

# Statistical properties of sketching algorithms

Daniel Ahfock<sup>1,2</sup>, William J. Astle<sup>1</sup>, and Sylvia Richardson<sup>1</sup>

<sup>1</sup>*MRC Biostatistics Unit, University of Cambridge*

<sup>2</sup>*School of Mathematics and Physics, University of Queensland*

## Abstract

Sketching is a probabilistic data compression technique that has been largely developed in the computer science community. Numerical operations on big datasets can be intolerably slow; sketching algorithms address this issue by generating a smaller surrogate dataset. Typically, inference proceeds on the compressed dataset. Sketching algorithms generally use random projections to compress the original dataset and this stochastic generation process makes them amenable to statistical analysis. We argue that the sketched data can be modelled as a random sample, thus placing this family of data compression methods firmly within an inferential framework. In particular, we focus on the Gaussian, Hadamard and Clarkson-Woodruff sketches, and their use in single pass sketching algorithms for linear regression with huge  $n$ . We explore the statistical properties of sketched regression algorithms and derive new distributional results for a large class of sketched estimators. A key result is a conditional central limit theorem for data oblivious sketches. An important finding is that the best choice of sketching algorithm in terms of mean square error is related to the signal to noise ratio in the source dataset. Finally, we demonstrate the theory and the limits of its applicability on two real datasets.

## 1 Introduction

Sketching is a general probabilistic data compression technique designed for Big Data applications (Cormode, 2011). Even routine calculations can be prohibitively computationally expensive on massive datasets. Computation time can be reduced to an acceptable level by allowing for some approximation error in the results. Sketching algorithms relax the computational task by generating a compressed version of the original dataset which then serves as a surrogate for calculations. The compressed dataset is referred to as a sketch, as it acts as a compact representation of the full dataset. Sketching algorithms use a randomised compression stage which makes them interesting from a statistical viewpoint. Sketching algorithms for linear regression have attracted significant attention in the numerical linear algebra and theoretical computer science communities (Woodruff, 2014; Mahoney, 2011). In this paper we investigate the statistical properties of sketched regression algorithms, a perspective which has received little attention up to now.

To describe sketched regression in more detail, we first assume the data consists of a  $n$ -length response vector  $\mathbf{y}$  and a  $n \times p$  matrix of covariates,  $\mathbf{X}$  which is of full rank. It is assumed throughout that  $n > p$ . The objective is to find the optimal least squares coefficients. Given sufficient

computational resources, these could be computed exactly as

$$\beta_F = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

The subscript  $F$  is used to indicate the connection to the full dataset. Only two quantities are needed in order to determine  $\beta_F$ , the Gram matrix  $\mathbf{X}^\top \mathbf{X}$ , and the marginal associations  $\mathbf{X}^\top \mathbf{y}$ . Calculation of  $\mathbf{X}^\top \mathbf{X}$  requires  $O(np^2)$  operations while computation of  $\mathbf{X}^\top \mathbf{y}$  needs only  $O(np)$  calculations. There are two broad methods for sketched regression, complete sketching and partial sketching. Complete sketching is based on approximating both  $\mathbf{X}^\top \mathbf{X}$  and  $\mathbf{X}^\top \mathbf{y}$ , whereas partial sketching only approximates the Gram matrix. Drineas et al. (2006) establish many important results for complete sketching, and Dhillon et al. (2013) and Pilanci and Wainwright (2016) give foundational results for partial sketching.

Sketching algorithms use random linear mappings to reduce the size of the dataset from  $n$  to  $k$  observations. The random linear mapping can be represented as a  $k \times n$  sketching matrix  $\mathbf{S}$ . Complete sketching generates a  $k$ -length sketched response vector  $\tilde{\mathbf{y}}$  and a  $k \times p$  matrix of sketched predictors  $\tilde{\mathbf{X}}$ . The sketched data are computed through the linear mappings  $\tilde{\mathbf{y}} = \mathbf{S}\mathbf{y}$  and  $\tilde{\mathbf{X}} = \mathbf{S}\mathbf{X}$ . Partial sketching only generates a  $k \times p$  matrix of sketched covariates  $\tilde{\mathbf{X}}$ . We again use the random mapping  $\tilde{\mathbf{X}} = \mathbf{S}\mathbf{X}$ .

The complete sketching estimator,  $\beta_S$ , is defined as the least squares coefficients using the sketched responses and predictors,

$$\beta_S = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{y}}. \tag{1}$$

The partial sketching estimator,  $\beta_P$ , is defined as

$$\beta_P = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{y}. \tag{2}$$

The key difference between (1) and (2) is that the partial sketched estimator  $\beta_P$  is constructed using the exact marginal associations  $\mathbf{X}^\top \mathbf{y}$ . Given the sketched data, computation of  $\beta_S$  or  $\beta_P$  requires only  $O(kp^2)$  operations, compared with the  $O(np^2)$  required for  $\beta_F$ .

The estimand within a sketching algorithm is the optimal coefficient vector  $\beta_F$ . Sketching algorithms have the property that given a fixed  $k$ , the approximation error  $\|\beta_S - \beta_F\|_2$  or  $\|\beta_P - \beta_F\|_2$  remains probabilistically bounded even as  $n \rightarrow \infty$ . Designing estimators for approximate computation with such properties is very difficult, and is a common goal in the development of techniques for Big Data (Bardenet and Maillard, 2015; Phillips, 2016). The favourable scaling properties of sketching algorithms are a critical factor in making them stand apart from simple subsampling approaches, where it can be difficult to establish universal worst case bounds for large  $n$  (Drineas et al., 2006; Ma et al., 2015). The fact that sketching algorithms provide finite  $k$  guarantees for arbitrarily large  $n$  is a major reason they have received so much attention in the computer science community.

There is a large literature concerned with designing appropriate distributions for the random sketching matrix  $\mathbf{S}$ . Our focus is on data-oblivious random projections, where the distribution of the sketching matrix is not a function of the source data  $[\mathbf{y}, \mathbf{X}]$ . An example is the Gaussian sketch, where each element is independently distributed as a  $N(0, 1/k)$  variate. We also consider

the Hadamard sketch and the Clarkson-Woodruff sketch, random projections that exploit structure and sparsity for computational efficiency.

Most existing results on the accuracy of sketching are universal worst case bounds (Woodruff, 2014; Mahoney and Drineas, 2016). This is typical for randomised algorithms, however a more detailed error analysis can provide important insights (Halko et al., 2011). We investigate the statistical properties of  $\beta_P$  and  $\beta_S$  when using data oblivious sketches. An important finding is that the signal to noise ratio in the source dataset strongly influences the relative efficiency of complete to partial sketching. The statistical analysis also allows the construction of exact confidence intervals for the Gaussian sketch, and asymptotic confidence intervals for other random projections, paving the way for their wider use in the statistical community interested in Big Data methods.

We start by reviewing the existing literature on sketching algorithms before investigating the statistical properties in more detail. At its core, sketched regression is a randomised algorithm for approximate computation of  $\beta_F$ . Repeated application of the sketching algorithm on the same dataset will produce different results. The first stage in our analysis is to establish the distributional properties of the sketched estimators with the source dataset fixed. This gives a clear statistical picture of the behaviour of the randomised algorithm. An important result is a conditional central limit theorem for the sketched dataset that connects the Hadamard and Clarkson-Woodruff projections to the Gaussian sketch. The regularity conditions have a intuitive interpretation in terms of the geometry of the source dataset. Our conditional analysis of the randomised algorithms is then extended to cover situations where sketching is used for approximate statistical inference. Given a statistical model for the response  $\mathbf{y} = \mathbf{X}\beta_0 + \epsilon$ , for a vector of population parameters  $\beta_0$ , and error terms  $\epsilon$ , we can determine properties of  $\beta_P$  and  $\beta_S$  by integrating over the conditional distributions of the sketched estimators that take  $\mathbf{y}$  as fixed.

## 2 Background and related work

### 2.1 Preliminaries

Before proceeding, it is worth mentioning alternatives to sketching, in particular iterative methods for calculating the least squares coefficients  $\beta_F$ . These include coordinate descent or stochastic gradient methods. Iterative methods are guaranteed to converge to  $\beta_F$  under very mild conditions. These iterative techniques assume that the entire dataset can be stored in memory in a single location, or require regular communication if the full dataset is distributed across multiple sites. Sketching algorithms are not burdened by these memory and communication costs, with the drawback of no convergence guarantees to  $\beta_F$ . Connections to iterative methods are postponed until the discussion, the focus for now is on the single pass estimators  $\beta_S$  and  $\beta_P$ .

The purpose of this section is to review the existing theoretical framework for sketching algorithms. Sketching algorithms are largely motivated through worst case guarantees. We recap how these bounds can be developed before studying the statistical properties of the sketched estimators.

It will be helpful to define a number of quantities related to the full dataset before moving on. Let  $TSS_F = \mathbf{y}^\top \mathbf{y}$ ,  $RSS_F = \|\mathbf{y} - \mathbf{X}\beta_F\|_2^2$ ,  $MSS_F = \|\mathbf{X}\beta_F\|_2^2$  and  $R_F^2 = MSS_F/TSS_F$ . These terms summarise the goodness of fit of the model. The total, residual and model sum of squares are given by  $TSS_F$ ,  $RSS_F$  and  $MSS_F$  respectively, with  $TSS_F = MSS_F + RSS_F$ . The proportion of

variance explained by the model is given by  $R_F^2$ . These values will be important in characterising the behaviour of  $\beta_S$  and  $\beta_P$ .

## 2.2 Worst case bounds

A key concept in the construction of sketching algorithms is the notion of an  $\epsilon$ -subspace embedding (Woodruff, 2014; Meng and Mahoney, 2013; Yang et al., 2015a).

**Definition 1.**  *$\epsilon$ -subspace embedding.*

For a given  $n \times d$  matrix  $\mathbf{A}$ , we call a  $k \times n$  matrix  $\mathbf{S}$  an  $\epsilon$ -subspace embedding for  $\mathbf{A}$ , if for all vectors  $\mathbf{z} \in \mathbb{R}^d$

$$(1 - \epsilon)\|\mathbf{Az}\|_2^2 \leq \|\mathbf{SAz}\|_2^2 \leq (1 + \epsilon)\|\mathbf{Az}\|_2^2.$$

Speaking broadly, an  $\epsilon$ -subspace preserves the linear structure of the columns of the original dataset up to some multiplicative  $(1 \pm \epsilon)$  factor. In particular, if  $\epsilon$  is small, the linear mapping  $\mathbf{S}$  approximately preserves the covariance structure of the source dataset. Most theoretical arguments for sketching algorithms are predicated on the idea that the sketching matrix  $\mathbf{S}$  is an  $\epsilon$ -subspace embedding for the source dataset. The general notion is that it is possible to use a linear mapping  $\mathbf{S}$  that reduces the sample size from  $n$  to  $k$  whilst preserving much of the linear information in the full dataset.

The issue of how to generate  $\epsilon$ -subspace embeddings is deferred until section 2.3, the present focus will be on the utility of  $\epsilon$ -subspace embeddings for linear regression problems. For now, assume that we have some method for generating  $\epsilon$ -subspace embeddings for the source data matrix  $\mathbf{A}$ . It will be convenient to refer to  $\tilde{\mathbf{A}} = \mathbf{SA}$  as an  $\epsilon$ -subspace embedding of  $\mathbf{A}$  if  $\mathbf{S}$  is an  $\epsilon$ -subspace embedding for  $\mathbf{A}$ . As regression is the focus from this point forward, we will define the source data matrix as  $\mathbf{A} = [\mathbf{y}, \mathbf{X}]$ , the sketched data matrix as  $\tilde{\mathbf{A}} = [\tilde{\mathbf{y}}, \tilde{\mathbf{X}}]$  and set  $d = p + 1$ .

The complete sketched estimator  $\beta_S$  is given by the least squares coefficients using the sketched responses  $\tilde{\mathbf{y}}$  and the sketched predictors  $\tilde{\mathbf{X}}$ ,

$$\beta_S = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta\|_2^2.$$

An  $\epsilon$ -subspace embedding is useful as it relates the sketched optimisation problem to the full dataset optimisation problem. If  $\tilde{\mathbf{A}} = [\tilde{\mathbf{y}}, \tilde{\mathbf{X}}]$  is an  $\epsilon$ -subspace embedding of  $\mathbf{A} = [\mathbf{y}, \mathbf{X}]$ , it must hold that for all  $\beta \in \mathbb{R}^p$ ,

$$(1 - \epsilon)\|\mathbf{y} - \mathbf{X}\beta\|_2^2 \leq \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta\|_2^2 \leq (1 + \epsilon)\|\mathbf{y} - \mathbf{X}\beta\|_2^2.$$

If  $\epsilon$  is small, minimising the sum of squared residuals on the sketched dataset is similar to minimising the sum of squared residuals on the full dataset. If this is the case, it can be expected that  $\beta_S$  will be close to  $\beta_F$ . It is possible to establish the concrete bounds, that if  $\tilde{\mathbf{A}}$  is an  $\epsilon$ -subspace embedding

of  $\tilde{\mathbf{A}}$  (Sarlos, 2006),

$$\|\beta_S - \beta_F\|_2^2 \leq \frac{\epsilon^2}{\sigma_{\min}^2(\mathbf{X})} RSS_F, \quad (3)$$

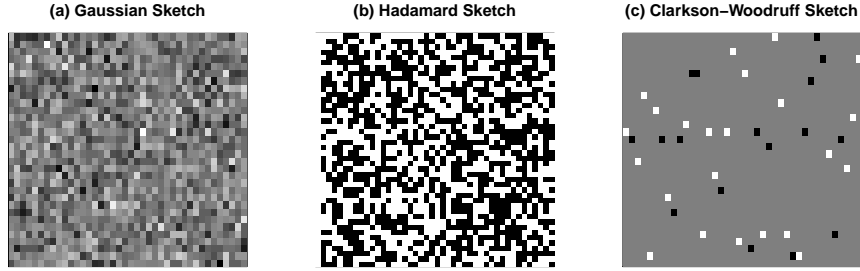
where  $\sigma_{\min}(\mathbf{X})$  represents the smallest singular value of the design matrix  $\mathbf{X}$ . A very similar argument can be used to motivate the partial sketched estimator  $\beta_P$ . Existing bounds for the partial sketch focus on the prediction error  $\|\mathbf{X}\beta_P - \mathbf{X}\beta_F\|_2^2$  (Becker et al., 2015; Pilanci and Wainwright, 2016). To make a direct comparison to (3) we establish a bound on the coefficient error

**Theorem 1.** *Suppose that  $\tilde{\mathbf{X}}$  is an  $\epsilon$ -subspace embedding of  $\mathbf{X}$  with  $\epsilon < 0.5$ . Then the following bound holds,*

$$\|\beta_P - \beta_F\|_2^2 \leq \frac{4\epsilon^2}{\sigma_{\min}^2(\mathbf{X})} MSS_F. \quad (4)$$

For proof see the supplementary material. The mild requirement that  $\epsilon < 0.5$  is imposed so that the bound matches the functional form of the complete sketching bound (3). Comparing the partial sketching bound to (3), we see that the tightness of the bound is controlled by the model sum of squares as opposed to the residual sum of squares. The sensitivity of partial sketching to the model sum of squares as opposed to the residual sum of squares has been noted in previous on partial sketching (Dhillon et al., 2013; Pilanci and Wainwright, 2016; Becker et al., 2015). This suggests that the signal to noise ratio in the source dataset will be important when selecting which sketched estimator to use. A naive conclusion is that complete sketching is preferred when  $RSS_F < 4MSS_F$ , or equivalently  $R_F^2 > 0.25$ . Such a result is hardly prescriptive, as the worst case bound is not necessarily indicative of expected performance. A second point of interest is that if the  $k \times n$  matrix  $\mathbf{S}$  is an  $\epsilon$ -subspace embedding for  $\mathbf{A} = [\mathbf{y}, \mathbf{X}]$ , it is also an  $\epsilon$ -subspace embedding for  $\mathbf{X}$ . This suggests that it is reasonable to compute both  $\beta_P$  and  $\beta_S$  from a single sketch, although it is not clear how to combine the estimators into a single point estimator. These issues will be explored in more depth by examining the statistical properties of both complete and partial sketching. Before moving on to the statistical analysis we review some of the existing methods for generating  $\epsilon$ -subspace embeddings.

There are two general categories of distributions for the random matrix  $\mathbf{S}$ , data aware random projections and data oblivious random projections. A data aware random projection uses information in the source data  $\mathbf{y}, \mathbf{X}$  to generate  $\mathbf{S}$ . In contrast, a data oblivious random projection can be sampled without knowledge of  $\mathbf{y}$  or  $\mathbf{X}$ . Data aware random projections are closely connected to finite population sampling methods in the statistics literature, and this is discussed in more detail in Section 2.4. Data oblivious random projections are more closely related to dimension reduction techniques such as multidimensional scaling. Our main focus is on data oblivious random projections. Data oblivious projections are designed to produce  $\epsilon$ -subspace embeddings for an arbitrary source data matrix with high probability.



**Figure 1:** Sampled sketching matrices  $\mathbf{S}$  for  $k = 32, n = 36$ . Elements in the sketching matrix are coloured based on the value. One and negative one are coloured as black and white respectively. Intermediate values are in shades of grey.

### 2.3 Data oblivious sketches

The Gaussian sketch was one of the first projections proposed for sketched regression (Sarlos, 2006). Recall that a Gaussian sketch is formed by independently sampling each element of  $\mathbf{S}$  from a  $N(0, 1/k)$  distribution. The drawback of the Gaussian sketch is that computation of the sketched data is quite demanding, taking  $O(ndk)$  operations. As such, there has been work on designing more computationally efficient random projections. The Hadamard sketch and the Clarkson-Woodruff sketch are two examples of more efficient methods for generating  $\epsilon$ -subspace embeddings.

The Hadamard sketch is a structured random matrix (Ailon and Chazelle, 2009). The sketching matrix is formed as  $\mathbf{S} = \Phi \mathbf{H} \mathbf{D} / \sqrt{k}$ , where  $\Phi$  is a  $k \times n$  matrix and  $\mathbf{H}$  and  $\mathbf{D}$  are both  $n \times n$  matrices. The fixed matrix  $\mathbf{H}$  is a Hadamard matrix of order  $n$ . A Hadamard matrix is a square matrix with elements that are either  $+1$  or  $-1$  and orthogonal rows. Hadamard matrices do not exist for all integers  $n$ , the source dataset can be padded with zeroes so that a conformable Hadamard matrix is available. The random matrix  $\mathbf{D}$  is a diagonal matrix where each nonzero element is an independent Rademacher random variable. The random matrix  $\Phi$  subsamples  $k$  rows of  $\mathbf{H}$  with replacement. The structure of the Hadamard sketch allows for fast matrix multiplication, reducing calculation of the sketched dataset to  $O(nd \log k)$  operations.

The Clarkson-Woodruff sketch is a sparse random matrix (Clarkson and Woodruff, 2013). The projection can be represented as the product of two independent random matrices,  $\mathbf{S} = \mathbf{\Gamma} \mathbf{D}$ , where  $\mathbf{\Gamma}$  is a random  $k \times n$  matrix and  $\mathbf{D}$  is a random  $n \times n$  matrix. The matrix  $\mathbf{\Gamma}$  is formed by choosing one element in each column independently and setting the entry to  $+1$ . The matrix  $\mathbf{D}$  is a diagonal matrix where each nonzero element is an independent Rademacher random variable. This results in a sparse  $\mathbf{S}$ , where there is only one nonzero entry per column. The sparsity of the Clarkson-Woodruff sketch speeds up matrix multiplication, dropping the complexity of generating the sketched dataset to  $O(nd)$ .

Figure 1 shows examples of the three sketches for  $k = 32, n = 36$ . The sketches are discussed in more detail in the supplementary material.

Data oblivious sketches are designed to give an  $\epsilon$ -subspace embedding for an arbitrary source dataset with at least probability  $(1 - \delta)$ . Sketching algorithms are appealing for large  $n$  problems as the required  $k$  to attain the  $(\delta, \epsilon)$  bound is independent of  $n$  for the Gaussian and Clarkson-Woodruff

Algorithm	Sketching time	Required sketch size $k$
Gaussian sketch	$O(ndk)$	$O[\{d + \log(1/\delta)\}/\epsilon^2]$
Hadamard sketch	$O(nd \log k)$	$O[(\sqrt{d} + \sqrt{\log n})^2 \{\log(d/\delta)\}/\epsilon^2]$
Clarkson-Woodruff Sketch	$O(nd)$	$O\{d^2/(\delta\epsilon^2)\}$

**Table 1:** Properties of different data oblivious random projections (Woodruff, 2014). The third column refers to the necessary sketch size  $k$  to obtain an  $\epsilon$ -subspace embedding for an arbitrary  $n \times d$  source dataset with at least probability  $(1 - \delta)$ .

sketches, and very weakly dependent on  $n$  for the Hadamard sketch. Table 1 summarises existing results on the necessary  $k$  to attain the  $(\epsilon, \delta)$  bound. Probabilistic worst case bounds for sketched regression are formed by noting that if a sketch produces an  $\epsilon$ -subspace embedding with probability at least  $(1 - \delta)$ , then the bounds in Section 2.2 must hold with probability at least  $(1 - \delta)$ . Woodruff (2014) gives an excellent survey of work in this area.

## 2.4 Data aware sketches

As mentioned, data aware random projections can also be used to generate  $\epsilon$ -subspace embeddings. Data aware sketching is closely related to finite population subsampling methods (Ma and Sun, 2015), in particular classic Hansen-Hurwitz estimators (Hansen and Hurwitz, 1943). Suppose we sample  $k$  observations from the original dataset with replacement using observation sampling weights  $\pi_1, \dots, \pi_n$ . Let the data aware sketching matrix be constructed as  $\mathbf{S} = \mathbf{R}\mathbf{W}$ . The  $k \times n$  matrix  $\mathbf{W}$  subsamples  $k$  rows of the source dataset with replacement. Each row of  $\mathbf{W}$  contains a single nonzero entry. Element  $W_{ij}$  is equal to one if the  $j$ th original observation is sampled in the  $i$ th sampling round for  $i = 1, \dots, k$  and  $j \in \{1, \dots, n\}$ . The  $k \times k$  diagonal matrix  $\mathbf{R}$  rescales the subsampled rows. The  $i$ th diagonal element of  $\mathbf{R}$  is set to  $1/(k\pi_j)^{1/2}$  if  $W_{ij}$  is equal to one, that is if row  $j$  in the source dataset is subsampled by the  $i$ th row of the subsampling matrix  $\mathbf{W}$ . Using a data aware sketch, the sketched dataset is defined as

$$\begin{aligned}\tilde{\mathbf{A}} &= \mathbf{S}\mathbf{A} \\ &= \mathbf{R}\Phi\mathbf{A}.\end{aligned}$$

The sketched dataset has the property that  $E_S(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \mid \mathbf{A}^T \mathbf{A}) = \mathbf{A}^T \mathbf{A}$ . The subsampling and rescaling can be interpreted as a Hansen-Hurwitz estimator of the full dataset sufficient statistics  $\mathbf{A}^T \mathbf{A}$ .

Data aware sketching algorithms use the leverage scores of the observations to define the sampling weights  $\pi_1, \dots, \pi_n$  (Mahoney, 2011; Woodruff, 2014). Let the singular value decomposition of the source data matrix  $\mathbf{A}$  be given by  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ , where  $\mathbf{U}$  is the  $n \times d$  matrix of left singular vectors,  $\mathbf{D}$  is a  $d \times d$  matrix with the singular values of  $\mathbf{A}$  on the diagonal, and  $\mathbf{V}$  is the  $d \times d$  matrix of right singular vectors. Let  $\mathbf{u}_i^T$  give the  $i$ th row in  $\mathbf{U}$ . The leverage score for observation  $i$  is defined as  $\|\mathbf{u}_i\|_2^2$ .

Suppose the original dataset is centred, so each column of  $\mathbf{A}$  has mean zero. The leverage scores then have a particularly intuitive interpretation in terms of the principal components decomposition of the source dataset. The row vector  $\mathbf{u}_i^T \mathbf{D}$  gives the coordinates of observation  $i$  on the principal

component axes. The elements of the vector  $\mathbf{u}_i$  give the coordinates of observation  $i$  in a scaled system where the variance along each principal coordinate axis is set to be one. The leverage score  $\|\mathbf{u}_i\|_2^2$  gives a measure of the distance from the origin in the principal coordinate system. This geometric perspective will also be of use when analysing data oblivious random projections.

Data oblivious random projections operate in a different manner to data aware random projections, as the sketched dataset is not a rescaled subset of the original instances. Data oblivious random projections generate a pseudo-dataset of  $k$  observations using the source dataset as a component in the generative process. In Section 4 we will show that the leverage scores have an important role in describing the asymptotic behaviour of data oblivious random projections. We first establish some exact distributional results for the estimators  $\beta_S$  and  $\beta_P$  under the Gaussian sketch in Section 3.1. In Section 4 we establish corresponding asymptotic results for the Hadamard and Clarkson-Woodruff projections under regularity conditions on the statistical leverage scores.

### 3 Gaussian sketching

#### 3.1 Complete sketching

The Gaussian sketch is mathematically tractable, and it is possible to establish a number of exact finite sample results regarding the performance of the sketched estimators. In this section we will develop the distribution of  $\beta_S$  when using a Gaussian sketch. As mentioned previously, all results treat  $\mathbf{y}$  and  $\mathbf{X}$  as fixed. The variability in  $\beta_S$  is solely due to the use of the random sketching matrix  $\mathbf{S}$ . Let  $(\tilde{y}_j, \tilde{\mathbf{x}}_j^\top)^\top$  refer to the  $j$ th row in the sketched data matrix  $\tilde{\mathbf{A}} = [\tilde{\mathbf{y}}, \tilde{\mathbf{X}}]$  for  $j = 1, \dots, k$ . Similarly, let  $\mathbf{s}_j^\top$  denote the  $j$ th row in the sketching matrix  $\mathbf{S}$ . The sketched dataset consists of  $k$  random units  $(\tilde{y}_j, \tilde{\mathbf{x}}_j^\top)$ , ( $j = 1, \dots, k$ ). The  $j$ th sketched response is given by  $\tilde{y}_j = \mathbf{s}_j^\top \mathbf{y}$ , and the  $j$ th sketched predictor is calculated as  $\tilde{\mathbf{x}}_j = \mathbf{s}_j^\top \mathbf{X}$  ( $j = 1, \dots, k$ ). The  $k$  sketched instances are independently distributed, because rows of the sketching matrix are independent.

We take an indirect route to find the distribution of  $\beta_S$ , by focusing on the distribution of the sketched data  $\tilde{\mathbf{A}} = [\tilde{\mathbf{y}}, \tilde{\mathbf{X}}]$  conditional on the original dataset  $\mathbf{A} = [\mathbf{y}, \mathbf{X}]$ . The initial step is to decompose the joint distribution on the sketched responses and predictors as the product of a marginal and conditional distribution. Specifically,

$$p(\tilde{\mathbf{y}}, \tilde{\mathbf{X}} \mid \mathbf{y}, \mathbf{X}) = p(\tilde{\mathbf{y}} \mid \tilde{\mathbf{X}}, \mathbf{y}, \mathbf{X})p(\tilde{\mathbf{X}} \mid \mathbf{y}, \mathbf{X}).$$

It can be shown that  $p(\tilde{\mathbf{y}} \mid \tilde{\mathbf{X}}, \mathbf{y}, \mathbf{X})p(\tilde{\mathbf{X}} \mid \mathbf{y}, \mathbf{X})$  has the structure of a hierarchical Gaussian linear model. We first show that the sketched dataset has a multivariate normal distribution, conditional on the source dataset. This follows as the sketched dataset can be expressed as a linear combination of Gaussian random variables. Specifically, row  $j$  in the sketched dataset is  $(\tilde{y}_j, \tilde{\mathbf{x}}_j^\top) = \mathbf{s}_j^\top \mathbf{A}$ . The random vector  $(\tilde{y}_j, \tilde{\mathbf{x}}_j^\top)^\top$  is given by the linear combination

$$\begin{bmatrix} \tilde{y}_j \\ \tilde{\mathbf{x}}_j \end{bmatrix} = \mathbf{A}^\top \mathbf{s}_j.$$

Conditional on  $\mathbf{A} = [\mathbf{y}, \mathbf{X}]$ ,  $\mathbf{A}^\top \mathbf{s}_j$  is a linear combination of independent Gaussians as  $\mathbf{s}_j \sim$



$N(\mathbf{0}, \mathbf{I}_d/k)$ . As affine transformations of Gaussians are also multivariate normal,  $(\tilde{y}_j, \tilde{\mathbf{x}}_j^\top)$  must then be jointly normally distributed, conditional on the source data  $\mathbf{A} = [\mathbf{y}, \mathbf{X}]$ . It is easily shown that the joint distribution of the sketched responses and predictors is then

$$\begin{bmatrix} \tilde{y}_j \\ \tilde{\mathbf{x}}_j \end{bmatrix} \Big| \mathbf{y}, \mathbf{X} \sim N \left( \begin{bmatrix} 0 \\ \mathbf{0} \end{bmatrix}, \frac{1}{k} \begin{bmatrix} \mathbf{y}^\top \mathbf{y} & \mathbf{y}^\top \mathbf{X} \\ \mathbf{X}^\top \mathbf{y} & \mathbf{X}^\top \mathbf{X} \end{bmatrix} \right), \quad (j = 1, \dots, k).$$

Standard results on the multivariate normal distribution give that the conditional distribution of  $\tilde{y}_j$  given  $\tilde{\mathbf{x}}_j$  is also normal. A routine calculation shows that the conditional mean is related to  $\boldsymbol{\beta}_F$ , that is  $E_S(\tilde{y}_j | \tilde{\mathbf{x}}_j, \mathbf{y}, \mathbf{X}) = \tilde{\mathbf{x}}_j^\top \boldsymbol{\beta}_F$ . The subscript  $S$  is used on the expectation operator to emphasise that only random quantity is the sketching matrix. The conditional variance is related to the prediction error on the source dataset  $RSS_F$ ,

$$\begin{aligned} \text{var}_S(\tilde{y}_j | \tilde{\mathbf{x}}_j, \mathbf{y}, \mathbf{X}) &= \frac{1}{k} \{ \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \} \\ &= \frac{1}{k} RSS_F. \end{aligned}$$

The subscript  $S$  is again used to recognise that the source of the variance is the random sketching matrix, the source dataset is fixed. The step in the second line follows from sum of squares partitions in linear models (Searle, 1997, Chapter 3). Therefore, the conditional distribution of  $\tilde{y}_j$  given the sketched predictors  $\tilde{\mathbf{x}}_j$  and the source dataset  $[\mathbf{y}, \mathbf{X}]$  is

$$\tilde{y}_j | \tilde{\mathbf{x}}_j, \mathbf{y}, \mathbf{X} \sim N \left( \tilde{\mathbf{x}}_j^\top \boldsymbol{\beta}_F, \frac{RSS_F}{k} \right) \quad (j = 1, \dots, k).$$

This is the exact form of a standard Gaussian linear model. The distribution  $p(\tilde{\mathbf{X}} | \mathbf{y}, \mathbf{X})$  is easily obtained as the marginal distribution of  $\tilde{\mathbf{x}}_j$  is also multivariate normal,

$$\tilde{\mathbf{x}}_j \sim N(\mathbf{0}, \mathbf{X}^\top \mathbf{X}/k), \quad (j = 1, \dots, k).$$

The sketching process can be described using the following hierarchical model,

$$\begin{aligned} \tilde{\mathbf{y}} | \tilde{\mathbf{X}}, \mathbf{y}, \mathbf{X} &\sim N \left( \tilde{\mathbf{X}} \boldsymbol{\beta}_F, \frac{RSS_F}{k} \mathbf{I}_k \right), \\ \tilde{\mathbf{X}} | \mathbf{y}, \mathbf{X} &\sim MN \left( \mathbf{0}, \mathbf{I}_k, \frac{1}{k} \mathbf{X}^\top \mathbf{X} \right). \end{aligned}$$

A Gaussian sketch effectively simulates a series of observations from a Gaussian linear model parametrised in terms of  $\boldsymbol{\beta}_F$  and  $RSS_F$ , where the design matrix has a matrix normal distribution. We now turn to the distribution of  $\boldsymbol{\beta}_S$ . The distribution of  $\boldsymbol{\beta}_S$  conditional on the sketched predictors follows immediately from standard results on linear models (Searle, 1997, Chapter 3).

$$\boldsymbol{\beta}_S | \tilde{\mathbf{X}}, \mathbf{y}, \mathbf{X} \sim N \left( \boldsymbol{\beta}_F, \frac{RSS_F}{k} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \right). \quad (5)$$

To obtain the marginal distribution of  $\boldsymbol{\beta}_S$  it is necessary to integrate over the random sketched

design matrix  $\widetilde{\mathbf{X}}$ . From properties of the normal distribution (Eaton, 2007), it is possible to show  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}}) | \mathbf{y}, \mathbf{X} \sim \text{Wishart}(k, \mathbf{X}^\top \mathbf{X}/k)$ . As such,

$$(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} | \mathbf{y}, \mathbf{X} \sim \text{Inverse-Wishart}(k, k(\mathbf{X}^\top \mathbf{X})^{-1}).$$

As seen in equation (5),  $\beta_S$  is normally distributed when conditioned on the random Inverse-Wishart matrix  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}$ . The marginal distribution of  $\beta_S$  can then be described using the Normal Inverse-Wishart distribution (Gelman et al., 2014, p.73). The following theorem characterises the distribution of  $\beta_S$  under the Gaussian sketch.

**Theorem 2.** *Suppose  $\beta_S$  is computed using a Gaussian sketch and  $k > p + 1$ . The conditional distribution of  $\beta_S$  is*

$$(i) \beta_S | \widetilde{\mathbf{X}}, \mathbf{y}, \mathbf{X} \sim N\left(\beta_F, \frac{RSS_F}{k} (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}\right).$$

The marginal distribution of  $\beta_S$  is

$$(ii) \beta_S | \mathbf{y}, \mathbf{X} \sim \text{Student}\left(\beta_F, \frac{RSS_F}{k-p+1} (\mathbf{X}^\top \mathbf{X})^{-1}, k-p+1\right).$$

For proof see the supplementary material.

An immediate application of result (i) is the ability to generate exact confidence intervals for the elements of  $\beta_S$ , methodology that does not appear to be present in the existing literature. Let  $\beta_S^{(i)}$  give the  $i$ th element of  $\beta_S$  and let  $\beta_F^{(i)}$  give the  $i$ th element of  $\beta_F$ . Let  $RSS_S$  denote the sketched residual sum of squares,  $RSS_S = \|\widetilde{\mathbf{y}} - \widetilde{\mathbf{X}}\beta_S\|_2^2$ . To construct a  $100(1-\alpha)\%$  confidence interval, let  $w_{ii} = (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}_{ii}$ , and  $t_{\text{crit}}$  denote the  $100(1-\alpha/2)$ th percentile of the  $t$ -distribution with  $k-p$  degrees of freedom. Then from standard results on Gaussian linear models (Searle, 1997),

$$\beta_S^{(i)} \pm t_{\text{crit}} \times (w_{ii} RSS_S / (k-p))^{1/2} \tag{6}$$

gives an exact  $100(1-\alpha)\%$  confidence interval for  $\beta_S^{(i)}$ . Again assuming that  $k > p + 1$ , it should be noted that the variance of  $\beta_S$ ,

$$\text{var}(\beta_S | \mathbf{y}, \mathbf{X}) = \frac{RSS_F}{(k-p+1)} (\mathbf{X}^\top \mathbf{X})^{-1} \tag{7}$$

is not dependent on the compression ratio  $k/n$ . Although  $RSS_F$  can be expected to grow linearly with  $n$ , this will generally be counterbalanced by  $(\mathbf{X}^\top \mathbf{X})^{-1}$  decreasing linearly with  $n$ . The distribution of the approximation error  $\|\beta_S - \beta_F\|_2$  will largely be controlled by the target dimension  $k$ . This speaks to the defining characteristic of sketching algorithms, that given a fixed  $k$ , the stochastic approximation error does not necessarily increase with size of the original dataset  $n$ .

### 3.2 Partial sketching

Partial sketching was first proposed by Dhillon et al. (2013) using uniform subsampling, and later studied for general sketches by Pilanci and Wainwright (2016). Existing results on partial sketching

highlight that the model sum of squares influences the approximation error of the partial sketched estimator  $\beta_P$ . It is simple to see that the variance of the partial sketched estimator will not be a function of the residual sum of squares. From the normal equations it holds that  $\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \beta_F$ . Using this property, we see that conditional on  $\mathbf{y}, \mathbf{X}$ , the variance of the random linear combination  $\beta_P = (\mathbf{X}^\top \mathbf{S}^\top \mathbf{S} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = (\mathbf{X}^\top \mathbf{S}^\top \mathbf{S} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \beta_F$  will be a function of the covariates  $\mathbf{X}$  and the fitted values  $\mathbf{X} \beta_F$ . The residual vector has no influence on the variance of the partial sketching estimator, and as such the variance of  $\beta_P$  will not be related to the residual sum of squares. This suggests that when the noise level is high, partial sketching may become preferable to complete sketching. This idea has been touched on in the existing literature, but specific guidelines are lacking (Becker et al., 2015; Dhillon et al., 2013). A statistical analysis can provide some insight into this issue.

The hierarchical model for complete sketching gave an intuitive statistical perspective on the mechanics of the algorithm. Partial sketching seems to lack a similar conceptual device. The least squares coefficients can be represented as the solution to the linear system of the equations  $\mathbf{X}^\top \mathbf{X} \mathbf{b} = \mathbf{X}^\top \mathbf{y}$ . Partial sketching simply returns the solution,  $\mathbf{b}$ , to the approximate linear system  $\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \mathbf{b} = \mathbf{X}^\top \mathbf{y}$ . Lacking a convenient representation for the estimator, we must proceed in a more pedestrian manner. The mean square error of the estimator  $\beta_P$  can be determined using only mean and variance information, and this will be the goal for now. The key observation is that  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mid \mathbf{y}, \mathbf{X} \sim \text{InvWishart}(k, k(\mathbf{X}^\top \mathbf{X})^{-1})$ . Conditional on  $\mathbf{y}, \mathbf{X}$ , the estimator  $\beta_P = (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{y}$  is a linear combination of the elements of an Inverse-Wishart random matrix. However, this is a non-standard distribution and it is difficult to directly express the distribution function of  $\beta_P$ . Despite this, it is straightforward to determine the mean and variance of  $\beta_P$ . From properties of the Inverse-Wishart distribution, it can be seen that the partial sketched estimator is biased, with mean

$$E_S(\beta_P \mid \mathbf{y}, \mathbf{X}) = \frac{k}{(k-p-1)} \beta_F,$$

where it is assumed that  $k > p + 3$ . This motivates an alternative unbiased estimator

$$\beta_P^* = \frac{(k-p-1)}{k} (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{y}.$$

Determining the variance of  $\beta_P$  and the unbiased  $\beta_P^*$  is a more lengthy computation (see supplementary material). The variance of the biased estimator  $\beta_P$  is

$$\text{var}(\beta_P \mid \mathbf{y}, \mathbf{X}) = \frac{k^2}{(k-p)(k-p-1)(k-p-3)} \left( MSS_F (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{k-p+1}{k-p-1} \beta_F \beta_F^\top \right). \quad (8)$$

The variance of the unbiased estimator  $\beta_P^*$  is

$$\text{var}(\beta_P^* \mid \mathbf{y}, \mathbf{X}) = \frac{(k-p-1)}{(k-p)(k-p-3)} \left( MSS_F (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{k-p+1}{k-p-1} \beta_F \beta_F^\top \right). \quad (9)$$

The variances of  $\beta_P$  and  $\beta_P^*$  have a similar structure to the variance of  $\beta_S$ . The main point of difference is that the variance of  $\beta_S$  depends on the residual sum of squares, whereas the variance of  $\beta_P$  and  $\beta_P^*$  depends on the model sum of squares.

As mentioned the explicit form of the sampling distribution is hard to obtain, but by making a connection with method of moments estimation it is possible to establish asymptotic normality of both  $\beta_P$  and  $\beta_P^*$  as  $k$  tends to infinity. This motivates the construction of approximate confidence intervals. As the exact variance is unknown we propose the following estimator

$$\text{var}(\beta_P^* | \mathbf{y}, \mathbf{X}) \approx \frac{(k-p-1)}{(k-p)(k-p-3)} \left( \left( \frac{k-p-1}{k} \right) MSS_S(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} + \beta_P^* \beta_P^{*\top} \right). \quad (10)$$

### 3.3 Relative efficiency

The relative efficacy of complete and partial sketching is also of interest. As the plug in estimator  $\beta_P$  has a higher mean square error than  $\beta_P^*$ , it will not be considered in this section. The performance of the complete sketching estimator  $\beta_S$  and the unbiased partial sketched estimator  $\beta_P^*$  will be compared in terms of mean squared error. As both  $\beta_F$  and  $\beta_P^*$  are unbiased, the mean squared error can be computed using their respective covariance matrices, that is

$$\begin{aligned} E_S (\|\beta_S - \beta_F\|_2^2 | \mathbf{y}, \mathbf{X}) &= \text{tr}(\text{var}(\beta_S)), \\ E_S (\|\beta_P^* - \beta_F\|_2^2 | \mathbf{y}, \mathbf{X}) &= \text{tr}(\text{var}(\beta_P^*)). \end{aligned}$$

Comparing (7) and (9), the variance of  $\beta_P^*$  is dependent on  $MSS_F$ , whereas the variance of  $\beta_S$  is dependent on  $RSS_F$ . This suggests that the signal to noise ratio in the source dataset will be an influential factor in determining which estimator is more efficient. When  $R_F^2$  is close to one complete sketching can be orders of magnitude more efficient than partial sketching, and when  $R_F^2$  is close to zero, partial sketching can be orders of magnitude more efficient than complete sketching.

### 3.4 Combined estimator

So far we have assumed that an analyst much choose between one of the two methods. Obtaining both  $\beta_P^*$  and  $\beta_S$  from a single sketch is computationally cheap, and may be an attractive strategy. The most demanding operation with the sketched data is calculating  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}$ . Given this quantity it is economical to compute both  $\beta_S$  and  $\beta_P^*$ . Becker et al. (2015) mention they are presently investigating such a strategy, but do not give any details. Our motivation for a combined estimator is driven by the fact even when using a single sketch  $(\widetilde{\mathbf{y}}, \widetilde{\mathbf{X}})$ , the two estimators are uncorrelated, that is  $\text{cov}(\beta_P^*, \beta_S) = \mathbf{0}$ . This is established by taking iterated expectations, and using the hierarchical model established in Section 3.1 (see supplementary material). A simple strategy is then to take a weighted combination of  $\beta_S$  and  $\beta_P^*$ . A combined estimator  $\beta_C$  can be defined as

$$\beta_C = \phi \beta_S + (1 - \phi) \beta_P^*,$$

for some  $0 \leq \phi \leq 1$ . The value of  $\phi$  that minimises the mean square error is

$$\phi_{\text{opt}} = \frac{\text{tr}(\text{var}(\beta_P^*))}{\text{tr}(\text{var}(\beta_P^*)) + \text{tr}(\text{var}(\beta_S))}.$$

Use of the weighted estimator is expected to be most beneficial when the signal to noise ratio

is moderate, that is  $R_F^2 \approx 0.5$ . When the signal to noise ratio is either very high or very low, there is little gain from using the weighted estimator as either the complete or partial estimator will dominate.

### 3.5 One-step correction

As noted by a referee, the combined estimator is related to another strategy in the sketching literature for improving  $\beta_S$ . Dhillon et al. (2013) and Pilanci and Wainwright (2016) propose a refinement procedure using gradient information from the source dataset. The one-step corrected estimator,  $\beta_H$ , is defined as

$$\begin{aligned}\beta_H &= \beta_S + (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \beta_S) \\ &= (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X}) \beta_S + (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}\quad (11)$$

The one-step estimator can be interpreted as a single step of the iterative Hessian sketch proposed by Pilanci and Wainwright (2016), initialised at  $\beta_S$ . The optimal least square solution  $\beta_F$  satisfies  $\mathbf{X}^\top (\mathbf{y} - \mathbf{X} \beta_F) = \mathbf{0}$  so

$$\begin{aligned}\beta_F &= \beta_F + (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \beta_F) \\ &= (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X}) \beta_F + (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}\quad (12)$$

Subtracting (12) from (11) gives the following expression for the error

$$\beta_H - \beta_F = (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X}) (\beta_S - \beta_F). \quad (13)$$

The expected squared error is then

$$E_S(\|\beta_H - \beta_F\|_2^2) = E_S \left\{ (\beta_S - \beta_F)^\top (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X})^\top (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X}) (\beta_S - \beta_F) \right\}.$$

We can then take iterated expectations using the hierarchical model in Section 3.1. The action of the sketch can be taken over  $\widetilde{\mathbf{X}}$  then over the conditional distribution  $\widetilde{\mathbf{y}}$  given  $\widetilde{\mathbf{X}}$ . Theorem 2 (i) gives the distribution of  $\beta_S$  conditional on  $\widetilde{\mathbf{X}}$ . We thus have

$$\begin{aligned}E_S(\|\beta_H - \beta_F\|_2^2) &= E_{\widetilde{\mathbf{X}}} \left[ E_{\widetilde{\mathbf{y}}|\widetilde{\mathbf{X}}} \left\{ (\beta_S - \beta_F)^\top (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X})^\top (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X}) (\beta_S - \beta_F) \mid \widetilde{\mathbf{X}} \right\} \right] \\ &= E_{\widetilde{\mathbf{X}}} \left[ \text{tr} \left( \text{var}(\beta_S \mid \widetilde{\mathbf{X}}) (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X})^\top (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X}) \right) \right] \\ &= E_{\widetilde{\mathbf{X}}} \left\{ \text{tr} \left( \frac{RSS_F}{k} (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X})^\top (\mathbf{I} - (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \mathbf{X}^\top \mathbf{X}) \right) \right\}.\end{aligned}\quad (14)$$

The key term in (14) is the random matrix  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}$ . Now as  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} \sim \text{Inverse-Wishart}(k, k(\mathbf{X}^\top \mathbf{X})^{-1})$ , it is possible to evaluate the expectation in (14) using moments of the Inverse-Wishart distribution. The exact expression involves  $E(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}$ ,  $E(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-2}$  and  $E(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-3}$ . Formulae for the required moments are given in Letac and Massam (2004). The main conclusions are that the one-step esti-

mator  $\beta_H$  can have a larger mean square error than  $\beta_S$  when the sketch size to variables ratio  $k/p$  is close to one. As  $k$  increases the one-step estimator becomes more efficient than both  $\beta_S$  and  $\beta_C$  with the optimal weight  $\phi_{\text{opt}}$ . The relative efficiency of  $\beta_C$  to  $\beta_S$  is at most two. The relative efficiency of  $\beta_H$  to  $\beta_S$  can be much larger, providing that  $k/p$  is sufficiently large. The exact relationship is a function of  $k$  and  $p$ . Direct comparisons between  $\beta_H$ ,  $\beta_S$  and  $\beta_P^*$  are difficult, as the one-step estimator  $\beta_H$  requires gradient information  $\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\beta_S)$ . Calculation of the gradient requires access to the full dataset. The single pass estimators  $\beta_S$  and  $\beta_P^*$  require only the sketched dataset and the summary statistic  $\mathbf{X}^\top\mathbf{y}$ . Additionally, the iterative correction can also be applied to  $\beta_P$  or  $\beta_P^*$ . We are currently investigating the properties of iterative sketching algorithms in more detail using the asymptotic results developed in this paper.

## 4 Asymptotics

### 4.1 Preliminaries

Finite sample distributions of random projection estimators can be mathematically intractable, and as such asymptotic analysis can be a powerful tool (Li et al., 2006; Diaconis and Freedman, 1984). It is a very difficult task to establish meaningful finite sample results for the Hadamard and Clarkson-Woodruff sketches, as they are discrete distributions over an enormous combinatorial space. The explicit finite sample distribution of the sketched estimators can be written as a sum over all these possible combinations, but such a representation is not very informative. Instead, it is useful to study the large  $n$  distribution of the estimators  $\beta_S$  and  $\beta_P$  to obtain an interpretable expression.

As  $\beta_F$  is the estimand in sketching algorithms, this requires conditioning on the source data in the asymptotic analysis. To elaborate, let  $\mathbf{A}_{(n)} = [\mathbf{y}_{(n)}, \mathbf{X}_{(n)}]$  represent the  $n \times d$  source data matrix of full column rank. Any source data matrix  $\mathbf{A}_{(n)}$  has a set of associated least squares coefficients, which will here be denoted  $\beta_F^{(n)}$ . The overall goal is to determine the asymptotic form of the distributions  $p(\beta_S | \mathbf{A}_{(n)})$  and  $p(\beta_P^* | \mathbf{A}_{(n)})$  for some arbitrary large dataset  $\mathbf{A}_{(n)}$ .

To take limits, we employ a fixed sequence of  $n \times d$  datasets, all of rank  $d$ . In the regression scenario this amounts to assuming that  $\mathbf{X}_{(n)}$  is of full column rank and that  $\mathbf{y}_{(n)}$  is not a perfect linear combination of the columns of  $\mathbf{X}_{(n)}$  for all  $n$ . Conditioning on  $\mathbf{A}_{(n)}$  is effectively the same as treating the full dataset as an arbitrary sequence of constants  $A_{ij}$  for  $i = 1, \dots, n, j = 1, \dots, d$ . This is analogous to large sample results for regression models where the design matrix is treated as arbitrary set of constants, and the random variables of interest are the error terms, for example see Van Der Vaart (1998, Section 2.5). Here the source dataset is treated as a sequence of constants and the random variables of interest are the elements of the sketching matrix.

The asymptotic analysis is carried out in two stages. The initial step is to establish asymptotic normality of the sketched dataset. The regularity condition for the central limit theorem highlights the influential role of the leverage scores of the observations in the source dataset. This is then followed by an analysis of the limiting distribution of  $\beta_S$ , and  $\beta_P^*$ . There is some related work by Ma et al. (2015) who develop Taylor series approximations for the bias and variance of data aware sketched regression estimators, where the asymptotic expansion is taken in the sketch size  $k$ . Our work is different as we study data oblivious random projections and build our asymptotic results from a conditional central limit theorem for the sketched data matrix. The conditional central limit

theorem is established for fixed  $k$  and  $d$ , taking the number of source observations to  $n$  to infinity.

## 4.2 Sketching central limit theorem

A central limit theorem for sparse sketching matrices with independent entries is given in Li et al. (2006). The Clarkson-Woodruff sketch and the Hadamard sketch have dependent entries, and as such we use a different method of proof. Under some regularity conditions the Hadamard and Clarkson-Woodruff sketches produce sketched data that asymptotically has the same matrix normal distribution as under the Gaussian sketch. Using a Gaussian sketch, conditional on  $\mathbf{A}$ ,

$$\tilde{\mathbf{A}} \sim MN(\mathbf{0}, \mathbf{I}_k, \mathbf{A}^\top \mathbf{A}/k). \quad (15)$$

Each row is statistically independent, and marginally normally distributed with covariance matrix  $\mathbf{A}^\top \mathbf{A}/k$ . Although asymptotic normality may not be particularly surprising seeing as the sketched data are linear combinations of random vectors, the proof is not immediate due to the dependence in the Hadamard and Clarkson-Woodruff sketches. The difficulties caused by the dependence are most easily illustrated for the Clarkson-Woodruff sketch.

---

### Algorithm 1 Clarkson-Woodruff sketch

---

$\tilde{\mathbf{A}} \leftarrow \mathbf{0}$       Initialise sketched dataset as  $k \times d$  matrix of zeroes  
For  $i = 1$  to  $i = n$   
    Sample  $z \sim \text{Uniform}(1, \dots, k)$       Sample random index  
    Sample  $r \sim \text{Uniform}(-1, +1)$       Sample random sign  
     $\tilde{\mathbf{A}}_z \leftarrow r \times \mathbf{A}_i + \tilde{\mathbf{A}}_z$       Multiply by  $r$  and add to row  $z$  in sketch  
Output  $\tilde{\mathbf{A}}$       Output sketched dataset

---

The behaviour of the Clarkson-Woodruff sketch can be represented as a many to less mapping. Each row in the source dataset is assigned to a single row in the sketched dataset. The Clarkson-Woodruff sketch has an alternative streaming construction that highlights this property, given in Algorithm 1. As each row in the source dataset only contributes to a single row in the sketched dataset, it might be expected that this results in some statistical dependence amongst the rows of the sketched dataset. Additionally, although it seems each row in the sketched dataset will be marginally normally distributed, it is not clear if joint asymptotic normality over all rows will hold. Similar conundrums arise when examining the Hadamard sketch in detail.

The  $k \times d$  random matrix  $\tilde{\mathbf{A}}$  is the output of a stochastic process governed by the fixed  $n \times d$  source dataset  $\mathbf{A}_{(n)}$  and the distribution of the random  $k \times n$  sketching matrix  $\mathbf{S}$ . The sketched dataset is a linear combination of random vectors, the number of which increases with  $n$ . As such, we can expect  $\tilde{\mathbf{A}}$  to demonstrate some stable limiting behaviour as  $n$  grows larger. Under an assumption on the limiting leverage scores of the source data matrix, we can establish a central limit theorem for the sketched dataset. Recall the singular value decomposition of the source dataset  $\mathbf{A}_{(n)} = \mathbf{U}_{(n)} \mathbf{D}_{(n)} \mathbf{V}_{(n)}^\top$ . The leverage score of observation  $i$  in the source dataset is defined as  $\|\mathbf{u}_{(n)i}\|_2^2$  where  $\mathbf{u}_{(n)i}^\top$  gives row  $i$  in  $\mathbf{U}_{(n)}$ . The leverage scores of the observations in the source data matrix have been identified an important structural property of sketching algorithms (Mahoney and Drineas,

2016). Assumption 1 highlights their role in establishing asymptotic normality of the sketched data matrix.

**Assumption 1.** *Let the singular value decomposition of the  $n \times d$  source dataset be given by  $\mathbf{A}_{(n)} = \mathbf{U}_{(n)}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^\top$ . Let  $\mathbf{u}_{(n)i}^\top$  give the  $i$ th row in  $\mathbf{U}_{(n)}$ . Assume that the maximum leverage score tends to zero, that is*

$$\lim_{n \rightarrow \infty} \max_{i=1, \dots, n} \|\mathbf{u}_{(n)i}\|_2^2 = 0.$$

Theorem 3 gives the sketching central limit theorem.

**Theorem 3.** *Consider a fixed sequence of arbitrary  $n \times d$  data matrices  $\mathbf{A}_{(n)}$ , where  $d$  is fixed. Let  $\mathbf{A}_{(n)} = \mathbf{U}_{(n)}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^\top$  represent the singular value decomposition of  $\mathbf{A}_{(n)}$ . Let  $\mathbf{S}$  be a  $k \times n$  Hadamard or Clarkson-Woodruff sketching matrix where  $k$  is also fixed. Suppose that Assumption 1 on the maximum leverage score is satisfied. Then as  $n$  tends to infinity with  $k$  and  $d$  fixed, we have the following convergence in distribution*

$$[\tilde{\mathbf{A}}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1} \mid \mathbf{A}_{(n)}] \rightarrow \text{MN}(\mathbf{0}, \mathbf{I}_k, \mathbf{I}_d/k).$$

The proof is given in the supplementary material. Heuristically, for large  $n$  we expect the matrix normal result (15) to approximately hold for the Hadamard and Clarkson-Woodruff sketches. The significance of Assumption 1 is perhaps best explained by making a connection to a version of the Lindeberg-Feller theorem for triangular arrays of uniformly bounded random variables.

**Theorem 4** (Billingsley, 1995, Chapter 5). *For each  $n \in \mathbb{N}$ , let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of independent random variables with  $E(Z_{ni}) = 0$  and  $\text{var}(Z_{ni}) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Let  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  and assume that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Suppose that we can form a sequence of upper bounds  $(K_n)_{n \in \mathbb{N}}$  such that for each  $n$ ,*

$$|Z_{ni}| \leq K_n \text{ almost surely for } i = 1, \dots, r_n.$$

*Then if  $K_n/s_n \rightarrow 0$  as  $n \rightarrow \infty$  we have the following convergence in distribution*

$$\frac{1}{s_n} \sum_{i=1}^{r_n} Z_{ni} \rightarrow N(0, 1)$$

In Theorem 4, the condition that  $K_n/s_n \rightarrow 0$  ensures that no random variable in a particular row of the array has too much pull over the sum  $\sum_{i=1}^{r_n} Z_{ni}$ . A triangular array of random variables satisfying the conditions in Theorem 4 is often said to be uniformly asymptotically negligible, in that no single term has undue influence over the random sum. We can make an analogy to the leverage score condition in the sketching central limit theorem (Theorem 3). The sum of the statistical leverage scores is always equal to the rank of the source dataset. As we have assumed that each dataset in the sequence is of rank  $d$ , we have that  $\sum_{i=1}^n \|\mathbf{u}_{(n)i}\|_2^2 = d$  for all  $n$ . As  $n$  grows we need the maximum contribution from a single term in the sum to tend to zero. The limiting leverage scores must satisfy an asymptotic negligibility condition, so that each individual observation provides a vanishingly small contribution to the total sum of the leverage scores.



As mentioned in the discussion of data aware sketching, the leverage scores have a particularly intuitive interpretation in terms of the principal components decomposition of the source dataset. The row vector  $\mathbf{u}_{(n)i}^\top \mathbf{D}_{(n)}$  gives the coordinates of observation  $i$  on the principal component axes. The elements of the vector  $\mathbf{u}_{(n)i}$  give the coordinates of observation  $i$  in a scaled system where the variance along each principal coordinate axis is set to be one. Treating the source dataset as a point cloud in Euclidean space, Assumption 1 essentially implies that there are no extreme outliers as  $n$  tends to infinity. Each observation must have a negligible contribution to the total variance along each principal component axis.

### 4.3 Sketching estimators

The central limit theorem for the sketched data suggests that the results about  $\beta_S$  and  $\beta_P$  for the Gaussian sketch will also approximately hold for the Hadamard and Clarkson-Woodruff sketches for large  $n$ . In order to establish convergence of the estimators it helps to adopt an extra assumption on the sequence of source datasets.

**Assumption 2.**

$$\lim_{n \rightarrow \infty} n^{-1} \begin{bmatrix} \mathbf{y}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{y}_{(n)}^\top \mathbf{X}_{(n)} \\ \mathbf{X}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \end{bmatrix} = \mathbf{Q} \quad \text{for some positive-definite matrix } \mathbf{Q}.$$

It is worth discussing the significance of the limiting matrix  $\mathbf{Q}$ . A useful comparison can be made to asymptotic theory for regression models, where a common assumption is that the design matrix satisfies the limit condition  $n^{-1} \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \rightarrow \mathbf{B}$ , where  $\mathbf{B}$  is some positive definite matrix (White, 1984; Greene, 1997). The development of asymptotic results is often eased by treating the covariates as a random sample, although this requires positing a realistic probability model for the covariates, which may be difficult. Treating the covariates as an arbitrary fixed sequence relaxes this assumption and covers more general scenarios. Although it is possible to establish asymptotic results when  $n^{-1} \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)}$  is not required to converge to any fixed matrix, proofs can become very technical (Fahrmeir and Tutz, 1994, Appendix A.2). Imposing a limiting value for  $n^{-1} \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)}$  simplifies arguments and can be seen as a compromise between making strong and weak assumptions about the covariates (Fahrmeir and Tutz, 1994, p.46). There is an analogous motivation for Assumption 2, the limiting matrix  $\mathbf{Q}$  is present to avoid specifying a probability model for the source dataset, without overcomplicating the mathematical analysis.

Setting up a limit theorem requires a little extra care with notation. As we have a sequence of datasets  $\mathbf{A}_{(n)}$ , there is a corresponding sequence of optimal least squares coefficients  $\beta_F^{(n)}$ . Similarly, there is a sequence of squared residual errors  $RSS_F^{(n)}$  and model sum of squares  $MSS_F^{(n)}$ . As the sequence of datasets are fixed,  $\beta_F^{(n)}$ ,  $RSS_F^{(n)}$  and  $MSS_F^{(n)}$  are a deterministic sequence.

Under Assumptions 1 and 2, it is possible to establish an asymptotic result for  $\beta_S$  and  $\beta_P$ .

**Theorem 5.** *Suppose that Assumptions 1 and 2 hold,  $k \geq p$ , and  $\beta_S$  is computed using a Hadamard or Clarkson-Woodruff sketch. Let  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^+$  denote the Moore-Penrose pseudo-inverse of  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})$ .*

Let

$$\tilde{\mathbf{V}}_{(n)} = \frac{RSS_F^{(n)}}{k} \left( \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \right)^+ \quad \text{and} \quad \mathbf{V}_{(n)} = \frac{RSS_F^{(n)}}{k-p+1} \left( \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \right)^{-1}.$$

Then as  $n \rightarrow \infty$ , convergence in distribution holds for

$$\begin{aligned} (i) & [\mathbf{V}_{(n)}^{-1/2} (\boldsymbol{\beta}_S - \boldsymbol{\beta}_F^{(n)}) \mid \mathbf{A}_{(n)}] \rightarrow \text{Student}(\mathbf{0}, \mathbf{I}_p, k-p+1), \\ (ii) & [\tilde{\mathbf{V}}_{(n)}^{-1/2} (\boldsymbol{\beta}_S - \boldsymbol{\beta}_F^{(n)}) \mid \mathbf{A}_{(n)}] \rightarrow N(\mathbf{0}, \mathbf{I}_p). \end{aligned}$$

For large  $n$ , we expect  $\boldsymbol{\beta}_S$  to be approximately distributed as per Theorem 5 for both the Hadamard and Clarkson-Woodruff sketches.

It is harder to establish a comparable limit theorem for  $\boldsymbol{\beta}_P^*$ , due to the non-standard distribution of  $\boldsymbol{\beta}_P^*$  when using a Gaussian sketch. There is no typical normalised distribution to target. Instead, we wish to show asymptotic equivalence in moments. The partially sketched estimator under the Hadamard and Clarkson-Woodruff sketches should have similar mean and variance properties to the Gaussian partially sketched estimator. An extra assumption has to be made to show convergence in moments. A sufficient condition is a stability condition on the singular values of the sketched data matrix.

**Assumption 3.** Let  $\mathbf{G}$  be the Gram matrix of the scaled sketched dataset,  $\mathbf{G} = n^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$ . Assume that the sequence of source datasets is such that  $E_S \left( \frac{1}{\sigma_{\min}^2(\mathbf{G})} \right)^2$  is finite for large enough  $n$ .

This additional regularity condition enables a formal limit theorem regarding the moments of  $\boldsymbol{\beta}_P^*$ .

**Theorem 6.** Suppose that Assumptions 1, 2 and 3 hold,  $k > p + 3$ , and  $\boldsymbol{\beta}_P^*$  is computed using a Hadamard or Clarkson-Woodruff sketch. Let

$$\mathbf{V}_{(n)} = \frac{(k-p-1)}{(k-p)(k-p-3)} \left( MSS_F^{(n)} (\mathbf{X}_{(n)}^\top \mathbf{X}_{(n)})^{-1} + \frac{k-p+1}{k-p-1} \boldsymbol{\beta}_F^{(n)} \boldsymbol{\beta}_F^{(n)\top} \right).$$

Then as  $n \rightarrow \infty$ ,

$$\begin{aligned} (i) & E_S(\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F^{(n)} \mid \mathbf{A}_{(n)}) \rightarrow \mathbf{0}. \\ (ii) & \text{var}_S \left( \mathbf{V}_{(n)}^{-1/2} (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F^{(n)}) \mid \mathbf{A}_{(n)} \right) \rightarrow \mathbf{I}_d \end{aligned}$$

Once again, the heavy notation may obscure the essence of the result. The subscript  $S$  is used to emphasise that the only source of randomness is the sketching matrix, and that the source dataset is fixed. The theorem suggests that the bias and variance of  $\boldsymbol{\beta}_P^*$  under the Clarkson-Woodruff and Hadamard sketches should be approximately equal to that under the Gaussian sketch. Specifically, we expect equations (8) and (9) to be good approximations for the variance of the sketched estimators using the Hadamard or Clarkson-Woodruff sketches.

The results here are meant to be useful heuristics to assess the uncertainty attached to the output of the randomised approximation algorithm. There is a need to communicate and quantify

the approximation error of sketching algorithms to end users (Lopes et al., 2018; Dobriban and Liu, 2018), and the asymptotic results developed in this section can be of use.

## 5 Unconditional results

So far we have treated the source dataset as fixed to isolate the approximation error introduced by the random projection. When sketching is used for statistical inference, we can extend the hierarchical model of Section 3.1 to include a source of variation at the population level. We take the design matrix  $\mathbf{X}$  as fixed and treat the response  $\mathbf{y}$  as random. We take the data generating process to be  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \varepsilon$ , where  $\varepsilon$  is a vector of  $n$  independently and identically distributed random variables with mean zero and variance  $\sigma^2$ . Let  $\gamma^2$  represent the average mean function sum of squares, so  $\gamma^2 = \|\mathbf{X}\boldsymbol{\beta}_0\|_2^2/n$ . At the population level, the ordinary least squares estimator satisfies (Searle, 1997),

$$\begin{aligned} E_y(\boldsymbol{\beta}_F) &= \boldsymbol{\beta}_0, \\ \text{var}_y(\boldsymbol{\beta}_F) &= \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}, \\ E_y(RSS_F) &= (n-p)\sigma^2, \\ E_y(MSS_F) &= p\sigma^2 + n\gamma^2. \end{aligned}$$

Taking iterated expectations, we can see that the Gaussian sketch gives an unbiased estimator of the population parameter  $\boldsymbol{\beta}_0$ ,

$$\begin{aligned} E_y(\boldsymbol{\beta}_S) &= E_y\{E_S(\boldsymbol{\beta}_S \mid \mathbf{y}, \mathbf{X})\} \\ &= E_y(\boldsymbol{\beta}_F) \\ &= \boldsymbol{\beta}_0 \end{aligned}$$

The unconditional variance of the Gaussian sketch can be obtained using the law of total variance,

$$\begin{aligned} \text{var}_y(\boldsymbol{\beta}_S) &= E_y\{\text{var}_S(\boldsymbol{\beta}_S \mid \mathbf{y}, \mathbf{X})\} + \text{var}_y\{E_S(\boldsymbol{\beta}_S \mid \mathbf{y}, \mathbf{X})\} \\ &= E_y\left(\frac{RSS_F}{(k-p+1)}(\mathbf{X}^\top \mathbf{X})^{-1}\right) + 0 \\ &= \frac{(n-p)\sigma^2}{(k-p+1)}(\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned} \tag{16}$$

We can also determine the unconditional properties of the partial sketch estimator  $\boldsymbol{\beta}_P^*$ . The estimator is also unbiased for  $\boldsymbol{\beta}_0$ ,

$$\begin{aligned} E_y(\boldsymbol{\beta}_P^*) &= E_y\{E_S(\boldsymbol{\beta}_P^* \mid \mathbf{y}, \mathbf{X})\} \\ &= E_y(\boldsymbol{\beta}_F) \\ &= \boldsymbol{\beta}_0. \end{aligned}$$

The unconditional variance of  $\beta_P^*$  is

$$\text{var}_y(\beta_P^*) = E_y \{ \text{var}_S(\beta_P^* | \mathbf{y}, \mathbf{X}) \} + \text{var}_y \{ E_S(\beta_P^* | \mathbf{y}, \mathbf{X}) \} \quad (17)$$

$$= E_y \left\{ \frac{(k-p-1)}{(k-p)(k-p-3)} \left( MSS_F(\mathbf{X}^\top \mathbf{X})^{-1} + \frac{k-p+1}{k-p-1} \beta_F \beta_F^\top \right) \right\} + 0 \quad (18)$$

$$= \frac{(k-p-1)}{(k-p)(k-p-3)} \left\{ (p\sigma^2 + n\gamma^2)(\mathbf{X}^\top \mathbf{X})^{-1} + \left( \frac{k-p+1}{k-p-1} \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} + \frac{k-p+1}{k-p-1} \beta_0 \beta_0^\top \right) \right\}. \quad (19)$$

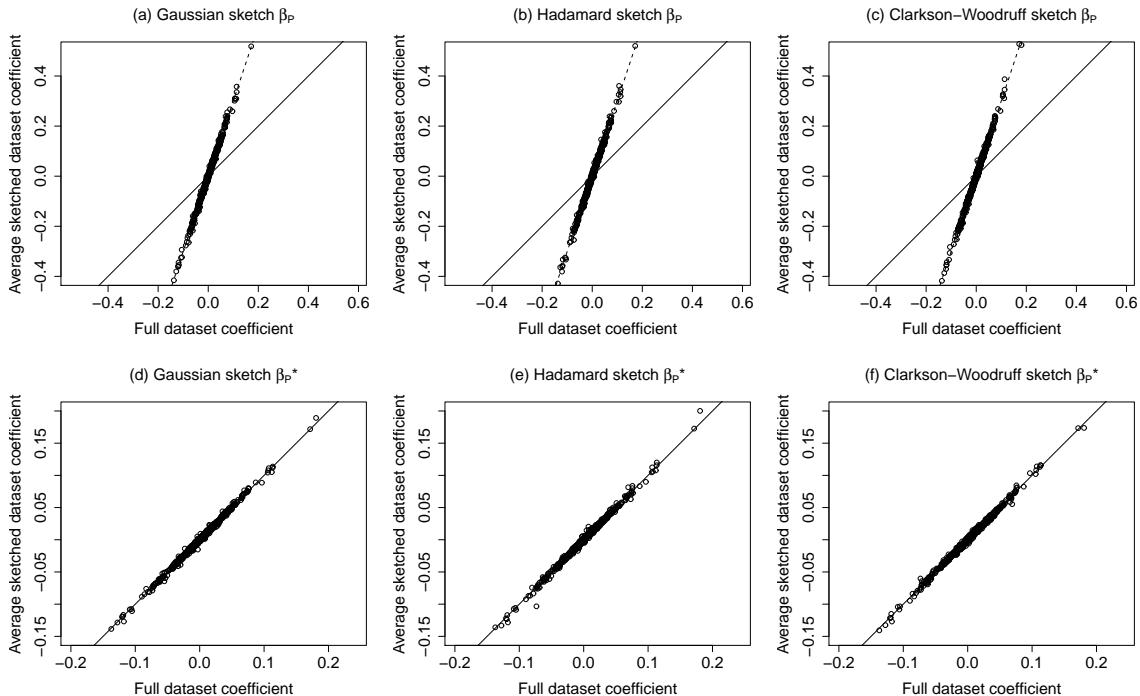
The most significant terms in the unconditional variance of  $\beta_S$  are  $n\sigma^2$  and  $(\mathbf{X}^\top \mathbf{X})^{-1}$ . The dominating terms in the unconditional variance of  $\beta_P^*$  are  $(\mathbf{X}^\top \mathbf{X})^{-1}$  and  $n\gamma^2 = \|\mathbf{X}\beta_0\|_2^2$ . We reach similar conclusions to the conditional analysis, in that we expect  $\beta_S$  to be more efficient when the signal to noise ratio is high, and  $\beta_P^*$  to be more efficient when the signal to noise ratio is low. Under Assumptions 1, 2 and 3, the variance expression give asymptotic approximations for the Hadamard and Clarkson-Woodruff projections. These results can be extended to account for more complicated error models on  $\varepsilon$  if it is still possible to determine  $E_y(\beta_F)$ ,  $\text{var}_y(\beta_F)$ ,  $E_y(RSS_F)$  and  $E_y(MSS_F)$ . In independent work, Chi and Ipsen (2018), also study the error rates of sketched regression, and additionally consider cases where the sketched design matrix does not have same rank as the full design matrix.

## 6 Data application

### 6.1 Human leukocyte antigen dataset

We compared the performance of the sketching estimators on a real genetic dataset taken from the UK Biobank database. We use a small extract from the data in Astle et al. (2016). The selected response variable was mean red cell volume (MCV), taken from the full blood count assay and adjusted for various technical and environmental covariates. Genome-wide imputed genotype data in expected allele dose format were available on  $n = 132353$  study subjects (Howie et al., 2009). We consider 1000 genetic variants in the Human leukocyte antigen (HLA) region of chromosome 6, selected so that no pair of variants had Pearson correlation of allelic scores greater than 0.8. The region was chosen as many associations were discovered in a genome-wide scan using univariable models; these associations were with variants with different allele frequencies, suggesting multiple distinct causal variants in the region. The aim is to perform a multivariable regression analysis to obtain variant effect size estimates that are conditional on the other variants in the region.

An early theoretical finding was that the partial sketched estimator  $\beta_P$  was biased. One thousand sketches were taken to estimate the bias  $E_S(\beta_P - \beta_F)$  with  $k = 1500$ . We also computed the bias corrected estimator  $\beta_P^*$  in each replication. Figure 2 plots the average value of the estimators against the true value of the least squares coefficient using the full dataset. The top row (a)-(c) shows results for  $\beta_P$ , and the bottom row (d)-(f) shows results for  $\beta_P^*$ . The first, second and third columns display the results for the Gaussian, Hadamard and Clarkson-Woodruff sketches respectively. The solid line in each panel is the identity line. The dashed line in the top row shows the theoretical bias, having slope  $k/(k-p-1)$ .



**Figure 2:** Bias of sketching estimators on the HLA dataset. Mean estimates are plotted against true values. In this scenario  $n = 132353, p = 1000, k = 1500$ . Solid line is the identity line and dashed line represents the theoretical bias factor.

The results in the top row show that  $\beta_P$  is biased for each of the random projections. The bias closely matches the theoretical factor. The bottom row shows that the adjusted estimator  $\beta_P^*$  appears to be unbiased, with the mean values falling closely along the identity line.

We also compared the complete and partially sketched estimators on mean square error and the coverage of confidence intervals at  $k = 1500$  and  $k = 10000$ . We also compared the data oblivious sketches to simple uniform subsampling with replacement. Simple random sampling is often referred to as the uniform sketch in the literature. We did not consider a combined estimator as the small  $R_F^2$  value would mean give an optimal complete sketching weight of close to zero. Table 2 reports the mean square error for each of the estimators. The signal to noise ratio is quite low for this dataset with  $R_F^2 = 0.02$ . We expect that partial sketching will be much more efficient than complete sketching on this dataset given the low signal to noise ratio. The simulation results support this idea, with  $\beta_P^*$  having a mean square error roughly sixty times smaller than  $\beta_S$  at both values of  $k$ . Results are very similar for each of the random projections, suggesting that the asymptotic approximations are reasonable for this dataset. For  $k = 1500$ , the mean square error of  $\beta_P$  is approximately ten times that of  $\beta_P^*$ . For  $k = 10000$ , there is less of a difference, as the ratio  $k/(k - p - 1)$  is closer to one. The bias adjusted estimator  $\beta_P^*$  has significant advantages over  $\beta_P$  when  $k/(k - p - 1)$  is larger than one.

Table 3 summarises the coverage of 95% confidence intervals for the sketched estimators. We report the overall proportion of intervals that contained the true value of the least squares estimate  $\beta_F$  over the two hundred and fifty sketches and  $p = 1000$  coefficients. The observed coverage is

	$k = 1500$			$k = 10000$		
	$\beta_S$	$\beta_P$	$\beta_P^*$	$\beta_S$	$\beta_P$	$\beta_P^*$
Gaussian	238 (3)	39 (0.7)	3.8 (0.08)	13.3 (0.17)	0.28 (0.004)	0.21 (0.002)
Hadamard	238 (4)	39 (0.7)	3.8 (0.07)	12.5 (0.16)	0.26 (0.003)	0.20 (0.002)
Clarkson-Woodruff	241 (3)	38 (0.8)	4.0 (0.05)	13.2 (0.16)	0.28 (0.004)	0.21 (0.002)
Uniform	375 (15)	105 (7.6)	10.7 (0.55)	13.8 (0.20)	0.38 (0.007)	0.29 (0.005)

**Table 2:** Mean square error of sketched estimators on HLA dataset. Standard errors are in brackets.

	$k = 1500$		$k = 10000$	
	$\beta_S$	$\beta_P^*$	$\beta_S$	$\beta_P^*$
Gaussian	0.950	0.953	0.950	0.951
Hadamard	0.949	0.949	0.954	0.954
Clarkson-Woodruff	0.947	0.952	0.951	0.950

**Table 3:** Coverage of confidence intervals on the HLA dataset. The largest standard error is 0.002

close the nominal level of 0.95 at both levels of  $k$ . The different random projections give very similar results, suggesting that the use of asymptotic approximations is again reasonable on this dataset. The intervals for the Hadamard sketch appear to be slightly conservative at  $k = 10000$ .

Table 4 reports the average sketching time for the data oblivious sketches. We computed ten sketches using each projection. The Gaussian sketch is an order of magnitude slower than the Hadamard projection and two orders of magnitude slower than the Clarkson-Wooduff sketch. The Gaussian sketch also scales more poorly as  $k$  increases, as is expected from Table 1.

## 6.2 Flights dataset

The sketching algorithms were also evaluated on the New York flights dataset available in the R package `nycflights13` (Wickham, 2014). Arrival delay was taken as the response, and departure delay, distance, departure time, origin and month and day were chosen to be the covariates. Rows of the dataset with missing data were omitted, leaving  $n = 327346$  and  $d = 47$ . The goal was to compare the accuracy of the various sketches on real data rather than to build a statistical model for the flights dataset. We compared the mean square error of the estimators and the coverage of confidence intervals for  $k = 5000$ . In contrast to the HLA dataset, the flights dataset has a very high  $R_F^2$  value of 0.99. We took five hundred sketches to compare complete and partial sketching.

Table 5 reports the mean square error of  $\beta_S, \beta_P$  and  $\beta_P^*$ . As expected, complete sketching has a much smaller mean square error than partial sketching. Table 6 summarises the coverage rates of

	$k = 1500$	$k = 10000$
Gaussian	522	3479
Hadamard	57	65
Clarkson-Woodruff	5.3	5.4

**Table 4:** Timings for sketching the HLA dataset in seconds. We report the average time to compute the sketched dataset  $\tilde{\mathbf{A}} = \mathbf{S}\mathbf{A}$ .

	$\beta_S$	$\beta_P$	$\beta_P^*$
Gaussian	60 (2)	14900 (400)	14900 (400)
Hadamard	63 (2)	14800 (500)	13900 (400)
Clarkson-Woodruff	66 (2)	15000 (500)	13800 (400)
Uniform	64 (2)	14600 (500)	14600 (400)

**Table 5:** Mean square error of sketched estimators on flights dataset with  $k = 5000$ . Standard errors are in brackets.

	$\beta_S$	$\beta_P^*$
Gaussian	0.948	0.951
Hadamard	0.950	0.948
Clarkson-Woodruff	0.948	0.947

**Table 6:** Coverage of 95% confidence intervals on the flights dataset with  $k = 5000$ . The largest standard error is 0.004

the 95% confidence intervals. We report the overall proportion of intervals that contained the true value of the least squares estimate over the five hundred sketches and  $p = 46$  coefficients.

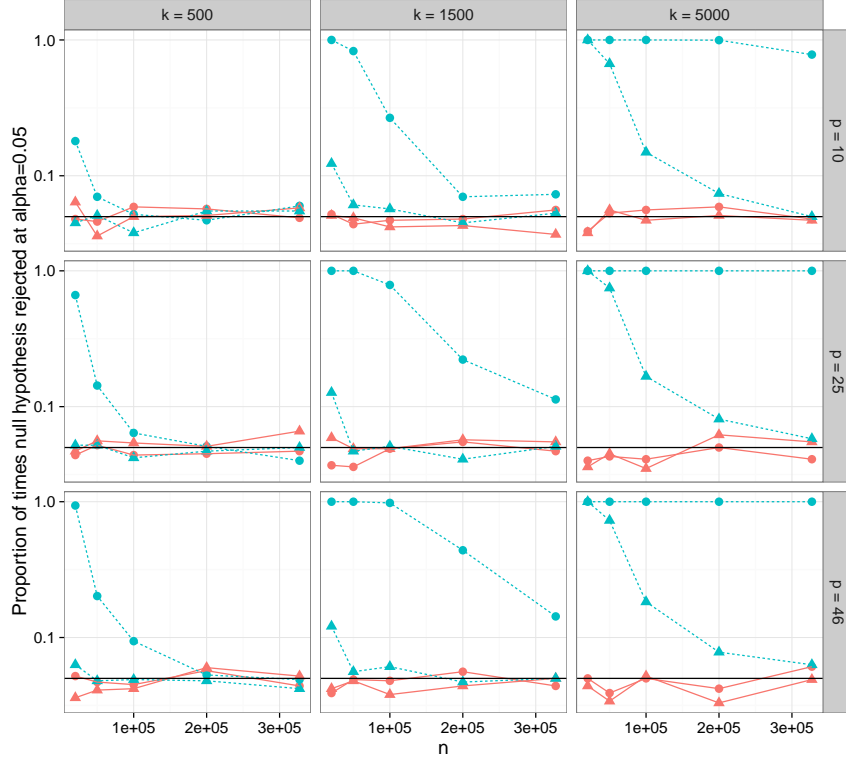
Table 7 reports the average sketching time for the data oblivious random projections. We generated ten sketches with each method. The Gaussian sketch is again considerably slower to apply than the Hadamard and Clarkson-Woodruff projections.

We also assessed the finite sample behaviour of the normal approximation in Theorem 3 at different levels of  $k$  and  $p$ . We dropped some predictors from the full flights dataset to give smaller datasets with  $p = 10$  and  $p = 25$  covariates. We then took subsamples of different sizes from each of the datasets. A single subsample was taken at each value of  $n$ , so the same subsampled dataset was being sketched each time. One thousand sketches were taken of each dataset at different values of  $k$ . We tested the joint multivariate normality of  $[\tilde{\mathbf{y}}, \tilde{\mathbf{X}}]$  and the normality of the sketched residual  $\tilde{\mathbf{e}} = \mathbf{S}(\mathbf{y} - \mathbf{X}\beta_F)$ . The squared Mahalanobis distance of the sketched observations was compared to the theoretical  $\chi^2$ -distribution. As  $n$  increases the rejection rate is expected to fall to the type one error rate of 0.05. Figure 3 plots the proportion of times the null hypothesis of normality is rejected against the size of the source dataset.

The Hadamard sketch appears to have a much faster rate of convergence than the Clarkson-Woodruff sketch. When using a Hadamard sketch, each row in the sketched dataset is a linear combination of  $n$  observations from the source dataset. When using a Clarkson-Woodruff sketch, each row in the sketched dataset is expected to be a combination of only  $n/k$  observations from the source dataset. As such,  $n/k$  must be large for the normal approximation to hold. As expected,

$k = 5000$	
Gaussian	404
Hadamard	5.8
Clarkson-Woodruff	0.2

**Table 7:** Timings for sketching the flights dataset in seconds. We report the average time to compute the sketched dataset  $\tilde{\mathbf{A}} = \mathbf{S}\mathbf{A}$ .



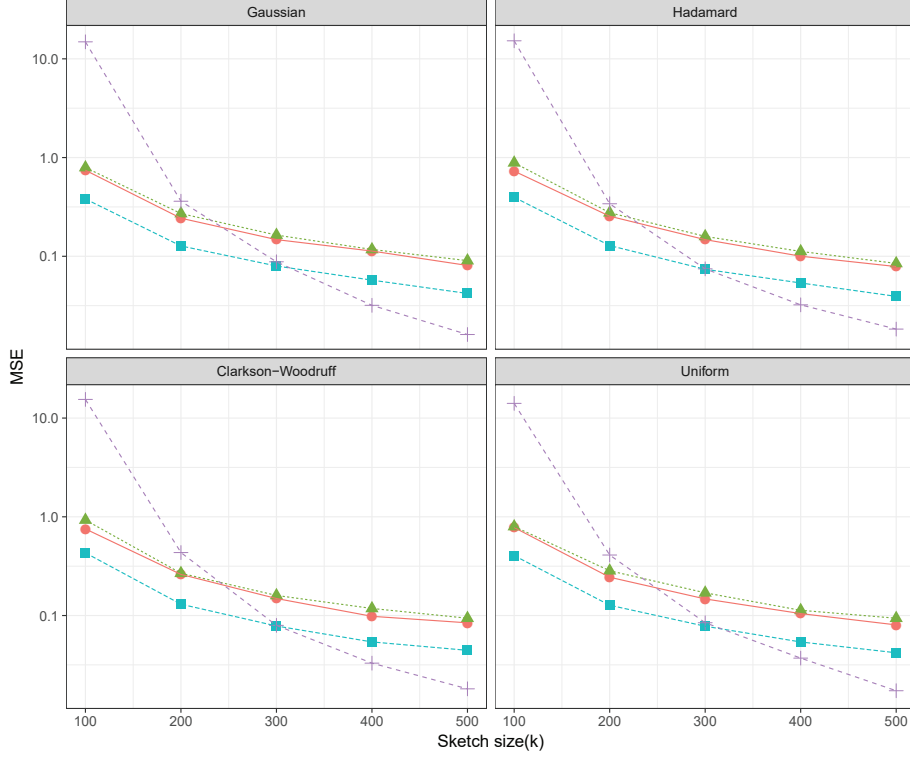
**Figure 3:** Proportion of times null hypothesis of normality is rejected against size of the source dataset ( $n$ ) for the Hadamard (solid line) and Clarkson-Woodruff sketches (dashed line). Results for tests of the sketched residual vector  $\tilde{\mathbf{e}} = \mathbf{S}(\mathbf{y} - \mathbf{X}\beta_F)$  are plotted as circles ( $\circ$ ), and results for tests of the entire sketched dataset  $[\tilde{\mathbf{y}}, \tilde{\mathbf{X}}]$  are plotted as triangles ( $\Delta$ ). The horizontal line gives the type 1 error of 0.05. The  $y$ -axis is on a log scale.

the rejection rate for the Clarkson-Woodruff sketch increases with  $k$ , but remains stable for the Hadamard sketch. In Fig. 3 the rejection rate for the Clarkson-Woodruff sketch increases with  $p$ . The Hadamard sketch seems to be less sensitive to the number of covariates. The extra  $\log k$  computation cost associated with the Hadamard sketch (Table 1) appears to have the benefit of accelerated convergence to normality. Even though joint normality may not be holding for the Clarkson-Woodruff sketch for the flights dataset, the coverage of the confidence intervals is still very good. As  $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta_F + \tilde{\mathbf{e}}$ , normality of the sketched residual is perhaps sufficient in justifying the approximate confidence intervals using Theorem 5 (ii). The sketched residual converges much more quickly than the full sketched data matrix, which perhaps explains the good coverage properties of the confidence intervals for  $\beta_S$  in Table 6.

### 6.3 Synthetic data

We also generated a synthetic dataset with  $n = 10000$ ,  $p = 50$  and an  $R_F^2$  of close to 0.5. The dataset consisted of responses  $y_i$  and covariates  $\mathbf{x}_i$ , ( $i = 1, \dots, n$ ). Covariates  $\mathbf{x}_i$  were drawn from a multivariate normal distribution with mean zero and covariance matrix  $\Sigma$ , with elements  $\Sigma_{ij} = 0.5^{|i-j|}$ . Responses were simulated independently using the standard linear model  $y_i = \mathbf{x}_i^T \beta_0 + \epsilon_i$  ( $i = 1, \dots, n$ ),





**Figure 4:** Comparison of sketching estimators on synthetic dataset with  $R_F^2 \approx 0.5$ . The  $y$ -axis is on a log scale. The average squared error for the sketching estimator is plotted against sketch size. Results are shown for  $\beta_S$  ( $\circ$ ),  $\beta_P^*$  ( $\Delta$ ), the weighted combined estimator  $\beta_C$  ( $\square$ ) and the one-step estimator  $\beta_H$  ( $+$ ).

where  $\epsilon_i$  is distributed as  $N(0, 0.45)$ . Each element of  $\beta_0$  was sampled independently from a  $N(0, 0.01)$  distribution. We compared the single pass estimators  $\beta_S$ ,  $\beta_P^*$  to the combined estimator  $\beta_C$  with the optimal weight  $\phi_{\text{opt}}$ , and the one-step estimator  $\beta_H$ . We applied the Gaussian, Hadamard, Clarkson-Woodruff and uniform subsampling sketches. We computed one hundred sketches at a range of sketch sizes  $k$ . We calculated the conditional sketching error  $\|\hat{\beta} - \beta_F\|_2^2$  for each sketched estimator  $\hat{\beta}$  in each replicate. Figure 4 plots the average error for the estimators  $\beta_S$ ,  $\beta_P^*$ ,  $\beta_C$  and  $\beta_H$  against the sketch size  $k$ . As expected, the combined estimator  $\beta_C$  has a mean square error that is roughly half that of  $\beta_S$  or  $\beta_P^*$  at all sketch sizes  $k$ . When  $k/p$  is small, the one-step estimator  $\beta_H$  has a higher mean square error than the single pass estimator  $\beta_S$ . As the ratio  $k/p$  increases, the one-step estimator  $\beta_H$  becomes more efficient than the weighted estimator  $\beta_C$ . This phenomenon can be studied in more detail using the moment results in Letac and Massam (2004). The results are similar for each of the data oblivious projections, suggesting that the asymptotic approximations are reasonable for this dataset. The uniform projection behaves similarly to the Gaussian projection, this is expected given that the covariates were simulated from a multivariate normal distribution.

## 7 Discussion

Sketching algorithms have emerged in the computer science community as a powerful device for the analysis of massive datasets (Mahoney and Drineas, 2016). Sketched regression algorithms use random projections to reduce the size of the original dataset, the sketched dataset is then used to estimate the optimal least squares coefficients. Most existing theory for sketched regression is from an algorithmic worst case perspective, and connects with random matrix theory and computational geometry (Raskutti and Mahoney, 2014; Thanei et al., 2017). In this paper we have provided a complementary statistical perspective and derived new tools for assessing the uncertainty attached to sketched estimators, as well as guidelines for choosing between competing sketching algorithms.

The sketching central limit theorem was essential in establishing the asymptotic behaviour of the Clarkson-Woodruff and Hadamard projections. The regularity condition on the limiting leverage scores of the source dataset connects with both the existing computer science literature on sketching, and classic central limit theorems from the statistics literature. The field of randomised algorithms is clearly at the interface of computer science and statistics, and it is pleasing to see some overlap in the fundamental theory underpinning a Big Data algorithm. It is also possible to use other methods to develop uncertainty tools for randomised algorithms. (Lopes et al., 2018) use the nonparametric bootstrap, and (Dobriban and Liu, 2018) use asymptotic results in random matrix theory. Together, these provide a practical suite of tools for end users.

Iterative methods, in particular stochastic gradient descent, have not been mentioned so far. For large  $n$  regression problems, stochastic gradient descent will produce iterates that converge to  $\beta_F$  under very mild conditions. Comparisons between single pass sketching and stochastic gradient methods are difficult, as the two techniques are not formulated for the exact same purpose. Single pass sketching algorithms are designed to return an approximate solution in finite time with probabilistically controlled error, whereas stochastic gradient methods are designed to converge to the exact solution asymptotically. It is perhaps more appropriate to compare stochastic gradient descent to iterative sketching methods, as iterative sketching algorithms also come with convergence guarantees to  $\beta_F$  (Pilanci and Wainwright, 2016; Gower and Richtik, 2015). Iterative sketching methods make use of approximate second order information that can lead to a potential improvement compared to first order stochastic gradient methods (Roosta-Khorasani and Mahoney, 2016). Our focus has been on characterising the approximation error attached to single pass sketching estimators.

There has been recent work in adapting sketching methods for statistical inference in large datasets, building from the worst case bounds in the computer science literature. Geppert et al. (2017) and Bardenet and Maillard (2015) investigate sketching algorithms for Bayesian regression, and derive bounds on the difference between the sketched posterior distribution and the full data posterior distribution. Yang et al. (2015b) consider sketched penalised regression, and give bounds between the sketched solution and the full data solution similar to the results in Section 2.2. Only complete sketching is considered in the aforementioned work. The results on the advantages of partial sketching in this paper could motivate adaptations that make use of the exact marginal associations  $\mathbf{X}^\top \mathbf{y}$ .

Sketching ideas have been used to develop methods for approximate non-linear regression (Avron et al., 2014; Banerjee et al., 2013). A related branch of work uses random projections to reduce the number

of predictors in regression and classification problems (Shah and Meinshausen, 2013; Cannings and Samworth, 2015; Guhaniyogi and Dunson, 2015).

## Acknowledgement

This work has been conducted using the UK Biobank resource under applications number 13745. Many thanks to Rajen Shah for helpful discussions.

## References

- Ailon, N. and Chazelle, B. (2009) The fast Johnson Lindenstrauss transform and approximate nearest neighbors. *SIAM Journal on Computing*, **39**, 302–322.
- Anderson, I. (1997) *Combinatorial Designs and Tournaments*. Oxford lecture series in mathematics and its applications. Clarendon Press.
- Astle, W. J., Elding, H., Jiang, T., Allen, D., Ruklisa, D., Mann, A. L., Mead, D., Bouman, H., Riveros-Mckay, F., Kostadima, M. A. et al. (2016) The allelic landscape of human blood cell trait variation and links to common complex disease. *Cell*, **167**, 1415–1429.
- Avron, H., Nguyen, H. and Woodruff, D. (2014) Subspace embeddings for the polynomial kernel. In *Advances in Neural Information Processing Systems*, 2258–2266.
- Banerjee, A., Dunson, D. B. and Tokdar, S. T. (2013) Efficient Gaussian process regression for large datasets. *Biometrika*, **100**, 75–89.
- Bardenet, R. and Maillard, O.-A. (2015) A note on replacing uniform subsampling by random projections in MCMC for linear regression of tall datasets. *HAL preprint 01248841*.
- Becker, S., Kawas, B., Petrik, M. and Ramamurthy, K. (2015) Robust partially-compressed least-squares. *arXiv: 1510.04905v1*.
- Billingsley, P. (1968) *Convergence of Probability Measures*. Wiley.
- Billingsley, P. (1995) *Probability and Measure*. Wiley Series in Probability and Statistics. Wiley.
- Cannings, T. I. and Samworth, R. J. (2015) Random projection ensemble classification. *arXiv: 1504.04595*.
- Chi, J. T. and Ipsen, I. C. F. (2018) Randomized least squares regression: Combining model- and algorithm-induced uncertainties. *arXiv e-prints*, arXiv:1808.05924.
- Clarkson, K. L. and Woodruff, D. P. (2013) Low rank approximation and regression in input sparsity time. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, 81–90. ACM.
- Cormode, G. (2011) Sketch techniques for approximate query processing. *Foundations and Trends in Databases*.

- Dhillon, P., Lu, Y., Foster, D. P. and Ungar, L. (2013) New subsampling algorithms for fast least squares regression. In *Advances in Neural Information Processing Systems*, 360–368.
- Diaconis, P. and Freedman, D. (1984) Asymptotics of graphical projection pursuit. *Annals of Statistics*, **12**, 793–815.
- Dobriban, E. and Liu, S. (2018) A new theory for sketching in linear regression. *arXiv e-prints*, arXiv:1810.06089.
- Drineas, P., Mahoney, M. W. and Muthukrishnan, S. (2006) Sampling algorithms for  $\ell_2$  regression and applications. In *Proceedings of the seventeenth annual ACM-SIAM symposium on Discrete algorithm*, 1127–1136. Society for Industrial and Applied Mathematics.
- Eaton, M. (2007) *Multivariate Statistics: A Vector Space Approach*. Institute of Mathematical Statistics.
- Fahrmeir, L. and Tutz, G. (1994) *Multivariate Statistical Modelling based on Generalized Linear Models*. Springer series in statistics. Springer-Verlag.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A. and Rubin, D. B. (2014) *Bayesian Data Analysis*. Boca Raton: Chapman & Hall, 3 edn.
- Geppert, L. N., Ickstadt, K., Munteanu, A., Quedenfeld, J. and Sohler, C. (2017) Random projections for Bayesian regression. *Statistics and Computing*, **27**, 79–101.
- Gower, R. M. and Richtrik, P. (2015) Randomized iterative methods for linear systems. *arXiv:1506.03296*.
- Greene, W. (1997) *Econometric Analysis*. Prentice-Hall international editions. Prentice Hall.
- Guhaniyogi, R. and Dunson, D. B. (2015) Bayesian compressed regression. *Journal of the American Statistical Association*, **110**, 1500–1514.
- Halko, N., Martinsson, P. G. and Tropp, J. A. (2011) Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Review*, **53**, 217–288.
- Hansen, M. H. and Hurwitz, W. N. (1943) On the theory of sampling from finite populations. *The Annals of Mathematical Statistics*, **14**, 333–362. URL <https://doi.org/10.1214/aoms/1177731356>.
- Howie, B. N., Donnelly, P. and Marchini, J. (2009) A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet*, **5**, e1000529.
- Letac, G. and Massam, H. (2004) All invariant moments of the wishart distribution. *Scandinavian Journal of Statistics*, **31**, 295–318.
- Li, P., Hastie, T. J. and Church, K. W. (2006) Very sparse random projections. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, 287–296. ACM.
- Loeve, M. (1977) *Probability Theory*. Springer.

- Lopes, M. E., Wang, S. and Mahoney, M. W. (2018) Error Estimation for Randomized Least-Squares Algorithms via the Bootstrap. *arXiv preprint*, arXiv:1803.08021.
- Ma, P., Mahoney, M. W. and Yu, B. (2015) A statistical perspective on algorithmic leveraging. *Journal of Machine Learning Research*, 861–911.
- Ma, P. and Sun, X. (2015) Leveraging for big data regression. *Wiley Interdisciplinary Reviews: Computational Statistics*, **7**, 70–76.
- Mahoney, M. (2011) Randomized algorithms for matrices and data. *Foundations and Trends in Machine Learning*, **3**, 123–224.
- Mahoney, M. and Drineas, P. (2016) Structural properties underlying high-quality randomized numerical linear algebra algorithms. In *Handbook of Big Data* (eds. P. Buhlmann, P. Drineas, M. Kane and M. van de Laan), 137–154. Chapman and Hall.
- Meng, X. and Mahoney, M. M. (2013) Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, 91–100. ACM.
- Phillips, J. M. (2016) Coresets and Sketches. *arXiv: 1601.00617*.
- Pilanci, M. and Wainwright, M. J. (2016) Iterative Hessian sketch: Fast and accurate solution approximation for constrained least-squares. *Journal of Machine Learning Research*, **17**, 1842–1879.
- Pruss, A. R. and Szynal, D. (2000) On the central limit theorem for negatively correlated random variables with negatively correlated squares. *Stochastic Processes and their Applications*, **87**, 299 – 309.
- Raskutti, G. and Mahoney, M. (2014) A statistical perspective on randomized sketching for ordinary least-squares. *arXiv: 1406.5986*.
- Roosta-Khorasani, F. and Mahoney, M. W. (2016) Sub-sampled Newton methods i: Globally convergent algorithms. *arXiv: 1601.04737*.
- Sarlos, T. (2006) Improved approximation algorithms for large matrices via random projections. In *47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06)*, 143–152. IEEE.
- Searle, S. R. (1997) *Linear Models*. New Jersey: Wiley-Interscience.
- Shah, R. D. and Meinshausen, N. (2013) Min-wise hashing for large-scale regression and classification with sparse data. *arXiv: 1308.1269*.
- Shorack, G. R. (2000) *Probability for Statisticians*. Springer Texts in Statistics. Springer.
- Svante, J. (1988) Some pairwise independent sequences for which the central limit theorem fails. *Stochastics: An International Journal of Probability and Stochastic Processes*, **23**, 439–448.

- Thanei, G.-A., Heinze, C. and Meinshausen, N. (2017) Random projections for large-scale regression. *arXiv: 1701.05325*.
- Van Der Vaart, A. (1998) *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics, 3. Cambridge University Press.
- White, H. (1984) *Asymptotic Theory for Econometricians*. Economic Theory, Econometrics and Mathematical Economics Series. Academic Press.
- Wickham, H. (2014) *nycflights13: Data about flights departing NYC in 2013*. Rstudio. R package version 0.1.
- Woodruff, D. P. (2014) Sketching as a tool for numerical linear algebra. *Foundations and Trends in Theoretical Computer Science*, **10**, 1–157.
- Yang, J., Meng, X. and Mahoney, M. W. (2015a) Implementing randomized matrix algorithms in parallel and distributed environments. *arXiv: 1502.03032*.
- Yang, T., Zhang, L., Lin, Q. and Jin, R. (2015b) Fast sparse least-squares regression with non-asymptotic guarantees. *arXiv: 1507.05185*.

## Supplementary Information

### A Sketching examples

As examples, we demonstrate the construction of a Hadamard sketch and a Clarkson-Woodruff sketch, for  $k = 3$ ,  $n = 4$ .

The Hadamard sketch matrix is formed as  $\mathbf{S} = \Phi \mathbf{H} \mathbf{D} / \sqrt{k}$ , where  $\Phi$  is a  $k \times n$  matrix and  $\mathbf{H}$  and  $\mathbf{D}$  are both  $n \times n$  matrices. The fixed matrix  $\mathbf{H}$  is a Hadamard matrix of order  $n$ . The random matrix  $\mathbf{D}$  is a diagonal matrix where each nonzero element is an independent Rademacher random variable. The random matrix  $\Phi$  subsamples  $k$  rows of  $\mathbf{H}$  with replacement. The display below shows an example of the random projection. The first matrix in the display represents  $\Phi \mathbf{H}$ , a subsample of three rows from a  $4 \times 4$  Hadamard matrix. In step 2, the diagonal matrix  $\mathbf{D}$  is generated, with random Rademacher random variables along the diagonal. The diagonal elements are shown above the matrix. In step 3 the matrix multiplication  $\Phi \mathbf{H} \mathbf{D}$  is performed. This outputs the sketching matrix  $\mathbf{S}$ .

$$\begin{array}{c}
 \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{pmatrix} \xrightarrow{\text{step 2}} \begin{array}{cccc} +1 & -1 & +1 & +1 \\ D_{11} & D_{22} & D_{33} & D_{44} \\ \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{pmatrix} \xrightarrow{\text{step 3}} \begin{array}{cccc} +1 & -1 & +1 & +1 \\ \times & \times & \times & \times \\ \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{pmatrix} \xrightarrow{\text{output}} \begin{pmatrix} 1 & 1 & 1 & -1 \\ 1 & 1 & -1 & 1 \\ 1 & 1 & 1 & -1 \end{pmatrix}
 \end{array}
 \end{array}$$

The Clarkson-Woodruff sketch is a sparse random matrix. The projection can be represented as the product of two independent random matrices,  $\mathbf{S} = \mathbf{\Gamma} \mathbf{D}$ , where  $\mathbf{\Gamma}$  is a random  $k \times n$  matrix and  $\mathbf{D}$  is a random  $n \times n$  matrix. The matrix  $\mathbf{\Gamma}$  is formed by choosing one element in each column

independently and setting the entry to +1. The matrix  $\mathbf{D}$  is a diagonal matrix where each nonzero element is an independent Rademacher random variable. This results in a sparse  $\mathbf{S}$ , where there is only one nonzero entry per column. The display below shows an example of the random projection. The first matrix in the display represents  $\mathbf{\Gamma}$ , a random matrix where a single element in each column is set to one. In step 2, the diagonal matrix  $\mathbf{D}$  is generated, with random Rademacher random variables along the diagonal. The diagonal elements are shown above the matrix. In step 3 the matrix multiplication  $\mathbf{\Gamma D}$  is performed. This outputs the sketching matrix  $\mathbf{S}$ .

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \xrightarrow{\text{step 2}} \begin{matrix} -1 & +1 & -1 & +1 \\ D_{11} & D_{22} & D_{33} & D_{44} \\ \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \end{matrix} \xrightarrow{\text{step 3}} \begin{matrix} +1 & -1 & +1 & +1 \\ \times & \times & \times & \times \\ \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{pmatrix} \end{matrix} \xrightarrow{\text{output}} \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & -1 & 1 & 0 \end{pmatrix}$$

## B Proof of Theorem 1

**Theorem 7.** *Suppose that  $\widetilde{\mathbf{X}}$  is an  $\epsilon$ -subspace embedding of  $\mathbf{X}$  with  $0 < \epsilon < 0.5$ . Then the following bound holds,*

$$\|\beta_P - \beta_F\|_2^2 \leq \frac{4\epsilon^2}{\sigma_{\min}^2(\mathbf{X})} MSS_F.$$

Let the singular value decomposition of  $\mathbf{X}$  be given by  $\mathbf{X} = \mathbf{U D V}^\top$ . The singular value decomposition will help to simplify expressions in later working. If the sketching matrix  $\mathbf{S}$  is an  $\epsilon$ -subspace embedding for the source dataset with  $0 < \epsilon < 1$ , then  $\mathbf{U}^\top \mathbf{S}^\top \mathbf{S U}$  is necessarily invertible. The expression for  $\beta_P$  can then be simplified to

$$\begin{aligned} \beta_P &= \mathbf{V D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S U})^{-1} \mathbf{D}^{-1} \mathbf{V}^\top \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{V D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S U})^{-1} \mathbf{D}^{-1} \mathbf{V}^\top \mathbf{V D U}^\top \mathbf{y} \\ &= \mathbf{V D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S U})^{-1} \mathbf{U}^\top \mathbf{y}. \end{aligned}$$

Similarly,  $\beta_F$  can be written as  $\beta_F = \mathbf{V D}^{-1} \mathbf{U}^\top \mathbf{y}$ . The Euclidean norm of the approximation error can thus be expressed as

$$\begin{aligned} \|\beta_P - \beta_F\|_2 &= \|\mathbf{V D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S U})^{-1} \mathbf{U}^\top \mathbf{y} - \mathbf{V D}^{-1} \mathbf{U}^\top \mathbf{y}\|_2 \\ &= \|\{\mathbf{V D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S U})^{-1} - \mathbf{V D}^{-1}\} \mathbf{U}^\top \mathbf{y}\|_2 \\ &= \|\{\mathbf{V D}^{-1} [(\mathbf{U}^\top \mathbf{S}^\top \mathbf{S U})^{-1} - \mathbf{I}_p]\} \mathbf{U}^\top \mathbf{y}\|_2. \end{aligned}$$

The model sum of squares can be written as

$$\begin{aligned}
MSS_F &= \|\mathbf{X}\boldsymbol{\beta}_F\|_2^2 \\
&= \|\mathbf{X}\mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top\mathbf{y}\|_2^2 \\
&= \|\mathbf{U}\mathbf{D}\mathbf{V}^\top\mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top\mathbf{y}\|_2^2 \\
&= \|\mathbf{U}\mathbf{U}^\top\mathbf{y}\|_2^2 \\
&= \|\mathbf{U}^\top\mathbf{y}\|_2^2.
\end{aligned} \tag{S.1}$$

The final line uses the fact that  $\mathbf{U}^\top\mathbf{U} = \mathbf{I}_p$ . Using the matrix norm induced by the Euclidean norm and the usual Euclidean norm for vectors we can form an upper bound on the error.

$$\begin{aligned}
\|\boldsymbol{\beta}_P - \boldsymbol{\beta}_F\|_2 &\leq \|\mathbf{V}\mathbf{D}^{-1}\{(\mathbf{U}^\top\mathbf{S}^\top\mathbf{S}\mathbf{U})^{-1} - \mathbf{I}_p\}\|_2\|\mathbf{U}^\top\mathbf{y}\|_2 \\
&\leq \|\mathbf{V}\mathbf{D}^{-1}\|_2\|\mathbf{U}^\top\mathbf{y}\|_2\|(\mathbf{U}^\top\mathbf{S}^\top\mathbf{S}\mathbf{U})^{-1} - \mathbf{I}_p\|_2 \\
&= \frac{MSS_F^{1/2}}{\sigma_{\min}(\mathbf{X})}\|(\mathbf{U}^\top\mathbf{S}^\top\mathbf{S}\mathbf{U})^{-1} - \mathbf{I}_p\|_2.
\end{aligned} \tag{S.2}$$

It remains to upper bound the maximum singular value of the matrix  $(\mathbf{U}^\top\mathbf{S}^\top\mathbf{S}\mathbf{U})^{-1} - \mathbf{I}_p$ . Let  $\mathbf{M} = \mathbf{U}^\top\mathbf{S}^\top\mathbf{S}\mathbf{U}$ . The maximum absolute value of the singular values of  $(\mathbf{U}^\top\mathbf{S}^\top\mathbf{S}\mathbf{U})^{-1} - \mathbf{I}_p$  will be given by  $\max(|1/\sigma_{\min}(\mathbf{M}) - 1|, |1/\sigma_{\max}(\mathbf{M}) - 1|)$ , where  $\sigma_{\min}(\mathbf{M})$  is the minimum singular value of  $\mathbf{M}$ , and  $\sigma_{\max}(\mathbf{M})$  is the maximum singular value of  $\mathbf{M}$ . If  $\mathbf{S}$  is an  $\epsilon$ -subspace embedding for the source covariate matrix  $\mathbf{X}$  then it must hold that  $\sigma_{\min}(\mathbf{M}) \geq 1 - \epsilon$ , and  $\sigma_{\max}(\mathbf{M}) \leq 1 + \epsilon$  (Woodruff, 2014, p.11). As such,  $\max(|1/\sigma_{\min}(\mathbf{M}) - 1|, |1/\sigma_{\max}(\mathbf{M}) - 1|) \leq |1/(1 - \epsilon) - 1|$ . It is simple to show that over the interval  $0 \leq \epsilon \leq 0.5$ ,  $|1/(1 - \epsilon) - 1| \leq 2\epsilon$ . This results in an upper bound on the singular value of interest,

$$\begin{aligned}
\|(\mathbf{U}^\top\mathbf{S}^\top\mathbf{S}\mathbf{U})^{-1} - \mathbf{I}_p\|_2 &\leq |1/(1 - \epsilon) - 1| \\
&\leq 2\epsilon.
\end{aligned}$$

Substituting this back into (S.2) gives that under the condition that  $\epsilon < 0.5$

$$\|\boldsymbol{\beta}_P - \boldsymbol{\beta}_F\|_2 \leq \frac{MSS_F^{1/2}}{\sigma_{\min}(\mathbf{X})} \times 2\epsilon.$$

Squaring both sides gives the final result, that if  $\epsilon < 0.5$

$$\|\boldsymbol{\beta}_P - \boldsymbol{\beta}_F\|_2^2 \leq \frac{4\epsilon^2}{\sigma_{\min}^2(\mathbf{X})} MSS_F.$$

## C Proof of Theorem 2 (Hierarchical model for the Gaussian sketch)

We use the following lemma about the Normal Inverse-Wishart distribution in many of our results (Gelman et al., 2014, p.73).

**Lemma 1.** *Suppose that  $\boldsymbol{\Sigma}$  is a random  $d \times d$  matrix and  $\mathbf{y}$  is a  $d$ -dimensional random vector from*



the following hierarchical model

$$\begin{aligned}\mathbf{y}|\Sigma &\sim N(\boldsymbol{\mu}, \Sigma/\kappa), \\ \Sigma &\sim \text{Inv-Wishart}(\Lambda, \nu),\end{aligned}$$

where  $\Lambda$  is a  $d \times d$  scale matrix,  $\nu$  is a scalar giving degrees of freedom, and  $\kappa$  is a scaling constant. Then marginally,

$$\mathbf{y} \sim \text{Student}(\boldsymbol{\mu}, \Lambda/(\kappa(\nu - d + 1)), \nu - d + 1).$$

Theorem 2 (ii) follows from setting  $\boldsymbol{\mu} = \boldsymbol{\beta}_F$ ,  $\Sigma = (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1}$ ,  $\kappa = k/RSS_F$ ,  $\Lambda = k(\mathbf{X}^\top \mathbf{X})^{-1}$ ,  $\nu = k$  and  $d = p$ . Theorem 2 (i) follows from standard results on linear models, for example see Searle (1997, Chapter 3).

## D Variance for partial sketching

Using a Gaussian sketch of size  $k$  where  $k > p + 3$ , the standard partial sketching estimator  $\boldsymbol{\beta}_P$  has variance

$$\text{var}(\boldsymbol{\beta}_P) = \frac{k^2}{(k-p)(k-p-1)(k-p-3)} \left( MSS_F(\mathbf{X}^\top \mathbf{X})^{-1} + \frac{(k-p+1)}{(k-p-1)} \boldsymbol{\beta}_F \boldsymbol{\beta}_F^\top \right). \quad (\text{S.3})$$

The bias corrected partial sketching estimator  $\boldsymbol{\beta}_P^*$  has variance

$$\text{var}(\boldsymbol{\beta}_P^*) = \frac{(k-p-1)}{(k-p)(k-p-3)} \left( MSS_F(\mathbf{X}^\top \mathbf{X})^{-1} + \frac{(k-p+1)}{(k-p-1)} \boldsymbol{\beta}_F \boldsymbol{\beta}_F^\top \right). \quad (\text{S.4})$$

We now prove (S.3) and (S.4).

Let the singular value decomposition of  $\mathbf{X}$  be given by  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ . The singular value decomposition will help to simplify expressions in later working. The sketched Gram matrix has the form  $\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} = \mathbf{V}\mathbf{D}\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U}\mathbf{D}\mathbf{V}^\top$ . As  $\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U} \sim \text{Wishart}(k, \mathbf{I}_p/k)$ , the matrix  $\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U}$  is almost surely invertible. The inverse Gram matrix can then be written as

$$\begin{aligned}(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^{-1} &= [\mathbf{D}\mathbf{V}^\top]^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U})^{-1} [\mathbf{V}\mathbf{D}]^{-1} \\ &= \mathbf{V}\mathbf{D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U})^{-1} \mathbf{D}^{-1} \mathbf{V}^\top.\end{aligned}$$

The expression for  $\boldsymbol{\beta}_P$  can then be simplified to

$$\begin{aligned}\boldsymbol{\beta}_P &= \mathbf{V}\mathbf{D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U})^{-1} \mathbf{D}^{-1} \mathbf{V}^\top \mathbf{X}^\top \mathbf{y} \\ &= \mathbf{V}\mathbf{D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U})^{-1} \mathbf{D}^{-1} \mathbf{V}^\top \mathbf{V}\mathbf{D}\mathbf{U}^\top \mathbf{y} \\ &= \mathbf{V}\mathbf{D}^{-1} (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U})^{-1} \mathbf{U}^\top \mathbf{y}.\end{aligned}$$

Let  $\mathbf{M} = (\mathbf{U}^\top \mathbf{S}^\top \mathbf{S}\mathbf{U})^{-1}$ . We know that  $\mathbf{M} \sim \text{Inverse-Wishart}(k, k\mathbf{I}_p)$ . Properties of the Inverse-

Wishart distribution give that that for  $i = 1, \dots, p$ ,

$$\text{var}(M_{ii}) = \frac{2k^2}{(k-p-1)^2(k-p-3)}. \quad (\text{S.5})$$

Additionally, for  $i, j = 1, \dots, p$ , where  $j \neq i$

$$\text{var}(M_{ij}) = \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)}. \quad (\text{S.6})$$

Finally we have that for  $i, j = 1, \dots, p$ ,  $i \neq j$ ,

$$\text{cov}(M_{ij}, M_{ji}) = \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)}, \quad (\text{S.7})$$

$$\text{cov}(M_{ii}, M_{jj}) = \frac{2k^2}{(k-p)(k-p-1)^2(k-p-3)}. \quad (\text{S.8})$$

All other covariances  $\text{cov}(M_{ij}, M_{br})$  are equal to zero unless they reduce to the cases in (S.7) or (S.8). Let  $\mathbf{z} = \mathbf{U}^\top \mathbf{y}$ . Let  $\mathbf{W} = \text{cov}(\mathbf{M}\mathbf{U}^\top \mathbf{y}) = \text{cov}(\mathbf{M}\mathbf{z})$ . The elements of  $\mathbf{W}$  can be determined using the properties in equations (S.5) to (S.8). Starting with the diagonal entries,

$$\begin{aligned} W_{ii} &= \text{var} \left( \sum_{j=1}^p M_{ij} z_j \right) \\ &= \sum_{j=1}^p z_j^2 \text{var}(M_{ij}) + \sum_{j=1}^p \sum_{w \neq j}^p z_j z_w \text{cov}(M_{ij}, M_{iw}). \end{aligned}$$

As  $\text{cov}(M_{ij}, M_{iw})$  is equal to zero for all  $w \neq j$  this simplifies to

$$\begin{aligned} W_{ii} &= \text{var} \left( \sum_{j=1}^p M_{ij} z_j \right) \\ &= \sum_{j=1}^p z_j^2 \text{var}(M_{ij}). \end{aligned}$$

It is helpful to split the sum into two pieces, a single term for  $j = i$  and then a sum over the remaining

indices. Grouping terms leads to an expression involving the model sum of squares  $MSS_F$ .

$$\begin{aligned}
W_{ii} &= z_i^2 \frac{2k^2}{(k-p-1)^2(k-p-3)} + \sum_{j=1, j \neq i}^p z_j^2 \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} \\
&= z_i^2 \frac{2k^2(k-p)}{(k-p)(k-p-1)^2(k-p-3)} + \sum_{j=1, j \neq i}^p z_j^2 \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} \\
&= z_i^2 \frac{2k^2(k-p-1) + 2k^2}{(k-p)(k-p-1)^2(k-p-3)} + \sum_{j=1, j \neq i}^p z_j^2 \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} \\
&= \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} \sum_{j=1}^p z_j^2 + \frac{k^2(k-p-1) + 2k^2}{(k-p)(k-p-1)^2(k-p-3)} z_i^2 \\
&= \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} MSS_F + \frac{k^2(k-p+1)}{(k-p)(k-p-1)^2(k-p-3)} z_i^2.
\end{aligned}$$

In the second line the first term is modified to have the same denominator as the remainder sum. In the third line we add and subtract by  $2k^2$  so that the numerator in the first term matches the numerator in the remainder sum. This allows the  $z_j$  terms to be grouped into a sum over the full set of indexes  $j = 1, \dots, p$  in the third line. The fourth line uses the fact that  $\sum_{j=1}^p z_j^2 = \mathbf{z}^\top \mathbf{z} = MSS_F$ . This was shown in the proof of Theorem 1 (S.1). For the off diagonal entries  $W_{ib}$  where  $b \neq i$ ,

$$\begin{aligned}
W_{ib} &= \text{cov} \left( \sum_{j=1}^p M_{ij} z_j, \sum_{r=1}^p M_{br} z_r \right) \\
&= \sum_{j=1}^p \sum_{r=1}^p z_j z_r \text{cov}(M_{ij}, M_{br}).
\end{aligned}$$

Now  $\text{cov}(M_{ij}, M_{br})$  is only nonzero for  $\text{cov}(M_{ib}, M_{bi})$  and  $\text{cov}(M_{ii}, M_{bb})$ . Using (S.7) and (S.8) we obtain

$$\begin{aligned}
W_{ib} &= z_i z_b \text{cov}(M_{ib}, M_{bi}) + z_i z_b \text{cov}(M_{ii}, M_{bb}) \\
&= \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} z_i z_b + \frac{2k^2}{(k-p)(k-p-1)^2(k-p-3)} z_i z_b \\
&= \frac{k^2(k-p+1)}{(k-p)(k-p-1)^2(k-p-3)} z_i z_b.
\end{aligned}$$

The entire covariance matrix  $\mathbf{W}$  can therefore be written compactly as

$$\begin{aligned}
\mathbf{W} &= \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} (MSS_F \mathbf{I}_p) + \frac{k^2(k-p+1)}{(k-p)(k-p-1)^2(k-p-3)} \mathbf{z} \mathbf{z}^\top \\
&= \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)} \left( MSS_F \mathbf{I}_p + \frac{(k-p+1)}{(k-p-1)} \mathbf{z} \mathbf{z}^\top \right) \\
&= \frac{k^2}{(k-p)(k-p-1)(k-p-3)} \left( MSS_F \mathbf{I}_p + \frac{(k-p+1)}{(k-p-1)} \mathbf{z} \mathbf{z}^\top \right).
\end{aligned}$$

Now  $\beta_P = \mathbf{V} \mathbf{D}^{-1} \mathbf{M} \mathbf{z}$ . Therefore  $\text{var}(\beta_P) = \mathbf{V} \mathbf{D}^{-1} \text{var}(\mathbf{M} \mathbf{z}) \mathbf{D}^{-1} \mathbf{V}^\top = \mathbf{V} \mathbf{D}^{-1} \mathbf{W} \mathbf{D}^{-1} \mathbf{V}^\top$ . The

variance of  $\beta_P$  is then a linear function of  $\mathbf{W}$ ,

$$\begin{aligned}
\text{var}(\beta_P) &= \mathbf{V}\mathbf{D}^{-1}\mathbf{W}\mathbf{D}^{-1}\mathbf{V}^\top \\
&= \mathbf{V}\mathbf{D}^{-1}\frac{k^2}{(k-p)(k-p-1)(k-p-3)}\left(MSS_F\mathbf{I}_p + \frac{(k-p+1)}{(k-p-1)}\mathbf{z}\mathbf{z}^\top\right)\mathbf{D}^{-1}\mathbf{V}^\top \\
&= \frac{k^2}{(k-p)(k-p-1)(k-p-3)}MSS_F(\mathbf{V}\mathbf{D}^{-2}\mathbf{V}^\top) + \\
&\quad \frac{k^2(k-p+1)}{(k-p)(k-p-1)^2(k-p-3)}\mathbf{V}\mathbf{D}^{-1}\mathbf{z}\mathbf{z}^\top\mathbf{D}^{-1}\mathbf{V}^\top. \tag{S.9}
\end{aligned}$$

Recall that  $\mathbf{z} = \mathbf{U}^\top\mathbf{y}$  and

$$\beta_F = (\mathbf{X}^\top\mathbf{X})^{-1}\mathbf{X}^\top\mathbf{y} \tag{S.10}$$

$$= \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top\mathbf{y} \tag{S.11}$$

$$= \mathbf{V}\mathbf{D}^{-1}\mathbf{z}. \tag{S.12}$$

The term  $\mathbf{V}\mathbf{D}^{-1}\mathbf{z}$  appears in (S.9). Substituting (S.12) into (S.9) gives

$$\begin{aligned}
\text{var}(\beta_P) &= \frac{k^2}{(k-p)(k-p-1)(k-p-3)}MSS_F(\mathbf{V}\mathbf{D}^{-2}\mathbf{V}^\top) + \\
&\quad \frac{k^2(k-p+1)}{(k-p)(k-p-1)^2(k-p-3)}\beta_F\beta_F^\top.
\end{aligned}$$

A final simplification can be made by noting that  $(\mathbf{X}^\top\mathbf{X})^{-1} = \mathbf{V}\mathbf{D}^{-2}\mathbf{V}^\top$  giving

$$\begin{aligned}
\text{var}(\beta_P) &= \frac{k^2}{(k-p)(k-p-1)(k-p-3)}MSS_F(\mathbf{X}^\top\mathbf{X})^{-1} + \frac{k^2(k-p+1)}{(k-p)(k-p-1)^2(k-p-3)}\beta_F\beta_F^\top \\
&= \frac{k^2}{(k-p)(k-p-1)(k-p-3)}\left(MSS_F(\mathbf{X}^\top\mathbf{X})^{-1} + \frac{(k-p+1)}{(k-p-1)}\beta_F\beta_F^\top\right).
\end{aligned}$$

The variance of  $\beta_P^* = [(k-p-1)/k]\beta_P$  is then

$$\begin{aligned}
\text{var}(\beta_P^*) &= \left(\frac{k-p-1}{k}\right)^2 \frac{k^2(k-p-1)}{(k-p)(k-p-1)^2(k-p-3)}\left(MSS_F(\mathbf{X}^\top\mathbf{X})^{-1} + \frac{(k-p+1)}{(k-p-1)}\beta_F\beta_F^\top\right) \\
&= \frac{(k-p-1)}{(k-p)(k-p-3)}\left(MSS_F(\mathbf{X}^\top\mathbf{X})^{-1} + \frac{(k-p+1)}{(k-p-1)}\beta_F\beta_F^\top\right). \tag{S.13}
\end{aligned}$$

## E Combined estimator results

We first show that  $\beta_P^*$  and  $\beta_S$  are uncorrelated. We again avoid explicitly conditioning on the source dataset  $[\mathbf{y}, \mathbf{X}]$  in every step, it is always treated as fixed. The covariance between  $\beta_P^*$  and  $\beta_S$  computed from the same sketch can be shown to be zero. Using the definition of covariance, and

taking iterated expectations

$$\begin{aligned}\text{cov}(\boldsymbol{\beta}_P^*, \boldsymbol{\beta}_S) &= \mathbb{E}_S \{ (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F)(\boldsymbol{\beta}_S - \boldsymbol{\beta}_F)^\top \} \\ &= \mathbb{E}_{\tilde{\mathbf{X}}} \left[ \mathbb{E}_{\tilde{\mathbf{y}}} \left\{ (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F)(\boldsymbol{\beta}_S - \boldsymbol{\beta}_F)^\top \mid \tilde{\mathbf{X}} \right\} \right].\end{aligned}$$

Recall the hierarchical model for complete sketching,

$$\tilde{\mathbf{y}} \mid \tilde{\mathbf{X}} \sim N \left( \tilde{\mathbf{X}}\boldsymbol{\beta}_F, \frac{RSS_F}{k} \mathbf{I}_k \right).$$

Equivalently,

$$\tilde{\mathbf{y}} \mid \tilde{\mathbf{X}} = \tilde{\mathbf{X}}\boldsymbol{\beta}_F + \tilde{\mathbf{e}},$$

where  $\tilde{\mathbf{e}} \mid \tilde{\mathbf{X}} \sim N(\mathbf{0}, \frac{RSS_F}{k} \mathbf{I}_k)$ . So

$$\boldsymbol{\beta}_S \mid \tilde{\mathbf{X}}, \mathbf{y}, \mathbf{X} = \boldsymbol{\beta}_F + (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{e}}.$$

Substituting back into the expression for the covariance,

$$\begin{aligned}\text{cov}(\boldsymbol{\beta}_P^*, \boldsymbol{\beta}_S) &= \mathbb{E}_{\tilde{\mathbf{X}}} \left\{ \mathbb{E}_{\tilde{\mathbf{e}} \mid \tilde{\mathbf{X}}} \left[ (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F)(\boldsymbol{\beta}_F + (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{e}} - \boldsymbol{\beta}_F)^\top \mid \tilde{\mathbf{X}} \right] \right\} \\ &= \mathbb{E}_{\tilde{\mathbf{X}}} \left\{ \left[ (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F)(\boldsymbol{\beta}_F + (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \mathbb{E}_{\tilde{\mathbf{e}} \mid \tilde{\mathbf{X}}}[\tilde{\mathbf{e}} \mid \tilde{\mathbf{X}}] - \boldsymbol{\beta}_F)^\top \mid \tilde{\mathbf{X}} \right] \right\} \\ &= \mathbb{E}_{\tilde{\mathbf{X}}} \left\{ \left[ (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F)(\boldsymbol{\beta}_F - \boldsymbol{\beta}_F)^\top \mid \tilde{\mathbf{X}} \right] \right\} \\ &= \mathbb{E}_{\tilde{\mathbf{X}}} \left\{ \left[ (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F)\mathbf{0}^\top \mid \tilde{\mathbf{X}} \right] \right\} \\ &= \mathbf{0}_{p \times p}.\end{aligned}$$

Simple calculus shows that the value which minimises the expected mean square error  $\mathbb{E}_S(\|\boldsymbol{\beta}_C - \boldsymbol{\beta}_F\|_2^2 \mid \mathbf{y}, \mathbf{X})$  is

$$\phi_{\text{opt}} = \frac{\text{tr}(\text{var}(\boldsymbol{\beta}_P^*))}{\text{tr}(\text{var}(\boldsymbol{\beta}_P^*)) + \text{tr}(\text{var}(\boldsymbol{\beta}_S))}.$$

## F Proof of Theorem 4 (central limit theorem under asymptotic negligibility condition)

A triangular array of random variables is a useful structure for studying weak convergence. To establish a triangular array, define for every  $n \in \mathbb{N}$  a collection of random variables  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$ . There are  $r_n$  random variables in row  $n$  of the array. Suppose that  $r_n = n$ . Visually we can represent

the first three rows of the array as

$$\begin{array}{ccc} Z_{11} & & \\ Z_{21} & Z_{22} & \\ Z_{31} & Z_{32} & Z_{33} \end{array}$$

**Theorem** (Billingsley, 1995, Chapter 5, Section 27). *For each  $n \in \mathbb{N}$ , let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of independent random variables with  $\mathbb{E}(Z_{ni}) = 0$  and  $\text{var}(Z_{ni}) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Let  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  and assume that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Suppose that we can form a sequence of upper bounds  $(K_n)_{n \in \mathbb{N}}$  such that*

$$|Z_{ni}| \leq K_n \text{ almost surely for } i = 1, \dots, r_n.$$

*Then if  $K_n/s_n \rightarrow 0$  as  $n \rightarrow \infty$  we have the convergence in distribution*

$$\frac{1}{s_n} \sum_{i=1}^{r_n} Z_{ni} \xrightarrow{d} N(0, 1)$$

Lindeberg's condition is a critical component in establishing asymptotic normality. We state Lindeberg's condition for triangular arrays of random variables.

**Definition 2** (Lindeberg's condition). *For each  $n \in \mathbb{N}$ , let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of random variables with  $\mathbb{E}(Z_{ni}) = 0$  and  $\text{var}(Z_{ni}) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Let  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  and suppose that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . The random variables are said to satisfy Lindeberg's condition if for all  $\eta > 0$ ,*

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^2} \sum_{i=1}^{r_n} \mathbb{E}(Z_{ni}^2 \mathbb{1}_{\{|Z_{ni}| > \eta s_n\}}) = 0. \quad (\text{S.14})$$

The triangular array of random variables does not have to have independent random variables in each row in order to satisfy the condition. The general form of the Lindeberg-Feller central limit theorem shows that a triangular array of independent random variables satisfying Lindeberg's condition is asymptotically normal after suitable scaling.

**Theorem 8** (Lindeberg-Feller). *For each  $n \in \mathbb{N}$ , let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of random variables with  $\mathbb{E}(Z_{ni}) = 0$  and  $\text{var}(Z_{ni}) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Let  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  and suppose that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Suppose the triangular array of random variables satisfies Lindeberg's condition (Definition 2). Then*

$$\frac{1}{s_n} \sum_{i=1}^{r_n} Z_{ni} \xrightarrow{d} N(0, 1)$$

For a proof see Loeve (1977). It can be difficult to show Lindeberg's condition directly. A stronger condition that implies the Lindeberg condition is the Lyapunov condition.

**Definition 3** (Lyapunov's condition). For each  $n \in \mathbb{N}$ , let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of random variables with  $\mathbb{E}(Z_{ni}) = 0$  and  $\text{var}(Z_{ni}) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Let  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  and suppose that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . The triangular array of random variables is said to satisfy Lyapunov's condition if there exists a  $\delta > 0$  such that

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta}) = 0. \quad (\text{S.15})$$

The Lyapunov condition implies the Lindeberg condition. We state this in a Lemma for later reference.

**Lemma 2.** *The Lyapunov condition implies the Lindeberg condition.*

To see this assume the Lyapunov condition is satisfied and fix  $\eta > 0$ . Now  $|Z_{ni}| \geq \eta s_n$  implies that  $1 \leq |Z_{ni}/(\eta s_n)|^\delta$ . We can then form an upper bound on the sequence of partial sums that appear in Lindeberg's condition.

$$\begin{aligned} \frac{1}{s_n^2} \sum_{i=1}^{r_n} \mathbb{E}(Z_{ni}^2 \mathbb{1}_{\{|Z_{ni}| > \eta s_n\}}) &\leq \frac{1}{s_n^2} \sum_{i=1}^{r_n} \mathbb{E}(Z_{ni}^2 |Z_{ni}/(\eta s_n)|^\delta \mathbb{1}_{\{|Z_{ni}| > \eta s_n\}}) \\ &= \frac{1}{s_n^2} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^2 |Z_{ni}/(\eta s_n)|^\delta \mathbb{1}_{\{|Z_{ni}| > \eta s_n\}}) \\ &= \frac{1}{s_n^2} \frac{1}{(\eta s_n)^\delta} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta} \mathbb{1}_{\{|Z_{ni}| > \eta s_n\}}) \\ &= \frac{1}{\eta^\delta} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta}). \end{aligned}$$

Assuming that Lyapunov's condition holds we can establish zero as an upper bound

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{s_n^2} \sum_{i=1}^{r_n} \mathbb{E}(Z_{ni}^2 \mathbb{1}_{\{|Z_{ni}| > \eta s_n\}}) &\leq \lim_{n \rightarrow \infty} \frac{1}{\eta^\delta} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta}) \\ &= \frac{1}{\eta^\delta} \lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta}) \\ &= 0. \end{aligned}$$

As the partial sums are lower bounded by zero, the Lyapunov condition implies the Lindeberg condition.

We now present a useful Lemma for showing the Lyapunov condition. The result is from Billingsley (1995) and applies to triangular arrays of uniformly bounded random variables.

**Lemma 3** (Billingsley, 1995). For each  $n \in \mathbb{N}$ , let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of random variables with  $\mathbb{E}(Z_{ni}) = 0$  and  $\text{var}(Z_{ni}) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Let  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  and suppose that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Suppose that we can form a sequence of upper bounds  $(K_n)_{n \in \mathbb{N}}$  such that

$$|Z_{ni}| \leq K_n \text{ almost surely for } i = 1, \dots, r_n.$$

Then if  $K_n/s_n \rightarrow 0$  as  $n \rightarrow \infty$  the Lyapunov condition holds for the triangular array of random variables.

Lemma 3 is useful as it does not impose a constant uniform bound on the random variables. In the special case where  $|Z_{ni}| \leq M$  almost surely for some constant  $M$  for all  $n \in \mathbb{N}$  and all  $i = 1, \dots, r_n$  we have that Lyapunov's condition is satisfied providing that  $s_n \rightarrow \infty$ . Lemma 3 allows for the bound  $K_n$  to increase with  $n$  as long as the rate of growth is slower than the rate of growth of  $s_n$ . Lyapunov's condition holds providing that  $K_n = o(s_n)$ .

The proof of Lemma 3 is given below. Again fix some  $\delta > 0$ . If  $|Z_{ni}| \leq K_n$  almost surely for  $i = 1, \dots, r_n$  it must hold that  $|Z_{ni}|^\delta \leq K_n^\delta$  as  $|Z_{ni}|, K_n$  and  $\delta$  are all positive. As such  $|Z_{ni}|^{2+\delta} = |Z_{ni}|^2 |Z_{ni}|^\delta \leq |Z_{ni}|^2 K_n^\delta$ . We can then form an upper bound on the sequence of partial sums that appear in Lyapunov's condition.

$$\begin{aligned} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta}) &\leq \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^2) K_n^\delta \\ &= \frac{K_n^\delta}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}|Z_{ni}|^2 \\ &= \frac{K_n^\delta}{s_n^{2+\delta}} s_n^2 \\ &= \left( \frac{K_n}{s_n} \right)^\delta. \end{aligned} \tag{S.16}$$

Now assuming that  $K_n = o(s_n)$  we have that  $K_n/s_n \rightarrow 0$  as  $n \rightarrow \infty$ . We then also have that

$$\begin{aligned} \lim_{n \rightarrow \infty} \left( \frac{K_n}{s_n} \right)^\delta &= \left( \lim_{n \rightarrow \infty} \frac{K_n}{s_n} \right)^\delta \\ &= 0, \end{aligned}$$

as the exponentiation by  $\delta > 0$  is a continuous function. Now taking limits on both sides of the inequality (S.16):

$$\lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta}) \leq \lim_{n \rightarrow \infty} \left( \frac{K_n}{s_n} \right)^\delta \tag{S.17}$$

$$= 0. \tag{S.18}$$

We also have the lower bound

$$0 \leq \lim_{n \rightarrow \infty} \frac{1}{s_n^{2+\delta}} \sum_{i=1}^{r_n} \mathbb{E}(|Z_{ni}|^{2+\delta}).$$

By the squeeze theorem we then have that  $K_n = o(s_n)$  is sufficient for Lyapunov's condition to hold.

The triangular array of independent random variables in Theorem 4 satisfies Lyapunov's condition by Lemma 3. As the Lyapunov condition implies the Lindeberg condition (Lemma 2) the general Lindeberg-Feller central limit theorem (Theorem 8) gives asymptotic normality of the scaled



row sums, thus proving Theorem 4.

## G Proof of Theorem 3 (Sketching central limit theorem)

**Assumption 1** Let the singular value decomposition of the  $n \times d$  source dataset be given by  $\mathbf{A}_{(n)} = \mathbf{U}_{(n)}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^\top$ . Let  $\mathbf{u}_{(n)i}^\top$  give the  $i$ th row in  $\mathbf{U}_{(n)}$ . Assume that the maximum leverage score tends to zero, that is

$$\lim_{n \rightarrow \infty} \max_{i=1, \dots, n} \|\mathbf{u}_{(n)i}\|_2^2 = 0.$$

Theorem 3 gives the sketching central limit theorem.

**Theorem.** Consider a fixed sequence of arbitrary  $n \times d$  data matrices  $\mathbf{A}_{(n)}$ , where  $d$  is fixed. Let  $\mathbf{A}_{(n)} = \mathbf{U}_{(n)}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^\top$  represent the singular value decomposition of  $\mathbf{A}_{(n)}$ . Let  $\mathbf{S}$  be a  $k \times n$  Hadamard or Clarkson-Woodruff sketching matrix where  $k$  is also fixed. Suppose that Assumption 1 on the maximum leverage score is satisfied. Then as  $n$  tends to infinity with  $k$  and  $d$  fixed,

$$[\tilde{\mathbf{A}}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1} \mid \mathbf{A}_{(n)}] \xrightarrow{d} \text{MN}(\mathbf{0}, \mathbf{I}_k, \mathbf{I}_d/k).$$

To prove the sketching central limit theorem it helps to restate Lemma 3. This helps to show the importance of the leverage scores in establishing asymptotic normality. Lemma 3 provided a sufficient condition for showing that Lindeberg's condition holds. We can restate Lemma 3 in terms of a normalised triangular array.

**Theorem 9** (Billingsley, 1995). For each  $n \in \mathbb{N}$  let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of random variables with  $\mathbb{E}[Z_{ni}] = 0$  and  $\text{var}(Z_{ni}^2) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Define  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  each  $n$ . Suppose that the rows of the triangular array are standardised such that  $s_n^2 = 1$  for all  $n$ . Suppose that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Suppose we have a sequence of upper bounds  $(Z_n)$  such that  $|Z_{ni}| \leq K_n$  almost surely for all  $i = 1, \dots, r_n$ . Then a sufficient condition for Lyapunov's condition to hold is  $K_n \rightarrow 0$  as  $n \rightarrow \infty$ .

The standardisation of the triangular array gives an intuitive condition for Lyapunov's and hence Lindeberg's condition to hold. We require that  $K_n \rightarrow 0$  as  $n \rightarrow \infty$ . We require that the upper bound tends to zero. All the random variables in the row must converge almost surely to zero. Almost sure convergence is stronger than convergence in probability and rules out pathological cases where a single random variable in a row can take a large value with small probability. Assumption 1 on the leverage scores in the sketching central limit theorem enforces a bounded growth condition that relates to Theorem 9.

Let  $n \in \mathbb{N}$  index the sequence of source datasets of increasing size. We assume that the source dataset consists of  $r_n$  observations where  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . For now we can take  $r_n = n$  to ease interpretation. We take the singular value decomposition of each dataset  $\mathbf{A}_{(n)} = \mathbf{U}_{(n)}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^\top$ . All results in this section treat the source dataset  $\mathbf{A}_{(n)}$  as fixed, only the sketching matrix is random.

We consider the sequence of whitened sketched datasets

$$\begin{aligned}\tilde{\mathbf{A}}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1} &= (\mathbf{S}\mathbf{A})\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1} \\ &= \mathbf{S}\mathbf{U}_{(n)}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^{\top}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1} \\ &= \mathbf{S}\mathbf{U}_{(n)}.\end{aligned}$$

The whitened sketched dataset  $\tilde{\mathbf{A}}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1}$  has a  $MN(\mathbf{0}, \mathbf{I}_k, \mathbf{I}_d/k)$  distribution when  $\mathbf{S}$  is a Gaussian sketch. We need to show that as  $n$  tends to infinity,  $\mathbf{S}\mathbf{U}_{(n)}$  converges in distribution to a  $MN(\mathbf{0}, \mathbf{I}_k, \mathbf{I}_d/k)$  random matrix for both the Clarkson-Woodruff and Hadamard sketches.

Let  $\mathbf{u}_{(n)i}^{\top}$  denote row  $i$  of the matrix of left singular vectors  $\mathbf{U}_{(n)}$ . We write  $\mathbf{u}_{(n)i}^{\top}$  so that that we can form a triangular array of left singular vectors. Taking  $r_n = n$ , the first three rows of the triangular array can be written as

$$\begin{array}{ccc}\mathbf{u}_{(1)1} & & \\ \mathbf{u}_{(2)1} & \mathbf{u}_{(2)2} & \\ \mathbf{u}_{(3)1} & \mathbf{u}_{(3)2} & \mathbf{u}_{(3)3}\end{array}$$

An important property is that for all  $n$ , the sum of the norms of the leverage scores always equals the number of variables in the source dataset  $d$ .

$$\sum_{i=1}^{r_n} \|\mathbf{u}_{(n)i}\|_2^2 = d. \quad (\text{S.19})$$

As  $n$  increases, the typical norm of each vector  $\mathbf{u}_{(n)i}$ ,  $i \in \{1, \dots, r_n\}$  is expected to decrease. For completeness we restate Assumption 1 in terms of the triangular array formulation.

**Assumption 1** Let the singular value decomposition of the  $r_n \times d$  source dataset be given by  $\mathbf{A}_{(n)} = \mathbf{U}_{(n)}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^{\top}$ . Let  $\mathbf{u}_{(n)i}^{\top}$  give the  $i$ th row in  $\mathbf{U}_{(n)}$  for  $i = 1, \dots, r_n$ . Assume that the maximum leverage score tends to zero, that is

$$\lim_{n \rightarrow \infty} \max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2^2 = 0.$$

This increasing collection of smaller quantities is similar to the behaviour of the triangular array of random variables in Theorem 9. The standardisation property in equation (S.19), namely that  $\sum_{i=1}^{r_n} \|\mathbf{u}_{(n)i}\|_2^2 = d$  for all  $n$  is similar to the assumption that  $s_n = 1$  in each row of the triangular array of random variables in Theorem 9. Assumption 1 on the leverage scores, where the maximum individual norm tends to zero is similar to the assumption that  $K_n \rightarrow 0$  in Theorem 9. This will be made more explicit in the proofs. Before moving on we make a note that assumption 1 also implies that the maximum square root of the leverage scores also tends to zero. As

$$\max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2 = \left( \max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2^2 \right)^{1/2} \quad (\text{S.20})$$

We have that

$$\begin{aligned}
\lim_{n \rightarrow \infty} \max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2 &= \lim_{n \rightarrow \infty} \left( \max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2^2 \right)^{1/2} \\
&= \left( \lim_{n \rightarrow \infty} \max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2^2 \right)^{1/2} \\
&= 0.
\end{aligned} \tag{S.21}$$

To establish joint asymptotic normality of the sketched data matrix we use the Cramér-Wold device.

**Lemma 4** (Cramér-Wold device). *Let  $(\mathbf{Z}_n)_{n \in \mathbb{N}}$  be a sequence of random vectors in  $\mathbb{R}^v$ . Let  $\mathbf{Z}$  denote another random vector also in  $\mathbb{R}^v$ . The sequence of random vectors  $(\mathbf{Z}_n)$  converges in distribution to  $\mathbf{Z}$  as  $n$  tends to infinity if and only if the sequence of random variables  $(\boldsymbol{\lambda}^\top \mathbf{Z}_n)_{n \in \mathbb{N}}$  converges in distribution to  $\boldsymbol{\lambda}^\top \mathbf{Z}$  for all unit vectors  $\boldsymbol{\lambda} \in \mathbb{R}^v$ .*

A proof is given in Shorack (2000, Chapter 13, Section 3). Let  $\mathbf{z}_n$  represent the  $kd$  length vector formed by stacking transposed rows of the whitened sketched dataset  $\tilde{\mathbf{U}} = \mathbf{S}\mathbf{U}_{(n)}$ . Let  $\tilde{\mathbf{u}}_j^\top$  give row  $j$  in  $\tilde{\mathbf{U}}$  for  $j = 1, \dots, k$ . Formally,

$$\mathbf{z}_n = \begin{bmatrix} \tilde{\mathbf{u}}_1 \\ \tilde{\mathbf{u}}_2 \\ \vdots \\ \tilde{\mathbf{u}}_k \end{bmatrix}. \tag{S.22}$$

Let us define the random matrix  $k \times d$  random matrix  $\mathbf{W}$  as having the matrix normal distribution

$$\mathbf{W} \sim MN(\mathbf{0}, \mathbf{I}_k, \mathbf{I}_d/k)$$

Let  $\mathbf{w}_i^\top$  refer to row  $i$  in  $\mathbf{W}$  for  $i = 1, \dots, k$ . Let  $\mathbf{z}_L$  refer to the stacked transposed rows of  $\mathbf{W}$ , so

$$\mathbf{z}_L = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \\ \vdots \\ \mathbf{w}_k \end{bmatrix}. \tag{S.23}$$

Let  $\boldsymbol{\lambda}$  be an arbitrary unit vector in  $\mathbb{R}^{k \times d}$ . It will be useful to also partition the vector  $\boldsymbol{\lambda}$  into  $k$  sub-vectors,

$$\boldsymbol{\lambda} = \begin{bmatrix} \boldsymbol{\lambda}_1 \\ \boldsymbol{\lambda}_2 \\ \vdots \\ \boldsymbol{\lambda}_k \end{bmatrix}, \tag{S.24}$$

where  $\boldsymbol{\lambda}_j$  is a  $d$ -dimensional vector for  $j = 1, \dots, k$ . For any unit vector  $\boldsymbol{\lambda} \in \mathbb{R}^{k \times d}$ ,  $\boldsymbol{\lambda}^\top \mathbf{z}_L$  is

distributed as  $N(0, 1/k)$ . We will aim to show that the distribution of the whitened sketched data  $\mathbf{S}\mathbf{A}_{(n)}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1}$  converges to that of  $\mathbf{W}$  through the Cramér-Wold device. We must show that for any fixed  $k \times d$  length unit vector  $\boldsymbol{\lambda}$ ,  $\boldsymbol{\lambda}^\top \mathbf{z}_n$  converges in distribution to  $N(0, 1/k)$  as  $n \rightarrow \infty$ .

We will rely on a central limit theorem for jointly symmetric, pairwise independent random variables (Pruss and Szynal, 2000). A collection of random variables  $(Z_1, \dots, Z_n)$  is said to be jointly symmetric if  $(Z_1, \dots, Z_n)$  has the same distribution as  $(q_1 Z_1, \dots, q_n Z_n)$ , where  $q_i \in \{+1, -1\}$  for  $i = 1, \dots, n$ . Given a set of random variables  $Y_1, \dots, Y_n$ , a jointly symmetric collection  $Z_1, \dots, Z_n$  can be formed by sampling  $n$  independent Rademacher random variables  $h_1, \dots, h_n$ , and setting  $Z_i = h_i Y_i$  (Pruss and Szynal, 2000). It is possible to establish a central limit theorem for jointly symmetric, pairwise independent random variables.

**Theorem 10** (Pruss and Szynal (2000), Theorem 1, Corollary 2). *For each  $n \in \mathbb{N}$ , let  $Z_{n1}, Z_{n2}, \dots, Z_{nr_n}$  be a sequence of jointly symmetric pairwise independent random variables with  $\mathbb{E}(Z_{ni}) = 0$  and  $\text{var}(Z_{ni}) = \sigma_{ni}^2$  for  $i = 1, \dots, r_n$ . Let  $s_n^2 = \sum_{i=1}^{r_n} \sigma_{ni}^2$  and assume that  $r_n \rightarrow \infty$  as  $n \rightarrow \infty$ . Suppose the triangular array of random variables satisfies Lindeberg's condition. Then as  $n \rightarrow \infty$ ,  $s_n^{-1} \sum_{i=1}^{r_n} Z_{ni}$  converges in distribution to  $N(0, 1)$ .*

Not all triangular arrays with pairwise independent random variables in each row satisfy a central limit theorem. The joint symmetry property is very important (Pruss and Szynal, 2000; Svante, 1988).

To use Theorem 10 we need to show that the triangular array of random variables satisfies Lindeberg's condition. As discussed this can be very difficult to establish directly. If the triangular array of random variables can be appropriately bounded, we can use Theorem 9 to show that Lyapunov's condition holds, and subsequently that Lindeberg's condition holds.

This is the approach we take in proving the sketching central limit theorem. The Cramér-Wold device is used to reduce the study of multivariate convergence to univariate convergence. We can then form a triangular array of random variables such that elements in each row are jointly symmetric and pairwise independent. We then show that triangular array satisfies Lindeberg's condition using Theorem 9. Assumption 1 on the maximum leverage score enforces the necessary cap on the rate of growth. Theorem 10 is then used to establish asymptotic normality.

### G.1 Clarkson-Woodruff sketch

The Clarkson-Woodruff sketch can be represented as the product of two independent random matrices,  $\mathbf{S} = \boldsymbol{\Gamma}\mathbf{D}$ , where  $\boldsymbol{\Gamma}$  is a random  $k \times n$  matrix and  $\mathbf{D}$  is a random  $n \times n$  matrix. The diagonal matrix  $\mathbf{D}$  contains  $n$  independent Rademacher random variables on the diagonal. Let  $h_i \in \{+1, -1\}$  be the random sign in element  $D_{ii}$ . The matrix  $\boldsymbol{\Gamma}$  is formed by choosing one element in each column independently and setting the entry to  $+1$ . Element  $\Gamma_{ij}$  is equal to  $+1$  if we add observation  $i$  in the original dataset to sketched observation  $j$ . The signs in row  $i$  are flipped if  $h_i$  is equal to negative one. Each observation in the original dataset is assigned to one sketched observation as each column of  $\boldsymbol{\Gamma}$  contains a single  $+1$  entry. Using a Clarkson-Woodruff sketch row  $j$  in the sketched data matrix

can be represented as

$$\tilde{\mathbf{u}}_j^\top = \sum_{i=1}^n h_i \Gamma_{ij} \mathbf{u}_{(n)i}^\top,$$

where  $h_i$  represents the random sign flip applied to row  $i$  of the original data matrix, and  $\Gamma_{ij}$  is the indicator variable which is equal to one if row  $i$  of the original data is added to row  $j$  of the sketched dataset.

Let us consider the linear combination  $\boldsymbol{\lambda}^\top \mathbf{z}$ , where  $\boldsymbol{\lambda}$  and  $\mathbf{z}$  are defined as in (S.22) and (S.24) respectively. The sum over the  $k$  rows in the sketched dataset can be rearranged into a sum over the  $n$  rows in the source dataset,

$$\begin{aligned} \boldsymbol{\lambda}^\top \mathbf{z}_n &= \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \tilde{\mathbf{u}}_j \\ &= \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \sum_{i=1}^n h_i \Gamma_{ij} \mathbf{u}_{(n)i} \\ &= \sum_{i=1}^n h_i \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i}. \end{aligned} \tag{S.25}$$

The scalar  $\boldsymbol{\lambda}^\top \mathbf{z}_n$  is equal to the sum of  $n$  independent random variables. Independence holds as the signs flips  $h_i$  on each observation are independent, and each column of  $\boldsymbol{\Gamma}$  is independent.

In the language of Theorem 9 we can form a triangular array of random variables setting

$$Z_{ni} = h_i \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i}. \tag{S.26}$$

for  $i = 1, \dots, n$  and  $n \in \mathbb{N}$ . The linear combination in (S.25) then be expressed as a row sum over the triangular array defined in (S.26):

$$\boldsymbol{\lambda}^\top \mathbf{z}_n = \sum_{i=1}^n Z_{ni}. \tag{S.27}$$

Our goal of showing that  $\boldsymbol{\lambda}^\top \mathbf{z}_n$  converges in distribution to a  $N(0, 1/k)$  random variable is achieved if we can show that  $\sum_{i=1}^n Z_{ni}$  converges in distribution to a  $N(0, 1/k)$  random variable.

It is worth making a connection to Theorem 10, because of the random sign flips  $h_i$  appearing in (S.26), we have a sequence of mutually independent jointly symmetric random variables. Mutually independent random variables are also necessarily pairwise independent. Theorem 10 can be used to establish asymptotic normality of the sum in (S.27) and hence the linear combination  $\boldsymbol{\lambda}^\top \mathbf{z}_n$ . To show that the triangular array of random variables defined in (S.26) satisfies Lindeberg's condition we use Theorem 9. Set  $s_n^2 = \sum_{i=1}^n \text{var}(Z_{ni})$ . We first determine  $s_n^2$ . We then form the necessary sequence of upper bounds  $K_n$  such that  $|Z_{ni}| \leq K_n$  almost surely for  $i = 1, \dots, n$ . The variance of

a single term in the sum (S.25) is

$$\text{var}(Z_{ni}) = \text{var} \left( h_i \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \right) \quad (\text{S.28})$$

$$= \sum_{j=1}^k \frac{1}{k} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \mathbf{u}_{(n)i}^\top \boldsymbol{\lambda}_j. \quad (\text{S.29})$$

The row-wise variance totals  $s_n^2$  are then

$$\begin{aligned} s_n^2 &= \sum_{i=1}^n \text{var}(Z_{ni}) \\ &= \sum_{i=1}^n \text{var} \left( h_i \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \right) \\ &= \frac{1}{k} \sum_{i=1}^n \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \mathbf{u}_{(n)i}^\top \boldsymbol{\lambda}_j \\ &= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \left( \sum_{i=1}^n \mathbf{u}_{(n)i} \mathbf{u}_{(n)i}^\top \right) \boldsymbol{\lambda}_j \\ &= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{U}_{(n)}^\top \mathbf{U}_{(n)} \boldsymbol{\lambda}_j \\ &= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{I}_d \boldsymbol{\lambda}_j. \\ &= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \boldsymbol{\lambda}_j \\ &= \frac{1}{k}. \end{aligned}$$

The fact that  $\mathbf{U}_{(n)}^\top \mathbf{U}_{(n)} = \mathbf{I}_d$  for all  $n$  serves as a useful normalisation to give stable limiting behaviour. The step in the last line follows as we have taken  $\boldsymbol{\lambda}$  to be a unit vector. We have  $s_n = 1/k$  for all  $n$  in the triangular array. We now establish a sequence of upper bounds ( $K_n$ ). As the random variables in the construction of construction of the sketch are bounded, we can bound the random variables in the triangular array using the leverage scores of the sequence of source dataset. Now as the random sign  $h_i \in \{+1, -1\}$

$$\begin{aligned} |Z_{ni}| &= |h_i \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i}| \\ &= \left| \left( \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j^\top \right) \mathbf{u}_{(n)i} \right|. \end{aligned} \quad (\text{S.30})$$

Now by the Cauchy-Schwarz inequality

$$\left| \left( \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j^\top \right) \mathbf{u}_{(n)i} \right| \leq \left\| \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j \right\|_2 \|\mathbf{u}_{(n)i}\|_2 \quad (\text{S.31})$$

Now as  $\Gamma_{ij} = 1$  for a single  $j \in \{1, \dots, k\}$  and is zero otherwise we have that

$$\begin{aligned} \left\| \sum_{j=1}^k \Gamma_{ij} \boldsymbol{\lambda}_j \right\|_2 &\leq \max_{j=1, \dots, k} \|\boldsymbol{\lambda}_j\|_2 \\ &\leq 1. \end{aligned} \quad (\text{S.32})$$

The last line follows as we have taken  $\boldsymbol{\lambda}$  to be a unit vector. Substituting (S.32) and (S.31) into (S.30) we arrive at

$$|Z_{ni}| \leq \|\mathbf{u}_{(n)i}\|_2.$$

We can then form the sequence of upper bounds  $K_n$ ,

$$K_n = \max_{i=1, \dots, n} \|\mathbf{u}_{(n)i}\|_2.$$

We have that  $|Z_{ni}| \leq K_n$  almost surely for  $i = 1, \dots, n$  and  $n \in \mathbb{N}$ . Assumption 1 controls the limiting behaviour of  $K_n = \max_{i=1, \dots, n} \|\mathbf{u}_{(n)i}\|_2$  (recall equation (S.21)). Taking limits and using Assumption 1 shows that  $K_n \rightarrow 0$ ,

$$\begin{aligned} \lim_{n \rightarrow \infty} K_n &= \lim_{n \rightarrow \infty} \max_{i=1, \dots, n} \|\mathbf{u}_{(n)i}\|_2 \\ &= 0. \end{aligned}$$

By theorem 9 we have that the triangular array of random variables in (S.26) satisfies Lindeberg's condition. As such the conditions of Theorem 10 are satisfied, giving that  $\boldsymbol{\lambda}^\top \mathbf{z}_n$  converges in distribution to  $N(0, 1/k)$ . Finally, the Cramér-Wold device gives that the whitened sketched dataset has a limiting matrix normal distribution, that is  $\tilde{\mathbf{A}} \mathbf{V}_{(n)} \mathbf{D}_{(n)}^{-1}$  converges in distribution to a  $MN(\mathbf{0}_{k \times d}, \mathbf{I}_k, \mathbf{I}_d/k)$  random matrix.

## G.2 Hadamard sketch

Recall that the Hadamard sketch is defined through  $\mathbf{S} = \boldsymbol{\Phi} \mathbf{H} \mathbf{D} / \sqrt{k}$ . Here  $\mathbf{H}$  is a Hadamard matrix. Hadamard matrices are square matrices with  $2^n$  rows for some integer  $n$ . To take limits we have to define our sequence of source datasets  $(\mathbf{A}_{(n)} = \mathbf{U}_{(n)} \mathbf{D}_{(n)} \mathbf{V}_{(n)}^\top)$  as having  $r_n = 2^n$  rows for  $n \in \mathbb{N}^+$ . In practice when taking a Hadamard sketch we pad the original dataset with zeros if the original number of observations is not a power of two. To rigorously establish asymptotic normality for the Hadamard sketch we have to take  $r_n = 2^n$ . The first three rows of the triangular array of left

singular vectors now looks like

$$\begin{aligned} & \mathbf{u}_{(1)1} \\ & \mathbf{u}_{(2)1} \quad \mathbf{u}_{(2)2} \\ & \mathbf{u}_{(3)1} \quad \mathbf{u}_{(3)2} \quad \mathbf{u}_{(3)3} \quad \mathbf{u}_{(n)4}. \end{aligned}$$

The intuition is the same as with the Clarkson-Woodruff sketch, as we move down the rows  $n$  we expect the norms of  $\mathbf{u}_{(n)i}$ ,  $i \in \{1, \dots, 2^n\}$  to decrease. This follows from the implicit row-wise normalisation property

$$\sum_{i=1}^{r_n} \|\mathbf{u}_{(n)i}\|_2^2 = d.$$

The indexing change to  $r_n = 2^n$  instead of  $r_n = n$  has very little impact on the underlying arguments.

There are two independent sources of randomness in a Hadamard sketch, the  $r_n = 2^n$  independent random Rademacher variables in the diagonal matrix  $\mathbf{D}$ , and the random matrix  $\Phi$  which subsamples  $k$  rows with replacement from the Hadamard matrix  $\mathbf{H}$ . Hadamard matrices have a number of properties that we will use (Anderson, 1997, section 3.2).

- (P1) The first column contains all ones.
- (P2) Every column other than the first contains an equal number of  $+1$  and  $-1$  entries.
- (P3) Consider any two different columns  $i$  and  $s$ , where  $i, s \in \{2, \dots, r_n\}$ ,  $i \neq s$ . Columns  $i$  and  $s$  will have  $+1$  together in a quarter of the rows, and  $-1$  together in a quarter of the rows. Furthermore, a quarter of the rows will have  $+1$  in column  $i$  and  $-1$  in column  $s$ . Similarly, a quarter of the rows will have  $-1$  in column  $i$  and  $+1$  in column  $s$ .

Let  $\mathbf{M}$  represent the random  $k \times n$  matrix from the subsampling operation  $\mathbf{M} = \Phi \mathbf{H}$ . Let  $m_{ji}$  refer to the element in row  $j$  and column  $i$  of  $\mathbf{M}$ . Each element in  $\mathbf{M}$  is equal to  $+1$  or  $-1$ . Let  $h_i \in \{+1, -1\}$  be the random sign in element  $D_{ii}$ . We now represent the Hadamard sketch as  $\mathbf{S} = \mathbf{M}\mathbf{D}/\sqrt{k}$ .

The structure of the Hadamard matrix gives the random matrix  $\mathbf{M}$  some useful properties. Consider an arbitrary row  $j$  in  $\mathbf{M}$ . By (P1) listed above regarding the first column of  $\mathbf{M}$ ,  $m_{j1} = 1$  with probability one. For the other columns,  $m_{ji} = 1$  with probability half, and  $m_{ji} = -1$  with probability half for  $i = 2, \dots, r_n$  by (P2). By (P3) listed above, we have pairwise independence between elements in row  $j$  of  $\mathbf{M}$ , that is  $p(m_{ji}|m_{js}) = p(m_{ji})$  for  $i, s \in \{1, \dots, r_n\}$ ,  $i \neq s$ . As rows of  $\mathbf{M}$  are sampled independently, each column of  $\mathbf{M}$  is pairwise independent.

Row  $j$  in the sketched dataset is given by

$$\tilde{\mathbf{u}}_j^\top = \frac{1}{\sqrt{k}} \sum_{i=1}^{r_n} m_{ji} h_i \mathbf{u}_{(n)i}^\top.$$

Let us again consider the linear combination  $\boldsymbol{\lambda}^\top \mathbf{z}_n$ , where  $\boldsymbol{\lambda}$  and  $\mathbf{z}_n$  are defined as in (S.22) and (S.24) respectively. The sum over the  $k$  rows in the sketched dataset can be rearranged into a sum



over the  $r_n = 2^n$  rows in the source dataset,

$$\begin{aligned}
\boldsymbol{\lambda}^\top \mathbf{z}_n &= \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \tilde{\mathbf{u}}_j \\
&= \frac{1}{\sqrt{k}} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \sum_{i=1}^{r_n} m_{ji} h_i \mathbf{u}_{(n)i} \\
&= \frac{1}{\sqrt{k}} \sum_{i=1}^{r_n} h_i \left( \sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j^\top \right) \mathbf{u}_{(n)i}.
\end{aligned} \tag{S.33}$$

In the language of Theorem 9 we can form a triangular array of random variables setting

$$Z_{ni} = \frac{1}{\sqrt{k}} h_i \left( \sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j^\top \right) \mathbf{u}_{(n)i}. \tag{S.34}$$

for  $i = 1, \dots, r_n$  and  $n \in \mathbb{N}$ . The linear combination in (S.33) can then be expressed as a row sum of the triangular array defined by (S.34)

$$\boldsymbol{\lambda}^\top \mathbf{z}_n = \sum_{i=1}^{r_n} Z_{ni}. \tag{S.35}$$

Our goal of showing that  $\boldsymbol{\lambda}^\top \mathbf{z}_n$  converges in distribution to a  $N(0, 1/k)$  random variable is achieved if we can show that  $\sum_{i=1}^{r_n} Z_{ni}$  converges in distribution to a  $N(0, 1/k)$  random variable.

The sequence of random variables in each row of the triangular array  $Z_{n1}, \dots, Z_{nr_n}$  are not mutually independent over  $i = 1, \dots, r_n$ . This is because the columns of  $\mathbf{M}$  are not mutually independent. However, as the columns of  $\mathbf{M}$  are pairwise independent, the random sums  $\sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j^\top$  appearing in (S.34) are also pairwise independent. Again making a connection to Theorem 10, the independent sign flips  $h_i$  appearing in (S.34) ensure that the random variables in each row of the triangular array are jointly symmetric and pairwise independent.

Theorem 10 can be used to establish asymptotic normality of the sum in (S.35) and hence the linear combination  $\boldsymbol{\lambda}^\top \mathbf{z}_n$ . To show that the triangular array of random variables defined in (S.34) satisfies Lindeberg's condition we use Theorem 9. Set  $s_n^2 = \sum_{i=1}^{r_n} \text{var}(Z_{ni})$ . We first determine  $s_n^2$ . We then form the necessary sequence of upper bounds  $K_n$  such that  $|Z_{ni}| \leq K_n$  almost surely for  $i = 1, \dots, r_n$ .

We start by considering the variance of a single term in the triangular array  $\text{var}(Z_{ni})$ . We have that

$$\text{var}(Z_{ni}) = \frac{1}{k} \text{var} \left( h_i \left( \sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j^\top \right) \mathbf{u}_{(n)i} \right) \tag{S.36}$$

It is important to consider the covariance between the elements of the sum over  $j = 1, \dots, k$ . For

$i \neq 1$  and  $j, v \in \{1, \dots, k\}$ ,  $j \neq v$  the covariance is zero

$$\begin{aligned} \text{cov}(h_i m_{ji} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i}, h_i m_{vi} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)i}) &= \mathbb{E}[h_i^2 m_{ji} m_{vi} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)i}] \\ &= \mathbb{E}[m_{ji} m_{vi}] \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)i} \\ &= 0. \end{aligned}$$

We use (P2) to conclude that  $\mathbb{E}[m_{ji} m_{vi}] = 0$ . Therefore for  $i = 2, \dots, r_n$

$$\begin{aligned} \text{var}(Z_{ni}) &= \frac{1}{k} \text{var} \left( \sum_{j=1}^k h_i m_{ji} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \right) \\ &= \frac{1}{k} \sum_{j=1}^k \text{var}(h_i m_{ji} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i}) \\ &= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \mathbf{u}_{(n)i}^\top \boldsymbol{\lambda}_j. \end{aligned} \tag{S.37}$$

Results are different for  $i = 1$  as the first column of the Hadamard matrix is all ones (P1). For  $j, v \in \{1, \dots, k\}$ ,  $j \neq v$  the covariance is

$$\begin{aligned} \text{cov}(h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1}, h_1 m_{v1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1}) &= \mathbb{E}[h_1^2 m_{j1} m_{v1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1}] \\ &= \mathbb{E}[m_{j1} m_{v1}] \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1} \\ &= \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1}. \end{aligned}$$

From (P1)  $m_{j1} = m_{v1} = 1$ . Now using the Cauchy-Schwarz inequality,

$$\begin{aligned} |\text{cov}(h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1}, h_1 m_{v1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1})| &= |\boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1}| \\ &\leq \|\boldsymbol{\lambda}_j\|_2 \|\mathbf{u}_{(n)1}\|_2 \|\boldsymbol{\lambda}_v\|_2 \|\mathbf{u}_{(n)1}\|_2 \\ &\leq \|\mathbf{u}_{(n)1}\|_2 \|\mathbf{u}_{(n)1}\|_2 \\ &= \|\mathbf{u}_{(n)1}\|_2^2 \end{aligned}$$

The second last last uses the fact that  $\boldsymbol{\lambda}$  is a unit vector and we must have  $\|\boldsymbol{\lambda}_j\|_2 \leq 1$ ,  $\|\boldsymbol{\lambda}_v\|_2 \leq 1$  for any  $j, k$ . From assumption 1, the right hand side of the previous inequality tends to zero as  $n$  tends to infinity. As such we conclude that  $|\text{cov}(h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1}, h_1 m_{v1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1})|$  is  $o(1)$ . Some covariance

terms appear in the expression for  $\text{var}(Z_{n1})$

$$\begin{aligned}
\text{var}(Z_{n1}) &= \frac{1}{k} \text{var} \left( \sum_{j=1}^k h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1} \right) \\
&= \frac{1}{k} \sum_{j=1}^k \text{var} (h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1}) + \frac{1}{k} 2 \sum_{j=1}^{k-1} \sum_{v=j+1}^k \text{cov} (h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1}, h_1 m_{v1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1}) \\
&= \frac{1}{k} \sum_{j=1}^k \text{var} (h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1}) + \frac{1}{k} 2 \sum_{j=1}^{k-1} \sum_{v=j+1}^k \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1} \boldsymbol{\lambda}_v^\top \mathbf{u}_{(n)1} \\
&= \frac{1}{k} \sum_{j=1}^k \text{var} (h_1 m_{j1} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1}) + o(1) \\
&= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)1} \mathbf{u}_{(n)1}^\top \boldsymbol{\lambda}_j + o(1) \tag{S.38}
\end{aligned}$$

The trailing term can be grouped into an  $o(1)$  term as the sketch size  $k$  is fixed in our analysis. Using (S.37) and (S.38) we can then determine the row-wise variance totals  $s_n^2$ :

$$\begin{aligned}
s_n^2 &= \frac{1}{k} \sum_{i=1}^{r_n} \text{var}(Z_{ni}) \\
&= \frac{1}{k} \sum_{i=1}^{r_n} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i} \mathbf{u}_{(n)i}^\top \boldsymbol{\lambda}_j + o(1) \\
&= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \left( \sum_{i=1}^{r_n} \mathbf{u}_{(n)i} \mathbf{u}_{(n)i}^\top \right) \boldsymbol{\lambda}_j + o(1) \\
&= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{U}_{(n)}^\top \mathbf{U}_{(n)} \boldsymbol{\lambda}_j + o(1) \\
&= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \mathbf{I}_d \boldsymbol{\lambda}_j + o(1) \\
&= \frac{1}{k} \sum_{j=1}^k \boldsymbol{\lambda}_j^\top \boldsymbol{\lambda}_j + o(1) \\
&= \frac{1}{k} + o(1).
\end{aligned}$$

The step in the last line follows as we have taken  $\boldsymbol{\lambda}$  to be a unit vector. The fact that  $\mathbf{U}_{(n)}^\top \mathbf{U}_{(n)} = \mathbf{I}_d$  for all  $n$  serves as a useful normalisation to give stable limiting behaviour. We are working with a triangular array where the rows are nearly standardised. Asymptotically in  $n$ ,  $s_n^2 \rightarrow 1/k$ .

We now establish a sequence of upper bounds  $(K_n)$ . As the random variables in the construction of construction of the Hadamard sketch are bounded, we can bound the random variables in the triangular array (S.34) using the leverage scores of the sequence of source datasets. Now as the

random sign  $h_i \in \{+1, -1\}$  we have that for all  $i = 1, \dots, r_n$ :

$$\begin{aligned} |Z_{ni}| &= \frac{1}{\sqrt{k}} |h_i \sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i}| \\ &= \frac{1}{\sqrt{k}} |\sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j^\top \mathbf{u}_{(n)i}|. \end{aligned}$$

Now using the Cauchy-Schwarz inequality,

$$\frac{1}{\sqrt{k}} \left| \left( \sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j^\top \right) \mathbf{u}_{(n)i} \right| \leq \frac{1}{\sqrt{k}} \left\| \left( \sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j \right) \right\|_2 \|\mathbf{u}_{(n)i}\|_2. \quad (\text{S.39})$$

Using the triangle inequality,

$$\left\| \left( \sum_{j=1}^k m_{ji} \boldsymbol{\lambda}_j \right) \right\|_2 \leq \sum_{j=1}^k \|m_{ji} \boldsymbol{\lambda}_j\|_2. \quad (\text{S.40})$$

Now as  $m_{ji} \in \{+1, -1\}$  for all  $j = 1, \dots, k$ ,

$$\sum_{j=1}^k \|m_{ji} \boldsymbol{\lambda}_j\|_2 = \sum_{j=1}^k \|\boldsymbol{\lambda}_j\|_2. \quad (\text{S.41})$$

As  $\boldsymbol{\lambda}$  is a unit vector we can easily form the bound

$$\sum_{j=1}^k \|\boldsymbol{\lambda}_j\|_2 \leq k. \quad (\text{S.42})$$

Substituting (S.41) and (S.42) into (S.39) leads to the upper bound for  $i = 1, \dots, r_n$ :

$$\begin{aligned} |Z_{ni}| &\leq \frac{1}{\sqrt{k}} k \|\mathbf{u}_{(n)i}\|_2 \\ &= \sqrt{k} \|\mathbf{u}_{(n)i}\|_2. \end{aligned} \quad (\text{S.43})$$

We can then form the sequence of upper bounds  $K_n$ :

$$K_n = \sqrt{k} \max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2.$$

We have that  $|Z_{ni}| \leq K_n$  almost surely for  $i = 1, \dots, r_n$  and  $n \in \mathbb{N}$ . Assumption 1 (recall equation (S.21)) gives the limiting behaviour of  $K_n$ . As the sketch size  $k$  is fixed in our analysis,

$$\begin{aligned} \lim_{n \rightarrow \infty} K_n &= \sqrt{k} \lim_{n \rightarrow \infty} \max_{i=1, \dots, r_n} \|\mathbf{u}_{(n)i}\|_2 \\ &= 0. \end{aligned}$$

We have that  $K_n \rightarrow 0$  as  $n \rightarrow \infty$ . As  $s_n^2 = 1/k + o(1)$  we have an asymptotically standardised array,

and  $K_n/s_n \rightarrow 0$ . We can use Theorem 9 to conclude that the triangular array of random variables defined in (S.34) satisfies Lindeberg's condition. As such, the conditions of Theorem 10 are satisfied. We conclude that the row sums in (S.35) converge in distribution to  $N(0, 1/k)$ . Finally, the Cramér-Wold device gives that the whitened sketched dataset has a limiting matrix normal distribution. That is the sequence of random matrices  $\widetilde{\mathbf{A}}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1}$  converges in distribution to a  $MN(\mathbf{0}, \mathbf{I}_k, \mathbf{I}_d/k)$  random matrix.

## H Proof of Theorem 5 (Complete sketching asymptotics)

**Assumption 2:**

$$\lim_{n \rightarrow \infty} n^{-1} \begin{bmatrix} \mathbf{y}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{y}_{(n)}^\top \mathbf{X}_{(n)} \\ \mathbf{X}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \end{bmatrix} = \mathbf{Q} \quad \text{for some positive-definite matrix } \mathbf{Q}.$$

**Theorem.** *Suppose that Assumptions 1 and 2 hold,  $k \geq p$ , and  $\beta_S$  is computed using a Hadamard or Clarkson-Woodruff sketch. Let  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^+$  denote the Moore-Penrose pseudo-inverse of  $(\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})$ . Let*

$$\widetilde{\mathbf{H}}_{(n)} = \frac{RSS_F^{(n)}}{k} (\widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}})^+ \quad \text{and} \quad \mathbf{H}_{(n)} = \frac{RSS_F^{(n)}}{k-p+1} (\mathbf{X}_{(n)}^\top \mathbf{X}_{(n)})^{-1}.$$

Then as  $n \rightarrow \infty$ , convergence in distribution holds for

$$\begin{aligned} (i) [\mathbf{H}_{(n)}^{-1/2} (\beta_S - \beta_F^{(n)}) | \mathbf{A}_{(n)}] &\rightarrow \text{Student}(\mathbf{0}, \mathbf{I}_p, k-p+1), \\ (ii) [\widetilde{\mathbf{H}}_{(n)}^{-1/2} (\beta_S - \beta_F^{(n)}) | \mathbf{A}_{(n)}] &\rightarrow N(\mathbf{0}, \mathbf{I}_p). \end{aligned}$$

Notation is slightly heavier in the proof compared to the main text for the sake of clarity. Again we do not explicitly condition on the source dataset  $\mathbf{A}_{(n)}$ , the source dataset is always fixed, and the only randomness is from the sketching matrix. The sketched data will be denoted  $\widetilde{\mathbf{y}}_{(n)}$  and  $\widetilde{\mathbf{X}}_{(n)}$  to denote the dependence on the  $n \times d$  source dataset. So  $\widetilde{\mathbf{y}}_{(n)} = \mathbf{S}\mathbf{y}_{(n)}$  and  $\widetilde{\mathbf{X}}_{(n)} = \mathbf{S}\mathbf{X}_{(n)}$ . The dimension of the sketched dataset does not change.

Assumption 2 is of assistance in establishing the limit theorem. Let

$$\mathbf{Q}_{(n)} = n^{-1} \begin{bmatrix} \mathbf{y}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{y}_{(n)}^\top \mathbf{X}_{(n)} \\ \mathbf{X}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \end{bmatrix}$$

The matrix  $\mathbf{Q}_{(n)}$  contains the sufficient statistics needed to fit a Gaussian linear model,  $\mathbf{y}_{(n)}^\top \mathbf{y}_{(n)}$ ,  $\mathbf{X}_{(n)}^\top \mathbf{y}_{(n)}$  and  $\mathbf{X}_{(n)}^\top \mathbf{X}_{(n)}$  given the source dataset  $\mathbf{A}_{(n)} = [\mathbf{y}_{(n)}, \mathbf{X}_{(n)}]$ . Assumption 2 states the averaged sufficient statistic matrix converges to a limiting matrix  $\mathbf{Q}$ . It will be helpful to partition the limiting matrix  $\mathbf{Q}$  as

$$\mathbf{Q} = \lim_{n \rightarrow \infty} n^{-1} \begin{bmatrix} \mathbf{y}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{y}_{(n)}^\top \mathbf{X}_{(n)} \\ \mathbf{X}_{(n)}^\top \mathbf{y}_{(n)} & \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \end{bmatrix} = \begin{bmatrix} s & \mathbf{m}^\top \\ \mathbf{m} & \mathbf{G} \end{bmatrix}, \quad (\text{S.44})$$

where  $s$  is a scalar,  $\mathbf{G}$  is a  $p \times p$  matrix and  $\mathbf{m}$  is a  $p$ -length column vector. The matrix  $\mathbf{G}$  is the limiting

averaged Gram matrix of the predictors. The vector  $\mathbf{m}$  is the limit of the predictor response inner products  $n^{-1}\mathbf{X}_{(n)}^\top \mathbf{y}_{(n)}$ , and the scalar  $s$  is the limit of the mean total sum of squares  $n^{-1}\mathbf{y}_{(n)}^\top \mathbf{y}_{(n)}$ .

As mentioned, the assumption of a sequence of source datasets also gives a sequence of optimal least squares coefficients and residual errors. Let  $\sigma_F^{2(n)} = RSS_F/n$ . Define the limiting least squares coefficient estimate as  $\boldsymbol{\beta} = \lim_{n \rightarrow \infty} \boldsymbol{\beta}_F^{(n)}$  and the limiting residual error as  $\sigma^2 = \lim_{n \rightarrow \infty} \sigma_F^{2(n)}$ . Both  $\boldsymbol{\beta}$  and  $\sigma^2$  can be expressed as functions of the matrix  $\mathbf{Q}$ . Specifically,

$$\boldsymbol{\beta} = \mathbf{G}^{-1}\mathbf{m}, \tag{S.45}$$

$$\sigma^2 = s - \mathbf{m}^\top \mathbf{G}^{-1}\mathbf{m}. \tag{S.46}$$

From Assumption 2, we have that  $n^{-1}\mathbf{V}_{(n)}\mathbf{D}_{(n)}^2\mathbf{V}_{(n)}^\top \rightarrow \mathbf{Q}$ . As such we have that  $n^{-1/2}\mathbf{D}_{(n)}\mathbf{V}_{(n)}^\top \rightarrow \mathbf{Q}^{1/2}$ . From the sketching central limit theorem the whitened sketched data converges to a matrix normal distribution

$$[\tilde{\mathbf{y}}_{(n)}, \tilde{\mathbf{X}}_{(n)}]\mathbf{V}_{(n)}\mathbf{D}_{(n)}^{-1} \xrightarrow{d} MN(\mathbf{0}_{k \times d}, \mathbf{I}_k, \mathbf{I}_d/k)$$

The benefit of adding Assumption 2 is that using Slutsky's theorem we have the additional convergence result

$$n^{-1/2}[\tilde{\mathbf{y}}_{(n)}, \tilde{\mathbf{X}}_{(n)}] \xrightarrow{d} MN(\mathbf{0}, \mathbf{I}_k, \mathbf{Q}/k).$$

To prove results (i) and (ii) we use the continuous mapping theorem (Van Der Vaart, 1998, p. 7) in conjunction with the previous convergence result. It will be helpful to define the random variables  $\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L$  as having the above limiting matrix normal distribution

$$[\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L] \sim MN(\mathbf{0}_{k \times d}, \mathbf{I}_k, \mathbf{Q}/k).$$

This is so we can say that

$$n^{-1/2}[\tilde{\mathbf{y}}_{(n)}, \tilde{\mathbf{X}}_{(n)}] \xrightarrow{d} [\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L].$$

**Lemma 5** (Continuous Mapping Theorem). *Let  $(\mathbf{Z}_n)_{n \in \mathbb{N}}$  indicate a sequence of random vectors in  $\mathbb{R}^d$  and  $\mathbf{Z}$  indicate another random vector in  $\mathbb{R}^d$ . Suppose the function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^m$  is continuous at every point of a set  $C$  such that  $P(\mathbf{Z} \in C) = 1$ . Then if  $\mathbf{Z}_n \xrightarrow{d} \mathbf{Z}$  then  $g(\mathbf{Z}_n) \xrightarrow{d} g(\mathbf{Z})$ .*

In Lemma 5, the function  $g : \mathbb{R}^d \rightarrow \mathbb{R}^m$  does not change with  $n$ , and the dimensions  $d$  and  $m$  are fixed when taking limits. The sketched estimator  $\boldsymbol{\beta}_S$  can be defined as a function of the sketched data that is continuous over the set where  $\tilde{\mathbf{X}}_{(n)}$  is of full rank. Formally we could say that  $\boldsymbol{\beta}_S = g(n^{-1/2}\tilde{\mathbf{y}}_{(n)}, n^{-1/2}\tilde{\mathbf{X}}_{(n)})$ . As  $\tilde{\mathbf{X}}_L$  is of rank  $p$  almost surely, and  $\tilde{\mathbf{X}}_{(n)} \xrightarrow{d} \tilde{\mathbf{X}}_L$  we can apply the continuous mapping theorem to determine the limiting distribution of the  $\boldsymbol{\beta}_S$ . The random matrix  $[\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L]$  can be described using a hierarchical model completely analogous in structure to the

hierarchical model established for the Gaussian sketch in Section 3.1 of the main text. Specifically,

$$\begin{aligned}\tilde{\mathbf{y}}_L \mid \tilde{\mathbf{X}}_L &\sim N\left(\tilde{\mathbf{X}}_L \boldsymbol{\beta}, \frac{1}{k} \sigma^2 \mathbf{I}_k\right), \\ \tilde{\mathbf{X}}_L &\sim MN\left(\mathbf{0}, \mathbf{I}_k, \frac{1}{k} \mathbf{G}\right).\end{aligned}$$

From Theorem 2 in the main text, and recalling that the function  $g$  outputs  $\boldsymbol{\beta}_S$ , we have that

$$g(\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L) \sim \text{Student}\left(\boldsymbol{\beta}, \frac{\sigma^2}{k-p+1} \mathbf{G}^{-1}, k-p+1\right).$$

As such, for the Hadamard and Clarkson-Woodruff sketches,

$$[\boldsymbol{\beta}_S \mid \mathbf{y}_{(n)}, \mathbf{X}_{(n)}] \xrightarrow{d} \text{Student}\left(\boldsymbol{\beta}, \frac{\sigma^2}{k-p+1} \mathbf{G}^{-1}, k-p+1\right).$$

Let

$$\mathbf{H}_{(n)} = \sigma_F^{2(n)} / (k-p+1) \left(n^{-1} \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)}\right)^{-1}.$$

Now as  $n^{-1} \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \rightarrow \mathbf{G}$ ,  $\sigma_F^{2(n)} \rightarrow \sigma^2$ , and  $\boldsymbol{\beta}_F^{(n)} \rightarrow \boldsymbol{\beta}$ , Slutsky's theorem can be used to arrive at (i),

$$\mathbf{H}_{(n)}^{-1/2} (\boldsymbol{\beta}_S - \boldsymbol{\beta}_F) \xrightarrow{d} \text{Student}(\mathbf{0}, \mathbf{I}_p, k-p+1).$$

For result (ii), let us define the function

$$\begin{aligned}f(n^{-1/2} \tilde{\mathbf{y}}_{(n)}, n^{-1/2} \tilde{\mathbf{X}}_{(n)}) &= \left[ n \left( \tilde{\mathbf{X}}_{(n)}^\top \tilde{\mathbf{X}}_{(n)} \right)^+ \right]^{-1/2} \left( \tilde{\mathbf{X}}_{(n)}^+ \tilde{\mathbf{y}}_{(n)} - \boldsymbol{\beta} \right) \\ &= \left[ n \left( \tilde{\mathbf{X}}_{(n)}^\top \tilde{\mathbf{X}}_{(n)} \right)^+ \right]^{-1/2} (\boldsymbol{\beta}_S - \boldsymbol{\beta}).\end{aligned}$$

This function transforms the  $\boldsymbol{\beta}_S$  so that the output is uncorrelated. This function is also continuous over the set where  $\tilde{\mathbf{X}}_{(n)}$  is of rank  $p$ . Again using the fact that  $\tilde{\mathbf{X}}_L$  has rank  $p$  almost surely, it follows from the continuous mapping theorem that  $f(n^{-1/2} \tilde{\mathbf{y}}_{(n)}, n^{-1/2} \tilde{\mathbf{X}}_{(n)}) \xrightarrow{d} f(\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L)$ . Result (ii) in Theorem 2 also applies to the hierarchical model for  $\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L$ , and gives the distribution of the transformed  $\boldsymbol{\beta}_S$  under the Gaussian sketch. The distribution of  $f(\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L)$  will be

$$f(\tilde{\mathbf{y}}_L, \tilde{\mathbf{X}}_L) \sim N\left(\mathbf{0}, \frac{\sigma^2}{k} \mathbf{I}_p\right).$$

As such, for the Clarkson-Woodruff and Hadamard sketches,

$$\left[ n \left( \tilde{\mathbf{X}}_{(n)}^\top \tilde{\mathbf{X}}_{(n)} \right)^+ \right]^{-1/2} (\boldsymbol{\beta}_S - \boldsymbol{\beta}) \xrightarrow{d} N\left(\mathbf{0}, \frac{\sigma^2}{k} \mathbf{I}_p\right).$$

Now let

$$\widetilde{\mathbf{H}}_{(n)} = n\sigma_F^{2(n)}/k \left( \widetilde{\mathbf{X}}^\top \widetilde{\mathbf{X}} \right)^+$$

Now as  $\sigma_F^{2(n)} \rightarrow \sigma^2$ , and  $\beta_F^{(n)} \rightarrow \beta$ , Slutsky's theorem can be used to arrive at (ii)

$$\widetilde{\mathbf{H}}_{(n)}^{-1/2} (\beta_S - \beta_F^{(n)}) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}_p).$$

## I Proof of Theorem 6 (Partial sketching asymptotics)

**Theorem.** *Suppose that Assumptions 1, 2 and 3 hold,  $k > p + 3$ , and  $\beta_P^*$  is computed using a Hadamard or Clarkson-Woodruff sketch. Let*

$$\mathbf{H}_{(n)} = \frac{(k-p-1)}{(k-p)(k-p-3)} \left( MSS_F^{(n)} (\mathbf{X}_{(n)}^\top \mathbf{X}_{(n)})^{-1} + \frac{k-p+1}{k-p-1} \beta_F^{(n)} \beta_F^{(n)\top} \right).$$

Then as  $n \rightarrow \infty$ ,

$$\begin{aligned} (i) \quad E_S(\beta_P^* - \beta_F^{(n)} | \mathbf{A}_{(n)}) &\rightarrow \mathbf{0}. \\ (ii) \quad \text{var}_S \left( \mathbf{H}_{(n)}^{-1/2} (\beta_P^* - \beta_F^{(n)}) | \mathbf{A}_{(n)} \right) &\rightarrow \mathbf{I}_d \end{aligned}$$

Application of the continuous mapping theorem gives that the distribution of  $\beta_S$  and  $\beta_P^*$  under the Hadamard and Clarkson-Woodruff sketches converges to the distribution of the estimators under the Gaussian sketch. This does not necessarily guarantee convergence in moments. To establish a limit theorem for the bias and variance of the estimators, we need a uniform integrability condition on the sketched dataset. The sketched data will be denoted  $\widetilde{\mathbf{X}}_{(n)}$  to denote the dependence on the  $n \times p$  source covariate matrix. So  $\widetilde{\mathbf{X}}_{(n)} = \mathbf{S}\mathbf{X}_{(n)}$ . We again do not explicitly condition on the source dataset  $\mathbf{A}_{(n)} = [\mathbf{y}_{(n)}, \mathbf{X}_{(n)}]$  in the following working.

Let  $\mathbf{G}_{(n)} = n^{-1} \widetilde{\mathbf{X}}_{(n)}^\top \widetilde{\mathbf{X}}_{(n)}$ . From the continuous mapping theorem and Theorem 3, it is known that

$$\mathbf{G}_{(n)}^{-1} \xrightarrow{d} \mathbf{W},$$

where  $\mathbf{W}$  has an Inverse-Wishart( $k, k\mathbf{Q}^{-1}$ ) distribution and  $\mathbf{Q}$  is the limiting matrix from assumption 2. We would like to establish convergence in first and second moments, that is

$$\begin{aligned} \mathbb{E}(\mathbf{G}_{(n)}^{-1}) &\rightarrow \mathbb{E}(\mathbf{W}), \\ \text{var}(\mathbf{G}_{(n)}^{-1}) &\rightarrow \text{var}(\mathbf{W}). \end{aligned}$$

If convergence in first and second moments occurs, then we can show that (i) and (ii) will hold. If  $\mathbb{E}(\mathbf{G}_{(n)}^{-1}) \rightarrow \mathbb{E}(\mathbf{W})$ , we can say that

$$\mathbb{E}(\beta_P^* - \beta) \rightarrow \mathbf{0},$$



where  $\boldsymbol{\beta}$  is the limiting ordinary least squares estimator (S.45), that is a function of the limiting matrix  $\mathbf{Q}$  in Assumption 2. From here, using that  $\lim_{n \rightarrow \infty} \boldsymbol{\beta}_F^{(n)}$ , Slutsky's theorem can be used to arrive at (i)

$$\mathbb{E}(\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F^{(n)}) \rightarrow \mathbf{0}.$$

To show convergence of the variance of the sketched estimator (ii), we define

$$\begin{aligned} \mathbf{H}_{(n)} &= \frac{(k-p-1)}{(k-p)(k-p-3)} \left( MSS_F^{(n)} \left( \mathbf{X}_{(n)}^\top \mathbf{X}_{(n)} \right)^{-1} + \frac{(k-p+1)}{(k-p-1)} \boldsymbol{\beta}_F^{(n)} \boldsymbol{\beta}_F^{(n)\top} \right), \\ \mathbf{H} &= \frac{(k-p-1)}{(k-p)(k-p-3)} \left( (s - \sigma^2) \mathbf{G}^{-1} + \frac{(k-p+1)}{(k-p-1)} \boldsymbol{\beta} \boldsymbol{\beta}^\top \right). \end{aligned}$$

Where  $s$ ,  $\sigma^2$  and  $\mathbf{G}$  are functions of the limiting matrix  $\mathbf{Q}$ , as in (S.44), (S.45) and (S.46). If  $\text{var}(\mathbf{G}_{(n)}^{-1}) \rightarrow \text{var}(\mathbf{W})$  it follows that

$$\text{var}_S \left( \mathbf{H}^{-1/2} (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}) \right) \rightarrow \mathbf{I}_d.$$

As  $\mathbf{H}_{(n)}$  converges to  $\mathbf{H}$  and  $\boldsymbol{\beta}_F^{(n)}$  converges to  $\boldsymbol{\beta}$  asymptotically with  $n$ , an application of Slutsky's theorem gives (ii),

$$\text{var}_S \left( \mathbf{H}_{(n)}^{-1/2} (\boldsymbol{\beta}_P^* - \boldsymbol{\beta}_F^{(n)}) \right) \rightarrow \mathbf{I}_d$$

As such, if we can establish that  $\text{var}(\mathbf{G}_{(n)}^{-1}) \rightarrow \text{var}(\mathbf{W})$  we have proved (ii). The following theorem describes the necessary conditions for such convergence to occur.

**Theorem 11.** (Billingsley, 1968, Theorem 5.4) *Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  be a sequence of random vectors. Suppose  $\mathbf{X}_n$  converges in distribution to a random variable  $\mathbf{Z}$  as  $n$  tends to infinity. For the additional convergence of moments  $\mathbb{E}[\mathbf{X}_n] \rightarrow \mathbb{E}[\mathbf{Z}]$  and  $\text{var}[\mathbf{X}_n] \rightarrow \text{var}[\mathbf{Z}]$ , it must hold that for all conformable constant vectors  $\boldsymbol{\lambda}$*

$$\lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} |\boldsymbol{\lambda}^\top \mathbf{X}_n|^2 \mathbb{1}_{\{|\boldsymbol{\lambda}^\top \mathbf{X}_n|^2 \geq M\}} = 0.$$

The above condition can be difficult to verify directly. It can be shown that if asymptotically  $|\boldsymbol{\lambda}^\top \mathbf{X}_n|$  has a bounded fourth moment, then the integrability condition is satisfied (Van Der Vaart, 1998, section 2.5).

A linear combination of the elements of the random matrix  $\mathbf{G}_{(n)}^{-1}$  can be written as  $\text{trace}(\boldsymbol{\Lambda} \mathbf{G}_{(n)}^{-1})$  for a  $p \times p$  matrix of constants  $\boldsymbol{\Lambda}$ . It is easier to work with this form rather than stacking the elements of the random matrix to form a random vector. From theorem 11, it is sufficient to show that the expected value of  $|\text{trace}(\boldsymbol{\Lambda} \mathbf{G}_{(n)}^{-1})|^4$  is finite for large  $n$  to show the desired convergence in moments.

As  $\text{trace}(\boldsymbol{\Lambda} \mathbf{G}_{(n)}^{-1})$  equals the sum of the singular values of the matrix  $\boldsymbol{\Lambda} \mathbf{G}_{(n)}^{-1}$ , we can form an

upper bound on the value,

$$\begin{aligned}\text{trace}(\mathbf{\Lambda}\mathbf{G}_{(n)}^{-1}) &\leq p\|\mathbf{\Lambda}\mathbf{G}_{(n)}^{-1}\|_2 \\ &\leq p\|\mathbf{\Lambda}\|_2\|\mathbf{G}_{(n)}^{-1}\|_2\end{aligned}$$

Squaring both sides gives an upper bound on the quantity that must satisfy the uniform integrability condition,

$$|\text{trace}(\mathbf{\Lambda}\mathbf{G}_{(n)}^{-1})|^2 \leq p^2\|\mathbf{\Lambda}\|_2^2\|\mathbf{G}_{(n)}^{-1}\|_2^2.$$

Squaring again gives an upper bound on the fourth moment of the linear combination of interest

$$\begin{aligned}|\text{trace}(\mathbf{\Lambda}\mathbf{G}_{(n)}^{-1})|^4 &\leq p^4\|\mathbf{\Lambda}\|_2^4\|\mathbf{G}_{(n)}^{-1}\|_2^4 \\ &= p^4\|\mathbf{\Lambda}\|_2^4\left(\frac{1}{\sigma_{\min}^2(\mathbf{G}_{(n)})}\right)^2\end{aligned}$$

By Assumption 3, the expectation of the right hand side is finite. As such, the uniform integrability condition holds and we can conclude that

$$\begin{aligned}\mathbb{E}(\mathbf{G}_{(n)}^{-1}) &\rightarrow \mathbb{E}(\mathbf{W}), \\ \text{var}(\mathbf{G}_{(n)}^{-1}) &\rightarrow \text{var}(\mathbf{W}).\end{aligned}$$

As discussed at the beginning of the proof this is sufficient to show that (i) and (ii) hold.