparse achieve excellent residual accuracy for rank-deficient problems, as well as LSRN. Blendenpk is not designed for rank-deficient problems and either returns a large residual or encounters numerical issues, returning NaN.

	nrow	ncol	rank
lp_ship12l	5533	1151	1042
Franz1	2240	768	755
GL7d26	2798	305	273
cis-n4c6-b2	1330	210	190
lp_modszk1	1620	687	686
rel5	240	35	24
ch5-5-b1	200	25	24
n3c5-b2	120	45	36
ch4-4-b1	72	16	15
n3c5-b1	45	10	9
n3c4-b1	15	6	5
connectus	394792	512	<458
landmark	71952	2704	2671
cis-n4c6-b3	5940	1330	1140

Table 3.2: Dimensions for the problems tested in Table 3.1

Ski-LLS-sparse on these problems with the residual shown in Table 3.1.¹³

We see both of Ski-LLS-dense and Ski-LLS-sparse have excellent residual accuracy comparing to SVD-based LAPACK solver. The result also shows that Blendenpk fails to accurately solve rank-deficient (1.3.1).

In our large scale numerical study with the Florida matrix collection, the residuals are also compared and the solution of Ski-LLS-sparse is no-less accurate than the state-of-the-art sparse solvers LS_SPQR and LS_HSL(see later sections).

3.5 Numerical performance

3.5.1 Solvers compared and their parameters

Recall Ski-LLS treats dense and sparse A differently in (1.3.1). For dense A , we compare to the state-of-the-art sketching solver Blendenpk, that has been shown to be four times faster than LAPACK on dense, large scale and moderately over-determined full rank problems [5]¹⁴. The parameters for Blendenpk are $m = 2.2d$ ¹⁵, $\tau_r = 10^{-6}$, $it_{max} = 10^4$ and $wisdom = 1$. The same wisdom data file

¹³These problems are given in a sparse format. We convert them to a dense format for dense solvers. Thus the dense solvers cannot assume any entry is a priori zero.

¹⁴Available at <https://github.com/haimav/Blendenpk>. For the sake of fair comparison, we wrote a C interface and uses the same LSQR routine as Ski-LLS.

¹⁵Chosen by calibration, see Appendix B.

as Ski-LLS is used.

For sparse A , we compare to the following solvers

1. HSL_MI35 (LS_HSL), that uses an incomplete Cholesky factorization of A to compute a preconditioner of the problem (1.3.1), before using LSQR. ¹⁶ The solver has been shown to be competitive for sparse problems (1.3.1) [39, 38]. We use $\tau_r = 10^{-6}$ and $it_{max} = 10^4$.
2. SPQR_SOLVE (LS_SPQR), that uses SPQR from Suitesparse to compute a sparse QR factorization of A , which is exploited to solve (1.3.1) directly. ¹⁷ The solver has been shown to be competitive for sparse problems [28].
3. LSRN, that uses the framework of Algorithm 1, with S having i.i.d. $N(0, 1/\sqrt{m})$ entries in Step 1; SVD factorization from Intel LAPACK of the matrix SA in Step 2; the same LSQR routine as Ski-LLS in Step 4. ¹⁸ LSRN has been shown to be an effective solver for possibly rank-deficient dense and sparse (1.3.1) under parallel computing environment [78]. However, parallel computing is outside the scope of this study and we therefore run LSRN in a serial environment. Hence the performance of LSRN may improve under parallelism. The default parameters are chosen to be $m = 1.1d$, $\tau_r = 10^{-6}$, $it_{max} = 10^4$.

3.5.2 Running time performance on randomly generated full rank dense A

Our first experiment compares of Ski-LLS for dense (1.3.1) with Blendenpik. For each matrix of different size shown in the x-axis of Figure 3.12, 3.13, 3.14, we generate a coherent, semi-coherent and incoherent dense matrix as defined in (3.4.3), (3.4.2), (3.4.1). Blendenpik, Ski-LLS (dense version) and Ski-LLS without R-CPQR are to solve (1.3.1) with b being a vector of all ones. The running time t with the residual $\|Ax - b\|_2$ are recorded, where x is the solution returned by the solvers. The residuals are all the same up to six significant figures, indicating all three solvers give an accurate solution of (1.3.1).

We see that using hashing instead of sampling yields faster solvers by comparing Ski-LLS without R-CPQR to Blendenpik, especially when the matrix A is of the form (3.4.3). We also see that Ski-LLS with R-CPQR is as fast as Blendenpik on full rank dense problems while being able to solve rank-deficient problems (Table 3.1).

The default parameters for Ski-LLS is used, and the parameters for Blendenpik is mentioned before.

¹⁶See http://www.hsl.rl.ac.uk/specs/hsl_mi35.pdf for a full specification. For the sake of fair comparison, we wrote a C interface and uses the same LSQR routine as Ski-LLS. We also disable the pre-processing of the data as it was not done for the other solvers. We then found using no scaling and no ordering was more effective than the default scaling and ordering. Hence we chose no scaling and no ordering in the comparisons. As a result, the performance of HSL may improve, however [39] experimented with the use of different scaling and ordering, providing some evidence that the improvement will not be significant.

¹⁷Available at <https://people.engr.tamu.edu/davis/suitesparse.html>.

¹⁸Note that LSRN does not contain the Step 3.

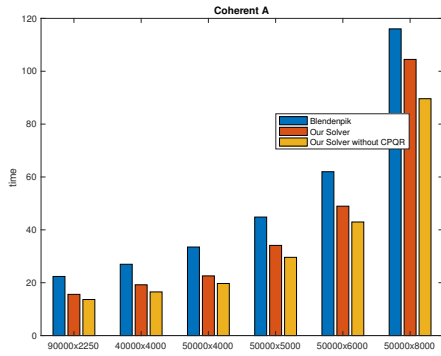


Figure 3.12: Time taken by solvers to compute the solution of problem (1.3.1) for A being coherent dense matrices of various sizes (x-axis)

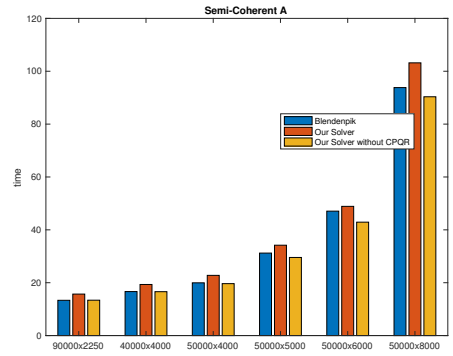


Figure 3.13: Time taken by solvers to compute the solution of problem (1.3.1) for A being semi-coherent dense matrices of various sizes (x-axis)

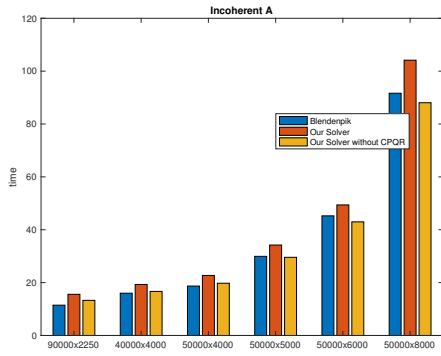


Figure 3.14: Time taken by solvers to compute the solution of problem (1.3.1) for A being incoherent dense matrices of various sizes (x-axis)

3.5.3 Running time performance on randomly generated full rank sparse A

Results, Sparse random matrices Figure 3.15, 3.16, 3.17 show the performance of Ski-LLS compared to LS_HSL and LS_SPQR on sparse random matrices of different types and sizes. We see Ski-LLS can be up to 10 times faster on this class of data matrix A . We also tested on LSRN but do not report the result because LSRN takes much longer than the other solvers for this class of data in the serial running environment.

Note that in this experiment, the solvers are compiled and run in a different machine then mentioned in Section 3.3.2, but all solvers are run on this machine in this experiment. The machine has 2.9 GHz Quad-Core Intel Core i7 CPU and 16MB RAM. Moreover, the parameter $s = 3$ is chosen for Ski-LLS ¹⁹. Furthermore, our solver was an old version without Step 3 implemented and uses the SPQR default ordering. Otherwise the settings are the same as the default settings for Ski-LLS, LS_HSL, LS_SPQR and LSRN.

3.5.4 Large scale benchmark of Ski-LLS-sparse on the Florida Matrix Collection

Performance profile Performance profile [30] has in recent years have become a popular and widely used tool for providing objective information when benchmarking software. In a typical performance profile here, we have the running time ratio against the fastest solver on the x-axis, reported in \log_2 scale. For each running time ratio a , we have the ratio of problems in the test set b on the y-axis such that for a particular solver, the running time ratio against the best solver is within a for b percent of the problems in the test set. For example, the intersect between the performance curve and the y-axis gives the ratio of the problems in the test set such that a particular solver is the fastest.

Running and testing environment specific to the benchmark experiment Given a problem A from the test set, let (r_1, r_2, r_3, r_4) be the residuals of solutions computed by the four solvers compared. And let $r = \min r_i$. A solver is declared as failed on this particular problem if one of the following two conditions holds

1. $r_i > (1 + \tau_r)r$ and $r_i > r + \tau_a$. So that the residual of the solution computed is neither relatively close nor close in absolute value to the residual of the best solution ²⁰.
2. The solver takes more than 800 wall clock seconds to compute a solution.

¹⁹According to Appendix C, the running time of Ski-LLS with $s = 2$ and $s = 3$ is similar.

²⁰Note that the residual of the best solution is in general only an upper bound of the minimal residual. However since it is too computational intensive to compute the minimal residual solution of all the problems in the Florida matrix collection, we use the residual of the best solution as a proxy.

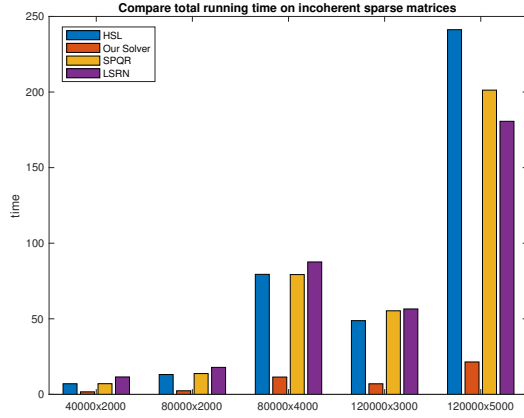


Figure 3.15: Running time comparison of SkiLLS with LS_HSL and LS_SPQR for randomly generated incoherent sparse matrices of different sizes.

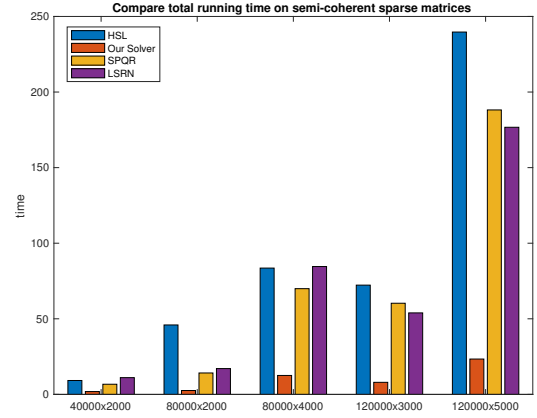


Figure 3.16: Running time comparison of SkiLLS with LS_HSL and LS_SPQR for randomly generated semi-coherent sparse matrices of different sizes.

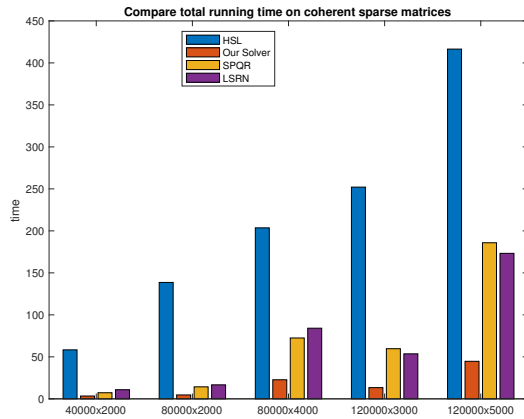


Figure 3.17: Running time comparison of SkiLLS with LS_HSL and LS_SPQR for randomly generated coherent sparse matrices of different sizes.

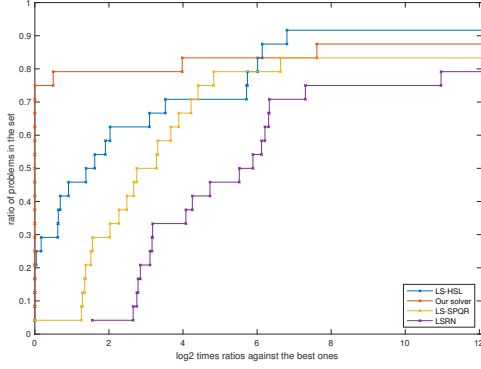


Figure 3.18: Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with $n \geq 30d$.

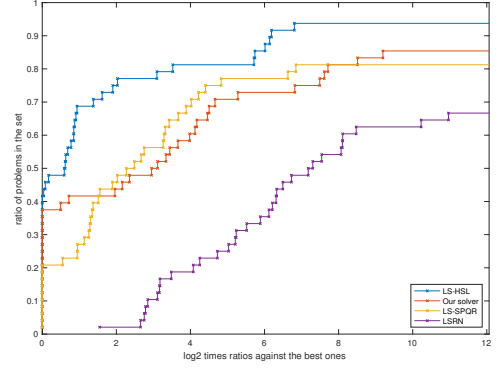


Figure 3.19: Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with $n \geq 10d$.

In the case that a solver is declared as failed, we set the running time of the solver to be 9999 seconds on the corresponding problem. This is because we want to include all the problems such that at least one of the solvers compared succeeded. As a result, a very large running time (ratio) could be due to either an inaccuracy of the solver or an inefficiency of the solver. We note that for all successful solvers, the running time is bounded above by 800 seconds so that there will be no confusion of whether a solver is successful or not.

The default parameters (as described in Section 3.3.2) for all solvers are used.

Results, highly over-determined matrices in the Florida Matrix Collection Figure 3.18 shows Ski-LLS is the fastest in 75% of problems in the Florida matrix collection with $n \geq 30d$.

What happens when the problem is moderately over-determined? Figure 3.19 shows LS_HSL is the fastest for the largest percentage of problems in the Florida matrix collection with $n \geq 10d$. Ski-LLS is still competitive and noticeably faster than LSRN.

Effect of condition number Many of the matrices in the Florida matrix condition has low condition numbers so that unpreconditioned LSQR converges in a few iterations. In those cases, it is disadvantageous to compare Ski-LLS to LS_HSL because we compute a better quality preconditioner through an complete orthogonal factorization.

Figure 3.20 show Ski-LLS is fastest in more than 50% of the moderately over-determined ($n \geq 10d$) problems if we only consider problems such that it takes LSQR more than 5 seconds to solve.

Effect of sparsity Figure 3.21 shows Ski-LLS is extremely competitive, being the fastest in all but one moderately over-determined problems with moderate sparsity ($\text{nnz}(A) \geq 0.01nd$).

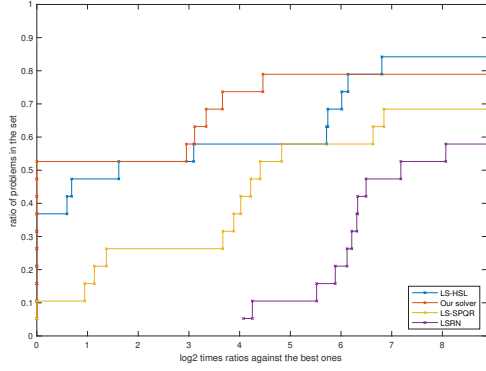


Figure 3.20: Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with $n \geq 10d$ and the unpreconditioned LSQR takes more than 5 seconds to solve.

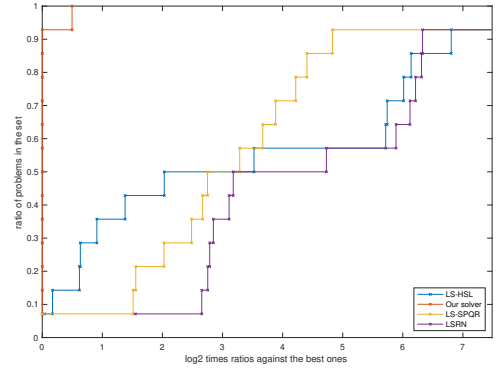


Figure 3.21: Performance profile comparison of Ski-LLS with LSRN, LS_HSL and LS_SPQR for all matrices $A \in \mathbb{R}^{n \times d}$ in the Florida matrix collection with $n \geq 10d$ and $\text{nnz}(A) \geq 0.01nd$.

Chapter 4

First order subspace method for general objectives

4.1 Introduction

This chapter expands the materials in [14, 13]. In Section 4.2, we first describe Algorithm 2, a generic algorithmic framework for solving (1.4.1) by taking successive steps computed from approximately minimising (random) reduced models. Our main result Theorem 4.2.1 provides complexity bound on the total number of iterations before Algorithm 2 drives the gradient of objective below ϵ , with high probability. Deducing from Theorem 4.2.1, we also show that the quantity $\min_{k \leq N} \|\nabla f(x_k)\|_2$ converges to zero with high probability, and the convergence of $\mathbb{E}[\min_{k \leq N} \|\nabla f(x_k)\|_2]$. The rate of these convergences depends on a function in one of the assumptions required by the framework.

In Section 4.3, we prove Theorem 4.2.1. The proof carefully counts the different types of iterations and uses a conditional form of the Chernoff bound whose proof we provide for completeness.

In Section 4.4, we describe Algorithm 3 that particularises Algorithm 2 by using random matrices to build random reduced models. We show how using random matrices that are oblivious JL embeddings satisfies the assumptions required for the convergence result in Theorem 4.2.1.

In Section 4.5, Algorithm 3 is further particularised to a quadratic-regularisation and a trust-region variant, depending on how the minimisation of the random reduced model is specified. Section 4.5 then uses Theorem 4.2.1 to show that both variants drive the full objective gradient $\nabla f(x_k)$ below ϵ in $\mathcal{O}(\frac{1}{\epsilon^2})$ iterations with high probability, matching the deterministic methods' iteration complexity.

In Section 4.6, we introduce non-linear least squares, a particular type of non-convex optimisation problems (1.4.1). We show how Algorithm 3 safe-guarded with trust-region leads to a subspace version of the well-known Gauss-Newton method, Randomised-Subspace Gauss Newton (R-SGN), with a convergence guarantee. We numerically illustrate the performance of R-SGN on non-linear least squares and logistic regression problems.

Related literature [17] proposes a generic algorithmic framework based on probabilistic models with an expected iteration complexity bound to generate a sufficiently small (true) gradient. Various strategies are discussed in [17] to generate such models both for derivative-based and derivative-free methods; however, subspace methods cannot be easily captured within the conditions and models used in [17]. In [45], a trust-region based method with probabilistically accurate models is proposed, with an iteration complexity bound of $\mathcal{O}(\epsilon^{-2})$ for the algorithm to drive the full gradient below ϵ , with high probability. [60] analyses a random subspace method with constant step size, where the sketching matrix S_k satisfies $\mathbb{E}[S_k^T S_k] = \mathbb{I}_d$ and $S_k S_k^T = \frac{d}{t} \mathbb{I}_d$. However, their convergence result requires the objective to be convex, or to satisfy the Polyak-Lojasiewicz inequality. Independently from our work/at the same time, [61] proposes a random subspace gradient descent method with linesearch; it uses Johnson-Lindenstrauss embedding properties in the analysis, similarly to our framework, but fewer ensembles are considered. However, their analysis only applies under various convexity assumptions of the objective.

4.2 General algorithmic framework and its convergence result

We consider the unconstrained optimisation problem

$$f^* = \min_{x \in \mathbb{R}^d} f(x). \quad (4.2.1)$$

4.2.1 Generic algorithmic framework and assumptions

We first describe a generic algorithmic framework that encompasses the main components of the unconstrained optimization schemes we analyse in this chapter. The scheme relies on building a local, reduced model of the objective function at each iteration, minimizing this model or reducing it in a sufficient manner and considering the step which is dependent on a stepsize parameter and which provides the model reduction (the stepsize parameter may be present in the model or independent of it). This step determines a new candidate point. The function value is then computed (accurately) at the new candidate point. If the function reduction provided by the candidate point is deemed sufficient, then the iteration is declared successful, the candidate point becomes the new iterate and the step size parameter is increased. Otherwise, the iteration is unsuccessful, the iterate is not updated and the step size parameter is reduced.

We summarize the main steps of the generic framework below.

Algorithm 2 Generic optimization framework based on randomly generated reduced models

Initialization

Choose a class of (possibly random) models $m_k(w_k(\hat{s})) = \hat{m}_k(\hat{s})$, where $\hat{s} \in \mathbb{R}^l$ with $l \leq d$ is the step parameter and $w_k : \mathbb{R}^l \rightarrow \mathbb{R}^d$ is the prolongation function. Choose constants $\gamma_1 \in (0, 1)$, $\gamma_2 = \gamma_1^{-c}$, for some $c \in \mathbb{N}^+$ (\mathbb{N}^+ refers to the set of positive natural numbers), $\theta \in (0, 1)$ and $\alpha_{\max} > 0$. Initialize the algorithm by setting $x_0 \in \mathbb{R}^d$, $\alpha_0 = \alpha_{\max} \gamma_1^p$ for some $p \in \mathbb{N}^+$ and $k = 0$.

1. Compute a reduced model and a step

Compute a local (possibly random) reduced model $\hat{m}_k(\hat{s})$ of f around x_k with $\hat{m}_k(0) = f(x_k)$. Compute a step parameter $\hat{s}_k(\alpha_k)$, where the parameter α_k is present in the reduced model or the step parameter computation. Compute a potential step $s_k = w_k(\hat{s}_k)$.

2. Check sufficient decrease

Compute $f(x_k + s_k)$ and check if sufficient decrease (parameterized by θ) is achieved in f with respect to $\hat{m}_k(0) - \hat{m}_k(\hat{s}_k(\alpha_k))$.

3. Update the parameter α_k and possibly take the potential step s_k

If sufficient decrease is achieved, set $x_{k+1} = x_k + s_k$ and $\alpha_{k+1} = \min\{\alpha_{\max}, \gamma_2 \alpha_k\}$ [this is referred to as a successful iteration].

Otherwise set $x_{k+1} = x_k$ and $\alpha_{k+1} = \gamma_1 \alpha_k$ [this is an unsuccessful iteration].

Increase the iteration count by setting $k = k + 1$ in both cases.

The generic framework and its assumptions we present here is similar to the framework presented in [17]. We extended their framework so that the proportionality constants for increase/decrease of the step size parameter are not required to be reciprocal, but reciprocal up to an integer power (see Assumption 2). Even though the framework and its assumptions are similar, our analysis and result are different and qualitatively improve upon their result (in the way that Theorem 4.2.1 implies their main result Theorem 2.1). Moreover, we show how to use random-embedding based sketching to build the reduced model in Section 4.4.

The connection between the generic framework and classical optimisation literature is detailed in [17]. Here we give a simple example. If we let $l = d$ and w_k be the identity function in Algorithm 2 so that $m_k(s) = \hat{m}_k(\hat{s})$; and if we let

$$m_k(s) = f(x_k) + \nabla f(x_k)^T s + \frac{1}{2} s^T B_k s,$$

where B_k is a Hessian approximation, and compute the step parameter (or in this case, since w_k is the identity function, the step) by seeking a solution of the problem

$$\min_{s \in \mathbb{R}^d} m_k(s) \quad \text{such that } \|s\|_2 \leq \alpha_k;$$

and if we define the sufficient decrease by

$$f(x_k) - f(x_k + s_k) \geq \theta [m_k(0) - m_k(s_k)],$$

then Algorithm 2 reduces to the (deterministic) trust-region method, see [85].

Because the model at each iteration is (possibly) random, x_k, s_k, α_k are in general random variables. We will use $\bar{x}_k, \bar{s}_k, \bar{\alpha}_k$ to denote their realizations. We define convergence in terms of a random variable N_ϵ , that can be a function of a positive scalar(s) ϵ , as well as the sequences $\{f(x_k)\}, \{\nabla f(x_k)\}, \{\nabla^2 f(x_k)\}$. For example,

$$N_\epsilon = \min \{k : \|\nabla f(x_k)\|_2 \leq \epsilon\} \quad (4.2.2)$$

will be used to represent convergence to a first-order local stationary point, as in [17]. We say that Algorithm 2 has not converged if $k \leq N_\epsilon$, and has converged otherwise. Furthermore, let us suppose that there is a class of iterations, hereafter will be referred to as **true iterations** such that Algorithm 2 satisfies the some conditions.

The convergence of Algorithm 2 relies on the following four assumptions. The first assumption states that given the current iterate (at any value), an iteration k is true at least with a fixed probability, and is independent of the truth values of all previous iterations.

Assumption 1. *There exists $\delta_S \in (0, 1)$ such that for any $\bar{x}_k \in \mathbb{R}^d$ and $k = 1, 2, \dots$*

$$\mathbb{P}(T_k | x_k = \bar{x}_k) \geq 1 - \delta_S,$$

where T_k is defined as

$$T_k = \begin{cases} 1, & \text{if iteration } k \text{ is true} \\ 0, & \text{otherwise.} \end{cases} \quad (4.2.3)$$

Moreover, $\mathbb{P}(T_0) \geq 1 - \delta_S$; and T_k is conditionally independent of T_0, T_1, \dots, T_{k-1} given $x_k = \bar{x}_k$.

The next assumption says that for α_k small enough, any true iteration before convergence is guaranteed to be successful.

Assumption 2. *For any $\epsilon > 0$, there exists an iteration-independent constant $\alpha_{low} > 0$ (that may depend on ϵ and the problem and algorithm parameters) such that if iteration k is true, $k < N_\epsilon$, and $\alpha_k \leq \alpha_{low}$ then iteration k is successful.*

The next assumption says that before convergence, true and successful iterations result in an objective decrease lower bounded by an (iteration-independent) function h , which is monotonically increasing in its two arguments, ϵ and α_k .

Assumption 3. *There exists a non-negative, non-decreasing function $h(z_1, z_2)$ such that, for any $\epsilon > 0$, if iteration k is true and successful with $k < N_\epsilon$, then*

$$f(x_k) - f(x_k + s_k) \geq h(\epsilon, \alpha_k), \quad (4.2.4)$$

where s_k is computed in step 1 of Algorithm 2. Moreover, $h(z_1, z_2) > 0$ if both $z_1 > 0$ and $z_2 > 0$.

The final assumption requires that the function values at successive iterations must form a non-increasing sequence throughout the algorithm.

Assumption 4. For any $k \in \mathbb{N}$, we have

$$f(x_k) \geq f(x_{k+1}). \quad (4.2.5)$$

The following Lemma is a simple consequence of Assumption 2.

Lemma 4.2.1. Let $\epsilon > 0$ and Assumption 2 hold with $\alpha_{low} > 0$. Then there exists $\tau_\alpha \in \mathbb{N}^+$, and $\alpha_{min} > 0$ such that

$$\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}, \quad (4.2.6)$$

$$\alpha_{min} \leq \alpha_{low},$$

$$\alpha_{min} \leq \frac{\alpha_0}{\gamma_2}, \quad (4.2.7)$$

where $\gamma_1, \gamma_2, \alpha_0$ are defined in Algorithm 2.

Proof. Let

$$\tau_\alpha = \left\lceil \log_{\gamma_1} \left(\min \left\{ \frac{\alpha_{low}}{\alpha_0}, \frac{1}{\gamma_2} \right\} \right) \right\rceil, \quad (4.2.8)$$

$$\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}. \quad (4.2.9)$$

We have that $\alpha_{min} \leq \alpha_0 \gamma_1^{\log_{\gamma_1} \left(\frac{\alpha_{low}}{\alpha_0} \right)} = \alpha_{low}$. Therefore by Assumption 2, if iteration k is true, $k < N_\epsilon$, and $\alpha_k \leq \alpha_{min}$ then iteration k is successful. Moreover, $\alpha_{min} \leq \alpha_0 \gamma_1^{\log_{\gamma_1} \left(\frac{1}{\gamma_2} \right)} = \frac{\alpha_0}{\gamma_2} = \alpha_0 \gamma_1^c$. It follows from $\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}$ that $\tau_\alpha \geq c$. Since $c \in \mathbb{N}^+$, we have $\tau_\alpha \in \mathbb{N}^+$ as well. \square

4.2.2 A probabilistic convergence result

Theorem 4.2.1 is our main result for Algorithm 2. It states a probabilistic bound on the total number of iterations N_ϵ needed to converge to ϵ -accuracy for the generic framework.

Theorem 4.2.1. Let Assumption 1, Assumption 2, Assumption 3 and Assumption 4 hold with $\delta_S \in (0, 1)$, $\alpha_{low} > 0$, $h : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}$ associated with α_{low} , for some $\tau_\alpha \in \mathbb{N}^+$. Let $\epsilon > 0$, f^* defined in (4.2.1). Run Algorithm 2 for N iterations. Suppose

$$\delta_S < \frac{c}{(c+1)^2} \quad (4.2.10)$$

Then for any $\delta_1 \in (0, 1)$ such that

$$g(\delta_S, \delta_1) > 0, \quad (4.2.11)$$

where

$$g(\delta_S, \delta_1) = \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1}. \quad (4.2.12)$$

If N satisfies

$$N \geq g(\delta_S, \delta_1) \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right], \quad (4.2.13)$$

we have that

$$\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N}.\footnote{1} \quad (4.2.14)$$

Remark 4. Note that $\frac{c}{(c+1)^2} \in (0, \frac{1}{4}]$ for $c \in \mathbb{N}^+$. Therefore (4.2.10) and (4.2.11) can only be satisfied for some c, δ_1 given that $\delta_S < \frac{1}{4}$. Thus our theory requires an iteration is true with probability at least $\frac{3}{4}$. Compared to the analysis in [17], which requires an iteration is true with probability at least $\frac{1}{2}$, our condition is stronger. This is due to the high probability nature of our result, while their convergence result is in expectation. Furthermore, we will see in Lemma 4.4.2 that we are able to impose arbitrarily small value of δ_S , thus satisfying the requirement, by choosing an appropriate dimension of the local reduced model $\hat{m}_k(\hat{s})$.

We show how our result leads to Theorem 2.1 in [17], which concerns $\mathbb{E}[N_\epsilon]$. We have, with N_0 defined as the RHS of (4.2.13),

$$\begin{aligned} \mathbb{E}[N_\epsilon] &= \int_0^\infty \mathbb{P}(N_\epsilon > M) dM \\ &= \int_0^{N_0} \mathbb{P}(N_\epsilon > M) dM + \int_{N_0}^\infty \mathbb{P}(N_\epsilon > M) dM \\ &\leq N_0 + \int_{N_0}^\infty \mathbb{P}(N_\epsilon > M) dM \\ &\leq N_0 + \int_{N_0}^\infty e^{-\frac{\delta_1^2}{2}(1-\delta_S)M} dM \\ &= N_0 + \frac{2}{\delta_1^2(1-\delta_S)} e^{-\frac{\delta_1^2}{2}(1-\delta_S)N_0}, \end{aligned}$$

where we used Theorem 4.2.1 to derive the last inequality. The result in [17] is in the form of $\mathbb{E}[N_\epsilon] \leq N_0$. Note that the discrepancy term is exponentially small in terms of N_0 and therefore the implication of our result is asymptotically the same as that in [17].

4.2.3 Corollaries of Theorem 4.2.1

Before we begin the proof, we state and prove three implications of Theorem 4.2.1, provided some mild assumptions on h with N_ϵ . These results show different flavours of Theorem 4.2.1.

The following expressions will be used, along with $g(\delta_S, \delta_1)$ is defined in (4.2.11)

$$q(\epsilon) = h(\epsilon, \gamma_1^c \alpha_{\min}), \quad (4.2.15)$$

$$D_1 = g(\delta_S, \delta_1)(f(x_0) - f^*), \quad (4.2.16)$$

$$D_2 = g(\delta_S, \delta_1) \frac{\tau_\alpha}{1+c}, \quad (4.2.17)$$

$$D_3 = \frac{\delta_1^2}{2}(1-\delta_S). \quad (4.2.18)$$

¹For the sake of clarity, we stress that N is a deterministic constant, namely, the total number of iterations that we run Algorithm 2. N_ϵ , the number of iterations needed before convergence, is a random variable.

From (4.2.15), (4.2.16), (4.2.17), (4.2.18), a sufficient condition for (4.2.13) to hold is

$$\begin{aligned} N &\geq g(\delta_S, \delta_1) \left[\frac{f(x_0) - f^*}{h(\epsilon, \gamma_1^c \alpha_{\min})} + \frac{\tau_\alpha}{1+c} \right] \\ &= \frac{D_1}{q(\epsilon)} + D_2; \end{aligned}$$

and (4.2.14) can be restated as

$$\mathbb{P}(N > N_\epsilon) \geq 1 - e^{-D_3 N}.$$

The first corollary gives the rate of change of $\min_{k \leq N} \|\nabla f(x_k)\|_2$ as $N \rightarrow \infty$. It will yield a rate of convergence by substituting in a specific expression of h (and hence q^{-1}).

Corollary 4.2.1. *Let Assumption 1, Assumption 2, Assumption 3, Assumption 4 hold.*

Let f^, q, D_1, D_2, D_3 be defined in (4.2.1), (4.2.15), (4.2.16), (4.2.17) and (4.2.18). Suppose (4.2.10) hold and let $\delta_1 \in (0, 1)$ satisfy (4.2.11). Then for any $N \in \mathbb{N}$ such that $q^{-1}\left(\frac{D_1}{N-D_2}\right)$ exists, we have*

$$\mathbb{P}\left(\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq q^{-1}\left(\frac{D_1}{N-D_2}\right)\right) \geq 1 - e^{-D_3 N}. \quad (4.2.19)$$

Proof. Let $N \in \mathbb{N}$ such that $q^{-1}\left(\frac{D_1}{N-D_2}\right)$ exists and let $\epsilon = q^{-1}\left(\frac{D_1}{N-D_2}\right)$. Then we have

$$\mathbb{P}\left(\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq q^{-1}\left(\frac{D_1}{N-D_2}\right)\right) = \mathbb{P}\left(\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq \epsilon\right) \geq \mathbb{P}(N > N_\epsilon), \quad (4.2.20)$$

where the inequality follows from the fact that $N \geq N_\epsilon$ implies $\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq \epsilon$. On the other hand, we have

$$\begin{aligned} N &= \frac{D_1}{\frac{D_1}{N-D_2}} + D_2 \\ &= \frac{D_1}{q(\epsilon)} + D_2. \end{aligned}$$

Therefore (4.2.13) holds; and applying Theorem 4.2.1, we have that $\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-D_3 N}$. Hence (4.2.20) gives the desired result. □

The next Corollary restates Theorem 4.2.1 for a fixed arbitrarily high success probability.

Corollary 4.2.2. *Let Assumption 1, Assumption 2, Assumption 3, Assumption 4 hold. Suppose (4.2.10) hold and let $\delta_1 \in (0, 1)$ satisfy (4.2.11). Then for any $\delta \in (0, 1)$, suppose*

$$N \geq \max\left\{\frac{D_1}{q(\epsilon)} + D_2, \frac{\log\left(\frac{1}{\delta}\right)}{D_3}\right\}, \quad (4.2.21)$$

where D_1, D_2, D_3, q are defined in (4.2.16), (4.2.17), (4.2.18) and (4.2.15). Then

$$\mathbb{P}\left(\min_{k \leq N} \|\nabla f(x_k)\|_2 < \epsilon\right) \geq 1 - \delta.$$

Proof. We have

$$\mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq \epsilon \right) \geq \mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-D_3 N} \geq 1 - \delta,$$

where the first inequality follows from definition of N_ϵ in (4.2.2), the second inequality follows from Theorem 4.2.1 (note that (4.2.21) implies (4.2.13)) and the last inequality follows from (4.2.21). \square

The next Corollary gives the rate of change of the expected value of $\min_{k \leq N} \|\nabla f(x_k)\|_2$ as N increases.

Corollary 4.2.3. *Let Assumption 1, Assumption 2, Assumption 3, Assumption 4 hold. Suppose (4.2.10) hold and let $\delta_1 \in (0, 1)$ satisfy (4.2.11). Then for any $N \in \mathbb{N}$ such that $q^{-1} \left(\frac{D_1}{N-D_2} \right)$ exists, where q, D_1, D_2 are defined in (4.2.15), (4.2.16), (4.2.17), we have*

$$\mathbb{E} \left[\min_{k \leq N} \|\nabla f(x_k)\|_2 \right] \leq q^{-1} \left(\frac{D_1}{N-D_2} \right) + \|\nabla f(x_0)\|_2 e^{-D_3 N},$$

where D_3 is defined in (4.2.18) and x_0 is chosen in Algorithm 2.

Proof. We have

$$\begin{aligned} & \mathbb{E} \left[\min_{k \leq N} \|\nabla f(x_k)\|_2 \right] \\ & \leq \mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq q^{-1} \left(\frac{D_1}{N-D_2} \right) \right) q^{-1} \left(\frac{D_1}{N-D_2} \right) \\ & + \mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 > q^{-1} \left(\frac{D_1}{N-D_2} \right) \right) \|\nabla f(x_0)\|_2 \\ & \leq q^{-1} \left(\frac{D_1}{N-D_2} \right) + e^{-D_3 N} \|\nabla f(x_0)\|_2, \end{aligned}$$

where for the first inequality, we split the integral in the definition of expectation

$$\begin{aligned} \mathbb{E} \left[\min_{k \leq N} \|\nabla f(x_k)\|_2 \right] &= \int_0^\infty \mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 = x \right) x dx \\ &= \int_0^{q^{-1} \left(\frac{D_1}{N-D_2} \right)} \mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 = x \right) x dx + \int_{q^{-1} \left(\frac{D_1}{N-D_2} \right)}^\infty \mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 = x \right) x dx \end{aligned}$$

; and used $\mathbb{P}(\min_{k \leq N} \|\nabla f(x_k)\|_2 = x) = 0$ for $x > \|\nabla f(x_0)\|_2$ which follows from $\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq \|\nabla f(x_0)\|_2$. For the second inequality,

we used $\mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 \leq q^{-1} \left(\frac{D_1}{N-D_2} \right) \right) \leq 1$ and $\mathbb{P} \left(\min_{k \leq N} \|\nabla f(x_k)\|_2 > q^{-1} \left(\frac{D_1}{N-D_2} \right) \right) \leq e^{-D_3 N}$ by (4.2.19). \square

4.3 Proof of Theorem 4.2.1

The proof of Theorem 4.2.1 involves a technical analysis of different types of iterations. An iteration can be true/false using Definition 4.4.1, successful/unsuccessful (Step 3 of Algorithm 2) and with an α_k above/below a certain value. The parameter α_k is important due to Assumption 2 and

Assumption 3 (that is, it influences the success of an iteration; and also the objective decrease in true and successful iterations).

Given that Algorithm 2 runs for N iterations, we use N with different subscripts to denote the total number of different types of iterations, detailed in Table 4.1. We note that they are all random variables because α_k , and whether an iteration is true/false, successful/unsuccessful all depend on the random model in Step 1 of Algorithm 2 and the previous (random) steps.

Symbol	Definition
N_T	Number of true iterations
N_F	Number of false iterations
N_{TS}	Number of true and successful iterations
N_S	Number of successful iterations
N_U	Number of unsuccessful iterations
N_{TU}	Number of true and unsuccessful iterations
$N_{T,\alpha_{min}}$	Number of true iterations such that $\alpha_k \leq \alpha_{min}$
$N_{S,\alpha_{min}}$	Number of successful iterations such that $\alpha_k \leq \alpha_{min}$
$N_{T,\alpha_{min}}$	Number of true iterations such that $\alpha_k > \alpha_{min}$
$N_{TS,\alpha_{min}}$	Number of true and successful iterations such that $\alpha_k > \alpha_{min}$
$N_{TU,\alpha_{min}}$	Number of true and unsuccessful iterations such that $\alpha_k > \alpha_{min}$
$N_{U,\alpha_{min}}$	Number of unsuccessful iterations such that $\alpha_k > \alpha_{min}$
$N_{S,\gamma_1^c \alpha_{min}}$	Number of successful iterations such that $\alpha_k > \gamma_1^c \alpha_{min}$
$N_{TS,\gamma_1^c \alpha_{min}}$	Number of true and successful iterations such that $\alpha_k > \gamma_1^c \alpha_{min}$
$N_{FS,\gamma_1^c \alpha_{min}}$	Number of false and successful iterations such that $\alpha_k > \gamma_1^c \alpha_{min}$

Table 4.1: List of random variables representing iteration counts given that Algorithm 2 has run for N iterations

The proof of Theorem 4.2.1 relies on the following three results relating the total number of different types of iterations.

The relationship between the total number of true iterations and the total number of iterations Lemma 4.3.1 shows that with high probability, a constant fraction of iterations of Algorithm 2 are true. This result is a conditional variant of the Chernoff bound [20].

Lemma 4.3.1. *Let Assumption 1 hold with $\delta_S \in (0, 1)$. Let Algorithm 2 run for N iterations. Then for any given $\delta_1 \in (0, 1)$,*

$$\mathbb{P}(N_T \leq (1 - \delta_S)(1 - \delta_1)N) \leq e^{-\frac{\delta_1^2}{2}(1 - \delta_S)N}, \quad (4.3.1)$$

where N_T is defined in Table 4.1.

The proof of Lemma 4.3.1 relies on the below technical result.

Lemma 4.3.2. *Let Assumption 1 hold with $\delta_S \in (0, 1)$. Let T_k be defined in (4.2.3). Then for any $\lambda > 0$ and $N \in \mathbb{N}$, we have*

$$\mathbb{E} \left[e^{-\lambda \sum_{k=0}^{N-1} T_k} \right] \leq \left[e^{(e^{-\lambda} - 1)(1 - \delta_S)} \right]^N.$$

Proof. Let $\lambda > 0$. We use induction on N . For $N = 1$, we want to show

$$\mathbb{E} [e^{-\lambda T_0}] \leq e^{(e^{-\lambda} - 1)(1 - \delta_S)}.$$

Let $g(x) = e^{-\lambda x}$. Note that

$$g(x) \leq g(0) + [g(1) - g(0)]x, \text{ for any } x \in [0, 1], \quad (4.3.2)$$

because $g(x)$ is convex. Substituting $x = T_0$, we have

$$e^{-\lambda T_0} \leq 1 + (e^{-\lambda} - 1)T_0.$$

Taking expectation, we have that

$$\mathbb{E} [e^{-\lambda T_0}] \leq 1 + (e^{-\lambda} - 1)\mathbb{E} [T_0]. \quad (4.3.3)$$

Moreover, we have

$$\mathbb{E} [T_0] \geq \mathbb{P} (T_0 = 1) \geq 1 - \delta_S,$$

where the first inequality comes from $T_0 \geq 0$ and the second inequality comes from Assumption 1.

Therefore, noting that $e^{-\lambda} - 1 < 0$, (4.3.3) gives

$$\mathbb{E} [e^{-\lambda T_0}] \leq 1 + (e^{-\lambda} - 1)(1 - \delta_S) \leq e^{(e^{-\lambda} - 1)(1 - \delta_S)}, \quad (4.3.4)$$

where the last inequality comes from $1 + y \leq e^y$ for $y \in \mathbb{R}$.

Having completed the initial step for the induction, let us assume

$$\mathbb{E} [e^{-\lambda \sum_{k=0}^{N-2} T_k}] \leq [e^{(e^{-\lambda} - 1)(1 - \delta_S)}]^{N-1}. \quad (4.3.5)$$

We have

$$\begin{aligned} & \mathbb{E} [e^{-\lambda \sum_{k=0}^{N-1} T_k}] \\ &= \mathbb{E} \left[\mathbb{E} [e^{-\lambda \sum_{k=0}^{N-1} T_k} | T_0, T_1, \dots, T_{N-2}, x_{N-1}] \right] \\ &= \mathbb{E} [e^{-\lambda \sum_{k=0}^{N-2} T_k} \mathbb{E} [e^{-\lambda T_{N-1}} | T_0, T_1, \dots, T_{N-2}, x_{N-1}]] \\ &= \mathbb{E} [e^{-\lambda \sum_{k=0}^{N-2} T_k} \mathbb{E} [e^{-\lambda T_{N-1}} | x_{N-1}]], \end{aligned} \quad (4.3.6)$$

where the first equality comes from the Tower property and the last equality follows from T_{N-1} is conditionally independent of T_0, T_1, \dots, T_{N-2} given x_{N-1} (see Assumption 1).

Substituting $x = T_{N-1}$ in (4.3.2), and taking conditional expectation, we have that

$$\mathbb{E} [e^{-\lambda T_{N-1}} | x_{N-1}] \leq 1 + (e^{-\lambda} - 1)\mathbb{E} [T_{N-1} | x_{N-1}].$$

On the other hand, we have that $\mathbb{E} [T_{N-1} | x_{N-1}] \geq \mathbb{P} [T_{N-1} = 1 | x_{N-1}] \geq 1 - \delta_S$, where we used $T_{N-1} \geq 0$ to derive the first inequality and $\mathbb{P} [T_{N-1} = 1 | x_{N-1} = \bar{x}_{N-1}] \geq 1 - \delta$ for any \bar{x}_{N-1} (see Assumption 1) to derive the second inequality. Hence

$$\mathbb{E} [e^{-\lambda T_{N-1}} | x_{N-1}] \leq e^{(e^{-\lambda} - 1)(1 - \delta_S)},$$

as in (4.3.4).

It then follows from (4.3.6) that

$$\mathbb{E} \left[e^{-\lambda \sum_{k=0}^{N-1} T_k} \right] \leq e^{(e^{-\lambda}-1)(1-\delta_S)} \mathbb{E} \left[e^{-\lambda \sum_{k=0}^{N-2} T_k} \right] \leq \left[e^{(e^{-\lambda}-1)(1-\delta_S)} \right]^N,$$

where we used (4.3.5) for the last inequality. \square

Proof of Lemma 4.3.1. Note that with N being the total number of iterations, we have $N_T = \sum_{k=0}^{N-1} T_k$, where T_k is defined in (4.2.3). Applying Markov inequality, we have that for any $\lambda > 0$,

$$\begin{aligned} \mathbb{P}(N_T \leq (1-\delta_S)(1-\delta_1)N) &= \mathbb{P}\left(e^{-\lambda N_T} \geq e^{-\lambda(1-\delta_S)(1-\delta_1)N}\right) \\ &\leq \mathbb{E} \left[e^{-\lambda N_T} \right] e^{\lambda(1-\delta_S)(1-\delta_1)N} \\ &= \mathbb{E} \left[e^{-\lambda \sum_{k=0}^{N-1} T_k} \right] e^{\lambda(1-\delta_S)(1-\delta_1)N} \\ &\leq e^{N(e^{-\lambda}-1)(1-\delta_S)+\lambda(1-\delta_S)(1-\delta_1)N}, \end{aligned} \quad (4.3.7)$$

where we used Lemma 4.3.2 to derive the last inequality.

Choosing $\lambda = -\log(1-\delta_1) > 0$, we have from (4.3.7)

$$\begin{aligned} \mathbb{P}(N_T \leq (1-\delta_S)(1-\delta_1)N) &\leq e^{N(1-\delta_S)[- \delta_1 - (1-\delta_1) \log(1-\delta_1)]} \\ &\leq e^{-\frac{\delta_1^2}{2}(1-\delta_S)N}, \end{aligned}$$

where we used $-\delta_1 - (1-\delta_1) \log(1-\delta_1) \leq -\delta_1^2/2$ for $\delta_1 \in (0, 1)$. \square

The relationship between the total number of true iterations with $\alpha_k \leq \alpha_{min}$ and the total number of iterations The next Lemma shows that we can have at most a constant fraction of iterations of Algorithm 2 that are true with $\alpha_k \leq \alpha_{min}$.

Lemma 4.3.3. *Let Assumption 2 hold with $\alpha_{low} > 0$ and $c \in \mathbb{N}^+$ and let α_{min} associated with α_{low} be defined in (4.2.6) with $\tau_\alpha \in \mathbb{N}^+$. Let $\epsilon > 0$, $N \in \mathbb{N}$ be the total number of iterations; and $N_{T, \overline{\alpha_{min}}}$ be defined in Table 4.1. Suppose $N \leq N_\epsilon$. Then*

$$N_{T, \overline{\alpha_{min}}} \leq \frac{N}{c+1}. \quad (4.3.8)$$

Proof. Let $k \leq N-1$ ². It follows from $N \leq N_\epsilon$ that $k < N_\epsilon$ and by definition of α_{min} (Lemma 4.2.1), iteration k is true with $\alpha_k \leq \alpha_{min}$ implies that iteration k is successful (with $\alpha_k \leq \alpha_{min}$).

Therefore we have

$$N_{T, \overline{\alpha_{min}}} \leq N_{S, \overline{\alpha_{min}}}. \quad (4.3.9)$$

If $N_{S, \overline{\alpha_{min}}} = 0$, then $N_{T, \overline{\alpha_{min}}} = 0$ and (4.3.8) holds. Otherwise let

$$\bar{k} = \max \{k \leq N-1 : \text{iteration } k \text{ is successful and } \alpha_k \leq \alpha_{min}\}. \quad (4.3.10)$$

²Note that $k = 0, 1, \dots, N-1$ if the total number of iterations is N .

Then for each $k \in \{0, 1, \dots, \bar{k}\}$, we have that either iteration k is successful and $\alpha_k \leq \alpha_{min}$, in which case $\alpha_{k+1} = \gamma_2 \alpha_k$ (note that (4.2.7) and $\alpha_k \leq \alpha_{min}$ ensure $\max\{\gamma_2 \alpha_k, \alpha_{max}\} = \gamma_2 \alpha_k$); or otherwise $\alpha_{k+1} \geq \gamma_1 \alpha_k$ (which is true for any iteration of *Algorithm 2*). Hence after $\bar{k} + 1$ iterations, we have

$$\begin{aligned} \alpha_{\bar{k}+1} &\geq \alpha_0 \gamma_2^{N_{S, \overline{\alpha_{min}}}} \gamma_1^{\bar{k}+1 - N_{S, \overline{\alpha_{min}}}} = \alpha_0 \left(\frac{\gamma_2}{\gamma_1} \right)^{N_{S, \overline{\alpha_{min}}}} \gamma_1^{\bar{k}+1} \\ &\geq \alpha_0 \left(\frac{\gamma_2}{\gamma_1} \right)^{N_{S, \overline{\alpha_{min}}}} \gamma_1^N, \end{aligned} \quad (4.3.11)$$

where we used $\bar{k} + 1 \leq N$ for the last inequality. On the other hand, we have

$$\alpha_{\bar{k}+1} = \gamma_2 \alpha_{\bar{k}} \leq \gamma_2 \alpha_{min},$$

where we used iteration \bar{k} is successful and $\alpha_{\bar{k}} \leq \alpha_{min}$ from (4.3.10). Therefore, combining the last displayed equation with (4.3.11), we have $\gamma_2 \alpha_{min} \geq \alpha_{\bar{k}+1} \geq \alpha_0 \left(\frac{\gamma_2}{\gamma_1} \right)^{N_{S, \overline{\alpha_{min}}}} \gamma_1^N$. Taking logarithm on both sides, we have

$$\log(\gamma_2 \alpha_{min}) \geq \log(\alpha_0) + N_{S, \overline{\alpha_{min}}} \log\left(\frac{\gamma_2}{\gamma_1}\right) + N \log(\gamma_1).$$

Rearranging, we have

$$N_{S, \overline{\alpha_{min}}} \leq p_0 N + p_1,$$

with $p_0 = \frac{\log(1/\gamma_1)}{\log(\gamma_2/\gamma_1)} = \frac{1}{c+1}$ and $p_1 = \frac{\log(\gamma_2 \alpha_{min}/\alpha_0)}{\log(\gamma_2/\gamma_1)} = \frac{c - \tau_\alpha}{c+1} \leq 0$ as $\tau_\alpha \geq c > 0$. Therefore we have $N_{S, \overline{\alpha_{min}}} \leq \frac{N}{c+1}$ and (4.3.9) then gives the desired result. \square

The relationship between the number of unsuccessful iterations and the number of successful iterations The next Lemma formalises the intuition that one cannot have too many unsuccessful iterations with $\alpha_k > \alpha_{min}$ compared to successful iterations with $\alpha_k > \gamma_1^c \alpha_{min}$, because unsuccessful iterations reduce α_k and only successful iterations with $\alpha_k > \gamma_1^c \alpha_{min}$ may compensate for these decreases. The conditions that $\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}$, $\gamma_2 = \frac{1}{\gamma_1^c}$ and $\alpha_{max} = \alpha_0 \gamma_1^p$ for some $\tau_\alpha, c, p \in \mathbb{N}^+$ are crucial in the (technical) proof.

Lemma 4.3.4. *Let Assumption 2 hold with $\alpha_{low} > 0$. Let α_{min} associated with α_{low} be defined in (4.2.6) with $\tau_\alpha \in \mathbb{N}^+$. Let $N \in \mathbb{N}$ be the total number of iterations of Algorithm 2 and $N_{U, \underline{\alpha_{min}}}$, $N_{S, \underline{\gamma_1^c \alpha_{min}}}$ be defined in Table 4.1. Then*

$$N_{U, \underline{\alpha_{min}}} \leq \tau_\alpha + c N_{S, \underline{\gamma_1^c \alpha_{min}}}.$$

Proof. Define

$$\beta_k = \log_{\gamma_1} \left(\frac{\alpha_k}{\alpha_0} \right). \quad (4.3.12)$$

Note that since $\alpha_{k+1} = \gamma_1 \alpha_k$ if iteration k is successful and $\alpha_{k+1} = \min\{\alpha_{max}, \gamma_2 \alpha_k\}$ otherwise, $\gamma_2 = \frac{1}{\gamma_1^c}$ and $\alpha_{max} = \alpha_0 \gamma_1^p$ with $c, p \in \mathbb{N}^+$, we have that $\beta_k \in \mathbb{Z}$. Moreover, we have that $\alpha_k = \alpha_0$

corresponds to $\beta_k = 0$, $\alpha_k = \alpha_{min}$ corresponds to $\beta_k = \tau_\alpha$ and $\alpha_k = \gamma^c \alpha_{min}$ corresponds to $\beta_k = \tau_\alpha + c$. Note also that on successful iterations, we have $\alpha_{k+1} \leq \gamma_2 \alpha_k = \gamma_1^{-c} \alpha_k$ (as $\alpha_{k+1} = \min\{\alpha_{max}, \gamma_2 \alpha_k\}$) so that $\beta_{k+1} \geq \beta_k - c$; and on unsuccessful iterations, we have $\beta_{k+1} = \beta_k + 1$.

Let $k_{start}^{(1)} = -1$; and define the following sets.

$$A^{(1)} = \left\{ k \in \left(k_{start}^{(1)}, N-1 \right] \cap \mathbb{N} : \beta_k = \tau_\alpha \right\}. \quad (4.3.13)$$

$$k_{end}^{(1)} = \begin{cases} \inf A^{(1)}, & \text{if } A^{(1)} \neq \emptyset \\ N, & \text{otherwise.} \end{cases}$$

$$\begin{aligned} M_1^{(1)} &= \left\{ k \in \left(k_{start}^{(1)}, k_{end}^{(1)} \right) : \text{iteration } k \text{ is unsuccessful with } \beta_k < \tau_\alpha \right\} \\ M_2^{(1)} &= \left\{ k \in \left(k_{start}^{(1)}, k_{end}^{(1)} \right) : \text{iteration } k \text{ is successful with } \beta_k < \tau_\alpha + c \right\}. \end{aligned} \quad (4.3.14)$$

Let $n_1^{(1)} = |M_1^{(1)}|$ and $n_2^{(1)} = |M_2^{(1)}|$, where $|\cdot|$ denotes the cardinality of a set.

If $k_{end}^{(1)} < N$, we have that $k_{end}^{(1)}$ is the first time β_k reaches τ_α . Because β_k starts at $0 < \tau_\alpha$ when $k = 0$; β_k increases by one on unsuccessful iterations and decreases by an integer on successful iterations (so that β_k remains an integer). So for $k \in \left(k_{start}^{(1)}, k_{end}^{(1)} \right)$, all iterates have $\beta_k < \tau_\alpha < \tau_\alpha + c$. It follows then the number of successful/unsuccessful iterations for $k \in \left(k_{start}^{(1)}, k_{end}^{(1)} \right)$ are precisely $n_1^{(1)}$ and $n_2^{(1)}$ respectively. Because β_k decreases by at most c on successful iterations, increases by one on unsuccessful iterations, starts at zero and $\beta_{k_{end}^{(1)}} \leq \tau_\alpha$, we have $0 + n_1^{(1)} - cn_2^{(1)} \leq \tau_\alpha$ (using $\beta_{k_{end}^{(1)}} \geq \beta_{k_{start}^{(1)}+1} + n_1^{(1)} - cn_2^{(1)}$). Rearranging gives

$$n_1^{(1)} \leq cn_2^{(1)} + \tau_\alpha. \quad (4.3.15)$$

If $k_{end}^{(1)} = N$, then we have that $\beta_k < \tau_\alpha$ for all $k \leq N-1$ and so $\beta_{k_{end}^{(1)}} \leq \tau_\alpha$. In this case we can derive (4.3.15) using the same argument. Moreover, since $k_{end}^{(1)} = N$, we have that

$$n_1^{(1)} = N_{U, \underline{\alpha_{min}}}, \quad (4.3.16)$$

$$n_1^{(2)} = N_{S, \underline{\gamma_1^c \alpha_{min}}}. \quad (4.3.17)$$

The desired result then follows.

Hence we only need to continue in the case where $k_{end}^{(1)} < N$, in which case let

$$\begin{aligned} B^{(1)} &= \left\{ k \in \left[k_{end}^{(1)}, N-1 \right] : \text{iteration } k \text{ is successful with } \beta_k < \tau_\alpha + c \right\} \\ k_{start}^{(2)} &= \begin{cases} \inf B^{(1)}, & \text{if } B^{(1)} \neq \emptyset \\ N, & \text{otherwise.} \end{cases} \end{aligned}$$

Note that there is no contribution to $N_{S, \underline{\gamma_1^c \alpha_{min}}}$ or $N_{U, \underline{\alpha_{min}}}$ for $k \in \left[k_{end}^{(1)}, k_{start}^{(2)} \right)$. There is no contribution to $N_{S, \underline{\gamma_1^c \alpha_{min}}}$ because $k_{start}^{(2)}$ is the first iteration (if any) that would make this contribution. Moreover, since $\beta_{k_{end}^{(1)}} = \tau_\alpha$ by definition of $k_{end}^{(1)}$, the first iteration with $\beta_k < \tau_\alpha$ for $k \geq k_{end}^{(1)}$ must be preceded by a successful iteration with $\beta_k < \tau_\alpha + c$ (note that in particular, since $k_{start}^{(2)}$ is the first such iteration, we have $\beta_{k_{start}^{(2)}} \geq \tau_\alpha$). Therefore there is no contribution to $N_{U, \underline{\alpha_{min}}}$

either for $k \in [k_{end}^{(1)}, k_{start}^{(2)}]$. Hence if $k_{start}^{(2)} = N$, we have (4.3.16), (4.3.17) and (4.3.15) gives the desired result.

Otherwise similarly to (4.3.13)–(4.3.14), let

$$\begin{aligned} A^{(2)} &= \left\{ k \in \left(k_{start}^{(2)}, N-1 \right] \cap \mathbb{N} : \beta_k = \tau_\alpha \right\}. \\ k_{end}^{(2)} &= \begin{cases} \inf A^{(2)}, & \text{if } A^{(2)} \neq \emptyset \\ N, & \text{otherwise.} \end{cases} \\ M_1^{(2)} &= \left\{ k \in \left(k_{start}^{(2)}, k_{end}^{(2)} \right) : \text{iteration } k \text{ is unsuccessful with } \beta_k < \tau_\alpha \right\} \\ M_2^{(2)} &= \left\{ k \in \left(k_{start}^{(2)}, k_{end}^{(2)} \right) : \text{iteration } k \text{ is successful with } \beta_k < \tau_\alpha + c \right\}. \end{aligned}$$

And let $n_1^{(2)} = |M_1^{(2)}|$ and $n_2^{(2)} = |M_2^{(2)}|$. Note that for $k \in (k_{start}^{(2)}, k_{end}^{(2)})$, we have $\tau_\alpha - c \leq \beta_{k_{start}^{(2)}+1}$ and $\beta_{k_{end}^{(2)}} \leq \tau_\alpha$ (the former is true as $\beta_{k_{start}^{(2)}} \geq l$ and iteration $k_{start}^{(2)}$ is successful). Therefore we have

$$\tau_\alpha - c + n_1^{(2)} - cn_2^{(2)} \leq \beta_{k_{start}^{(2)}+1} + n_1^{(2)} - cn_2^{(2)} \leq \beta_{k_{end}^{(2)}} \leq \tau_\alpha.$$

Rearranging gives

$$n_1^{(2)} \leq cn_2^{(2)} + \tau_\alpha - [\tau_\alpha - c] = cn_2^{(2)} + c, \quad (4.3.18)$$

Let $\hat{n}_1^{(1)}$ be the total number of iterations contributing to $N_{U, \alpha_{min}}$ with $k \in [k_{end}^{(1)}, k_{start}^{(2)}]$; and $\hat{n}_2^{(1)}$ be the total number of iterations contributing to $N_{S, \gamma_1^c \alpha_{min}}$ with $k \in [k_{end}^{(1)}, k_{start}^{(2)}]$. Since there is no contribution to either for $k \in [k_{end}^{(1)}, k_{start}^{(2)}]$ as argued before, and iteration $k_{start}^{(2)}$ by definition contributes to $N_{S, \gamma_1^c \alpha_{min}}$ by one, we have

$$\hat{n}_1^{(1)} = 0, \quad (4.3.19)$$

$$\hat{n}_2^{(1)} = 1. \quad (4.3.20)$$

Using (4.3.15), (4.3.18), (4.3.19) and (4.3.20), we have

$$n_1^{(1)} + \hat{n}_1^{(1)} + n_1^{(2)} \leq c \left(n_2^{(1)} + \hat{n}_2^{(1)} + n_2^{(2)} \right) + \tau_\alpha. \quad (4.3.21)$$

If $k_{end}^{(2)} = N$ the desired result follows. Otherwise define $B^{(2)}$ in terms of $k_{end}^{(2)}$, and $k_{start}^{(3)}$ in terms of $B^{(2)}$ similarly as before. If $k_{start}^{(3)} = N$, then we have the desired result as before. Otherwise repeat what we have done (define $A^{(3)}$, $k_{end}^{(3)}$, $M_1^{(3)}$, $M_2^{(3)}$ etc). Note that we will reach either $k_{end}^{(i)} = N$ for some $i \in \mathbb{N}$ or $k_{start}^{(i)} = N$ for some $i \in \mathbb{N}$, because if $k_{end}^{(i)} < N$ and $k_{start}^{(i)} < N$ for all i , we have that $k_{start}^{(i)} < k_{end}^{(i)} \leq k_{start}^{(i+1)}$ by definitions. So $k_{start}^{(i)}$ is strictly increasing, contradicting $k_{start}^{(i)} < N$ for all i . In the case with $k_{end}^{(i)} = N$ or $k_{start}^{(i)} = N$, the desired result will follow using our previous argument. \square

An intermediate result bounding the total number of iterations With Lemma 4.3.1, Lemma 4.3.3, Lemma 4.3.4, we show a bound on the total number of iterations of Algorithm 2 in terms of the number of true and successful iterations with α_k above a certain constant.

Lemma 4.3.5. *Let Assumption 1 and Assumption 2 hold with $\delta_S \in (0, 1)$, $c, \tau_\alpha \in \mathbb{N}^+$. Let N be the total number of iterations. Then for any $\delta_1 \in (0, 1)$ such that $g(\delta_S, \delta_1) > 0$, we have that $\mathbb{P}\left(N < g(\delta_S, \delta_1) \left[N_{TS, \underline{\alpha_0 \gamma_1^{c+\tau_\alpha}}} + \frac{\tau_\alpha}{1+c}\right]\right) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N}$ where $g(\delta_S, \delta_1)$ is defined in (4.2.12).*

Proof. We decompose the number of true iterations as

$$N_T = N_{T, \overline{\alpha_{min}}} + N_{T, \underline{\alpha_{min}}} = N_{T, \overline{\alpha_{min}}} + N_{TS, \underline{\alpha_{min}}} + N_{TU, \underline{\alpha_{min}}} \leq N_{T, \overline{\alpha_{min}}} + N_{TS, \underline{\alpha_{min}}} + N_{U, \underline{\alpha_{min}}}, \quad (4.3.22)$$

where $N_T, N_{T, \overline{\alpha_{min}}}, N_{T, \underline{\alpha_{min}}}, N_{TS, \underline{\alpha_{min}}}, N_{TU, \underline{\alpha_{min}}}, N_{U, \underline{\alpha_{min}}}$ are defined in Table 4.1.

From Lemma 4.3.4, we have

$$\begin{aligned} N_{U, \underline{\alpha_{min}}} &\leq \tau_\alpha + cN_{S, \underline{\gamma_1^c \alpha_{min}}} \\ &= \tau_\alpha + cN_{TS, \underline{\gamma_1^c \alpha_{min}}} + cN_{FS, \underline{\gamma_1^c \alpha_{min}}} \\ &\leq \tau_\alpha + cN_{TS, \underline{\gamma_1^c \alpha_{min}}} + cN_F \\ &\leq \tau_\alpha + cN_{TS, \underline{\gamma_1^c \alpha_{min}}} + c(N - N_T), \end{aligned}$$

It then follows from (4.3.22) that

$$N_T \leq N_{T, \overline{\alpha_{min}}} + N_{TS, \underline{\alpha_{min}}} + \tau_\alpha + cN_{TS, \underline{\gamma_1^c \alpha_{min}}} + c(N - N_T).$$

Rearranging, we have

$$N_T \leq \frac{N_{T, \overline{\alpha_{min}}}}{1+c} + \frac{1}{1+c} \left[N_{TS, \underline{\alpha_{min}}} + cN_{TS, \underline{\gamma_1^c \alpha_{min}}} \right] + \frac{\tau_\alpha + cN}{1+c}.$$

Using Lemma 4.3.3 to bound $N_{T, \overline{\alpha_{min}}}$; $N_{TS, \underline{\alpha_{min}}} \leq N_{TS, \underline{\gamma_1^c \alpha_{min}}}$; and $\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}$ gives

$$N_T \leq \left[1 - \frac{c}{(c+1)^2} \right] N + N_{TS, \underline{\alpha_0 \gamma_1^{c+\tau_\alpha}}} + \frac{\tau_\alpha}{1+c}. \quad (4.3.23)$$

Combining with Lemma 4.3.1; and rearranging gives the result. \square

The bound on true and successful iterations The next lemma bounds the total number of true and successful iterations with $\alpha_k > \alpha_0 \gamma_1^{c+\tau_\alpha}$.

Lemma 4.3.6. *Let Assumption 3 and Assumption 4 hold. Let $\epsilon > 0$ and $N \in \mathbb{N}$ be defined in Table 4.1. Suppose $N \leq N_\epsilon$. Then*

$$N_{TS, \underline{\alpha_0 \gamma_1^{c+\tau_\alpha}}} \leq \frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})}$$

where f^* is defined in (4.2.1), and x_0 is chosen in the initialization of Algorithm 2.

Proof. We have, using Assumption 4 and Assumption 3 respectively for the two inequalities

$$\begin{aligned}
f(x_0) - f(x_N) &= \sum_{k=0}^{N-1} f(x_k) - f(x_{k+1}) \\
&\geq \sum_{\substack{\text{Iteration } k \text{ is true and successful} \\ \text{with } \alpha_k \geq \alpha_0 \gamma_1^{c+\tau_\alpha}}} f(x_k) - f(x_{k+1}) \\
&\geq \sum_{\substack{\text{Iteration } k \text{ is true and successful} \\ \text{with } \alpha_k \geq \alpha_0 \gamma_1^{c+\tau_\alpha}}} h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) \\
&= N_{TS, \alpha_0 \gamma_1^{c+\tau_\alpha}} h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}). \tag{4.3.24}
\end{aligned}$$

Noting $f(x_N) \geq f^*$ and $h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) > 0$ by Assumption 3, rearranging (4.3.24) gives the result. \square

The final proof We are ready to prove Theorem 4.2.1 using Lemma 4.3.5 and Lemma 4.3.6.

Proof of Theorem 4.2.1. We have

$$N_\epsilon \geq N \implies \frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} \geq N_{TS, \alpha_0 \gamma_1^{c+\tau_\alpha}} \quad \text{by Lemma 4.3.6} \tag{4.3.25}$$

$$\stackrel{(4.2.13)}{\implies} N \geq g(\delta_S, \delta_1) \left[N_{TS, \alpha_0 \gamma_1^{c+\tau_\alpha}} + \frac{\tau_\alpha}{1+c} \right]. \tag{4.3.26}$$

Therefore by Lemma 4.3.5, we have $\mathbb{P}(N_\epsilon \geq N) \leq \mathbb{P}\left(N \geq g(\delta_S, \delta_1) \left[N_{TS, \alpha_0 \gamma_1^{c+\tau_\alpha}} + \frac{\tau_\alpha}{1+c} \right]\right) \leq e^{-\frac{\delta_1^2}{2}(1-\delta_S)N}$. \square

4.4 An algorithmic framework based on sketching

4.4.1 A generic random subspace method based on sketching

Algorithm 3 particularises Algorithm 2 by specifying the local reduced model as one generated by sketching using a random matrix; the step transformation function; and the criterion for sufficient decrease. We leave specification of the computation of the step parameter to the next section.

Algorithm 3 A generic random subspace method based on sketching

Initialization

Choose a matrix distribution \mathcal{S} of matrices $S \in \mathbb{R}^{l \times d}$. Let $\gamma_1, \gamma_2, \theta, \alpha_{max}, x_0, \alpha_0$ be defined in Algorithm 2 with $\hat{m}_k(\hat{s})$ and w_k specified below in (4.4.1) and (4.4.2).

1. Compute a reduced model and a step

In Step 1 of Algorithm 2, draw a random matrix $S_k \in \mathbb{R}^{l \times d}$ from \mathcal{S} , and let

$$\hat{m}_k(\hat{s}) = f(x_k) + \langle S_k \nabla f(x_k), \hat{s} \rangle + \frac{1}{2} \langle \hat{s}, S_k B_k S_k^T \hat{s} \rangle; \quad (4.4.1)$$

$$w_k(\hat{s}_k) = S_k^T \hat{s}_k, \quad (4.4.2)$$

where $B_k \in \mathbb{R}^{d \times d}$ is a user provided matrix.

Compute \hat{s}_k by approximately minimising $\hat{m}_k(\hat{s})$ such that at least $\hat{m}_k(\hat{s}_k) \leq \hat{m}_k(0)^3$ where α_k appears as a parameter, and set $s_k = w_k(\hat{s}_k)$ as in Algorithm 2.

2. Check sufficient decrease

In Step 2 of Algorithm 2, let sufficient decrease be defined by the condition

$$f(x_k) - f(x_k + s_k) \geq \theta [\hat{m}_k(0) - \hat{m}_k(\hat{s}_k(\alpha_k))]. \quad (4.4.3)$$

3. Update the parameter α_k and possibly take the potential step s_k

Follow Step 3 of Algorithm 2.

With the concrete criterion for sufficient decrease, we have that Assumption 4 is satisfied by Algorithm 3.

Lemma 4.4.1. *Algorithm 3 satisfies Assumption 4.*

Proof. If iteration k is successful, (4.4.3) with $\theta \geq 0$ and $\hat{m}_k(\hat{s}_k) \leq \hat{m}_k(0)$ (specified in Algorithm 3) give $f(x_k) - f(x_k + s_k) \geq 0$. If iteration k is unsuccessful, we have $s_k = 0$ and therefore $f(x_k) - f(x_k + s_k) = 0$. \square

Next, we define what a true iteration is for Algorithm 3 and show Assumption 1 is satisfied with \mathcal{S} being a variety of random ensembles.

Definition 4.4.1. *Iteration k is a true iteration if*

$$\|S_k \nabla f(x_k)\|_2^2 \geq (1 - \epsilon_S) \|\nabla f(x_k)\|_2^2, \quad (4.4.4)$$

$$\|S_k\|_2 \leq S_{max}, \quad (4.4.5)$$

where $S_k \in \mathbb{R}^{l \times d}$ is the random matrix drawn in Step 1 of Algorithm 3, and $\epsilon_S \in (0, 1)$, $S_{max} > 0$ are iteration-independent constants.

Remark 5. *In [17], true iterations are required to satisfy*

$$\|\nabla m_k(0) - \nabla f(x_k)\|_2 \leq \kappa \alpha_k \|\nabla m_k(0)\|_2,$$

where $\kappa > 0$ is a constant and α_k in their algorithm is bounded by α_{max} . The above equation implies

$$\|\nabla m_k(0)\|_2 \geq \frac{\|\nabla f(x_k)\|_2}{1 + \kappa\alpha_{max}},$$

which implies (4.4.4) with $1 - \epsilon_S = \frac{1}{1 + \kappa\alpha_{max}}$ and $\delta_S^{(1)} = p$. Since Assumption 6 is easily satisfied for a variety of random matrix distributions \mathcal{S} we see that their requirement is stronger than our (main) requirement for true iterations.

We first show that with this definition of the true iterations, Assumption 1 holds if the following two conditions on the random matrix distribution \mathcal{S} are met.

Assumption 5. *There exists $\epsilon_S, \delta_S^{(1)} \in (0, 1)$ such that for a(ny) fixed $y \in \{\nabla f(x) : x \in \mathbb{R}^d\}$, S_k drawn from \mathcal{S} satisfies*

$$\mathbb{P}\left(\|S_k y\|_2^2 \geq (1 - \epsilon_S) \|y\|_2^2\right) \geq 1 - \delta_S^{(1)}. \quad (4.4.6)$$

Assumption 6. *There exists $\delta_S^{(2)} \in [0, 1), S_{max} > 0$ such that for S_k randomly drawn from \mathcal{S} , we have*

$$\mathbb{P}(\|S_k\|_2 \leq S_{max}) \geq 1 - \delta_S^{(2)}.$$

Lemma 4.4.2. *Let Assumption 5 and Assumption 6 hold with $\epsilon_S, \delta_S^{(2)} \in (0, 1), \delta_S^{(1)} \in [0, 1), S_{max} > 0$. Suppose that $\delta_S^{(1)} + \delta_S^{(2)} < 1$. Let true iterations be defined in Definition 4.4.1. Then Algorithm 3 satisfies Assumption 1 with $\delta_S = \delta_S^{(1)} + \delta_S^{(2)}$.*

The proof of Lemma 4.4.2 makes use of the following elementary result in probability theory, whose proof is included for completeness.

Lemma 4.4.3. *Let $n \in \mathbb{N}^+$ and A_1, A_2, \dots, A_n be events. Then we have*

$$\mathbb{P}(A_1 \cap A_2 \dots \cap A_n) = 1 - \mathbb{P}(A_1^c) - \mathbb{P}(A_2^c) - \dots - \mathbb{P}(A_n^c).$$

Proof. We have

$$\begin{aligned} \mathbb{P}(A_1 \cap A_2 \dots \cap A_n) &= 1 - \mathbb{P}((A_1 \cap \dots \cap A_n)^c) \\ &= 1 - \mathbb{P}(A_1^c \cup \dots \cup A_n^c) \\ &\geq 1 - \sum_{k=1}^n \mathbb{P}(A_k^c). \end{aligned}$$

□

Proof of Lemma 4.4.2. Let $\bar{x}_k \in \mathbb{R}^d$ be given. Note that this determines $\nabla f(\bar{x}_k) \in \mathbb{R}^d$. Let $A_k^{(1)}$ be the event that (4.4.4) hold and $A_k^{(2)}$ be the event that (4.4.5) hold. Thus $T_k = A_k^{(1)} \cap A_k^{(2)}$. Note that given $x_k = \bar{x}_k$, T_k only depends on S_k , which is independent of all previous iterations. Hence T_k is conditionally independent of T_0, T_1, \dots, T_{k-1} given $x_k = \bar{x}_k$.

Next, we have for $k \geq 1$,

$$\mathbb{P}\left(A_k^{(1)} \cap A_k^{(2)} | x_k = \bar{x}_k\right) \geq 1 - \mathbb{P}\left(\left(A_k^{(1)}\right)^c | x_k = \bar{x}_k\right) - \mathbb{P}\left(\left(A_k^{(2)}\right)^c | x_k = \bar{x}_k\right), \quad (4.4.7)$$

by Lemma 4.4.3.

Note that

$$\begin{aligned} \mathbb{P}\left(A_k^{(1)} | x_k = \bar{x}_k\right) &= \mathbb{P}\left[A_k^{(1)} | x_k = \bar{x}_k, \nabla f(x_k) = \nabla f(\bar{x}_k)\right] \\ &= \mathbb{P}\left[A_k^{(1)} | \nabla f(x_k) = \nabla f(\bar{x}_k)\right] \\ &\geq 1 - \delta_S^{(1)}, \end{aligned} \quad (4.4.8)$$

where the first equality follows from the fact that $x_k = \bar{x}_k$ implies $\nabla f(x_k) = \nabla f(\bar{x}_k)$; the second equality follows from the fact that given $\nabla f(x_k) = \nabla f(\bar{x}_k)$, $A_k^{(1)}$ is independent of x_k ; and the inequality follows from applying Assumption 5 with $y = \nabla f(\bar{x}_k)$.

On the other hand, because $A_k^{(2)}$ is independent of x_k , we have that

$$\mathbb{P}\left(A_k^{(2)} | x_k = \bar{x}_k\right) = \mathbb{P}\left(A_k^{(2)}\right) \geq 1 - \delta_S^{(2)}, \quad (4.4.9)$$

where the inequality follows from Assumption 6. It follows from (4.4.7) using (4.4.8) and (4.4.9) that for $k \geq 1$,

$$\mathbb{P}\left(A_k^{(1)} \cap A_k^{(2)} | x_k = \bar{x}_k\right) \geq 1 - \delta_S^{(1)} - \delta_S^{(2)} = 1 - \delta_S.$$

For $k = 0$, we have $\mathbb{P}\left(A_0^{(1)}\right) \geq 1 - \delta_S^{(1)}$ by Assumption 5 with $y = \nabla f(x_0)$ and $\mathbb{P}\left(A_0^{(2)}\right) \geq 1 - \delta_S^{(2)}$ by Assumption 6. So $\mathbb{P}\left(A_0^{(1)} \cap A_0^{(2)}\right) \geq 1 - \delta_S$ by Lemma 4.4.3. □

Next, we give four distributions \mathcal{S} that satisfy Assumption 5 and Assumption 6, thus satisfying Assumption 1 and can be used in Algorithm 3. Other random ensembles are possible, for example, Subsampled Randomised Hadamard Transform (Definition 1.2.3), Hashed Randomised Hadamard Transform (Definition 2.4.2), and many more (see discussion of random ensembles in Chapter 2).

4.4.2 The random matrix distribution \mathcal{S} in Algorithm 3

4.4.2.1 Gaussian sketching matrices

(Scaled) Gaussian matrices have independent and identically distributed normal entries (see Definition 1.2.2). The next result, which is a consequence of the scaled Gaussian matrices being an oblivious JL embedding (Definition 2.2.4), shows that using scaled Gaussian matrices with Algorithm 3 satisfies Assumption 5. The proof is included for completeness but can also be found in [25].

Lemma 4.4.4. *Let $S \in \mathbb{R}^{l \times d}$ be a scaled Gaussian matrix so that each entry is $N(0, l^{-1})$. Then S satisfies Assumption 5 with any $\epsilon_S \in (0, 1)$ and $\delta_S^{(1)} = e^{-\epsilon_S^2 l/4}$.*

Proof. Since (4.4.6) is invariant to the scaling of y and is trivial for $y = 0$, we may assume without loss of generality that $\|y\|_2 = 1$.

Let $R = \sqrt{l}S$, so that each entry of R is distributed independently as $N(0, 1)$. Then because the sum of independent Gaussian random variables is distributed as a Gaussian random variable; $\|y\|_2 = 1$; and the fact that rows of S are independent; we have that the entries of Ry , denoted by z_i for $i \in [l]$, are independent $N(0, 1)$ random variables. Therefore, for any $-\infty < q < \frac{1}{2}$, we have that

$$\mathbb{E} \left[e^{q\|Ry\|_2^2} \right] = \mathbb{E} \left[e^{q\sum_{i=1}^l z_i^2} \right] = \prod_{i=1}^l \mathbb{E} \left[e^{qz_i^2} \right] = (1 - 2q)^{-l/2}, \quad (4.4.10)$$

where we used $\mathbb{E} \left[e^{qz_i^2} \right] = \frac{1}{1-2q}$ for $z_i \in N(0, 1)$ and $-\infty < q < \frac{1}{2}$.

Hence, by Markov inequality, we have that, for $q < 0$,

$$\mathbb{P} \left(\|Ry\|_2^2 \leq l(1 - \epsilon_S) \right) = \mathbb{P} \left(e^{q\|Ry\|_2^2} \geq e^{ql(1 - \epsilon_S)} \right) \leq \frac{\mathbb{E} \left[e^{q\|Ry\|_2^2} \right]}{e^{ql(1 - \epsilon_S)}} = (1 - 2q)^{-l/2} e^{-ql(1 - \epsilon_S)}, \quad (4.4.11)$$

where the last inequality comes from (4.4.10).

Noting that

$$(1 - 2q)^{-l/2} e^{-ql(1 - \epsilon_S)} = \exp \left[-l \left(\frac{1}{2} \log(1 - 2q) + q(1 - \epsilon_S) \right) \right], \quad (4.4.12)$$

which is minimised at $q_0 = -\frac{\epsilon_S}{2(1 - \epsilon_S)} < 0$, we choose $q = q_0$ and the right hand side of (4.4.11) becomes

$$e^{\frac{1}{2}l[\epsilon_S + \log(1 - \epsilon_S)]} \leq e^{-\frac{1}{4}l\epsilon_S^2}, \quad (4.4.13)$$

where we used $\log(1 - x) \leq -x - x^2/2$, valid for all $x \in [0, 1]$.

Hence we have

$$\begin{aligned} & \mathbb{P} \left(\|Sy\|_2^2 \leq (1 - \epsilon_S) \|y\|_2^2 \right) \\ &= \mathbb{P} \left(\|Sy\|_2^2 \leq (1 - \epsilon_S) \right) \quad \text{by } \|y\|_2 = 1 \\ &= \mathbb{P} \left(\|Ry\|_2^2 \leq l(1 - \epsilon_S) \right) \quad \text{by } S = \frac{1}{\sqrt{l}}R \\ &\leq e^{-\frac{l\epsilon_S^2}{4}} \quad \text{by (4.4.13) and (4.4.11).} \end{aligned}$$

□

In order to show using scaled Gaussian matrices satisfies Assumption 6, we make use of the following bound on the maximal singular value of scaled Gaussian matrices.

Lemma 4.4.5 ([26] Theorem 2.13). *Given $l, d \in \mathbb{N}$ with $l \leq d$, consider the $d \times l$ matrix Γ whose entries are independent $N(0, d^{-1})$. Then for any $\delta > 0$,⁴*

$$\mathbb{P} \left(\sigma_{\max}(\Gamma) \geq 1 + \sqrt{\frac{l}{d}} + \sqrt{\frac{2 \log(\frac{1}{\delta})}{d}} \right) < \delta, \quad (4.4.14)$$

⁴We set $t = \sqrt{\frac{2 \log(\frac{1}{\delta})}{d}}$ in the original theorem statement.

where $\sigma_{\max}(\cdot)$ denotes the largest singular value of its matrix argument.

The next lemma shows that using scaled Gaussian matrices satisfies Assumption 6.

Lemma 4.4.6. *Let $S \in \mathbb{R}^{l \times d}$ be a scaled Gaussian matrix. Then S satisfies Assumption 6 with any $\delta_S^{(2)} \in (0, 1)$ and*

$$S_{\max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}}.$$

Proof. We have $\|S\|_2 = \|S^T\|_2 = \sqrt{\frac{d}{l}} \left\| \sqrt{\frac{l}{d}} S^T \right\|_2$. Applying Lemma 4.4.5 with $\Gamma = \sqrt{\frac{l}{d}} S^T$, we have that

$$\mathbb{P} \left(\sigma_{\max} \left(\sqrt{\frac{l}{d}} S^T \right) \geq 1 + \sqrt{\frac{l}{d}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{d}} \right) < \delta_S^{(2)}.$$

Noting that $\|S\|_2 = \sqrt{\frac{d}{l}} \sigma_{\max}(\Gamma)$, and taking the event complement gives the result. \square

4.4.2.2 s -hashing matrices

Comparing to Gaussian matrices, s -hashing matrices, including the $s = 1$ case, (defined in Definition 1.2.5) are sparse so that it preserves the sparsity (if any) of the vector/matrix it acts on; and the corresponding linear algebra computation is faster. The next two lemmas show that using s -hashing matrices satisfies Assumption 5 and Assumption 6.

Lemma 4.4.7 ([58] Theorem 13, also see [23] Theorem 5 for a simpler proof). *Let $S \in \mathbb{R}^{l \times d}$ be an s -hashing matrix. Then S satisfies Assumption 5 for any $\epsilon_S \in (0, 1)$ and $\delta_S^{(1)} = e^{-\frac{\epsilon_S^2}{C_1^2}}$ given that $s = C_2 \epsilon_S l$, where C_1, C_2 are problem-independent constants.*

Lemma 4.4.8. *Let $S \in \mathbb{R}^{l \times d}$ be an s -hashing matrix. Then S satisfies Assumption 6 with $\delta_S^{(2)} = 0$ and $S_{\max} = \sqrt{\frac{d}{s}}$.*

Proof. Note that for any matrix $A \in \mathbb{R}^{l \times d}$, $\|A\|_2 \leq \sqrt{d} \|A\|_\infty$; and $\|S\|_\infty = \frac{1}{\sqrt{s}}$. The result follows from combining these two facts. \square

4.4.2.3 (Stable) 1-hashing matrices

In [19], a variant of 1-hashing matrix is proposed that satisfies Assumption 5 but with better S_{\max} bound. The construction is given as follows.

Definition 4.4.2. *Let $l < d \in \mathbb{N}^+$. A stable 1-hashing matrix $S \in \mathbb{R}^{l \times d}$ has one non-zero per column, whose value is ± 1 with equal probability, with the row indices of the non-zeros given by the sequence I constructed as the following. Repeat $[l]$ (that is, the set $\{1, 2, \dots, l\}$) for $\lceil d/l \rceil$ times to*

obtain a set D . Then randomly sample d elements from D without replacement to construct sequence I .⁵

Remark 6. Comparing to a 1-hashing matrix, a stable 1-hashing matrix still has 1 non-zero per column. However its construction guarantees that each row has at most $\lceil d/l \rceil$ non-zeros because the set D has at most $\lceil d/l \rceil$ repeated row indices and the sampling is done without replacement.

In order to show using stable 1-hashing matrices satisfies Assumption 5, we need to following result from [19].

Lemma 4.4.9 (Theorem 5.3 in [19]). *The matrix $S \in \mathbb{R}^{l \times d}$ defined in Definition 4.4.2 satisfies the following: given $0 < \epsilon, \delta < 1/2$, there exists $l = \mathcal{O}\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ such that for any $x \in \mathbb{R}^d$, we have that $\mathbb{P}(\|Sx\|_2 \geq (1 - \epsilon)\|x\|_2) > 1 - \delta$.*

Lemma 4.4.10. *Let $S \in \mathbb{R}^{l \times d}$ be a stable 1-hashing matrix. Let $\epsilon_S \in (0, 3/4)$ and suppose that $e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}} \in (0, 1/2)$, where C_3 is a problem-independent constant. Then S satisfies Assumption 5 with $\delta_S^{(1)} = e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}}$.*

Proof. Let $\bar{\epsilon} = \epsilon_S - 1/4 \in (0, 1/2)$. From Lemma 4.4.9, we have that there exists $C_3 > 0$ such that with $\delta_S^{(1)} = e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}}$, S satisfies $\mathbb{P}(\|Sx\|_2 \geq (1 - \bar{\epsilon})\|x\|_2) > 1 - \delta_S^{(1)}$. Note that $\|Sx\|_2 \geq (1 - \bar{\epsilon})\|x\|_2$ implies $\|Sx\|_2^2 \geq (1 - 2\bar{\epsilon} + \bar{\epsilon}^2)\|x\|_2^2$, which implies $\|Sx\|_2^2 \geq (1 - \bar{\epsilon} - 1/4)\|x\|_2^2$ because $\bar{\epsilon}^2 - \bar{\epsilon} \geq -1/4$ for $\bar{\epsilon} \in (0, 1/2)$. The desired result follows. \square

The next lemma shows that using stable 1-hashing matrices satisfies Assumption 6. Note that the bound S_{max} is smaller than that for 1-hashing matrices; and, assuming $l > s$, smaller than that for s -hashing matrices as well.

Lemma 4.4.11. *Let $S \in \mathbb{R}^{l \times d}$ be a stable 1-hashing matrix. Then S satisfies Assumption 6 with $\delta_S^{(2)} = 0$ and $S_{max} = \sqrt{\lceil d/l \rceil}$.*

Proof. Let D be defined in Definition 4.4.2. we have that

$$\|Sx\|_2^2 = \left(\sum_{1 \leq j \leq d, I(j)=1} \pm x_j\right)^2 + \left(\sum_{1 \leq j \leq d, I(j)=2} \pm x_j\right)^2 + \cdots + \left(\sum_{1 \leq j \leq d, I(j)=l} \pm x_j\right)^2 \quad (4.4.15)$$

$$\leq \left(\sum_{1 \leq j \leq d, I(j)=1} |x_j|\right)^2 + \left(\sum_{1 \leq j \leq d, I(j)=2} |x_j|\right)^2 + \cdots + \left(\sum_{1 \leq j \leq d, I(j)=l} |x_j|\right)^2 \quad (4.4.16)$$

$$\leq \lceil d/l \rceil \left(\sum_{1 \leq j \leq d, I(j)=1} x_j^2 + \sum_{1 \leq j \leq d, I(j)=2} x_j^2 + \cdots + \sum_{1 \leq j \leq d, I(j)=l} x_j^2 \right) \quad (4.4.17)$$

$$= \lceil d/l \rceil \|x\|_2^2, \quad (4.4.18)$$

⁵One may also conceptually think S as being constructed from taking the first d columns of a random column permutation of the matrix $T = [I_{l \times l}, I_{l \times l}, \dots, I_{l \times l}]$ where the identity matrix $I_{l \times l}$ is concatenated by columns $\lceil d/l \rceil$ times.

where the \pm on the first line results from the non-zero entries of S having random signs, and the last inequality is because for any vector $v \in \mathbb{R}^n$, $\|v\|_1^2 \leq n\|v\|_2^2$; and $I(j) = k$ is true for at most $\lceil d/l \rceil$ indices j . \square

4.4.2.4 Sampling matrices

(Scaled) Sampling matrices $S \in \mathbb{R}^{l \times d}$ (defined in Definition 1.2.1) randomly select rows of vector/matrix it acts on (and scale it). Next we show that sampling matrices satisfy Assumption 5. The following expression that represents the maximum non-uniformity (see Definition 2.2.6) of the objective gradient will be used

$$\nu = \max \left\{ \frac{\|y\|_\infty}{\|y\|_2}, y = \nabla f(x) \text{ for some } x \in \mathbb{R}^d \right\}. \quad (4.4.19)$$

The following concentration result will be useful.

Lemma 4.4.12 ([98]). *Consider a finite sequence of independent random numbers $\{X_k\}$ that satisfies $X_k \geq 0$ and $|X_k| \leq P$ almost surely. Let $\mu = \sum_k \mathbb{E}[X_k]$, then $\mathbb{P}(\sum_k X_k \leq (1 - \epsilon)\mu) \leq e^{-\frac{\epsilon^2 \mu}{2P}}$.*

Lemma 4.4.13. *Let $S \in \mathbb{R}^{l \times d}$ be a scaled sampling matrix. Let ν be defined in (4.4.19). Then S satisfies Assumption 5 for any $\epsilon_S \in (0, 1)$ with $\delta_S^{(1)} = e^{-\frac{\epsilon_S^2 l}{2d\nu^2}}$.*

Proof. Note that (4.4.6) is invariant to scaling of y and trivial for $y = 0$. Therefore we may assume $\|y\|_2 = 1$ without loss of generality.

We have $\|Sy\|_2^2 = \frac{l}{d} \sum_{k=1}^l [(Ry)_k]^2$ where $R \in \mathbb{R}^{l \times d}$ is an (un-scaled) sampling matrix⁶ and $(Ry)_k$ denotes the k^{th} entry of Ry . Let $X_k = [(Ry)_k]^2$. Note that because the rows of R are independent, X_k are independent. Moreover, because $(Ry)_k$ equals to some entry of y , and $\|y\|_\infty \leq \nu$ by definition of ν and $\|y\|_2 = 1$; we have $[(Ry)_k]^2 \leq \nu^2$. Finally, note that $\mathbb{E}[X_k] = \frac{1}{d} \|y\|_2^2 = \frac{1}{d}$; so that $\sum_k \mathbb{E}[X_k] = \frac{l}{d}$.

Therefore applying Lemma 4.4.12 with $\epsilon = \epsilon_S$ we have

$$\mathbb{P} \left(\sum_{k=1}^l [(Ry)_k]^2 \leq (1 - \epsilon_S) \frac{l}{d} \right) \geq e^{-\frac{\epsilon_S^2 l}{2d\nu^2}}.$$

Using $\|Sy\|_2^2 = \frac{l}{d} \sum_{k=1}^l [(Ry)_k]^2$ gives the result. \square

We note that the theoretical property for scaled sampling matrices is different to Gaussian/ s -hashing matrices in the sense that the required value of l depends on ν . Note that $\frac{1}{d} \leq \nu^2 \leq 1$ with both bounds attainable. Therefore in the worst case, for fixed value of $\epsilon_S, \delta_S^{(1)}$, l is required to be $\mathcal{O}(d)$ and no dimensionality reduction is achieved by sketching. This is not surprising given that sampling based random methods often require adaptively increasing the sampling size for convergence (reference). However note that for ‘nice’ objective functions such that $\nu^2 = \mathcal{O}(\frac{1}{d})$, sampling

⁶I.e. each row of R has a one at a random column.

matrices have the same theoretical property as Gaussian/ s -hashing matrices. The attractiveness of sampling lies in the fact that only a subset of entries of the gradient need to be evaluated.

Sampling matrices also have bounded 2-norms, thus Assumption 6 is satisfied.

Lemma 4.4.14. *Let $S \in \mathbb{R}^{l \times d}$ be a scaled sampling matrix. Then Assumption 6 is satisfied with $\delta_S^{(2)} = 0$ and $S_{max} = \sqrt{\frac{d}{l}}$.*

Proof. We have that $\|Sx\|_2^2 \leq \frac{d}{l} \|x\|_2^2$ for any $x \in \mathbb{R}^d$. \square

We summarises this section in Table 4.2, where we also give l in terms of ϵ_S and $\delta_S^{(1)}$ by rearranging the expressions for $\delta_S^{(1)}$. Note that for s -hashing matrices, s is required to be $C_2 \epsilon_S l$ (see Lemma 4.4.7), while for scaled sampling matrices, ν is defined in (4.4.19). One may be concerned about the exponential increase of the embedding dimension l as ϵ_S goes to zero. However, ϵ_S may in fact be taken as some $\mathcal{O}(1)$ constant that is smaller than 1 (or $3/4$ in the case of stable 1-hashing). The reason being that the iterative nature of Algorithm 3 mitigates the inaccuracies of the embedding. See, e.g., the complexity bound in Theorem 4.5.1.

	ϵ_S	$\delta_S^{(1)}$	l	$\delta_S^{(2)}$	S_{max}
Scaled Gaussian	$(0, 1)$	$e^{-\frac{\epsilon_S^2 l}{4}}$	$4\epsilon_S^{-2} \log\left(\frac{1}{\delta_S^{(1)}}\right)$	$(0, 1)$	$1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}}$
s -hashing	$(0, 1)$	$e^{-\frac{\epsilon_S^2 l}{C_1}}$	$C_1 \epsilon_S^{-2} \log\left(\frac{1}{\delta_S^{(1)}}\right)$	0	$\sqrt{\frac{d}{s}}$
Stable 1-hashing	$(0, \frac{3}{4})$	$e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}}$	$C_3(\epsilon_S - 1/4)^{-2} \log\left(\frac{1}{\delta_S^{(1)}}\right)$	0	$\sqrt{\lceil \frac{d}{l} \rceil}$
Scaled sampling	$(0, 1)$	$e^{-\frac{\epsilon_S^2 l}{2d\nu^2}}$	$2d\nu^2 \epsilon_S^{-2} \log\left(\frac{1}{\delta_S^{(1)}}\right)$	0	$\sqrt{\frac{d}{l}}$

Table 4.2: Summary of theoretical properties of using different random ensembles with Algorithm 3.

4.5 Random subspace quadratic regularisation and subspace trust region methods

In this section, we analyse two methods for computing the trial step \hat{s}_k given the sketching based model in Algorithm 3. We show that using both methods: quadratic regularisation and trust-region, satisfy Assumption 2 and Assumption 3. Using Theorem 4.2.1, we show that the iteration complexity for both methods is $\mathcal{O}(\epsilon^{-2})$ to bring the objective's gradient below ϵ .

First we show that Assumption 3 holds for Algorithm 3 if the following model reduction condition is met.

Assumption 7. *There exists a non-negative, non-decreasing function $\bar{h} : \mathbb{R}^2 \rightarrow \mathbb{R}$ such that on each true iteration k of Algorithm 3 we have*

$$\hat{m}_k(0) - \hat{m}_k(\hat{s}_k(\alpha_k)) \geq \bar{h}(\|S_k \nabla f(x_k)\|_2, \alpha_k),$$

where $S_k, \hat{m}_k, \alpha_k, \hat{s}_k$ are defined in Algorithm 3.

Lemma 4.5.1. *Let Assumption 7 hold with \bar{h} and true iterations defined in Definition 4.4.1. Then Algorithm 3 satisfies Assumption 3 with $h(\epsilon, \alpha_k) = \theta \bar{h}((1 - \epsilon_S)^{1/2} \epsilon, \alpha_k)$, where ϵ_S is defined in (4.4.4).*

Proof. Let k be a true and successful iteration with $k < N_\epsilon$ for some $\epsilon > 0$ where N_ϵ is defined in (4.2.2). Then, using the fact that the iteration is true, successful, Assumption 7 and $k < N_\epsilon$, we have

$$\begin{aligned} f(x_k) - f(x_k + s_k) &\geq \theta [\hat{m}_k(0) - \hat{m}_k(\hat{s}_k(\alpha_k))] \\ &\geq \theta \bar{h}(\|S_k \nabla f(x_k)\|_2, \alpha_k) \\ &\geq \theta \bar{h}((1 - \epsilon_S)^{1/2} \|\nabla f(x_k)\|_2, \alpha_k) \\ &\geq \theta \bar{h}((1 - \epsilon_S)^{1/2} \epsilon, \alpha_k). \end{aligned}$$

□

The next Lemma is a standard result and we include its proof for completeness. It is needed to show random subspace quadratic regularisation and trust region methods satisfy Assumption 2.

Lemma 4.5.2. *In Algorithm 3, suppose that $\|B_k\|_2 \leq B_{max}$ for all k where B_{max} is independent of k , and f is continuously differentiable with L -Lipschitz continuous gradient. Then for any $\hat{s}_k \in \mathbb{R}^l$ and $S_k \in \mathbb{R}^{l \times d}$, let $s_k = S_k^T \hat{s}_k \in \mathbb{R}^d$. We have that*

$$|f(x_k + s_k) - \hat{m}_k(\hat{s}_k)| \leq \left(\frac{L + B_{max}}{2} \right) \|S_k^T \hat{s}_k\|_2^2. \quad (4.5.1)$$

Proof. As f is continuously differentiable with L -Lipschitz gradient, we have from Corollary 8.4 in [16] that

$$|f(x_k + S_k \hat{s}_k) - \langle S_k \nabla f(x_k), \hat{s}_k \rangle| \leq \frac{L}{2} \|S_k^T \hat{s}_k\|_2^2. \quad (4.5.2)$$

The above equation and triangle inequality implies

$$\begin{aligned} |f(x_k + s_k) - \hat{m}_k(\hat{s}_k)| &= |f(x_k + s_k) - f(x_k) - \langle S_k \nabla f(x_k), \hat{s}_k \rangle - \frac{1}{2} \langle S_k^T \hat{s}_k, B_k S_k^T \hat{s}_k \rangle| \\ &\leq \left(\frac{L}{2} + \frac{1}{2} \|B_k\|_2 \right) \|S_k^T \hat{s}_k\|_2^2 \\ &\leq \frac{L + B_{max}}{2} \|S_k^T \hat{s}_k\|_2^2, \end{aligned} \quad (4.5.3)$$

where we used $\|B_k\|_2 \leq B_{max}$ to derive the last inequality.

□

4.5.1 Random subspace quadratic regularisation with sketching

Here we present Algorithm 4, a generic random subspace quadratic regularisation method with sketching, which is a particular form of Algorithm 3 where the step is computed using a quadratic regularisation framework (see Page 13 in Chapter 1). We show that in addition to Assumption 4 which is satisfied by Algorithm 3, Algorithm 4 satisfies Assumption 2 and Assumption 3.

Algorithm 4 A generic random subspace quadratic regularisation method with sketching
Initialization

Choose a matrix distribution \mathcal{S} of matrices $S \in \mathbb{R}^{l \times d}$. Choose constants $\gamma_1 \in (0, 1)$, $\gamma_2 = \gamma_1^{-c}$, for some $c \in \mathbb{N}^+$, $l \in \mathbb{N}^+$, $\theta \in (0, 1)$ and $\alpha_{\max}, B_{\max} > 0$. Initialize the algorithm by setting $x_0 \in \mathbb{R}^d$, $\alpha_0 = \alpha_{\max} \gamma_1^p$ for some $p \in \mathbb{N}^+$ and $k = 0$.

1. Compute a reduced model and a step

Draw a random matrix $S_k \in \mathbb{R}^{l \times d}$ from \mathcal{S} , and let

$$\hat{m}_k(\hat{s}) = f(x_k) + \langle S_k \nabla f(x_k), \hat{s} \rangle + \frac{1}{2} \langle \hat{s}, S_k B_k S_k^T \hat{s} \rangle \quad (4.5.4)$$

where $B_k \in \mathbb{R}^{d \times d}$ is a positive-semi-definite user provided matrix with $\|B_k\|_2 \leq B_{\max}$.

Compute \hat{s}_k by approximately minimising $\hat{l}_k(\hat{s}) = \hat{m}_k(\hat{s}) + \frac{1}{2\alpha_k} \|S_k^T \hat{s}_k\|_2^2$ such that the following two conditions hold

$$\|\nabla \hat{l}_k(\hat{s}_k)\|_2 \leq \kappa_T \|S_k^T \hat{s}_k\|_2, \quad (4.5.5)$$

$$\hat{l}_k(\hat{s}_k) \leq \hat{l}_k(0), \quad (4.5.6)$$

where $\kappa_T \geq 0$ is a user chosen constant. And set $s_k = S_k^T \hat{s}_k$.

2. Check sufficient decrease

Let sufficient decrease be defined by the condition

$$f(x_k) - f(x_k + s_k) \geq \theta [\hat{m}_k(0) - \hat{m}_k(\hat{s}_k(\alpha_k))].$$

3, Update the parameter α_k and possibly take the potential step s_k

If sufficient decrease is achieved, set $x_{k+1} = x_k + s_k$ and $\alpha_{k+1} = \min\{\alpha_{\max}, \gamma_2 \alpha_k\}$ [a successful iteration].

Otherwise set $x_{k+1} = x_k$ and $\alpha_{k+1} = \gamma_1 \alpha_k$ [an unsuccessful iteration].

Increase the iteration count by setting $k = k + 1$ in both cases.

We note that

$$\hat{m}_k(0) - \hat{m}_k(\hat{s}_k) = \hat{l}_k(0) - \hat{l}_k(\hat{s}_k) + \frac{1}{2\alpha_k} \|S_k^T \hat{s}_k\|_2^2 \geq \frac{1}{2\alpha_k} \|S_k^T \hat{s}_k\|_2^2, \quad (4.5.7)$$

where we have used (4.5.6). Lemma 4.5.3 shows Algorithm 4 satisfies Assumption 2.

Lemma 4.5.3. *Let f be continuously differentiable with L -Lipschitz continuous gradient. Then Algorithm 4 satisfies Assumption 2 with*

$$\alpha_{\text{low}} = \frac{1 - \theta}{L + B_{\max}}.$$

Proof. Let $\epsilon > 0$ and $k < N_\epsilon$, and assume iteration k is true with $\alpha_k \leq \alpha_{low}$, define

$$\rho_k = \frac{f(x_k) - f(x_k + s_k)}{\hat{m}_k(0) - \hat{m}_k(\hat{s}_k)}.$$

We have

$$\begin{aligned} |1 - \rho_k| &= \frac{|f(x_k + s_k) - \hat{m}_k(s_k)|}{|\hat{m}_k(0) - \hat{m}_k(\hat{s}_k)|} \quad \text{by } \hat{m}_k(0) = f(x_k) \\ &\leq \frac{\left(\frac{L+B_{max}}{2}\right) \|S_k^T \hat{s}_k\|_2^2}{\frac{1}{2\alpha_k} \|S_k^T \hat{s}_k\|_2} \\ &\leq 1 - \theta, \end{aligned}$$

where the first inequality follows from Lemma 4.5.2 and (4.5.7). The above equation implies that $\rho_k \geq \theta$ and therefore iteration k is successful.⁷ \square

The next Lemma shows Algorithm 4 satisfies Assumption 7, thus satisfying Assumption 3 by Lemma 4.5.1.

Lemma 4.5.4. *Algorithm 4 satisfies Assumption 7 with*

$$\bar{h}(z_1, z_2) = \frac{z_1^2}{2\alpha_{max} (S_{max} (B_{max} + z_2^{-1}) + \kappa_T)^2}, \quad (4.5.8)$$

where S_{max} is defined in (4.4.5).

Proof. Let iteration k be true. Using the definition of \hat{l}_k , we have

$$\nabla \hat{l}_k(\hat{s}_k) = S_k \nabla f(x_k) + S_k B_k S_k^T \hat{s}_k + \frac{1}{\alpha_k} S_k S_k^T \hat{s}_k,$$

It follows that

$$\begin{aligned} \|S_k \nabla f(x_k)\|_2 &= \left\| -S_k \left(B_k + \frac{1}{\alpha_k} \right) S_k^T \hat{s}_k + \nabla \hat{l}_k(\hat{s}_k) \right\|_2 \\ &\leq \left(S_{max} \left(B_{max} + \frac{1}{\alpha_k} \right) \right) \|S_k^T \hat{s}_k\|_2 + \|\nabla \hat{l}_k(\hat{s}_k)\|_2 \\ &\leq \left(S_{max} \left(B_{max} + \frac{1}{\alpha_k} \right) + \kappa_T \right) \|S_k^T \hat{s}_k\|_2, \end{aligned} \quad (4.5.9)$$

where we used $\|S_k\|_2 \leq S_{max}$ on true iterations and $\|B_k\|_2 \leq B_{max}$ to derive the first inequality and (4.5.5) to derive the last inequality.

Therefore, using (4.5.7) and (4.5.9), we have

$$\begin{aligned} \hat{m}_k(0) - \hat{m}_k(\hat{s}_k) &\geq \frac{1}{2\alpha_k} \|S_k^T \hat{s}_k\|_2^2 \\ &\geq \frac{1}{2\alpha_k} \left(\frac{1}{S_{max} \left(B_{max} + \frac{1}{\alpha_k} \right) + \kappa_T} \right)^2 \|S_k \nabla f(x_k)\|_2^2 \\ &\geq \frac{1}{2\alpha_{max}} \left(\frac{1}{S_{max} \left(B_{max} + \frac{1}{\alpha_k} \right) + \kappa_T} \right)^2 \|S_k \nabla f(x_k)\|_2^2, \end{aligned} \quad (4.5.10)$$

⁷For ρ_k to be well-defined, we need the denominator to be strictly positive. But this is shown in (4.5.10).

satisfying Assumption 7.

□

4.5.2 Iteration complexity of random subspace quadratic regularisation methods

Here we derive complexity results for three concrete implementations of Algorithm 4 that use different random ensembles. Many other random ensembles are possible. As a reminder, the below expression, introduced earlier in this chapter, will be needed.

$$\tau_\alpha = \left\lceil \log_{\gamma_1} \left(\min \left\{ \frac{\alpha_{low}}{\alpha_0}, \frac{1}{\gamma_2} \right\} \right) \right\rceil \quad (4.5.11)$$

Applying Lemma 4.5.1, Lemma 4.5.3, Lemma 4.5.4 for Algorithm 4, we have that Assumption 2 and Assumption 3 are satisfied with

$$\begin{aligned} \alpha_{low} &= \frac{1 - \theta}{L + B_{max}} \\ h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) &= \theta \bar{h} \left((1 - \epsilon_S)^{1/2} \epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha} \right) \\ &= \frac{\theta(1 - \epsilon_S)\epsilon^2}{2\alpha_{max} (S_{max} (B_{max} + \alpha_0^{-1} \gamma_1^{-c-\tau_\alpha}) + \kappa_T)^2} \end{aligned} \quad (4.5.12)$$

Moreover, Assumption 4 for Algorithm 4 is satisfied by applying Lemma 4.4.1. The following three subsections give complexity results of Algorithm 4 with different random ensembles. We suggest the reader to refer back to Table 4.2 for a summary of their theoretical properties.

4.5.2.1 Using scaled Gaussian matrices

Algorithm 4 with scaled Gaussian matrices have a (high-probability) iteration complexity of $\mathcal{O}(\frac{d}{l}\epsilon^{-2})$ to drive $\nabla f(x_k)$ below ϵ and l can be chosen as a (problem dimension-independent) constant (see Table 4.2).

Theorem 4.5.1. *Suppose f is continuously differentiable with L -Lipschitz continuous gradient. Let $\delta_S^{(2)}, \epsilon_S, \delta_1 > 0$, $l \in \mathbb{N}^+$ such that*

$$\delta_S < \frac{c}{(c+1)^2}, \quad \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} > 0,$$

where $\delta_S = e^{-l\epsilon_S^2/4} + \delta_S^{(2)}$. Run Algorithm 4 with \mathcal{S} being the distribution of scaled Gaussian matrices, for N iterations with

$$N \geq \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

where

$$h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) = \frac{\theta(1 - \epsilon_S)\epsilon^2}{2\alpha_{max} \left(\left[1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}} \right] (B_{max} + \alpha_0^{-1} \gamma_1^{-c-\tau_\alpha}) + \kappa_T \right)^2}$$

and τ_α is given in (4.5.11). Then, we have

$$\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N},$$

where N_ϵ is defined in (4.2.2).

Proof. We note that Algorithm 4 is a particular form of Algorithm 2, therefore Theorem 4.2.1 applies. Moreover Assumption 2, Assumption 3 and Assumption 4 are satisfied. Applying Lemma 4.4.2, Lemma 4.4.4 and Lemma 4.4.6 for scaled Gaussian matrices, Assumption 1 is satisfied with

$$S_{max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}}$$

$$\delta_S = e^{-\epsilon_S^2 l/4} + \delta_S^{(2)}.$$

Applying Theorem 4.2.1 and substituting the expression of S_{max} above in (4.5.12) gives the desired result. \square

4.5.2.2 Using stable 1-hashing matrices

Algorithm 4 with stable 1-hashing matrices have a (high-probability) iteration complexity of $\mathcal{O}\left(\frac{d}{l}\epsilon^{-2}\right)$ to drive $\nabla f(x_k)$ below ϵ and l can be chosen as a (problem dimension-independent) constant (see Table 4.2).

Theorem 4.5.2. *Suppose f is continuously differentiable with L -Lipschitz continuous gradient. Let $\delta_1 > 0$, $\epsilon_S \in (0, 3/4)$, $l \in \mathbb{N}^+$ such that*

$$\delta_S < \frac{c}{(c+1)^2}, \quad \left[(1-\delta_S)(1-\delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} > 0,$$

where $\delta_S = e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}}$ and C_3 is defined in Lemma 4.4.10. Run Algorithm 4 with \mathcal{S} being the distribution of stable 1-hashing matrices, for N iterations with

$$N \geq \left[(1-\delta_S)(1-\delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

where

$$h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) = \frac{\theta(1-\epsilon_S)\epsilon^2}{2\alpha_{max} \left(\sqrt{\lceil d/l \rceil} (B_{max} + \alpha_0^{-1} \gamma_1^{-c-\tau_\alpha}) + \kappa_T \right)^2}$$

and τ_α is given in (4.5.11). Then, we have

$$\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N},$$

where N_ϵ is defined in (4.2.2).

Proof. Applying Lemma 4.4.2, Lemma 4.4.10 and Lemma 4.4.11 for stable 1-hashing matrices, Assumption 1 is satisfied with

$$S_{max} = \sqrt{\lceil d/l \rceil}$$

$$\delta_S = e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}}.$$

Applying Theorem 4.2.1 and substituting the expression of S_{max} above in (4.5.12) gives the desired result. \square

4.5.2.3 Using sampling matrices

Algorithm 4 with scaled sampling matrices have a (high-probability) iteration complexity of $\mathcal{O}(\frac{d}{l}\epsilon^{-2})$ to drive $\nabla f(x_k)$ below ϵ . However, unlike in the previous two cases, here l depends on the problem dimension d and a problem specific constant ν (see Table 4.2). If $\nu = \mathcal{O}(1/d)$, then l can be chosen as a problem dimension-independent constant.

Theorem 4.5.3. *Suppose f is continuously differentiable with L -Lipschitz continuous gradient. Let $\delta_1 > 0$, $\epsilon_S \in (0, 1)$, $l \in \mathbb{N}^+$ such that*

$$\delta_S < \frac{c}{(c+1)^2}, \quad \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} > 0,$$

where $\delta_S = e^{-\frac{\epsilon_S^2 l}{2d\nu^2}}$ and ν is defined in (4.4.19). Run Algorithm 4 with \mathcal{S} being the distribution of scaled sampling matrices, for N iterations with

$$N \geq \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

where

$$h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) = \frac{\theta(1 - \epsilon_S)\epsilon^2}{2\alpha_{max} \left(\sqrt{d/l} (B_{max} + \alpha_0^{-1} \gamma_1^{-c-\tau_\alpha}) + \kappa_T \right)^2}$$

and τ_α is given in (4.5.11). Then, we have

$$\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N},$$

where N_ϵ is defined in (4.2.2).

Proof. Applying Lemma 4.4.2, Lemma 4.4.13 and Lemma 4.4.14 for scaled sampling matrices, Assumption 1 is satisfied with

$$S_{max} = \sqrt{d/l}$$

$$\delta_S = e^{-\frac{\epsilon_S^2 l}{2d\nu^2}}.$$

Applying Theorem 4.2.1 and substituting the expression of S_{max} above in (4.5.12) gives the desired result. \square

Remark 7. The dependency on ϵ in the iteration complexity matches that for the full-space quadratic regularisation method (Page 13). Note that for each ensemble considered, there is dimension-dependence in the bound of the form $\frac{d}{l}$. We may eliminate the dependence on d in the iteration complexity by fixing the ratio $\frac{d}{l}$ to be a constant.

4.5.3 Random subspace trust region methods with sketching

Here we present a generic random subspace trust region method with sketching, Algorithm 5, which is a particular form of Algorithm 3 where the step is computed using a trust region framework (see Page 12 in Chapter 1).

Algorithm 5 A generic random subspace trust region method with sketching

Initialization

Choose a matrix distribution \mathcal{S} of matrices $S \in \mathbb{R}^{l \times d}$. Choose constants $\gamma_1 \in (0, 1)$, $\gamma_2 = \gamma_1^{-c}$, for some $c \in \mathbb{N}^+$, $l \in \mathbb{N}^+$, $\theta \in (0, 1)$ and $\alpha_{\max}, B_{\max} > 0$. Initialize the algorithm by setting $x_0 \in \mathbb{R}^d$, $\alpha_0 = \alpha_{\max} \gamma_1^p$ for some $p \in \mathbb{N}^+$ and $k = 0$.

1. Compute a reduced model and a step

Draw a random matrix $S_k \in \mathbb{R}^{l \times d}$ from \mathcal{S} , and let

$$\hat{m}_k(\hat{s}) = f(x_k) + \langle S_k \nabla f(x_k), \hat{s} \rangle + \frac{1}{2} \langle \hat{s}, S_k B_k S_k^T \hat{s} \rangle \quad (4.5.13)$$

where $B_k \in \mathbb{R}^{d \times d}$ is a user provided matrix with $\|B_k\|_2 \leq B_{\max}$.

Compute \hat{s}_k by approximately minimising $\hat{m}_k(\hat{s})$ such that for some $C_7 > 0$,

$$\|\hat{s}_k\|_2 \leq \alpha_k \quad (4.5.14)$$

$$\hat{m}_k(0) - \hat{m}_k(\hat{s}_k) \geq C_7 \|S_k \nabla f(x_k)\|_2 \min \left\{ \alpha_k, \frac{\|S_k \nabla f(x_k)\|_2}{\|B_k\|_2} \right\}. \quad (4.5.15)$$

2. Check sufficient decrease

Let sufficient decrease be defined by the condition

$$f(x_k) - f(x_k + s_k) \geq \theta [\hat{m}_k(0) - \hat{m}_k(\hat{s}_k(\alpha_k))].$$

3. Update the parameter α_k and possibly take the potential step s_k

If sufficient decrease is achieved, set $x_{k+1} = x_k + s_k$ and $\alpha_{k+1} = \min \{\alpha_{\max}, \gamma_2 \alpha_k\}$ [successful iteration].

Otherwise set $x_{k+1} = x_k$ and $\alpha_{k+1} = \gamma_1 \alpha_k$. [unsuccessful iteration].

Increase the iteration count by setting $k = k + 1$ in both cases.

Remark 8. Lemma 4.3 in [85] shows there always exists $\hat{s}_k \in \mathbb{R}^l$ such that (4.5.15) holds. Specifically, define $g_k = S_k \nabla f(x_k)$. If $g_k = 0$, one may take $\hat{s}_k = 0$; and otherwise one may take \hat{s}_k to be the Cauchy point (that is, the point where the model \hat{m}_k is minimised in the negative model gradient direction within the trust region), which can be computed by $\hat{s}_k^c = -\tau_k \frac{\alpha_k}{\|g_k\|_2} g_k$, where $\tau_k = 1$ if $g_k^T B_k g_k \leq 0$; and $\tau_k = \min \left(\frac{\|g_k\|_2^3}{g_k^T B_k g_k \alpha_k}, 1 \right)$ otherwise.

Lemma 4.5.5 shows that Algorithm 5 satisfies Assumption 2.

Lemma 4.5.5. *Suppose f is continuously differentiable with L -Lipschitz continuous gradient. Then Algorithm 5 satisfies Assumption 2 with*

$$\alpha_{low} = (1 - \epsilon_S)^{1/2} \epsilon \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right). \quad (4.5.16)$$

Proof. Let $\epsilon > 0$ and $k < N_\epsilon$, and assume iteration k is true with $\alpha_k \leq \alpha_{low}$, define

$$\rho_k = \frac{f(x_k) - f(x_k + s_k)}{\hat{m}_k(0) - \hat{m}_k(\hat{s}_k)}.$$

Then we have

$$\begin{aligned} |1 - \rho_k| &= \left| \frac{f(x_k + s_k) - \hat{m}_k(\hat{s}_k)}{\hat{m}_k(0) - \hat{m}_k(\hat{s}_k)} \right| \\ &\leq \frac{(L + \frac{1}{2}B_{max}) \|S_k^T \hat{s}_k\|_2^2}{C_7 \|S_k \nabla f(x_k)\|_2 \min \left(\alpha_k, \frac{\|S_k \nabla f(x_k)\|_2}{\|B_k\|_2} \right)} \\ &\leq \frac{(L + \frac{1}{2}B_{max}) S_{max}^2 \alpha_k^2}{C_7 \|S_k \nabla f(x_k)\|_2 \min \left(\alpha_k, \frac{\|S_k \nabla f(x_k)\|_2}{\|B_k\|_2} \right)} \\ &\leq \frac{(L + \frac{1}{2}B_{max}) S_{max}^2 \alpha_k^2}{C_7 (1 - \epsilon_S)^{1/2} \epsilon \min \left(\alpha_k, \frac{(1 - \epsilon_S)^{1/2} \epsilon}{B_{max}} \right)} \\ &\leq 1 - \theta, \end{aligned}$$

where the first inequality follows from (4.5.15) and Lemma 4.5.2, the second inequality follows from (4.4.5) and $\|\hat{s}_k\|_2 \leq \alpha_k$, the third inequality follows from (4.4.4) and the fact that $\nabla f(x_k) > \epsilon$ for $k < N_\epsilon$, the last inequality follows from $\alpha_k \leq \alpha_{low}$ and (4.5.16). It follows then $\rho_k \geq \theta$ and iteration k is successful. ⁸ \square

The next lemma shows that Algorithm 5 satisfies Assumption 7, thus satisfying Assumption 3.

Lemma 4.5.6. *Algorithm 5 satisfies Assumption 7 with*

$$\bar{h}(z_1, z_2) = C_7 \min(z_1 z_2, z_1^2 / B_{max}).$$

Proof. Use (4.5.15) with $\|B_k\|_2 \leq B_{max}$. \square

4.5.4 Iteration complexity of random subspace trust region methods

Here we derive complexity results for three concrete implementations of Algorithm 5 that use different random ensembles. The exposition follows closely Section 4.5.2. And the complexity results are in the same order in $\epsilon, \frac{d}{l}$ but have different constants.

⁸Note that for k being a true iteration with $k < N_\epsilon$, (4.5.15) along with (4.4.4), $\alpha_k > 0$ gives $\hat{m}_k(0) - \hat{m}_k(\hat{s}_k) > 0$ so that ρ_k is well defined.

Applying Lemma 4.5.1, Lemma 4.5.5, Lemma 4.5.6 for Algorithm 5, we have that Assumption 2 and Assumption 3 are satisfied with

$$\begin{aligned}\alpha_{low} &= (1 - \epsilon_S)^{1/2} \epsilon \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right) \\ h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) &= \theta \bar{h} \left((1 - \epsilon_S)^{1/2} \epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha} \right) \\ &= \theta C_7 \min \left\{ (1 - \epsilon_S)^{1/2} \epsilon \alpha_0 \gamma_1^{c+\tau_\alpha}, (1 - \epsilon_S) \epsilon^2 / B_{max} \right\}\end{aligned}\quad (4.5.17)$$

Here, unlike in the analysis of Algorithm 4, α_{low} (and consequently τ_α) depends on ϵ . We make this dependency on ϵ explicit. Using the definition of τ_α in (4.2.8) and substituting in the expression for α_{low} , we have

$$\begin{aligned}\alpha_0 \gamma_1^{c+\tau_\alpha} &= \alpha_0 \gamma_1^c \gamma_1^{\left\lceil \log_{\gamma_1} \left(\min \left\{ (1 - \epsilon_S)^{1/2} \epsilon \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right) \alpha_0^{-1}, \gamma_2^{-1} \right\} \right) \right\rceil} \\ &\geq \alpha_0 \gamma_1^c \gamma_1 \min \left\{ (1 - \epsilon_S)^{1/2} \epsilon \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right) \alpha_0^{-1}, \gamma_2^{-1} \right\} \\ &= \gamma_1^{c+1} \min \left\{ (1 - \epsilon_S)^{1/2} \epsilon \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right), \alpha_0 \gamma_2^{-1} \right\},\end{aligned}$$

where we used $\lceil y \rceil \leq y + 1$ to derive the inequality. Therefore, (4.5.17) implies

$$\begin{aligned}h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) &\geq \theta C_7 \min \left\{ \gamma_1^{c+1} \min \left\{ (1 - \epsilon_S) \epsilon^2 \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right), (1 - \epsilon_S)^{1/2} \epsilon \alpha_0 \gamma_2^{-1} \right\}, \frac{(1 - \epsilon_S) \epsilon^2}{B_{max}} \right\} \\ &= \theta C_7 (1 - \epsilon_S) \epsilon^2 \min \left\{ \gamma_1^{c+1} \min \left\{ \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right), \frac{\alpha_0}{(1 - \epsilon_S)^{1/2} \epsilon \gamma_2} \right\}, \frac{1}{B_{max}} \right\} \\ &= \theta C_7 (1 - \epsilon_S) \epsilon^2 \gamma_1^{c+1} \min \left\{ \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2}B_{max})S_{max}^2}, \frac{1}{B_{max}} \right), \frac{\alpha_0}{(1 - \epsilon_S)^{1/2} \epsilon \gamma_2} \right\}.\end{aligned}\quad (4.5.18)$$

where the last equality follows from $\gamma_1^{c+1} < 1$. Moreover, Assumption 4 for Algorithm 5 is satisfied by applying Lemma 4.4.1. The following three subsections give complexity results of Algorithm 5 using different random ensembles with Algorithm 5. Again, we suggest the reader to refer back to Table 4.2 for a summary of their theoretical properties.

4.5.4.1 Using scaled Gaussian matrices

Algorithm 5 with scaled Gaussian matrices have a (high-probability) iteration complexity of $\mathcal{O}(\frac{d}{l} \epsilon^{-2})$ to drive $\nabla f(x_k)$ below ϵ and l can be chosen as a (problem dimension-independent) constant (see Table 4.2).

Theorem 4.5.4. *Suppose f is continuously differentiable with L -Lipschitz continuous gradient. Let $\delta_S^{(2)}, \epsilon_S, \delta_1 > 0, l \in \mathbb{N}^+$ such that*

$$\delta_S < \frac{c}{(c+1)^2}, \quad \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} > 0,$$

where $\delta_S = e^{-l\epsilon_S^2/4} + \delta_S^{(2)}$. Run Algorithm 5 with \mathcal{S} being the distribution of scaled Gaussian matrices, for N iterations with

$$N \geq \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

where

$$h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) \tag{4.5.19}$$

$$= \theta C_7 (1 - \epsilon_S) \epsilon^2 \gamma_1^{c+1} \min \left\{ \min \left(\frac{C_7(1 - \theta)}{(L + \frac{1}{2} B_{max}) \left[1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}} \right]^2}, \frac{1}{B_{max}} \right), \frac{\alpha_0}{(1 - \epsilon_S)^{1/2} \epsilon \gamma_2} \right\} \tag{4.5.20}$$

Then, we have

$$\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-\frac{\delta_1^2}{2}(1 - \delta_S)N},$$

where N_ϵ is defined in (4.2.2).

Proof. We note that Algorithm 5 is a particular version of Algorithm 2 therefore Theorem 4.2.1 applies. Applying Lemma 4.4.2, Lemma 4.4.4 and Lemma 4.4.6 for scaled Gaussian matrices, Assumption 1 is satisfied with

$$S_{max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}}$$

$$\delta_S = e^{-\epsilon_S^2 l/4} + \delta_S^{(2)}.$$

Applying Theorem 4.2.1 and substituting the expression of S_{max} in (4.5.18) gives the desired result. \square

4.5.4.2 Using stable 1-hashing matrices

Algorithm 5 with stable 1-hashing matrices have a (high-probability) iteration complexity of $\mathcal{O}(\frac{d}{l} \epsilon^{-2})$ to drive $\nabla f(x_k)$ below ϵ and l can be chosen as a (problem dimension-independent) constant (see Table 4.2).

Theorem 4.5.5. Suppose f is continuously differentiable with L -Lipschitz continuous gradient. Let $\delta_1 > 0$, $\epsilon_S \in (0, 3/4)$, $l \in \mathbb{N}^+$ such that

$$\delta_S < \frac{c}{(c+1)^2}, \quad \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} > 0,$$

where $\delta_S = e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}}$ and C_3 is defined in Lemma 4.4.10. Run Algorithm 5 with \mathcal{S} being the distribution of stable 1-hashing matrices, for N iterations with

$$N \geq \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

where

$$h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) = \theta C_7 (1 - \epsilon_S) \epsilon^2 \gamma_1^{c+1} \min \left\{ \min \left(\frac{C_7(1-\theta)}{(L + \frac{1}{2} B_{max}) \lceil d/l \rceil}, \frac{1}{B_{max}} \right), \frac{\alpha_0}{(1 - \epsilon_S)^{1/2} \epsilon \gamma_2} \right\}$$

Then, we have

$$\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N},$$

where N_ϵ is defined in (4.2.2).

Proof. Applying Lemma 4.4.2, Lemma 4.4.10 and Lemma 4.4.11 for stable 1-hashing matrices, Assumption 1 is satisfied with

$$S_{max} = \sqrt{\lceil d/l \rceil}$$

$$\delta_S = e^{-\frac{l(\epsilon_S - 1/4)^2}{C_3}}.$$

Applying Theorem 4.2.1 and substituting the expression of S_{max} in (4.5.18) gives the desired result. \square

4.5.4.3 Using sampling matrices

Algorithm 5 with scaled sampling matrices have a (high-probability) iteration complexity of $\mathcal{O}(\frac{d}{l}\epsilon^{-2})$ to drive $\nabla f(x_k)$ below ϵ . Similar to Algorithm 4 with scaled sampling matrices, here l depends on the problem dimension d and a problem specific constant ν (see Table 4.2). If $\nu = \mathcal{O}(1/d)$, then l can be chosen as a problem dimension-independent constant.

Theorem 4.5.6. *Suppose f is continuously differentiable with L -Lipschitz continuous gradient. Let $\delta_1 > 0$, $\epsilon_S \in (0, 1)$, $l \in \mathbb{N}^+$ such that*

$$\delta_S < \frac{c}{(c+1)^2}, \quad \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} > 0,$$

where $\delta_S = e^{-\frac{\epsilon^2 l}{2d\nu^2}}$ and ν is defined in (4.4.19). Run Algorithm 5 with \mathcal{S} being the distribution of scaled sampling matrices, for N iterations with

$$N \geq \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1} \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

where

$$h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha}) = \theta C_7 (1 - \epsilon_S) \epsilon^2 \gamma_1^{c+1} \min \left\{ \min \left(\frac{C_7(1-\theta)}{(L + \frac{1}{2} B_{max}) d/l}, \frac{1}{B_{max}} \right), \frac{\alpha_0}{(1 - \epsilon_S)^{1/2} \epsilon \gamma_2} \right\}$$

Then, we have

$$\mathbb{P}(N \geq N_\epsilon) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N},$$

where N_ϵ is defined in (4.2.2).

Proof. Applying Lemma 4.4.2, Lemma 4.4.13 and Lemma 4.4.14 for scaled sampling matrices, Assumption 1 is satisfied with

$$S_{max} = \sqrt{d/l}$$

$$\delta_S = e^{-\frac{\epsilon^2 l}{2dv^2}}.$$

Applying Theorem 4.2.1 and substituting the expression of S_{max} in (4.5.18) gives the desired result. \square

Remark 9. *Similar to Algorithm 4, Algorithm 5 matches the iteration complexity of the corresponding (full-space) trust region method; and the l/d dependency can be eliminated by setting l to be a constant fraction of d .*

Remark 10. *Although Algorithm 4 and Algorithm 5 with a(any) of the above three random ensembles only require l directional derivative evaluations of f per iteration, instead of d derivative evaluations required by the (full-space) methods, the iteration complexities are increased by a factor of d/l . Therefore, theoretically, Algorithm 4 and Algorithm 5 do not reduce the total number of gradient evaluations. However, the computational cost of the step \hat{s}_k is typically reduced from being proportional to d^2 to being proportional to l^2 (for example, if we are solving a non-linear least squares problem and choose $B_k = J_k^T J_k$) thus we still gain in having a smaller computational complexity. In practice, our theoretical analysis may not be tight and therefore we could gain in having both a smaller gradient evaluation complexity and a smaller computational complexity. See numerical illustrations in Section 4.6. In addition, by reducing the number of variables from d (which can be arbitrarily large) to l (which can be set as a constant, see Table 4.2), Algorithm 4 and Algorithm 5 reduce the memory requirement of the computation of the sub-problem at each iteration, comparing to the corresponding full-space methods.*

4.6 Randomised Subspace Gauss-Newton (R-SGN) for non-linear least squares

We consider the nonlinear least-squares problem (defined in (1.4.9))

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{2} \sum_{i=1}^n \|r_i(x)\|_2^2 = \frac{1}{2} \|r(x)\|_2^2$$

where $r = (r_1, \dots, r_n) : \mathbb{R}^d \rightarrow \mathbb{R}^n$ is a smooth vector of nonlinear (possibly nonconvex) residual functions. We define the Jacobian (matrix of first order derivatives) as

$$J(x) = \left(\frac{\partial r_i(x)}{\partial x_j} \right)_{ij} \in \mathbb{R}^{n \times d}$$

and can then compactly write the gradient as $\nabla f(x) = J(x)^T r(x)$. It can be shown e.g. in [85], that the gradient and Hessian of $f(x)$ is then given by

$$\begin{aligned}\nabla f(x) &= J(x)^T r(x), \\ \nabla^2 f(x) &= J(x)^T J(x) + \sum_{i=1}^n r_i(x) \nabla^2 r_i(x)\end{aligned}$$

The classical Gauss-Newton (GN) algorithm applies Newton's method to minimising f with only the first-order $J(x)^T J(x)$ term in the Hessian, dropping the second-order terms involving the Hessians of the residuals r_i . This is equivalent to linearising the residuals in (1.4.9) so that

$$r(x + s) \approx r(x) + J(x)s,$$

and minimising the resulting model in the step $s \in \mathbb{R}^d$. Thus, at every iterate x_k , Gauss-Newton approximately minimises the following convex quadratic local model

$$f(x_k) + \langle J(x_k)^T r(x_k), s \rangle + \frac{1}{2} \langle s, J(x_k)^T J(x_k) s \rangle$$

over $s \in \mathbb{R}^d$. In our approach, which we call Random Subspace Gauss-Newton (R-SGN), we reduce the dimensionality of this model by minimising in an l -dimensional randomised subspace $\mathcal{L} \subset \mathbb{R}^d$, with $l \ll d$, by approximately minimising the following reduced model

$$f(x_k) + \langle J_S(x_k)^T r(x_k), \hat{s} \rangle + \frac{1}{2} \langle \hat{s}, J_S(x_k)^T J_S(x_k) \hat{s} \rangle \quad (4.6.1)$$

over $\hat{s} \in \mathbb{R}^l$, where $J_S(x_k) = J(x_k) S_k^T \in \mathbb{R}^{n \times l}$ denotes the reduced Jacobian for $S_k \in \mathbb{R}^{l \times d}$ being a randomly generated sketching matrix. Note that with $B_k = J_k^T J_k$, Algorithm 3 framework can be applied directly to this subspace Gauss-Newton method; guaranteeing its convergence under assumptions of model minimisation and sketching matrices. Compared to the classical Gauss-Newton model, in addition to the speed-up gained due to the model dimension being reduced from d to l , this reduced model also offers the computational advantage that it only needs to evaluate l Jacobian actions, giving $J_S(x_k)$, instead of the full Jacobian matrix $J(x_k)$.

In its simplest form, when S_k is a scaled sampling matrix, J_S can be thought of as a random subselection of columns of the full Jacobian J , which leads to variants of our framework that are Block-Coordinate Gauss-Newton (BC-GN) methods. In this case, for example, if the Jacobian were being calculated by finite-differences of the residual r , only a small number of evaluations of r along coordinate directions would be needed; such a BC-GN variant has already been used for parameter estimation in climate modelling [96]. Note that theoretically, the convergence of BC-GN method requires an upper bound on $\frac{\|\nabla f(x_k)\|_\infty}{\|\nabla f(x_k)\|_2}$ for all $k \in \mathbb{N}$ (for more details, see the discussion of sampling matrices on page 90) and Theorem 4.5.6.

More generally, S_k can be generated from any matrix distribution that satisfies Assumption 5, Assumption 6, e.g/ scaled Gaussian matrices or s -hashing matrices. In our work jointly done with

Jaroslav Fowkes [14, 13], we showcase the numerical performance of R-SGN methods with different sketching matrices. In this thesis we provide some numerical illustrations; the code used to produce these illustrations is written by Jaroslav Fowkes, and the results below appear in [14, 13].

Large-scale CUTEst problems We look at the behaviour of R-SGN on three large-scale ($d \approx 5,000$ to $10,000$) non-linear least squares problems from the CUTEst collection [40]. The three problems are given in Table 4.3. we run R-SGN five times (and take the average performance) on each problem until we achieve a 10^{-1} decrease in the objective, or failing that, for a maximum of 20 iterations. Furthermore, we plot the objective decrease against cumulative Jacobian action evaluations⁹ for each random run with subspace-sizes of 1%, 5%, 10%, 50%, 100% of the full-space-sizes.

Name	d	n	Name	d	n	Name	d	n
ARTIF	5,000	5,000	BRATU2D	4,900	4,900	OSCIGRNE	10,000	10,000

Table 4.3: The 3 large-scale CUTEst test problems.

Let us start by looking at the performance of R-SGN with scaled sampling matrices S_k . In Figure 4.1, we see that the objective decrease against cumulative Jacobian action evaluations for ARTIF, BRATU2D and OSCIGRNE. On ARTIF, we see that while R-SGN exhibits comparable performance, Gauss-Newton is clearly superior from a Jacobian action budget perspective. On BRATU2D we see that R-SGN really struggles to achieve any meaningful decrease, as does Gauss-Newton initially but then switches to a quadratic regime and quickly converges. On OSCIGRNE, we see that R-SGN with subspace sizes of $0.05d, 0.1d, 0.5d$ sometimes performs very well (outperforming Gauss-Newton) but sometimes struggles, and on averages Gauss-Newton performs better.

Next, we compare the performance of R-SGN with scaled Gaussian sketching matrices S_k . In Figure 4.2, we can see the objective decrease against cumulative Jacobian action evaluations for ARTIF, BRATU2D and OSCIGRNE. On ARTIF, we see that R-SGN with a subspace size of $0.5d$ outperforms Gauss-Newton initially before stagnating. On BRATU2D we once again see that R-SGN struggles to achieve any meaningful decrease, as does Gauss-Newton initially but then switches to a quadratic regime and quickly converges. However, on OSCIGRNE we see that R-SGN with a subspace size of $0.5d$ consistently outperforms Gauss-Newton.

Finally, we compare the performance of R-SGN with 3-hashing sketching matrices S_k . In Figure 4.3, we can see the objective decrease against cumulative Jacobian action evaluations for ARTIF, BRATU2D and OSCIGRNE. On ARTIF, we again see that R-SGN with a subspace size of $0.5d$ outperforms Gauss-Newton initially before stagnating. On BRATU2D we once again see that R-SGN struggles to achieve any meaningful decrease, as does Gauss-Newton initially but then switches to a

⁹The total number of evaluations of Jacobian-vector product used by the algorithm.

quadratic regime and quickly converges. However, on OSCIGRNE we again see that R-SGN with a subspace size of $0.5d$ consistently outperforms Gauss-Newton.

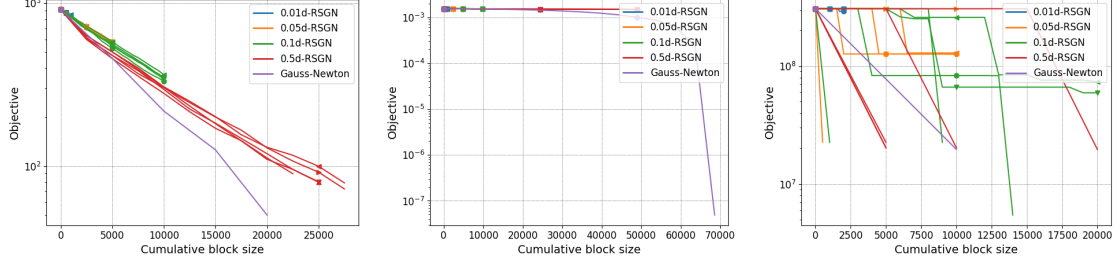


Figure 4.1: ARTIF (left), BRATU2D (middle) and OSCIGRNE (right) objective value against cumulative Jacobian action size for R-SGN with coordinate sampling.

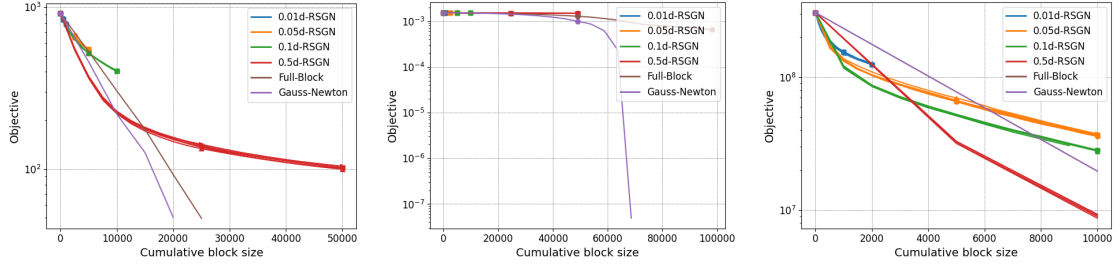


Figure 4.2: ARTIF (left), BRATU2D (middle) and OSCIGRNE (right) objective value against cumulative Jacobian action size for R-SGN with Gaussian sketching.

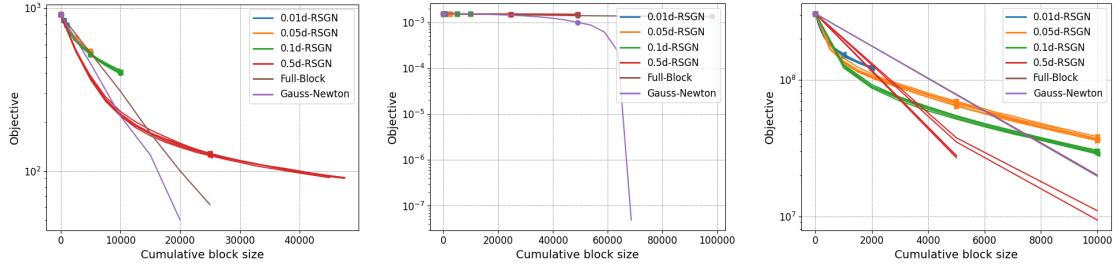


Figure 4.3: ARTIF (left), BRATU2D (middle) and OSCIGRNE (right) objective value against cumulative Jacobian action size for R-SGN with 3-hashing sketching.

Large scale machine learning problems Here we only use scaled sampling sketching matrices. We consider logistic regressions¹⁰, written in the form (1.4.9), by letting $r_i(x) = \ln(1 + \exp(-y_i a_i^T x))$, where $a_i \in \mathbb{R}^d$ are the observations and $y_i \in \{-1, 1\}$ are the class labels; we also include a quadratic regularization term $\lambda \|x\|_2^2$ by treating it as an additional residual.

We test on the CHEMOTHERAPY and GISETTE datasets from OpenML [99] for 100 iterations with $\lambda = 10^{-10}$, using subspace-sizes (or block-sizes, as here we are using the BC-GN variant by

¹⁰Here in order to fit in the non-linear least squares framework, we square the logistic losses r_i in the objective

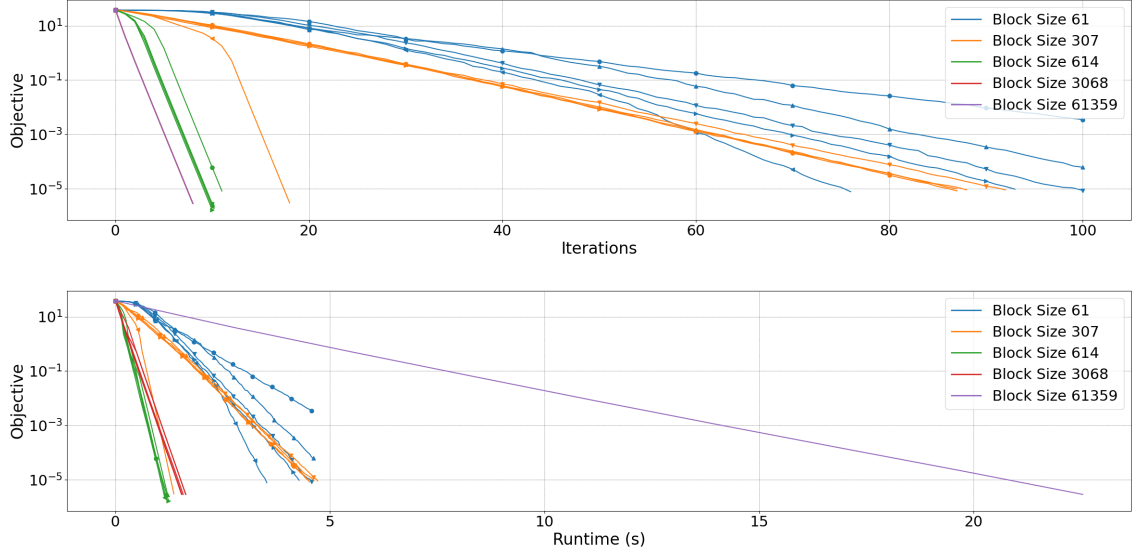


Figure 4.4: R-SGN on the CHEMOTHERAPY dataset

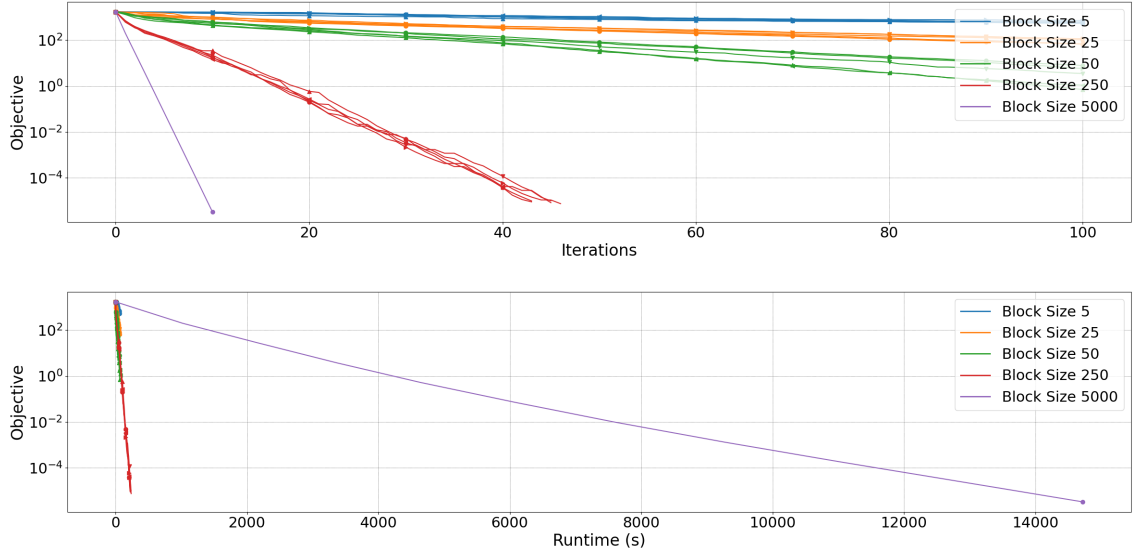


Figure 4.5: R-SGN on the GISETTE dataset

using the sampling sketching) of 0.1%, 0.5%, 1%, 5% and 100% of the full-space-sizes for the 61,359 dimensional CHEMOTHERAPY dataset and the 5,000 dimensional GISETTE dataset; in a similar testing setup to [42]. We perform five runs of the algorithm for each block size starting at $x_0 = 0$ (and take the average performance). We terminate once the objective $f(x_k)$ goes below 10^{-5} and plot $f(x_k)$ against iterations and runtime in each Figure. On the CHEMOTHERAPY dataset, we see from Figure 4.4 that we are able to get comparable performance to full Gauss-Newton ($d = 61,359$ in purple) using only 1% of the original block size ($l = 614$ in green) at 1/20th of the runtime. For the GISETTE dataset, we see from Figure 4.5 that similarly, we are able to get good performance

compared to GN ($d = 5,000$ in purple) using 5% of the original block size ($l = 250$ in red) at 1/60th of the runtime.

Chapter 5

Second order subspace methods for general objectives

5.1 Introduction

In this chapter we continue our investigation on subspace methods for the minimisation of general objectives. In the last chapter we saw that if the sketching matrix S_k stays bounded and is sufficiently accurate to capture the gradient of the objective at the current iterate with positive probability, convergence of the subspace methods occurs at essentially the same rate as classical full-space first order methods. It is known that if second order information of the objective function is available, cubic regularisation full-space methods achieve faster convergence rates for general non-convex objective [16]. In this chapter, we first show that the same can be obtained with subspace methods. Namely, when the sketching matrix S_k captures sufficiently accurate second order information, essentially the same faster rate of convergence can be achieved. We then show that this faster rate of convergence can be achieved also in the case of sparse second derivatives, without requiring the low rank/subspace embedding condition. Next, we show that a class of subspace methods converge to an approximate second order minimum in the subspace, in the sense that the subspace Hessian at the limit point is almost positive-semi-definite. Finally, we show that the second order subspace method with Gaussian sketching converges to an approximate second order minimum in the full space.

5.2 R-ARC: random subspace adaptive cubic regularisation method

First, we describe the random subspace cubic regularisation algorithm (R-ARC).

Algorithm 6 Random subspace cubic regularisation algorithm (R-ARC)

Initialization

Choose a matrix distribution \mathcal{S} of matrices $S \in \mathbb{R}^{l \times d}$. Choose constants $\gamma_1 \in (0, 1)$, $\gamma_2 > 1$, $\theta \in (0, 1)$, $\kappa_T, \kappa_S \geq 0$ and $\alpha_{\max} > 0$ such that $\gamma_2 = \frac{1}{\gamma_1^c}$, for some $c \in \mathbb{N}^+$. Initialize the algorithm by setting $x_0 \in \mathbb{R}^d$, $\alpha_0 = \alpha_{\max} \gamma^p$ for some $p \in \mathbb{N}^+$ and $k = 0$.

1. Compute a reduced model and a trial step

In Step 1 of Algorithm 2, draw a random matrix $S_k \in \mathbb{R}^{l \times d}$ from \mathcal{S} , and let

$$\begin{aligned} \hat{m}_k(\hat{s}) &= f(x_k) + \langle S_k \nabla f(x_k), \hat{s} \rangle + \frac{1}{2} \langle \hat{s}, S_k \nabla^2 f(x_k) S_k^T \hat{s} \rangle + \frac{1}{3\alpha_k} \|S_k^T \hat{s}\|_2^3 \\ &= \hat{q}_k(\hat{s}) + \frac{1}{3\alpha_k} \|S_k^T \hat{s}\|_2^3, \end{aligned} \quad (5.2.1)$$

where $\hat{q}_k(\hat{s})$ is the second order Taylor series of $f(x_k + S_k^T \hat{s}_k)$ around x_k ;

Compute \hat{s}_k by approximately minimising (5.2.1) such that

$$\hat{m}_k(\hat{s}_k) \leq \hat{m}_k(0) \quad (5.2.2)$$

$$\|\nabla \hat{m}_k(\hat{s}_k)\|_2 \leq \kappa_T \|S_k^T \hat{s}_k\|_2^2 \quad (5.2.3)$$

$$\nabla^2 \hat{m}_k(\hat{s}_k) \succeq -\kappa_S \|S_k^T \hat{s}_k\|_2, \quad (5.2.4)$$

where we may drop (5.2.4) if only convergence to a first order critical point is desired.

Compute a trial step

$$s_k = w_k(\hat{s}_k) = S_k^T \hat{s}_k, \quad (5.2.5)$$

2. Check sufficient decrease

In Step 2 of Algorithm 2, check sufficient decrease as defined by the condition

$$f(x_k) - f(x_k + s_k) \geq \theta [\hat{q}_k(0) - \hat{q}_k(\hat{s})], \quad (5.2.6)$$

3. Update the parameter α_k and possibly take the trial step s_k

If (5.2.6) holds, set $x_{k+1} = x_k + s_k$ and $\alpha_{k+1} = \min\{\alpha_{\max}, \gamma_2 \alpha_k\}$ [successful iteration].

Otherwise set $x_{k+1} = x_k$ and $\alpha_{k+1} = \gamma_1 \alpha_k$ [unsuccessful iteration]

Increase the iteration count by setting $k = k + 1$ in both cases.

Here note that Algorithm 6 is a specific form of Algorithm 2. Therefore the convergence result in Theorem 4.2.1 can be applied, provided that the four assumptions of the theorem can be shown to hold here. In the remaining sections of this chapter, we give different definitions of the two key terms in the convergence result Theorem 4.2.1: N_ϵ and true iterations. These lead to different requirements for the matrix distribution \mathcal{S} , and iteration complexities to drive $\|\nabla f(x_k)\|_2 < \epsilon$, $\lambda_{\min}(S_k \nabla^2 f(x_k) S_k^T) > -\epsilon_H$ and/or $\lambda_{\min}(\nabla^2 f(x_k)) > -\epsilon_H$

Compared to Algorithm 3, Algorithm 6 lets B_k be the Hessian at the iterate $\nabla^2 f(x_k)$, although we only need it in the form of $S_k \nabla^2 f(x_k) S_k^T$ so that the full Hessian never needs to be computed. Furthermore, we let \hat{s}_k , the reduced step, be computed by minimising a cubically regularised subspace model, corresponding to the classical approaches in [15], also see Section 1.4.2 on page 15. As in

the corresponding full-space method, the combination of the availability of second order information and the cubic regularisation term leads to improved iteration complexity for our subspace methods to drive $\nabla f(x_k) < \epsilon$ and convergence to a second order critical point.

Remark 11. *Two strategies for computing \hat{s}_k by minimising (5.2.1) are given in [15], either requiring a factorisation of $S_k \nabla^2 f(x_k) S_k^T$ (in a Newton-like algorithm), or repeated matrix-vector products involving $S_k \nabla^2 f(x_k) S_k^T$ (in a Lanczos-based algorithm). Although we note that the iteration complexity, and the evaluation complexity of f and its derivatives (which are the focuses of this chapter) are unaffected by the computation complexity of calculating \hat{s}_k , Algorithm 6 significantly reduces the computation of this inner problem by reducing the dimension of the Hessian from $d \times d$ to $l \times l$ comparing to the full-space counterpart. (In addition to reducing the gradient and Hessian evaluation complexity per iteration.)*

5.3 Fast convergence rate assuming subspace embedding of the Hessian matrix

Our first convergence result shows Algorithm 6 drives $\|\nabla f(x_k)\|_2$ below ϵ in $\mathcal{O}(\epsilon^{-3/2})$ iterations, given that \mathcal{S} has an embedding property (a necessary condition of which is that \mathcal{S} is an oblivious subspace embedding for matrices of rank $r + 1$, where r is the maximum rank of $\nabla^2 f(x_k)$ across all iterations).

Define N_ϵ and true iterations based on (one-sided) subspace embedding In order to prove convergence of Algorithm 6, we show that Assumption 1, Assumption 2, Assumption 3, Assumption 4 that are needed for Theorem 4.2.1 to hold are satisfied. To this end, we first define N_ϵ , the criterion for convergence, as $\min\{k : \|\nabla f(x_{k+1})\|_2 \leq \epsilon\}$. Then, we define the true iterations based on achieving an embedding of the Hessian and the gradient.

Definition 5.3.1. Let $\epsilon_S^{(2)} \in (0, 1)$, $S_{max} > 0$. Iteration k is $(\epsilon_S^{(2)}, S_{max})$ -true if

$$\|S_k M_k z_k\|_2^2 \geq (1 - \epsilon_S^{(2)}) \|M_k z_k\|_2^2, \quad \text{for all } z_k \in \mathbb{R}^{d+1} \quad (5.3.1)$$

$$\|S_k\|_2 \leq S_{max}, \quad (5.3.2)$$

where $M_k = [\nabla f(x_k) \quad \nabla^2 f(x_k)] \in \mathbb{R}^{d \times (d+1)}$. Note that all vectors are column vectors.

Remark 12. (5.3.1) implies

$$\|S_k \nabla f(x_k)\|_2^2 \geq (1 - \epsilon_S^{(2)}) \|\nabla f(x_k)\|_2^2, \quad (5.3.3)$$

by taking $z_k = [1, 0, \dots, 0]^T$. Thus Definition 5.3.1 is a stronger condition than Definition 4.4.1, the definition of true iterations for our convergence result of first order subspace methods.

5.3.1 Auxiliary results

In this subsection we provide some useful results needed to prove our assumptions in Theorem 4.2.1.

Lemma 5.3.1. *In Algorithm 6, if iteration k is successful, then*

$$f(x_{k+1}) \leq f(x_k) - \frac{\theta}{3\alpha_k} \|S_k^T \hat{s}_k\|_2^3$$

Proof. From the definition of successful iterations and (5.2.6)

$$\begin{aligned} f(x_{k+1}) &= f(x_k + s_k) \\ &\leq f(x_k) - \theta [\hat{m}_k(0) - \hat{m}_k(\hat{s}_k)] - \frac{\theta}{3\alpha_k} \|S_k^T \hat{s}_k\|_2^3 \\ &\leq f(x_k) - \frac{\theta}{3\alpha_k} \|S_k^T \hat{s}_k\|_2^3, \end{aligned} \tag{5.3.4}$$

where in the last inequality, we used (5.2.2). \square

The gradient of the model has the expression

$$\nabla \hat{m}_k(\hat{s}_k) = S_k \nabla f(x_k) + S_k \nabla^2 f(x_k) S_k^T \hat{s}_k + \frac{1}{\alpha_k} S_k S_k^T \hat{s}_k \|S_k^T \hat{s}_k\|_2. \tag{5.3.5}$$

The following lemma bounds the size of the step at true iterations.

Lemma 5.3.2. *Assume f is twice continuously differentiable with L_H -Lipschitz Hessian $\nabla^2 f$ and $k < N_\epsilon$. Suppose that iteration k is $(\epsilon_S^{(2)}, S_{max})$ -true. We have*

$$\|S_k^T \hat{s}_k\|_2^2 \geq \frac{\epsilon}{2} \min \left\{ \frac{2}{L_H}, \left(\frac{1}{\alpha_k} S_{max} + \kappa_T \right)^{-1} \sqrt{1 - \epsilon_S^{(2)}} \right\} \tag{5.3.6}$$

Proof. (5.3.5) and the triangle inequality give

$$\begin{aligned} \|S_k \nabla f(x_k) + S_k \nabla^2 f(x_k) S_k^T \hat{s}_k\|_2 &= \left\| \frac{1}{\alpha_k} S_k S_k^T \hat{s}_k \|S_k^T \hat{s}_k\|_2 - \nabla \hat{m}_k(\hat{s}_k) \right\|_2 \\ &\leq \frac{1}{\alpha_k} \|S_k\|_2 \|S_k^T \hat{s}_k\|_2^2 + \|\nabla \hat{m}_k(\hat{s}_k)\|_2 \\ &\leq \left(\frac{1}{\alpha_k} \|S_k\|_2 + \kappa_T \right) \|S_k^T \hat{s}_k\|_2^2 \quad \text{by (5.2.3)} \end{aligned} \tag{5.3.7}$$

$$\leq \left(\frac{1}{\alpha_k} S_{max} + \kappa_T \right) \|S_k^T \hat{s}_k\|_2^2, \tag{5.3.8}$$

where we used (5.3.2). On the other hand, we have that

$$\begin{aligned} &\|S_k \nabla f(x_k) + S_k \nabla^2 f(x_k) S_k^T \hat{s}_k\|_2 \\ &= \left\| S_k M_k [1, (S_k^T \hat{s}_k)^T]^T \right\|_2 \\ &\geq \sqrt{1 - \epsilon_S^{(2)}} \|\nabla f(x_k) + \nabla^2 f(x_k) s_k\|_2 \quad \text{by (5.3.1) with } z_k = [1, (S_k^T \hat{s}_k)^T]^T \\ &= \sqrt{1 - \epsilon_S^{(2)}} \|\nabla f(x_{k+1}) - [\nabla f(x_{k+1}) - \nabla f(x_k) - \nabla^2 f(x_k) s_k]\|_2 \end{aligned} \tag{5.3.9}$$

$$\geq \sqrt{1 - \epsilon_S^{(2)}} \|\nabla f(x_{k+1})\|_2 - \|[\nabla f(x_{k+1}) - \nabla f(x_k) - \nabla^2 f(x_k) s_k]\|_2 \tag{5.3.10}$$

Note that by Taylor's Theorem, because f is twice continuously differentiable with L_H -Lipschitz $\nabla^2 f$, we have that $\nabla f(x_k + s_k) = \nabla f(x_k) + \int_0^1 \nabla^2 f(x_k + ts_k) s_k dt$. Therefore, we have

$$\|\nabla f(x_{k+1}) - \nabla f(x_k) - \nabla^2 f(x_k) s_k\|_2 = \left\| \int_0^1 [\nabla^2 f(x_k + ts_k) - \nabla^2 f(x_k)] s_k dt \right\|_2 \quad (5.3.11)$$

$$\leq \int_0^1 \|s_k\|_2 \|\nabla^2 f(x_k + ts_k) - \nabla^2 f(x_k)\|_2 dt \quad (5.3.12)$$

$$\leq \|s_k\|_2 \int_0^1 L_H t \|s_k\|_2 dt \quad (5.3.13)$$

$$= \frac{1}{2} L_H \|s_k\|_2^2 \quad (5.3.14)$$

by Lipschitz continuity of $\nabla^2 f$. Next we discuss two cases,

1. If $L_H \|s_k\|_2^2 > \epsilon$, then we have the desired result in (5.3.6).
2. If $L_H \|s_k\|_2^2 \leq \epsilon$, then (5.3.10), and the fact that $\|\nabla f(x_{k+1})\|_2 \geq \epsilon$ by $k < N_\epsilon$, imply that

$$\|S_k \nabla f(x_k) + S_k \nabla^2 f(x_k) S_k^T \hat{s}_k\|_2 \geq \sqrt{1 - \epsilon_S^{(2)}} \frac{\epsilon}{2}.$$

Then (5.3.8) implies

$$\|S_k^T s_k\|_2^2 \geq \left(\frac{1}{\alpha_k} S_{max} + \kappa_T \right)^{-1} \sqrt{1 - \epsilon_S^{(2)}} \frac{\epsilon}{2}.$$

This again gives the desired result. □

5.3.2 Satisfying the assumptions of Theorem 4.2.1

Here we only address the case where \mathcal{S} is the distribution of scaled Gaussian matrices. But \mathcal{S} could also be the distribution of scaled sampling matrices, s -hashing matrices, SRHT matrices and HRHT matrices because those distributions also satisfy similar properties detailed below, namely, having a bounded two-norm with high probability (Lemma 5.3.3), and having a one-sided subspace embedding property (Lemma 5.3.4).

Concerning scaled Gaussian matrices, we have the following results.

Lemma 5.3.3 (Lemma 4.4.6). *Let $S \in \mathbb{R}^{l \times d}$ be a scaled Gaussian matrix (Definition 1.2.2). Then for any $\delta_S^{(2)} > 0$, S satisfies (5.3.2) with probability $1 - \delta_S^{(2)}$ and*

$$S_{max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}}.$$

Lemma 5.3.4 (Theorem 2.3 in [101]). *Let $\epsilon_S^{(2)} \in (0, 1)$ and $S \in \mathbb{R}^{l \times d}$ be a scaled Gaussian matrix. Then for any fixed $d \times (d+1)$ matrix M with rank at most $r+1$, with probability $1 - \delta_S^{(3)}$ we have that simultaneously for all $z \in \mathbb{R}^{d+1}$, $\|SMz\|_2^2 \geq (1 - \epsilon_S^{(2)}) \|Mz\|_2^2$, where*

$$\delta_S^{(3)} = e^{-\frac{\iota(\epsilon_S^{(2)})^2}{C_l} + r + 1} \quad (5.3.15)$$

and C_l is an absolute constant.

Satisfying Assumption 1 (page 71)

Lemma 5.3.5. Suppose that $\nabla^2 f(x_k)$ has rank at most $r \leq d$ for all k ; $S \in \mathbb{R}^{l \times d}$ is drawn as a scaled Gaussian matrix. Let $\epsilon_S^{(2)}, \delta_S^{(2)} \in (0, 1)$ such that $\delta_S^{(2)} + \delta_S^{(3)} < 1$ where $\delta_S^{(3)}$ is defined in (5.3.15).

Then Algorithm 6 satisfies Assumption 1 with $\delta_S = \delta_S^{(2)} + \delta_S^{(3)}$ and $S_{max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}}$, with true iterations defined in Definition 5.3.1.

Proof. Let $x_k = \bar{x}_k \in \mathbb{R}^d$ be given. This determines $\nabla f(x_k), \nabla^2 f(x_k)$ and hence M_k . As $\nabla^2 f(x_k)$ has rank at most r , M_k has rank at most $r + 1$. Consider the events

$$\begin{aligned} A_k^{(1)} &= \left\{ \|S_k M_k z\|_2^2 \geq (1 - \epsilon_S^{(2)}) \|M_k z\|_2^2, \quad \forall z \in \mathbb{R}^{d+1} \right\} \\ A_k^{(2)} &= \{\|S_k\|_2 \leq S_{max}\}. \end{aligned}$$

Note that iteration k is true if and only if $A_k^{(1)}$ and $A_k^{(2)}$ occur. It follows from Lemma 5.3.4 that $\mathbb{P}(A_k^{(1)} | x_k = \bar{x}_k) \geq 1 - \delta_S^{(3)}$; and from Lemma 5.3.3 that $\mathbb{P}(A_k^{(2)}) \geq 1 - \delta_S^{(2)}$. Since $A_k^{(2)}$ is independent of x_k , we have $\mathbb{P}(A_k^{(2)} | x_k = \bar{x}_k) = \mathbb{P}(A_k^{(2)}) \geq 1 - \delta_S^{(2)}$.

Hence, we have $\mathbb{P}(A_k^{(1)} \cap A_k^{(2)} | x_k = \bar{x}_k) \geq 1 - \mathbb{P}((A_k^{(1)})^c | x_k = \bar{x}_k) - \mathbb{P}((A_k^{(2)})^c | x_k = \bar{x}_k) \geq 1 - \delta_S^{(2)} - \delta_S^{(3)}$. A similar argument shows that $\mathbb{P}(A_0^{(1)} \cap A_0^{(2)}) \geq 1 - \delta_S^{(2)} - \delta_S^{(3)}$, as x_0 is fixed.

Moreover, given $x_k = \bar{x}_k$, $A_k^{(1)}$ and $A_k^{(2)}$ only depend on S_k , which is drawn randomly at iteration k . Hence given $x_k = \bar{x}_k$, $A_k^{(1)} \cap A_k^{(2)}$ is independent of whether the previous iterations are true or not. Hence Assumption 1 is true. \square

Satisfying Assumption 2 (page 71)

Lemma 5.3.6. Let f be twice continuously differentiable with L_H -Lipshitz continuous Hessian $\nabla^2 f$.

Algorithm 6 satisfies Assumption 2 with

$$\alpha_{low} = \frac{2(1 - \theta)}{L_H} \quad (5.3.16)$$

Proof. From (5.2.2), we have that

$$f(x_k) - \hat{q}_k(\hat{s}_k) \geq \frac{1}{3\alpha_k} \|S_k^T \hat{s}_k\|_2^3.$$

Using Lemma 5.3.2, in true iterations with $k < N_\epsilon$, we have that $\|S_k^T \hat{s}_k\|_2 > 0$. Therefore we can define ¹

$$\rho_k = \frac{f(x_k) - f(x_k + s_k)}{f(x_k) - \hat{q}_k(\hat{s}_k)}, \quad (5.3.17)$$

¹Note that Algorithm 6 does not use the ratio ρ_k in (5.3.17), but uses (5.2.6). This is because the denominator of (5.3.17) may be zero before termination, on account of sketching/subspace techniques being used.

with

$$|1 - \rho_k| = \frac{|f(x_k + s_k) - \hat{q}_k(\hat{s}_k)|}{|f(x_k) - \hat{q}_k(\hat{s}_k)|}.$$

The numerator can be bounded by

$$|f(x_k + s_k) - \hat{q}_k(\hat{s}_k)| \leq \frac{1}{6} L_H \|s_k\|_2^2,$$

by Corollary A.8.4 in [16]. Therefore, we have

$$|1 - \rho_k| \leq \frac{\frac{1}{6} L_H \|s_k\|_2^3}{\frac{1}{3\alpha_k} \|s_k\|_2^3} = \frac{1}{2} \alpha_k L_H \leq 1 - \theta \quad \text{by (5.3.16) and } \alpha_k \leq \alpha_{low}. \quad (5.3.18)$$

Thus $1 - \rho_k \leq |1 - \rho_k| \leq 1 - \theta$ so $\rho_k \geq \theta$ and iteration k is successful. \square

Satisfying Assumption 3 (page 71)

Lemma 5.3.7. *Let f be twice continuously differentiable with L_H -Lipschitz continuous Hessian. Algorithm 6 with true iterations defined in Definition 5.3.1 satisfies Assumption 3 with*

$$h(\epsilon, \alpha_k) = \frac{\theta}{3\alpha_{max}} \left(\frac{\epsilon}{2}\right)^{3/2} \min \left\{ \frac{2^{3/2}}{L_H^{3/2}}, \left(\frac{\sqrt{1 - \epsilon_S^{(2)}}}{\frac{1}{\alpha_k} S_{max} + \kappa_T} \right)^{3/2} \right\}. \quad (5.3.19)$$

Proof. For true and successful iterations with $k < N_\epsilon$, use Lemma 5.3.2 with Lemma 5.3.1 and $\alpha_k \leq \alpha_{max}$. \square

Satisfying Assumption 4 (page 72) The next lemma shows that the function value following Algorithm 6 is non-increasing.

Lemma 5.3.8. *Algorithm 6 satisfies Assumption 4.*

Proof. In Algorithm 6, we either have $x_{k+1} = x_k$ when the step is unsuccessful, in which case $f(x_k) = f(x_{k+1})$; or the step is successful, in which case we have $f(x_{k+1}) - f(x_k) \leq 0$ by Lemma 5.3.1. \square

5.3.3 Iteration complexity of Algorithm 6 to decrease $\nabla f(x_k)$ below ϵ

We have shown that Algorithm 6 satisfies Assumption 1, Assumption 2, Assumption 3 and Assumption 4. Noting that Algorithm 6 is a particular case of Algorithm 2, we apply Theorem 4.2.1 to arrive at the main result of this section.

Theorem 5.3.1. *Let \mathcal{S} be the distribution of scaled Gaussian matrices $S \in \mathbb{R}^{l \times d}$ defined in Definition 1.2.2. Suppose that f is bounded below by f^* , twice continuously differentiable with L_H -Lipschitz $\nabla^2 f$, $\nabla^2 f(x_k)$ has rank at most r for all k and let $\epsilon > 0$. Choose $l = 4C_l(\log 16 + r + 1)$; $\epsilon_S^{(2)} = \frac{1}{2}$; $\delta_S^{(2)} = \frac{1}{16}$; so that $\delta_S^{(3)} = e^{-\frac{\iota(\epsilon_S^{(2)})^2}{C_l} + r + 1} = \frac{1}{16}$; $\delta_S = \frac{1}{8}$; $S_{max} = 1 + \frac{\sqrt{d} + \sqrt{2 \log 16}}{\sqrt{4C_l(\log 16 + r + 1)}}$, where C_l is*

defined in (5.3.15). Run Algorithm 6 for N iterations. Suppose that $\delta_S < \frac{c}{(c+1)^2}$ (i.e. $\frac{c}{(c+1)^2} > \frac{1}{8}$). Then for any $\delta_1 \in (0, 1)$ with

$$g(\delta_1) > 0,$$

where

$$g(\delta_1) = \left[\frac{7}{8}(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1},$$

if $N \in \mathbb{N}$ satisfies

$$N \geq g(\delta_1) \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{4C_l(\log 16 + r + 1)}{1 + c} \right],$$

where $h(\epsilon, \alpha_k)$ is defined in (5.3.19) with $\epsilon_S^{(2)}$, S_{max} defined in the theorem statement, α_{low} is given in (5.3.16) and $\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}$ associated with α_{low} , for some $\tau_\alpha \in \mathbb{N}^+$. Then we have that

$$\mathbb{P} \left(\min_{k \leq N} \{ \|\nabla f(x_{k+1})\|_2 \} \leq \epsilon \right) \geq 1 - e^{-\frac{7\delta_1^2}{16}N}.$$

5.3.3.1 Discussion

Use other random ensembles than the scaled Gaussian matrices in Algorithm 6 Although Theorem 5.3.1 requires \mathcal{S} to be the distribution of scaled Gaussian matrices, qualitatively similar result, namely, convergence with a rate of $\mathcal{O}(\epsilon^{-3/2})$ with exponentially high probability, can be established for s -hashing matrices (defined in Definition 1.2.5), Subsampled Randomised Hadamard Transforms (defined in Definition 1.2.3) and Hashed Randomised Hadamard Transforms (defined in Definition 2.4.2). The proof for satisfying Assumption 1 needs to be modified, using the upper bounds for S_{max} and the subspace embedding properties of these ensembles instead. Consequently, the constants in Theorem 5.3.1 will change, but the convergence rate and the form of the result stays the same (as the results in Section 4.5.2).

Comparison with the adaptive cubic regularisation method with random models in [17]

We achieve the same $\mathcal{O}(\epsilon^{-3/2})$ convergence rate as [17], which is optimal for non-convex optimisations using second order models [16], and the same for deterministic adaptive cubic regularisation method. One main difference between our work and [17] is the definition of true iterations. Instead of Definition 5.3.1, they define true iterations as those iterations that satisfy

$$\|\nabla f(x_k) - \nabla m_k(s_k)\|_2 \leq \kappa_g \|s_k\|_2^2 \quad (5.3.20)$$

$$\|\nabla^2 f(x_k) - \nabla^2 m_k(s_k)\|_2 \leq \kappa_H \|s_k\|_2, \quad (5.3.21)$$

where $\kappa_g, \kappa_H > 0$ are constants.

This difference leads to different potential applications of the two frameworks. In their work, they proposed to use sampling with adaptive sample sizes for problems having the finite sum structure ($f = \sum_i f_i$) to construct the model m_k , or to use finite differences in the context of derivative free optimisation to construct the model m_k . However, without other assumptions, even just in

order to obtain condition (5.3.20), one may need a sample size that may be impractically large. In contrast, in our framework, the sketching size is fixed and even then, true iterations happen sufficiently frequently for scaled Gaussian matrices (and indeed for other random embeddings, see remark above). However, since the subspace dimension l is proportional to the rank of the Hessian matrix r , the Hessian matrix $\nabla^2 f$ is assumed to have a lower rank r than the full space dimension l , as otherwise Algorithm 6 does not save computation/gradient/Hessian evaluations compared to the deterministic version. Another difference is that our convergence result Theorem 5.3.1 is expressed in the high probability form, while the result in [17] is in expectation. Our result is stronger because it leads to an equivalent expectation result in [17], see Corollary 4.2.3 on Page 75.

Inexact local models constructed by subsampling for sums of functions have also been proposed for cubic regularization and other Newton-type methods in [59, 103, 104, 108]. Our emphasis here is related to reducing specifically the dimension of the variable domain (rather than the observational space).

5.4 Fast convergence rate assuming the sparsity of the Hessian matrix

This section is mostly conceptual and is an attempt to show the fast convergence rate of Algorithm 6 can be achieved without assuming subspace embedding of the Hessian matrix. Here, we maintain N_ϵ as $\min\{k : \|\nabla f(x_{k+1})\|_2 \leq \epsilon\}$, similarly to the last section. However, in the definition of true iterations, we replace the condition (5.3.1) on subspace embedding of the Hessian with the condition that the sketched Hessian $S_k \nabla^2 f(x_k)$ has a small norm. This may be achieved when the Hessian matrix has sparse rows and we choose S_k to be a scaled sampling matrix. We show that this new definition of true iterations still allows the same $\mathcal{O}(\epsilon^{-3/2})$ iteration complexity to drive the norm of the objective's gradient norm below ϵ . Specifically, true iterations are defined as follows.

Definition 5.4.1. Let $\epsilon_S \in (0, 1)$, $S_{max} > 0$. Iteration k is (ϵ_S, S_{max}) -true if

$$\|S_k \nabla^2 f(x_k)\|_2 \leq c_k \epsilon^{\frac{1}{2}}, \quad (5.4.1)$$

$$\|S_k \nabla f(x_k)\|_2^2 \geq (1 - \epsilon_S) \epsilon^2, \quad (5.4.2)$$

$$\|S_k\|_2 \leq S_{max}, \quad (5.4.3)$$

where $c_k = \sqrt{\frac{4(1-\epsilon_S)^{1/2} S_{max}}{3\alpha_{max}}}$ and α_{max} is a user-chosen constant in Algorithm 6.

Note that the desired accuracy ϵ appears in this particular definition of true iterations.

Consequently, the requirements on the objective and the sketching dimension l may be stronger for smaller ϵ . For simplicity, we assume $\kappa_T = 0$ (where κ_T is a user chosen parameter in (5.2.3) in Algorithm 6) in this section, namely $\nabla \hat{m}_k(\hat{s}_k) = 0$, and it follows from (5.2.3) that

$$S_k \nabla f(x_k) = \frac{1}{\alpha_k} S_k S_k^T \hat{s}_k \|S_k^T \hat{s}_k\|_2 - S_k \nabla^2 f(x_k) S_k^T \hat{s}_k. \quad (5.4.4)$$

The proofs that Assumption 2 and Assumption 4 are satisfied are identical to the previous section, while the following technical lemma helps us to satisfy Assumption 3.

Lemma 5.4.1. *Let $\epsilon > 0$. Let $\epsilon_S \in (0, 1)$, $\kappa_T = 0$. Suppose we have (5.4.1), (5.4.2) and (5.4.3). Then*

$$\|S_k^T \hat{s}_k\|_2 \geq \alpha_k \sqrt{\frac{(1 - \epsilon_S)^{1/2} \epsilon}{3S_{max}\alpha_{max}}}.$$

Proof. Let $b = \|S_k \nabla^2 f(x_k)\|_2$, $x = \|S_k^T \hat{s}_k\|_2$, then taking 2-norm of (5.4.4) with (5.4.2), $\|S_k\|_2 \leq S_{max}$ and the triangle inequality gives

$$\begin{aligned} (1 - \epsilon_S)^{\frac{1}{2}} \epsilon &\leq \|S_k \nabla f(x_k)\|_2 \leq \frac{S_{max}}{\alpha_k} x^2 + bx \\ \implies \frac{S_{max}}{\alpha_k} x^2 + bx - (1 - \epsilon_S)^{\frac{1}{2}} \epsilon &\geq 0 \\ \implies x^2 + \frac{\alpha_k b}{S_{max}} x - \frac{(1 - \epsilon_S)^{\frac{1}{2}} \epsilon \alpha_k}{S_{max}} &\geq 0 \\ \implies \left(x + \frac{\alpha_k b}{2S_{max}}\right)^2 &\geq \frac{(1 - \epsilon_S)^{\frac{1}{2}} \epsilon \alpha_k}{S_{max}} + \frac{\alpha_k^2 b^2}{4S_{max}^2} \\ \xRightarrow{x, b \geq 0} x &\geq \sqrt{\frac{(1 - \epsilon_S)^{\frac{1}{2}} \epsilon \alpha_k}{S_{max}} + \frac{\alpha_k^2 b^2}{4S_{max}^2}} - \frac{\alpha_k b}{2S_{max}}. \end{aligned}$$

Introduce $a = \frac{(1 - \epsilon_S)^{\frac{1}{2}} \epsilon \alpha_k}{S_{max}}$ and the function $y(b) = \frac{\alpha_k b}{2S_{max}}$, then the above gives

$$x \geq \sqrt{a + y(b)^2} - y(b). \quad (5.4.5)$$

We note, by taking derivative, that given $y(b) \geq 0$, the RHS of (5.4.5) is monotonically decreasing with $y(b)$. Therefore given $b \leq c_k \epsilon^{\frac{1}{2}}$ and thus $y(b) \leq y(c_k \epsilon^{\frac{1}{2}})$, we have

$$x \geq \sqrt{a + y(c_k \epsilon^{\frac{1}{2}})^2} - y(c_k \epsilon^{\frac{1}{2}}).$$

The choice of $c_k = \sqrt{\frac{4(1 - \epsilon_S)^{1/2} S_{max}}{3\alpha_{max}}} \leq \sqrt{\frac{4(1 - \epsilon_S)^{1/2} S_{max}}{3\alpha_k}}$ gives $a \geq 3y(c_k \epsilon^{\frac{1}{2}})^2$. And therefore we have $x \geq y(c_k \epsilon^{\frac{1}{2}})$. Noting that $x = \|S_k^T \hat{s}_k\|_2$ and substituting the expression for y and c_k gives the desired result. □

Lemma 5.4.2. *Following the framework of Algorithm 6 with $\kappa_T = 0$, let $\epsilon > 0$. Define true iterations as iterations that satisfy (5.4.1), (5.4.2), and (5.4.3). Then Assumption 3 is satisfied with*

$$h(\epsilon, \alpha_k) = \frac{\theta \alpha_k^2 \epsilon^{3/2}}{3} \left[\frac{(1 - \epsilon_S)^{1/2}}{3S_{max}\alpha_{max}} \right]^{3/2}. \quad (5.4.6)$$

Proof. A true and successful iteration k gives

$$f(x_k) - f(x_k + s_k) \geq \frac{\theta}{3\alpha_k} \|S_k^T \hat{s}_k\|_2^3$$

by Lemma 5.3.1 and combining with the conclusion of Lemma 5.4.1 gives the result. □

With Assumption 2, Assumption 3 and Assumption 4 satisfied, applying Theorem 4.2.1 gives the following result for Algorithm 6.

Theorem 5.4.1. *Let f be bounded below by f^* and twice continuously differentiable with L_H -Lipschitz continuous Hessian. Run Algorithm 6 for N iterations. Suppose Assumption 1 hold with $\delta_S \in (0, 1)$ and true iterations defined in Definition 5.4.1. Suppose $\delta_S < \frac{c}{(c+1)^2}$.*

Then for any $\delta_1 \in (0, 1)$ such that $g(\delta_S, \delta_1) > 0$ where $g(\delta_S, \delta_1)$ is defined in (4.2.12). If N satisfies

$$N \geq g(\delta_S, \delta_1) \left[\frac{f(x_0) - f^*}{h(\epsilon, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right], \quad (5.4.7)$$

where h is given in (5.4.6), α_{low} is given in (5.3.16) and $\alpha_{min}, \tau_\alpha$ are given in Lemma 4.2.1; then we have

$$\mathbb{P} \left(\min_{k \leq N} \{ \|\nabla f(x_{k+1})\|_2 \} \leq \epsilon \right) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N}.$$

Remark 13. *In order to satisfy Assumption 1, we require that at each iteration, with positive probability, (5.4.1), (5.4.2) and (5.4.3) hold. This maybe achieved for objective functions whose Hessian only has a few non-zero rows, with S being a scaled sampling matrix. Because if $\nabla^2 f(x_k)$ only has a few non-zero rows, we have that $S_k \nabla^2 f(x_k) = 0$ with positive probability, thus satisfying (5.4.1). Scaled sampling matrices also satisfy (5.4.2) and (5.4.3) (See Lemma 4.4.13 and Lemma 4.4.14).*

5.5 Convergence to second order (subspace) critical points

In this section, we show that Algorithm 6 converges to a (subspace) second order critical point of $f(x)$. Our convergence aim here is

$$N_\epsilon = N_{\epsilon_H}^{(2)} = \min \{k : \lambda_{min}(S_k \nabla^2 f(x_k) S_k^T) \geq -\epsilon_H\} \quad (5.5.1)$$

And we define S_{max} -true iterations as

Definition 5.5.1. *Let $S_{max} > 0$. An iteration k is true if $\|S_k\|_2 \leq S_{max}$.*

Compared to Section 5.3, here we have a less restrictive definition of true iterations. Consequently it is easy to show Assumption 1 is true.

Satisfying Assumption 1 For S being a scaled Gaussian matrix, Lemma 4.4.6 gives that Algorithm 6 satisfies Assumption 1 with with any $\delta_S \in (0, 1)$ and

$$S_{max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log \left(\frac{1}{\delta_S} \right)}{l}}.$$

Results for other random ensembles can be found in Section 4.4.2.

Satisfying Assumption 2

Lemma 5.5.1. *Let f be twice continuously differentiable with L_H -Lipschitz continuous Hessian. Algorithm 6 satisfies Assumption 2 with*

$$\alpha_{low} = \frac{2(1-\theta)}{L_H} \quad (5.5.2)$$

The proof is similar to Lemma 5.3.6, where the condition $\|S_k^T \hat{s}_k\|_2 > 0$ on true iterations before convergence is ensured by Lemma 5.5.2.

Satisfying Assumption 3 We can calculate

$$\nabla^2 \hat{m}_k(\hat{s}) = S_k \nabla^2 f(x_k) S_k^T + \frac{1}{\alpha_k} \left[\|S_k^T \hat{s}\|_2^{-1} (S_k S_k^T \hat{s}) (S_k S_k^T \hat{s})^T + \|S_k^T \hat{s}\|_2 S_k S_k^T \right]. \quad (5.5.3)$$

Therefore for any $y \in \mathbb{R}^l$, we have

$$y^T \nabla^2 \hat{m}_k(\hat{s}) y = y^T S_k \nabla^2 f(x_k) S_k^T y + \frac{1}{\alpha_k} \left[\|S_k^T \hat{s}\|_2^{-1} \left[(S_k S_k^T \hat{s})^T y \right]^2 + \|S_k^T \hat{s}\|_2 (S_k^T y)^2 \right]. \quad (5.5.4)$$

The following Lemma says that if the subspace Hessian has negative curvature, then the step size is bounded below by the size of the negative curvature. (But also depends on α_k .)

Lemma 5.5.2. *If $\lambda_{min}(S_k \nabla^2 f(x_k) S_k^T) < -\epsilon_H$; and $\|S_k\|_2 \leq S_{max}$, then*

$$\|S_k^T \hat{s}_k\|_2 \geq \epsilon_H \left[\frac{2S_{max}^2}{\alpha_k} + \kappa_S \right]^{-1}.$$

Proof. Let $y \in \mathbb{R}^l$. Using (5.2.4) and (5.5.4) we have that

$$y^T S_k \nabla^2 f(x_k) S_k^T y \geq -\frac{1}{\alpha_k} \left[\|S_k^T \hat{s}_k\|_2^{-1} \left[(S_k S_k^T \hat{s}_k)^T y \right]^2 + \|S_k^T \hat{s}_k\|_2 (S_k^T y)^2 \right] - \kappa_S \|S_k^T \hat{s}_k\|_2.$$

Given $\|S_k\|_2 \leq S_{max}$, we have that $S_k^T y \leq S_{max} \|y\|_2$. So we have

$$y^T S_k \nabla^2 f(x_k) S_k^T y \geq -\frac{1}{\alpha_k} \left[\|S_k^T \hat{s}_k\|_2^{-1} \left[(S_k S_k^T \hat{s}_k)^T y \right]^2 + \|S_k^T \hat{s}_k\|_2 S_{max}^2 \|y\|_2^2 \right] - \kappa_S \|S_k^T \hat{s}_k\|_2.$$

Minimising over $\|y\|_2 = 1$, noting that $\max_{\|y\|_2=1} \left((S_k S_k^T \hat{s}_k)^T y \right)^2 = \|S_k S_k^T \hat{s}_k\|_2^2$, we have

$$\begin{aligned} -\epsilon_H > \lambda_{min}(S_k \nabla^2 f(x_k) S_k^T) &\geq -\frac{1}{\alpha_k} \left[\|S_k^T \hat{s}_k\|_2^{-1} \|S_k S_k^T \hat{s}_k\|_2^2 + \|S_k^T \hat{s}_k\|_2 S_{max}^2 \right] - \kappa_S \|S_k^T \hat{s}_k\|_2 \\ &\geq -\frac{1}{\alpha_k} \left[\|S_k^T \hat{s}_k\|_2^{-1} S_{max}^2 \|S_k^T \hat{s}_k\|_2^2 + \|S_k^T \hat{s}_k\|_2 S_{max}^2 \right] - \kappa_S \|S_k^T \hat{s}_k\|_2 \\ &= -\frac{2S_{max}^2}{\alpha_k} \|S_k^T \hat{s}_k\|_2 - \kappa_S \|S_k^T \hat{s}_k\|_2. \end{aligned}$$

Rearranging gives the result. \square

Lemma 5.5.3. *Algorithm 6 satisfies Assumption 3 with*

$$h(\epsilon_H, \alpha_k) = \frac{\theta \epsilon_H^3}{3\alpha_k} \left[\frac{2S_{max}^2}{\alpha_k} + \kappa_S \right]^{-3} \quad (5.5.5)$$

Proof. Using Lemma 5.3.1, on successful iterations, we have $f(x_k) - f(x_k + s_k) \geq \frac{\theta}{3\alpha_k} \|S_k^T \hat{s}_k\|_2^3$. Consequently, $k \leq N_\epsilon$ (note the definition (5.5.1) of N_ϵ in this section) and Lemma 5.5.2 give the lower bound h that holds in true, successful and $k \leq N_\epsilon$ iterations. \square

Satisfying Assumption 4

Lemma 5.5.4. *Algorithm 6 satisfies Assumption 4.*

The proof of this lemma is identical to Lemma 5.3.8.

Convergence result Applying Theorem 4.2.1, we have a convergence result for Algorithm 6 to a point where the subspace Hessian has approximately non negative curvature. While the statement is for scaled Gaussian matrices, it is clear from the above proof that similar results apply to a wide range of sketching matrices.

Theorem 5.5.1. *Let $\epsilon_H > 0, l \in \mathbb{N}^+$. Let \mathcal{S} be the distribution of scaled Gaussian matrices $S \in \mathbb{R}^{l \times d}$. Suppose that f is bounded below by f^* , twice continuously differentiable with L_H -Lipschitz continuous Hessian $\nabla^2 f$. Choose $\delta_S = \frac{1}{16}$; so that $S_{max} = 1 + \frac{\sqrt{d} + \sqrt{2 \log 16}}{\sqrt{l}}$. Let h be defined in (5.5.5), α_{low} be given in (5.5.2) and $\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}$ associated with α_{low} , for some $\tau_\alpha \in \mathbb{N}^+$. Suppose that $\delta_S < \frac{c}{(c+1)^2}$.*

Then for any $\delta_1 \in (0, 1)$ with

$$g(\delta_S, \delta_1) > 0,$$

where

$$g(\delta_S, \delta_1) = \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1}$$

; if $N \in \mathbb{N}$ satisfies

$$N \geq g(\delta_S, \delta_1) \left[\frac{f(x_0) - f^*}{h(\epsilon_H, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

we have that

$$\mathbb{P}(\min\{k : \lambda_{min}(S_k^T \nabla^2 f(x_k) S_k) \geq -\epsilon_H\} \leq N) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N}.$$

Remark 14. *We see that the convergence rate to a (subspace) second order critical point is ϵ_H^{-3} .*

5.6 Convergence to second order (full space) critical points

In this section, we show that, if \mathcal{S} is the distribution of scaled Gaussian matrices, Algorithm 6 will converge to a (full-space) second order critical point, with a rate matching the classical full space algorithm.

We define

$$N_\epsilon = N_{\epsilon_H}^{(3)} = \min \{k : \lambda_{min}(\nabla^2 f(x_k)) \geq -\epsilon_H\} \quad (5.6.1)$$

The following definition of true iterations assumes that $\nabla^2 f(x_k)$ has rank $r \leq d$.

Definition 5.6.1. *Let $S_{max} > 0, \epsilon_S \in (0, 1)$. An iteration k is (ϵ_S, S_{max}) -true if the following two conditions hold*

$$1. \|S_k\|_2 \leq S_{max}.$$

2. There exists an eigen-decomposition of $\nabla^2 f(x_k) = \sum_{i=1}^r \lambda_i u_i u_i^T$ with $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r$, such that with $w_i = S_k u_i$,

$$1 - \epsilon_S \leq \|w_r\|_2^2 \leq 1 + \epsilon_S, \quad (5.6.2)$$

$$(w_i^T w_r)^2 \leq 16l^{-1}(1 + \epsilon_S) \quad \text{for all } i \neq r. \quad (5.6.3)$$

Note that since $S_k \in \mathbb{R}^{l \times d}$ is a (scaled) Gaussian matrix, and the set $\{u_i\}$ is orthonormal, we have that $\{w_i\}$ are independent Gaussian vectors, with entries being $N(0, l^{-1})$. (5.6.3) simply requires that those high-dimensional Gaussian vectors are approximately orthogonal, which is known to happen with high probability [100]. We proceed to show that the four assumptions needed for Theorem 4.2.1 hold, and then apply Theorem 4.2.1 for this particular definition of N_ϵ .

Satisfying Assumption 1 As before, we show each of the two conditions in Definition 5.6.1 hold with high probability, and then use the union bound (See the proof of Lemma 4.4.2) to show Assumption 1 is true. Note that the conditional independence between iterations is clear here because given x_k , whether the iteration is true or not only depends on the random matrix S_k and is independent of all the previous iterations.

For the first condition in Definition 5.6.1, We have that for S being a scaled Gaussian matrix, Lemma 4.4.6 gives that Algorithm 6 satisfies Assumption 1 with with any $\delta_S^{(2)} \in (0, 1)$ and

$$S_{max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log \left(1/\delta_S^{(2)} \right)}{l}}. \quad (5.6.4)$$

Lemma 5.6.1 shows that (5.6.2) holds with high probability.

Lemma 5.6.1. *Let $w_i \in \mathbb{R}^l$ with w_{ij} be independent $N(0, l^{-1})$. Let $\epsilon_S \in (0, 1)$. Then we have for some problem-independent constant C ,*

$$\mathbb{P} \left(\left| \|w_i\|_2^2 - 1 \right| \leq \epsilon_S \right) \geq 1 - 2e^{-\frac{l\epsilon_S^2}{C}}. \quad (5.6.5)$$

Proof. The proof is standard. One side of the bound is established in Lemma 4.4.4. Also see [25]. We note that $C \approx 4$. □

Next, we show that conditioning on (5.6.2) being true, (5.6.3) holds with high probability. We first study the case for a single fixed i instead of all i .

Lemma 5.6.2. *Let $\epsilon_S \in (0, 1)$ and suppose w_r satisfies (5.6.5). Then with (conditional) probability at least 0.9999, independent of w_r , we have $(w_i^T w_r)^2 \leq 16l^{-1}(1 + \epsilon_S)$.*

Proof. We have $(w_i^T w_r)^2 = \|w_r\|_2^2 \left(w_i^T \frac{w_r}{\|w_r\|_2} \right)^2$. The term inside the bracket is an $N(0, l^{-1})$ random variable independent of w_r , because sum of independent normal random variables is still normal. Note that for a normal random variable $N(0, \sigma^2)$, with probability at least 0.9999, its absolute value lies within $\pm 4\sigma$. Therefore we have that with probability at least 0.9999,

$$\left(w_i^T \frac{w_r}{\|w_r\|_2} \right)^2 \leq 16l^{-1}. \quad (5.6.6)$$

Combining with (5.6.5) gives the result. \square

Corollary 5.6.1 shows that conditioning on (5.6.2) being true, (5.6.3) is true with high probability.

Corollary 5.6.1. *With (conditional) probability at least $0.9999^{(r-1)}$, we have that $(w_i^T w_r)^2 \leq (1 + \epsilon_S)16l^{-1}$ for all $i \neq r$.*

Proof. Note that conditioning on $\|w_r\|_2^2$, $w_i^T w_r$ are independent events. Therefore we simply multiply the probability. \square

The following Lemma shows that the second condition in Definition 5.6.1 is true with high probability.

Lemma 5.6.3. *Let $\epsilon_S > 0$. Let $A_1 = \left\{ \left| \|w_r\|_2^2 - 1 \right| \leq \epsilon_S \right\}$, and $A_2 = \left\{ (w_i^T w_r)^2 \leq 16l^{-1}(1 + \epsilon_S) \right\}$, $\forall i \neq r$.*

Then with probability at least $(0.9999)^{r-1} \left(1 - 2e^{-\frac{l\epsilon_S^2}{C}} \right)$, we have that A_1 and A_2 hold simultaneously.

Proof. We have $\mathbb{P}(A_1 \cap A_2) = \mathbb{P}[A_2|A_1] \mathbb{P}(A_1)$. Using Lemma 5.6.1 and Corollary 5.6.1 gives the result. \square

Therefore, using (5.6.4), Lemma 5.6.3 and the union bound we have the following

Lemma 5.6.4. *Let $\epsilon_S > 0$, $l \in \mathbb{N}^+$, $\delta_S^{(2)} > 0$ such that*

$$\delta_S = (0.9999)^{r-1} \left(1 - 2e^{-\frac{l\epsilon_S^2}{C}} \right) + \delta_S^{(2)} < 1. \quad (5.6.7)$$

Then Algorithm 6 with (S_{max}, ϵ_S) -true iterations defined in Definition 5.6.1 satisfies Assumption 1 where $S_{max} = 1 + \sqrt{\frac{d}{l}} + \sqrt{\frac{2 \log(1/\delta_S^{(2)})}{l}}$.

Satisfying Assumption 2

Lemma 5.6.5. *Let f be twice continuously differentiable with L_H -Lipschitz continuous Hessian. Then Algorithm 6 with true iterations defined in Definition 5.6.1 satisfies Assumption 2 with*

$$\alpha_{low} = \frac{2(1 - \theta)}{L_H} \quad (5.6.8)$$

The proof is identical to Lemma 5.3.6. ²

²Except that we need $\|S_k^T \hat{s}_k\|_2 > 0$ in true iterations before convergence. But this is shown in (5.6.13).

Satisfying Assumption 3 Lemma 5.6.6 is a key ingredient. It shows that in true iterations, the subspace Hessian's negative curvature ($\lambda_{\min}(S_k \nabla^2 f(x_k) S_k^T)$) is proportional to the full Hessian's negative curvature ($\lambda_{\min}(\nabla^2 f(x_k))$).

Lemma 5.6.6. *Suppose iteration k is true with $\epsilon_S^{(1)} \in (0, 1)$ and $k < N_\epsilon$. Let $\kappa_H = \min\{0, \lambda_1/\lambda_r\}$. Suppose*

$$1 - \epsilon_S + 16 \frac{r-1}{l} \frac{1 + \epsilon_S}{1 - \epsilon_S} \frac{\lambda_1}{\lambda_r} \geq 0. \quad (5.6.9)$$

Then we have that

$$\lambda_{\min}(S_k \nabla^2 f(x_k) S_k^T) \leq -\epsilon_H m(\epsilon_S, r, l, \kappa_H),$$

where

$$m(\epsilon_S, r, l, \kappa_H) = \left(1 - \epsilon_S + 16 \frac{r-1}{l} \frac{1 + \epsilon_S}{1 - \epsilon_S} \kappa_H\right). \quad (5.6.10)$$

Proof. Using the eigen-decomposition of $\nabla^2 f(x_k)$, we have that $S_k \nabla^2 f(x_k) S_k^T = \sum_{i=1}^r \lambda_i w_i w_i^T$.

Use the Rayleigh quotient expression of minimal eigenvalue (with w_r being the trial vector):

$$\lambda_{\min}(S_k \nabla^2 f(x_k) S_k^T) \leq \frac{\sum_{i=1}^r \lambda_i (w_i^T w_r)^2}{w_r^T w_r}$$

We have

$$\begin{aligned} & \frac{\sum_{i=1}^r \lambda_i (w_i^T w_r)^2}{w_r^T w_r} \\ &= (w_r^T w_r) \lambda_r + \frac{\sum_{i=1}^{r-1} \lambda_i (w_i^T w_r)^2}{w_r^T w_r} \\ &\leq (1 - \epsilon_S) \lambda_r + \lambda_1 \frac{\sum_{i=1}^{r-1} (w_i^T w_r)^2}{w_r^T w_r} \\ &\leq (1 - \epsilon_S) \lambda_r + 16 \frac{r-1}{l} \frac{1 + \epsilon_S}{1 - \epsilon_S} \lambda_1, \end{aligned} \quad (5.6.11)$$

where the two inequalities follow from (5.6.2) and (5.6.3) because iteration k is true. Next we discuss two cases.

1. If $\lambda_1 < 0$, then $\kappa_H = 0$ because $\lambda_r < -\epsilon_H < 0$. Thus, $m(\epsilon_S, r, l, \kappa_H) = 1 - \epsilon_S$. The desired result follows from (5.6.11) by noting that the second term $16 \frac{r-1}{l} \frac{1 + \epsilon_S}{1 - \epsilon_S} \lambda_1 < 0$ and $\lambda_r < -\epsilon_H$.
2. If $\lambda_r \geq 0$, then $\kappa_H = \frac{\lambda_1}{\lambda_r}$ and from (5.6.11), we have

$$\begin{aligned} (1 - \epsilon_S) \lambda_r + 16 \frac{r-1}{l} \frac{1 + \epsilon_S}{1 - \epsilon_S} \lambda_1 &= \lambda_r \left(1 - \epsilon_S + 16 \frac{r-1}{l} \frac{1 + \epsilon_S}{1 - \epsilon_S} \frac{\lambda_1}{\lambda_r}\right) \\ &\leq -\epsilon_H m(\epsilon_S, r, l, \kappa_H), \end{aligned}$$

where we used (5.6.9) and $\lambda_r < -\epsilon_H$ to derive the inequality. And the desired result follows. \square

Remark 15. (5.6.9) always holds if $\lambda_1 \leq 0$ (where recall that λ_i are eigen-values of $\nabla^2 f(x_k)$ and $k < N_\epsilon$ implies $\lambda_r < -\epsilon_H < 0$). If $\lambda_1 > 0$, then (5.6.9) holds if we have $\kappa(\nabla^2 f(x_k)) \frac{r-1}{l} \leq \frac{(1 - \epsilon_S)^2}{16(1 + \epsilon_S)}$ where $\kappa(\nabla^2 f(x_k)) = |\frac{\lambda_1}{\lambda_r}|$ is the condition number of $\nabla^2 f(x_k)$.

We conclude that Assumption 3 is satisfied.

Lemma 5.6.7. *Algorithm 6 with \mathcal{S} being the distribution of scaled Gaussian matrices, true iteration defined in Definition 5.6.1 and N_ϵ defined in (5.6.1) satisfies Assumption 3 with*

$$h(\epsilon_H, \alpha_k) = \frac{\theta \epsilon_H^3 m(\epsilon_S, r, l, \kappa_H)^3}{3\alpha_k} \left[\frac{2S_{max}^2}{\alpha_k} + \kappa_S \right]^{-3}. \quad (5.6.12)$$

Proof. Let iteration k be true and successful with $k < N_\epsilon$. Lemma 5.6.6 gives that

$$\lambda_{min}(S_k \nabla^2 f(x_k) S_k^T) \leq -\epsilon_H m(\epsilon_S, r, l, \kappa_H).$$

Then we have

$$\|S_k^T \hat{s}_k\|_2 \geq \epsilon_H m(\epsilon_S, r, l, \kappa_H) \left[\frac{2S_{max}^2}{\alpha_k} + \kappa_S \right]^{-1}, \quad (5.6.13)$$

by applying Lemma 5.5.2 with $\epsilon_H = \epsilon_H m(\epsilon_S, r, l, \kappa_H)$. The desired result follows by applying Lemma 5.3.1. \square

Satisfying Assumption 4 The identical proof as the last section applies because Assumption 4 is not affected by the change of definitions of N_ϵ and true iterations.

Convergence of Algorithm 6 to a second order (full-space) critical point Applying Theorem 4.2.1, the next theorem shows that using Algorithm 6 with scaled Gaussian matrices achieves convergence to a second order critical point, with a rate matching the classical full-space method.

Theorem 5.6.1. *In Algorithm 6, let \mathcal{S} be the distribution of scaled Gaussian matrices. Let $\epsilon_H > 0$ and N_ϵ be defined in (5.6.1). Define true iterations in Definition 5.6.1. Suppose f is lower bounded by f^* and twice continuously differentiable with L_H -Lipschitz continuous Hessian $\nabla^2 f$.*

Choose $\delta_S^{(2)} = \frac{1}{16}$; so that $S_{max} = 1 + \frac{\sqrt{d} + \sqrt{2 \log 16}}{\sqrt{l}}$. Let δ_S be defined in (5.6.7). Let h be defined in (5.6.12), α_{low} be given in (5.6.8) and $\alpha_{min} = \alpha_0 \gamma_1^{\tau_\alpha}$ associated with α_{low} (See Lemma 4.2.1). Suppose that $\delta_S < \frac{c}{(c+1)^2}$. Then for any $\delta_1 \in (0, 1)$ with

$$g(\delta_S, \delta_1) > 0,$$

where

$$g(\delta_S, \delta_1) = \left[(1 - \delta_S)(1 - \delta_1) - 1 + \frac{c}{(c+1)^2} \right]^{-1};$$

if N satisfies

$$N \geq g(\delta_S, \delta_1) \left[\frac{f(x_0) - f^*}{h(\epsilon_H, \alpha_0 \gamma_1^{c+\tau_\alpha})} + \frac{\tau_\alpha}{1+c} \right],$$

we have that

$$\mathbb{P}(\min\{k : \lambda_{min}(\nabla^2 f(x_k)) \geq -\epsilon_H\} \leq N) \geq 1 - e^{-\frac{\delta_1^2}{2}(1-\delta_S)N}.$$

Proof. Applying Lemma 5.6.4, Lemma 5.6.5, Lemma 5.6.7, we have that the four assumptions in Theorem 4.2.1 are satisfied. Then applying Theorem 4.2.1 gives the desired result. \square

Chapter 6

Conclusion and future directions

In this thesis, we studied random embeddings and their application to optimisation problems and algorithms in order to achieve faster and more scalable solutions.

After introducing the necessary background related to random embeddings — Johnson-Lindenstrauss lemma, subspace and oblivious embeddings, and commonly used random ensembles — we analysed the subspace embedding property of hashing embeddings when the matrix whose column space is to be embedded has low coherence. We found that 1-hashing embeddings achieve the same theoretical dimensionality reduction property as the scaled Gaussian matrices if the coherence of the data is sufficiently low. This result motivated us to propose a new type of general random subspace embeddings – where the typically-used subsampling is replaced by hashing when combined with coherence-reducing transformations; this is the case of Subsampled- versus the novel Hashed-Randomised Hadamard Transform. Some open questions remain, that would be worthwhile pursuing and that would further enrich our understanding of this fascinating area of random matrices. For example, in our Theorem 2.3.1, we showed that 1-hashing matrices provide an oblivious subspace embedding of optimal size $m = \mathcal{O}(d)$ provided the input coherence is sufficiently low, of order $1/d$. Though the former, size requirement, cannot be improved in order, the latter, coherence one, probably can be improved to allow a larger class of input matrices to be embedded.

In chapter 3, we cascade our findings about sparse random matrices to the development of efficient solvers for large-scale linear least-squares problems, building on the success of the randomised Blendenpik algorithm and state-of-the-art numerical linear algebra techniques. We additionally present comprehensive benchmarking results of our proposed solver Ski-LLS against both random embedding-based, and deterministic, solvers. We found that our solver, SKi-LLS, which is available as an open source C++ code, outperforms not only sketching-based solvers but also state-of-the-art deterministic sparse solvers on certain subsets of the Florida collection of large-scale sparse matrices. Future development of our solver Ski-LLS may include incorporation and testing of other sparse ensembles such as the stable 1-hashing proposed in [19] (see also Section 4.4.2.3).

After considering reducing the dimensionality of the observational space in the linear least squares, we next turned to applying random embeddings to reduce the dimensionality of the variable/parameter space, leading to random subspace algorithmic variants of standard optimization algorithms for nonconvex problems. We showed that the $\mathcal{O}(\epsilon^{-2})$ convergence rate of first-order-type methods to obtain an approximately small gradient value, within ϵ , can be preserved, with high probability, when the gradient and the search direction are sketched/randomly projected. Various sketching matrices are allowed, of dimension independent of problem size, and can be used in a generic algorithmic framework that incorporates quadratic regularization and trust region variants. A current direction here is to particularise our general random subspace framework to linesearch methods, which in light of [17], is clearly possible, with similar complexity bounds being obtained.

When the second order information is also available, we investigated in Chapter 5, a Random subspace variant of Adaptive Cubic Regularization (R-ARC). We found that when the Hessian information is low rank, and Gaussian sketching is used to generate the subspace, the optimal complexity of order $\mathcal{O}(\epsilon^{-3/2})$ of the full-space algorithm is recovered with high probability, for generating a sufficiently small gradient. The complexity of achieving approximate second order criticality using R-ARC is also addressed, with similar outcomes in relation to the complexity of the full space variant provided again, that low-rank assumptions hold for the curvature information. Our focus in Chapters 4 and 5 was theoretical, but it has informed us about the potential and strength of fixed-size random projections to generating suitable subspaces for minimization, and strongly convergent ensuing algorithmic variants. Future work would be to numerically implement and test the more general variants (not just Gauss-Newton type) on large-scale general objectives, which would likely involve further, careful algorithm development.

Finally, we see potential in the techniques in this thesis to apply to other problem classes in numerical analysis, either more directly or with further development, such as to low rank matrix approximation, linear programming and nonlinear sum of functions arising in machine learning. Could we solve more large scale numerical analysis problems faster with random embeddings? Or perhaps we can find domain-tailored random embeddings for specific problems in machine learning, finance and other applications?

Bibliography

- [1] D. Achlioptas. Database-friendly random projections. In *Proceedings of the twentieth ACM SIGMOD-SIGACT-SIGART symposium on Principles of database systems*, pages 274–281, 2001.
- [2] N. Ailon and B. Chazelle. Approximate nearest neighbors and the fast Johnson-Lindenstrauss transform. In *STOC'06: Proceedings of the 38th Annual ACM Symposium on Theory of Computing*, pages 557–563. ACM, New York, 2006.
- [3] N. Ailon and E. Liberty. An almost optimal unrestricted fast Johnson-Lindenstrauss transform. *ACM Trans. Algorithms*, 9(3):Art. 21, 1–12, 2013.
- [4] E. Anderson, Z. Bai, C. Bischof, L. S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide*. Society for Industrial and Applied Mathematics, third edition, 1999.
- [5] H. Avron, P. Maymounkov, and S. Toledo. Blendenpik: supercharging Lapack's least-squares solver. *SIAM J. Sci. Comput.*, 32(3):1217–1236, 2010.
- [6] H. Avron, E. Ng, and S. Toledo. Using perturbed QR factorizations to solve linear least-squares problems. *SIAM J. Matrix Anal. Appl.*, 31(2):674–693, 2009.
- [7] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski. *Optimization with Sparsity-Inducing Penalties*, volume 4:1. Foundations and Trends in Machine Learning, 2011.
- [8] A. S. Berahas, R. Bollapragada, and J. Nocedal. An investigation of Newton-Sketch and subsampled Newton methods. *Optimization Methods and Software*, 2020.
- [9] A. Björck. *Numerical methods for least squares problems*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1996.
- [10] J. Bourgain, S. Dirksen, and J. Nelson. Toward a unified theory of sparse dimensionality reduction in Euclidean space. *Geom. Funct. Anal.*, 25(4):1009–1088, 2015.
- [11] Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford. Lower bounds for finding stationary points I. *Math. Program.*, 184(1-2, Ser. A):71–120, 2020.

- [12] C. Cartis, J. Fiala, and Z. Shao. Hashing embeddings of optimal dimension, with applications to linear least squares. *arXiv e-prints*, page arXiv:2105.11815, May 2021.
- [13] C. Cartis, J. Fiala, and Z. Shao. Randomised subspace methods for non-convex optimization, with applications to nonlinear least-squares. *arXiv e-prints, in preparation*, 2022.
- [14] C. Cartis, J. Fowkes, and Z. Shao. A randomised subspace gauss-newton method for nonlinear least-squares. In *Thirty-seventh International Conference on Machine Learning*, 2020. In Workshop on Beyond First Order Methods in ML Systems.
- [15] C. Cartis, N. I. Gould, and P. L. Toint. Adaptive cubic regularisation methods for unconstrained optimization. part i: motivation, convergence and numerical results. *Mathematical Programming*, 127(2):245–295, 2011.
- [16] C. Cartis, N. I. M. Gould, and P. L. Toint. *Evaluation complexity of algorithms for non-convex optimization*. MOS-SIAM series on Optimization. Society for Industrial and Applied Mathematics (SIAM), 2022.
- [17] C. Cartis and K. Scheinberg. Global convergence rate analysis of unconstrained optimization methods based on probabilistic models. *Mathematical Programming*, 169(2):337–375, 2018.
- [18] J. Cerdán, D. Guerrero, J. Marín, and J. Mas. Preconditioners for rank deficient least squares problems. *J. Comput. Appl. Math.*, 372:112621, 2020.
- [19] L. Chen, S. Zhou, and J. Ma. Stable sparse subspace embedding for dimensionality reduction. *Knowledge-Based Systems*, 195:105639, 2020.
- [20] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Ann. Math. Statistics*, 23:493–507, 1952.
- [21] K. L. Clarkson and D. P. Woodruff. Low-rank approximation and regression in input sparsity time. *J. ACM*, 63(6):Art. 54, 1–45, 2017.
- [22] M. B. Cohen. Nearly tight oblivious subspace embeddings by trace inequalities. In *Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 278–287. ACM, New York, 2016.
- [23] M. B. Cohen, T. S. Jayram, and J. Nelson. Simple analyses of the sparse Johnson-Lindenstrauss transform. In *1st Symposium on Simplicity in Algorithms*, volume 61 of *OASICs OpenAccess Ser. Inform.*, pages Art. No. 15, 9. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern, 2018.

- [24] Y. Dahiya, D. Konomis, and D. P. Woodruff. An empirical evaluation of sketching for numerical linear algebra. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '18, pages 1292–1300, New York, NY, USA, 2018. Association for Computing Machinery.
- [25] S. Dasgupta and A. Gupta. An elementary proof of a theorem of Johnson and Lindenstrauss. *Random Structures Algorithms*, 22(1):60–65, 2003.
- [26] K. R. Davidson and S. J. Szarek. Local operator theory, random matrices and Banach spaces. In *Handbook of the geometry of Banach spaces, Vol. I*, pages 317–366. North-Holland, Amsterdam, 2001.
- [27] T. A. Davis. *Direct methods for sparse linear systems*, volume 2 of *Fundamentals of Algorithms*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2006.
- [28] T. A. Davis. Algorithm 915, SuiteSparseQR: multifrontal multithreaded rank-revealing sparse QR factorization. *ACM Trans. Math. Software*, 38(1):Art. 1, 1–22, 2011.
- [29] T. A. Davis and Y. Hu. The University of Florida sparse matrix collection. *ACM Trans. Math. Software*, 38(1):Art. 1, 1–25, 2011.
- [30] E. D. Dolan and J. J. Moré. Benchmarking optimization software with performance profiles. *Mathematical programming*, 91(2):201–213, 2002.
- [31] D. L. Donoho and M. Elad. Optimally sparse representation in general (nonorthogonal) dictionaries via l^1 minimization. *Proc. Natl. Acad. Sci. USA*, 100(5):2197–2202, 2003.
- [32] P. Drineas, M. W. Mahoney, and S. Muthukrishnan. Sampling algorithms for l_2 regression and applications. In *Proceedings of the Seventeenth Annual ACM-SIAM Symposium on Discrete Algorithm*, SODA '06, page 1127–1136, USA, 2006. Society for Industrial and Applied Mathematics.
- [33] F. Facchinei, G. Scutari, and S. Sagratella. Parallel selective algorithms for nonconvex big data optimization. *IEEE Transactions on Signal Processing*, 63(7):1874–1889, 2015.
- [34] D. C.-L. Fong and M. Saunders. LSMR: An iterative algorithm for sparse least-squares problems. *SIAM J. Sci. Comput.*, 33(5):2950–2971, 2011.
- [35] C. Freksen, L. Kamma, and K. G. Larsen. Fully understanding the hashing trick. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS'18, pages 5394–5404, Red Hook, NY, USA, 2018. Curran Associates Inc.

- [36] A. Gnanasekaran and E. Darve. Hierarchical Orthogonal Factorization: Sparse Least Squares Problems. *arXiv e-prints*, page arXiv:2102.09878, Feb. 2021.
- [37] G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, third edition, 1996.
- [38] N. Gould and J. Scott. The state-of-the-art of preconditioners for sparse linear least-squares problems: the complete results. Technical report, STFC Rutherford Appleton Laboratory, 2015. Available at ftp://cuter.rl.ac.uk/pub/nimg/pubs/GoulScot16b_toms.pdf.
- [39] N. Gould and J. Scott. The state-of-the-art of preconditioners for sparse linear least-square problems. *ACM Trans. Math. Software*, 43(4):Art. 36, 1–35, 2017.
- [40] N. I. Gould, D. Orban, and P. L. Toint. CUTEst: a constrained and unconstrained testing environment with safe threads for mathematical optimization. *Computational Optimization and Applications*, 60(3):545–557, 2015.
- [41] R. Gower, D. Goldfarb, and P. Richtárik. Stochastic block BFGS: Squeezing more curvature out of data. In M. F. Balcan and K. Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1869–1878, New York, 2016. PMLR.
- [42] R. Gower, D. Koralev, F. Lieder, and P. Richtárik. RSN: Randomized subspace Newton. In *Advances in Neural Information Processing Systems*, pages 614–623, 2019.
- [43] R. M. Gower and P. Richtárik. Randomized iterative methods for linear systems. *SIAM J. Matrix Anal. Appl.*, 36(4):1660–1690, 2015.
- [44] R. M. Gower, P. Richtárik, and F. Bach. Stochastic quasi-gradient methods: variance reduction via Jacobian sketching. *Mathematical Programming*, 2020.
- [45] S. Gratton, C. W. Royer, L. N. Vicente, and Z. Zhang. Complexity and global rates of trust-region methods based on probabilistic models. *IMA Journal of Numerical Analysis*, 38(3):1579–1597, 2018.
- [46] S. Gratton, A. Sartenaer, and P. L. Toint. Recursive trust-region methods for multiscale nonlinear optimization. *SIAM Journal on Optimization*, 19(1):414–444, 2008.
- [47] A. Griewank. The modification of newton’s method for unconstrained optimization by bounding cubic terms. Technical report, Technical report NA/12, 1981.
- [48] D. Grishchenko, F. Iutzeler, and J. Malick. Proximal gradient methods with adaptive subspace sampling. *Mathematics of Operations Research*, 2021.

- [49] G. W. Howell and M. Baboulin. Iterative solution of sparse linear least squares using lu factorization. In *Proceedings of the International Conference on High Performance Computing in Asia-Pacific Region*, HPC Asia 2018, pages 47–53, New York, NY, USA, 2018. Association for Computing Machinery.
- [50] P. Indyk and R. Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *STOC '98 (Dallas, TX)*, pages 604–613. ACM, New York, 1999.
- [51] M. A. Iwen, D. Needell, E. Rebrova, and A. Zare. Lower Memory Oblivious (Tensor) Subspace Embeddings with Fewer Random Bits: Modewise Methods for Least Squares. *SIAM J. Matrix Anal. Appl.*, 42(1):376–416, 2021.
- [52] C. Iyer, H. Avron, G. Kollias, Y. Ineichen, C. Carothers, and P. Drineas. A randomized least squares solver for terabyte-sized dense overdetermined systems. *J. Comput. Sci.*, 36:100547, 2019.
- [53] C. Iyer, C. Carothers, and P. Drineas. Randomized sketching for large-scale sparse ridge regression problems. In *2016 7th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems (ScalA)*, pages 65–72, 2016.
- [54] M. Jagadeesan. Understanding sparse JL for feature hashing. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [55] C. Jin, R. Ge, P. Netrapalli, S. M. Kakade, and M. I. Jordan. How to escape saddle points efficiently. In *International Conference on Machine Learning*, pages 1724–1732. PMLR, 2017.
- [56] W. B. Johnson and J. Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. In *Conference in modern analysis and probability (New Haven, Conn., 1982)*, volume 26 of *Contemp. Math.*, pages 189–206. Amer. Math. Soc., Providence, RI, 1984.
- [57] N. Kahale. Least-squares regressions via randomized Hessians. *arXiv e-prints*, page arXiv:2006.01017, June 2020.
- [58] D. M. Kane and J. Nelson. Sparser Johnson-Lindenstrauss transforms. *J. ACM*, 61(1):Art. 4, 23, 2014.
- [59] J. M. Kohler and A. Lucchi. Sub-sampled cubic regularization for non-convex optimization. *arXiv e-prints*, May 2017.
- [60] D. Kozak, S. Becker, A. Doostan, and L. Tenorio. Stochastic subspace descent. *arXiv preprint arXiv:1904.01145*, 2019.

- [61] D. Kozak, S. Becker, A. Doostan, and L. Tenorio. A stochastic subspace approach to gradient-free optimization in high dimensions. *Computational Optimization and Applications*, 79(2):339–368, June 2021.
- [62] J. Lacotte and M. Pilanci. Faster Least Squares Optimization. *arXiv e-prints*, page arXiv:1911.02675, Nov. 2019.
- [63] J. Lacotte and M. Pilanci. Effective dimension adaptive sketching methods for faster regularized least-squares optimization. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 19377–19387. Curran Associates, Inc., 2020.
- [64] J. Lacotte and M. Pilanci. Optimal Randomized First-Order Methods for Least-Squares Problems. *arXiv e-prints*, page arXiv:2002.09488, Feb. 2020.
- [65] J. Lacotte, M. Pilanci, and M. Pavone. High-dimensional optimization in adaptive random subspaces. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [66] J. D. Lee, M. Simchowitz, M. I. Jordan, and B. Recht. Gradient descent only converges to minimizers. In *Conference on learning theory*, pages 1246–1257. PMLR, 2016.
- [67] S. Liu, T. Liu, A. Vakilian, Y. Wan, and D. P. Woodruff. Extending and Improving Learned CountSketch. *arXiv e-prints*, page arXiv:2007.09890, July 2020.
- [68] M. Locatelli and F. Schoen. *Global Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2013.
- [69] N. Loizou and P. Richtárik. Momentum and stochastic momentum for stochastic gradient, Newton, proximal point and subspace descent methods. *Comput. Optim. Appl.*, 77(3):653–710, 2020.
- [70] M. Lopes, S. Wang, and M. Mahoney. Error estimation for randomized least-squares algorithms via the bootstrap. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3217–3226. PMLR, 10–15 Jul 2018.
- [71] Z. Lu and L. Xiao. A randomized nonmonotone block proximal gradient method for a class of structured nonlinear programming. *SIAM Journal on Numerical Analysis*, 55(6):2930–2955, 2017.

- [72] H. Luo, A. Agarwal, N. Cesa-Bianchi, and J. Langford. Efficient second order online learning by sketching. In *Advances in Neural Information Processing Systems*, pages 902–910, 2016.
- [73] M. W. Mahoney. Randomized algorithms for matrices and data. *Found. Trends Mach. Learn.*, 3(2):123–224, Feb. 2011.
- [74] M. W. Mahoney, J. C. Duchi, and A. C. Gilbert, editors. *The mathematics of data*, volume 25 of *IAS/Park City Mathematics Series*. American Mathematical Society, Providence, RI; Institute for Advanced Study (IAS), Princeton, NJ, 2018. Papers based on the lectures presented at the 26th Annual Park City Mathematics Institute Summer Session, July 2016.
- [75] P.-G. Martinsson. Blocked rank-revealing QR factorizations: How randomized sampling can be used to avoid single-vector pivoting. *arXiv e-prints*, page arXiv:1505.08115, May 2015.
- [76] P.-G. Martinsson, G. Quintana Ortí, N. Heavner, and R. van de Geijn. Householder QR factorization with randomization for column pivoting (HQRFP). *SIAM J. Sci. Comput.*, 39(2):C96–C115, 2017.
- [77] X. Meng and M. W. Mahoney. Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression. In *STOC’13—Proceedings of the 2013 ACM Symposium on Theory of Computing*, pages 91–100. ACM, New York, 2013.
- [78] X. Meng, M. A. Saunders, and M. W. Mahoney. LSRN: a parallel iterative solver for strongly over- or underdetermined systems. *SIAM J. Sci. Comput.*, 36(2):C95–C118, 2014.
- [79] J. Nelson and H. L. Nguyen. OSNAP: faster numerical linear algebra algorithms via sparser subspace embeddings. In *2013 IEEE 54th Annual Symposium on Foundations of Computer Science—FOCS 2013*, pages 117–126. IEEE Computer Soc., Los Alamitos, CA, 2013.
- [80] J. Nelson and H. L. Nguyen. Sparsity lower bounds for dimensionality reducing maps. In *STOC’13—Proceedings of the 2013 ACM Symposium on Theory of Computing*, pages 101–110. ACM, New York, 2013.
- [81] J. Nelson and H. L. Nguyen. Lower bounds for oblivious subspace embeddings. In *Automata, languages, and programming. Part I*, volume 8572 of *Lecture Notes in Comput. Sci.*, pages 883–894. Springer, Heidelberg, 2014.
- [82] Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM J. Optim.*, 22(2):341–362, 2012.
- [83] Y. Nesterov. *Lectures on convex optimization*, volume 137 of *Springer Optimization and Its Applications*. Springer, Cham, 2018. Second edition of [MR2142598].

- [84] Y. Nesterov and B. T. Polyak. Cubic regularization of Newton method and its global performance. *Math. Program.*, 108(1, Ser. A):177–205, 2006.
- [85] J. Nocedal and S. J. Wright. *Numerical optimization*. Springer Series in Operations Research and Financial Engineering. Springer, New York, second edition, 2006.
- [86] C. C. Paige and M. A. Saunders. LSQR: an algorithm for sparse linear equations and sparse least squares. *ACM Trans. Math. Software*, 8(1):43–71, 1982.
- [87] A. Patrascu and I. Necoara. Efficient random coordinate descent algorithms for large-scale structured nonconvex optimization. *Journal of Global Optimization*, 61(1):19–46, 2015.
- [88] M. Pilanci and M. J. Wainwright. Newton sketch: A near linear-time optimization algorithm with linear-quadratic convergence. *SIAM Journal on Optimization*, 27(1):205–245, 2017.
- [89] M. J. D. Powell. On search directions for minimization algorithms. *Math. Programming*, 4:193–201, 1973.
- [90] P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, 156:433–484, 2015.
- [91] P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Math. Program.*, 144(1-2, Ser. A):1–38, 2014.
- [92] V. Rokhlin and M. Tygert. A fast randomized algorithm for overdetermined linear least-squares regression. *Proc. Natl. Acad. Sci. USA*, 105(36):13212–13217, 2008.
- [93] T. Sarlos. Improved approximation algorithms for large matrices via random projections. In *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS’06)*, pages 143–152, 2006.
- [94] J. Scott and M. Tuma. HSL_MI28: An efficient and robust limited-memory incomplete cholesky factorization code. *ACM Trans. Math. Softw.*, 40(4):Art. 36, 1–35, July 2014.
- [95] Z. Shao, C. Cartis, and F. Jan. A randomised subspace gauss-newton method for nonlinear least-squares. In *Thirty-seventh International Conference on Machine Learning*, 2020. In Workshop on Beyond First Order Methods in ML Systems.
- [96] S. F. Tett, K. Yamazaki, M. J. Mineter, C. Cartis, and N. Eizenberg. Calibrating climate models using inverse methods: case studies with HadAM3, HadAM3P and HadCM3. *Geoscientific Model Development*, 10:3567–3589, 2017.
- [97] J. A. Tropp. Improved analysis of the subsampled randomized Hadamard transform. *Adv. Adapt. Data Anal.*, 3(1-2):115–126, 2011.

- [98] J. A. Tropp. User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.*, 12(4):389–434, 2012.
- [99] J. Vanschoren, J. N. van Rijn, B. Bischl, and L. Torgo. OpenML: networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.
- [100] R. Vershynin. *High-dimensional probability*, volume 47 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge, 2018. An introduction with applications in data science, With a foreword by Sara van de Geer.
- [101] D. P. Woodruff. Sketching as a tool for numerical linear algebra. *Found. Trends Theor. Comput. Sci.*, 10(1-2):1–157, 2014.
- [102] S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151:3–34, 2015.
- [103] P. Xu, F. Roosta-Khorasan, and M. W. Mahoney. Newton-type methods for non-convex optimization under inexact Hessian information. *arXiv e-prints*, Aug 2017.
- [104] P. Xu, F. Roosta-Khorasan, and M. W. Mahoney. Second-order optimization for non-convex machine learning: An empirical study. *arXiv e-prints*, Aug 2017.
- [105] Y. Xu and W. Yin. Block stochastic gradient iteration for convex and nonconvex optimization. *SIAM Journal on Optimization*, 25(3):1686–1716, 2015.
- [106] Y. Xu and W. Yin. A globally convergent algorithm for nonconvex optimization based on block coordinate update. *Journal of Scientific Computing*, 72(2):700–734, 2017.
- [107] Y. Yang, M. Pesavento, Z.-Q. Luo, and B. Ottersten. Inexact block coordinate descent algorithms for nonsmooth nonconvex optimization. *IEEE Transactions on Signal Processing*, 68:947–961, 2020.
- [108] Z. Yao, P. Xu, F. Roosta-Khorasan, and M. W. Mahoney. Inexact non-convex Newton-type methods. *arXiv e-prints*, Feb 2018.
- [109] Z. Fu. Package snobfit. <http://reflectometry.org/danse/docs/snobfit>, 2009.
- [110] R. Zhu. Gradient-based sampling: An adaptive importance sampling for least-squares. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS’16, pages 406–414, Red Hook, NY, USA, 2016. Curran Associates Inc.