

EVERYTHING IS VECCHIA: UNIFYING COLUMN NYSTRÖM AND SPARSE VECCHIA APPROXIMATIONS

EAGAN KAMINETZ*

Abstract. This thesis examines two types of factored matrix approximations that are fast and accurate in different contexts: low-rank approximations of positive definite matrices and sparse factored approximations of positive definite matrix inverses. The thesis combines these two approaches using the theoretical framework of Vecchia approximation. The main algorithmic contribution is a greedy method for computing the low-rank portion of a low-rank plus sparse approximation that achieves near-optimal theoretical guarantees. The greedy method is applied to produce preconditioners for iteratively solving linear systems arising from Gaussian processes, where numerical experiments demonstrate significant convergence rate improvements over existing methods.

Key words. Vecchia approximation, Nyström approximation, kernel matrix, Gaussian process, Factorized sparse approximate inverse

AMS subject classifications. 65F55, 65C99, 68T05

1. Motivation. This thesis asks and answers the following question: what is a fast, accurate algorithm that provides a factored approximation of a positive definite matrix?

A factored approximation of a positive definite matrix is useful for scientific computing in at least two ways. The approximation can be used as a “direct solver”, when it is substituted for the original matrix when solving linear systems, calculating determinants, or calculating eigenvectors. Direct solvers are computationally efficient, and they induce error, but the error level is acceptable given a highly accurate approximation. As an alternative approach, a factored approximation can serve as a preconditioner for an iterative method. Iterative methods for linear solves and eigenvalue solves drive the error all the way to zero, but the rate of convergence depends on the approximation quality of the matrix preconditioner; for details see [section 2](#).

In the large-data regime, both direct solvers and preconditioned iterative solvers are much faster than dense factorization methods for scientific computing, and they stretch the limits of the calculations that can be performed. For example, given a $10^5 \times 10^5$ matrix in single precision, it is impossible to execute dense factorization methods using standard software (MATLAB or python) on a laptop-scale computer (64GB RAM). In contrast, preconditioned iterative solvers execute in minutes or hours, and direct solvers are even faster.

This thesis will examine a large, existing body of literature devoted to factored approximations of positive definite matrices. The literature is divided into low-rank approximations [\[4, 5, 6\]](#) and sparse factored approximations of the matrix inverse [\[11, 18\]](#), as well as combinations of the two approaches [\[24\]](#). All these approximation methods execute quickly and apply in the large-data regime, since they examine just a subset of entries in the matrix. Yet, they are based on different philosophies, since the low-rank approximations are accurate for matrices with quickly decaying eigenvalues and the sparse factored approximations are valid for matrices with off-diagonal decay in the inverse. It is unclear whether these approximations can be combined in an overarching framework.

The thesis considers a unified framework for low-rank plus sparse approximation, and it makes the following three contributions:

*University of California San Diego, La Jolla, CA (ekaminetz@ucsd.edu).

1. The thesis synthesizes the existing literature by showing how existing low-rank and sparse approximation methods [4, 5, 6, 11, 18, 24] can be understood as particular examples of Vecchia approximation [22] with different sparsity patterns. Therefore, the optimality guarantees of the Vecchia approximation [12, 23, 19] extend to all these methods.
2. Using the theory of Vecchia approximations, this thesis introduces a greedy method to build the low-rank part of a low-rank plus sparse approximation [24], and it proves near-optimal theoretical guarantees for this greedy method.
3. Numerical experiments show that the greedy method leads to significant convergence rate improvements when used to precondition iterative methods.

The rest of the thesis is organized as follows. [Subsection 1.1](#) introduces notation. [Section 2](#) provides background on approximate matrix factorizations. [Section 3](#) unifies approaches in [section 2](#) using the theoretical framework of Vecchia approximation. [Section 4](#) introduces the greedy algorithm. Last, [section 5](#) presents numerical experiments.

1.1. Notation. We denote scalars in lower case italics: m, n, r . We denote vectors in lower case boldface: \mathbf{u}, \mathbf{v} . We denote matrices in boldface capital letters: \mathbf{A}, \mathbf{B} . We use \mathbf{u}_i to refer to the i th entry of a vector $u \in \mathbb{C}^n$, and we use $\mathbf{A}_{i,j}$ to refer to the (i, j) entry of a matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$. Given index sets $S, T \subseteq \{1, \dots, n\}$, we use \mathbf{u}_S to refer to the subvector $(\mathbf{u}_i)_{i \in S}$, and we use $\mathbf{A}_{S,T}$ to refer to the submatrix $(\mathbf{A}_{i,j})_{i \in S, j \in T}$. Similarly, $\mathbf{A}_{i,:}$ and $\mathbf{A}_{:,i}$ indicate the i th row and column of \mathbf{A} .

We represent the vector of all ones by $\mathbf{1}$. We use \mathbf{e}_i to represent the i th standard basis vector, which has i th entry 1 and all other entries 0. We write $[n]$ to mean $\{1 \dots n\}$, and write $m : k$ to refer to $\{m \dots k\}$, especially when indexing vectors or matrices. We use $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ to indicate the multivariate Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. The conjugate transpose of a matrix $\mathbf{A} \in \mathbb{C}^{m \times n}$ is denoted \mathbf{A}^* , and the Moore-Penrose pseudoinverse is \mathbf{A}^+ .

The Nystrom approximation of a positive definite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ with respect to a test matrix $\mathbf{X} \in \mathbb{C}^{n \times r}$ is defined as

$$\mathbf{A}\langle \mathbf{X} \rangle = \mathbf{A}\mathbf{X}(\mathbf{X}^* \mathbf{A}\mathbf{X})^{-1} \mathbf{X}^* \mathbf{A}.$$

The column Nystrom approximation using an index set $R \subseteq \{1, \dots, n\}$ is defined as

$$\mathbf{A}\langle R \rangle = \mathbf{A}\langle \mathbf{I}(:, R) \rangle = \mathbf{A}_{:,R}(\mathbf{A}_{R,R})^{-1} \mathbf{A}_{R,:}.$$

The Schur complement with respect to index set R is defined as

$$\mathbf{A}/R = \mathbf{A} - \mathbf{A}\langle R \rangle.$$

2. Methods for factored approximations. This section describes several methods to construct factored approximations for positive definite matrices $\mathbf{A} \in \mathbb{C}^{n \times n}$. The main purpose of this section is to provide motivation and to set the notation.

2.1. Approximate Cholesky factorization. Every positive definite matrix \mathbf{A} has a unique Cholesky factorization, $\mathbf{A} = \mathbf{L}\mathbf{L}^*$, where \mathbf{L} is a lower triangular matrix. Since the inverse \mathbf{A}^{-1} is positive definite, it also admits a unique Cholesky factorization, $\mathbf{A}^{-1} = \mathbf{C}\mathbf{C}^*$. However, the factorizations $\mathbf{A} = \mathbf{L}\mathbf{L}^*$ and $\mathbf{A} = \mathbf{C}^{-*}\mathbf{C}^{-1}$ are distinct, because \mathbf{L} is lower triangular while \mathbf{C}^{-*} is upper triangular.

In practice, it is often too expensive to form a Cholesky factorization exactly. However, we can consider an approximate Cholesky factorization $\hat{\mathbf{A}} = \hat{\mathbf{L}}\hat{\mathbf{L}}^*$ or $\hat{\mathbf{A}}^{-1} = \hat{\mathbf{C}}\hat{\mathbf{C}}^*$, where the matrix $\hat{\mathbf{L}}$ or $\hat{\mathbf{C}}$ is lower triangular and sparse. When such an approximation is available, we can put this factorization into action to solve linear systems $\mathbf{A}\mathbf{x} = \mathbf{b}$ using a direct or indirect approach.

In the direct approach, we replace the original matrix \mathbf{A} with an approximate Cholesky factorization $\hat{\mathbf{A}}$ and then solve the modified linear system $\hat{\mathbf{A}}\hat{\mathbf{x}} = \mathbf{b}$. At this point, solving the system is extremely efficient. For example, given an approximate factorization $\hat{\mathbf{A}}^{-1} = \hat{\mathbf{C}}\hat{\mathbf{C}}^*$ we can solve the linear system using two matrix–vector products, which require at most $\mathcal{O}(n^2)$ operations and run even faster if $\hat{\mathbf{C}}$ is sparse. Alternatively, given an approximate factorization $\hat{\mathbf{A}} = \hat{\mathbf{L}}\hat{\mathbf{L}}^*$, we can solve the linear system using two instances of triangular substitution. The runtime for triangular substitution is linear in the number of nonzero entries, so it requires $\mathcal{O}(n^2)$ operations for a dense matrix and runs even faster when the Cholesky factor \mathbf{L} is sparse.

In an indirect method, we do not substitute $\hat{\mathbf{A}}$ for \mathbf{A} directly. Rather, we use the approximation $\hat{\mathbf{A}}$ as a preconditioner in the preconditioned conjugate gradient (PCG) algorithm. We run PCG for some number of iterations, producing better and better estimates of $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$ with each iteration. PCG requires just $\mathcal{O}(n^2)$ operations per iteration, and after k iterations it produces an approximation $\hat{\mathbf{x}}^{(k)}$ with error bounded according to [9, Ch. 11]:

$$\|\hat{\mathbf{x}}^{(k)} - \mathbf{x}\|_{\mathbf{K}} \leq 2\|\mathbf{x}\|_{\mathbf{K}} \exp\left(-\frac{2k}{\sqrt{\kappa(\mathbf{L}^{-1}\mathbf{K}\mathbf{L}^{-*})}}\right)$$

or

$$\|\hat{\mathbf{x}}^{(k)} - \mathbf{x}\|_{\mathbf{K}} \leq 2\|\mathbf{x}\|_{\mathbf{K}} \exp\left(-\frac{2k}{\sqrt{\kappa(\mathbf{C}^*\mathbf{K}\mathbf{C})}}\right)$$

The condition number $\kappa(\mathbf{M})$ of a positive definite matrix $\mathbf{M} \in \mathbb{C}^{n \times n}$ is defined as the ratio of the largest eigenvalue of \mathbf{M} to the smallest.

2.2. Column Nyström approximation. Suppose we are given a positive definite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ with rapidly decreasing eigenvalues. Then, we can consider a factored approximation

$$\hat{\mathbf{A}} = \mathbf{B}\mathbf{D}\mathbf{B}^*, \quad \text{where } \mathbf{B} \in \mathbb{C}^{n \times r} \text{ and } \mathbf{D} \in \mathbb{C}^{r \times r} \text{ with } r \ll n.$$

This is called a low-rank approximation, because $\text{rank}(\mathbf{B}\mathbf{D}\mathbf{B}^*) \leq r$. If such a low-rank approximation exists and we can find it quickly, then we can replace \mathbf{A} by the low-rank approximation $\hat{\mathbf{A}}$ to speed up linear algebra operations. For example, we can quickly multiply $\hat{\mathbf{A}}$ by a vector in a sequence of three multiplications using each matrix in the factorization. We can also efficiently compute the shifted inverse $(\hat{\mathbf{A}} + \lambda\mathbf{I})^{-1}$ for $\lambda > 0$ by using the Woodbury formula.

A popular and useful low-rank approximation is the *Nyström approximation*:

$$(2.1) \quad \mathbf{A}\langle \mathbf{X} \rangle := \mathbf{A}\mathbf{X}(\mathbf{X}^*\mathbf{A}\mathbf{X})^{-1}\mathbf{X}^*\mathbf{A}.$$

The Nyström approximation can be defined with respect to any test matrix $\mathbf{X} \in \mathbb{C}^{n \times r}$, and it has rank at most r . The approximation (2.1) can be justified as the matrix $\mathbf{H} = \mathbf{A}\mathbf{X}\mathbf{D}\mathbf{X}^*\mathbf{A}$ that is closest to \mathbf{A} in the spectral norm while maintaining a positive semidefinite residual $\mathbf{A} - \mathbf{H}$.

The Nyström approximation is defined for any $\mathbf{X} \in \mathbb{C}^{n \times r}$, so how should we choose \mathbf{X} ? In an ideal world, there is no better choice than making the columns the dominant r eigenvectors of \mathbf{A} , since computing $\mathbf{A}\langle\mathbf{X}\rangle$ via (2.1) then recovers the r -truncated eigendecomposition. In practice, running a few iterations of an iterative eigenvector-finding method can give a high-accuracy Nyström approximation when the eigenvalues decay quickly; see [21]. However, this approach is too expensive for many computational settings, since it requires $\mathcal{O}(n^2r)$ operations to execute.

A more efficient Nyström approximation is generated by the matrix $\mathbf{X} = \mathbf{I}(:, R)$, which is a block of columns from the identity indexed by a cardinality- r index set $R \subseteq \{1, \dots, n\}$, often called a *pivot set*. Then, the Nyström approximation takes a simplified form:

$$(2.2) \quad \mathbf{A}\langle R \rangle = \mathbf{A}\langle \mathbf{I}(:, R) \rangle = \mathbf{A}_{:,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R,:}.$$

We call (2.2) the *column Nyström* approximation, and it can be computed in just $\mathcal{O}(nr^2)$ operations by first forming and factoring the matrix $\mathbf{A}_{R,R}$ and then incorporating by the matrix $\mathbf{A}_{:,R}$. This approach does not require reading all entries of the matrix!

The column Nyström approximation can be represented as a block matrix approximation, after a suitable reordering of the indices. Given an index set $R = \{k_1, \dots, k_r\}$, we introduce the permutation matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ that satisfies $\mathbf{P}\mathbf{e}_i = \mathbf{e}_{k_i}$ for each $1 \leq i \leq r$. Then the column Nyström approximation of \mathbf{A} can be represented as

$$\mathbf{P}^* \mathbf{A}\langle R \rangle \mathbf{P} = \begin{bmatrix} \mathbf{A}_{R,R} & \mathbf{A}_{R,R^c} \\ \mathbf{A}_{R^c,R} & \mathbf{A}_{R^c,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R^c,R} \end{bmatrix}.$$

The residual matrix, which is called the Schur complement, can be represented as

$$\mathbf{P}^*(\mathbf{A}/R)\mathbf{P} = \mathbf{P}^*(\mathbf{A} - \mathbf{A}\langle R \rangle)\mathbf{P} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_{R^c,R^c} - \mathbf{A}_{R^c,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R^c,R} \end{bmatrix}$$

The index reordering makes clear that the selected rows and columns are represented exactly, while the remaining rows and columns are interpolated by a linear combination of the selected columns.

2.3. Vecchia approximation. An alternative way to approximate a positive definite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ or $\mathbf{A} \in \mathbb{C}^{n \times n}$ is based on approximating the distribution of a real- or complex-valued random Gaussian vector $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{A})$ [17]. If we introduce an exact inverse Cholesky factorization $\mathbf{A}^{-1} = \mathbf{C}\mathbf{C}^*$, then we observe $\mathbf{z} = \mathbf{C}^*\mathbf{y}$ is a white noise vector, $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, so the vectors \mathbf{z} and \mathbf{y} are linked by the regression formula

$$(2.3) \quad \mathbf{z}_i = \sum_{j=i}^n \overline{\mathbf{C}_{ji}} \mathbf{y}_j \quad \text{and consequently} \quad \mathbf{y}_i = - \sum_{j=i+1}^n \frac{\overline{\mathbf{C}_{ji}}}{\mathbf{C}_{ii}} \mathbf{y}_j + \frac{\mathbf{z}_i}{\mathbf{C}_{ii}}.$$

By this formula, when we linearly regress \mathbf{y}_i on $\mathbf{y}_{i+1}, \dots, \mathbf{y}_n$, the coefficients are $-\overline{\mathbf{C}_{ji}}/\mathbf{C}_{ii}$ and the residual is $\mathcal{N}(0, 1/\mathbf{C}_{ii}^2)$.

Next, we observe that many covariance matrices arising in applications have entries \mathbf{y}_i and \mathbf{y}_j that satisfy the conditional independence relationship

$$(2.4) \quad \mathbf{y}_i \perp \mathbf{y}_j \mid \mathbf{y}_{j+1} \dots \mathbf{y}_n.$$

By the regression formula (2.3), the conditional independence relationship is equivalent to a zero entry in the Cholesky factor $\mathbf{C}_{ji} = 0$. Given the tendency for many

conditional independence relationships to hold simultaneously, we consider an approximation $\hat{\mathbf{A}}^{-1} = \hat{\mathbf{C}}\hat{\mathbf{C}}^*$ where the approximate Cholesky factor is sparse

$$\hat{\mathbf{C}}_{ij} = 0, \quad \text{if } j \notin S_i,$$

and the sparsity pattern is determined by index sets $\{S_i\}_{i=1}^n$ which contain a small, bounded number of elements subject to $\{i\} \subseteq S_i \subseteq \{i, \dots, n\}$.

Taking one step further, we can explicitly approximate the Gaussian distribution $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{A})$ using the sparsity pattern $\{S_i\}_{i=1}^n$. Using the chain rule of conditional probabilities, the *exact* Gaussian density function p is

$$p(\mathbf{y}) = \prod_{i=1}^n p_{i|\{i+1, \dots, n\}}(\mathbf{y}_i | \mathbf{y}_{\{i+1, \dots, n\}}),$$

where $p_{i|\{i+1, \dots, n\}}$ is the conditional density of \mathbf{y}_i given variables $\mathbf{y}_{i+1}, \dots, \mathbf{y}_n$. If we only consider the conditional dependencies specified in the sparsity pattern $\{S_i\}_{i=1}^n$, then the Gaussian density is naturally approximated as

$$(2.5) \quad \hat{p}(\mathbf{y}) = \prod_{i=1}^n p_{i|S_i \setminus \{i\}}(\mathbf{y}_i | \mathbf{y}_{S_i \setminus \{i\}}).$$

Each conditional density $p_{i|S_i \setminus \{i\}}(\mathbf{y}_i | \mathbf{y}_{S_i \setminus \{i\}})$ is representing a Gaussian conditional distribution, which can be written as

$$\mathbf{y}_i \sim \mathcal{N}\left(-\frac{\mathbf{e}_1^*(\mathbf{A}_{S_i, S_i})^{-1}}{\mathbf{e}_1^*(\mathbf{A}_{S_i, S_i})^{-1}\mathbf{e}_1} \begin{bmatrix} 0 \\ \mathbf{y}_{S_i \setminus \{i\}} \end{bmatrix}, \frac{1}{\mathbf{e}_1^*(\mathbf{A}_{S_i, S_i})^{-1}\mathbf{e}_1}\right).$$

Taken as a whole, the density \hat{p} is representing a multivariate Gaussian distribution $\mathcal{N}(\mathbf{0}, (\hat{\mathbf{C}}\hat{\mathbf{C}}^*)^{-1})$, where the inverse Cholesky factor has a closed-form, computationally tractable expression:

$$(2.6) \quad \hat{\mathbf{C}}_{S_i, i} = \frac{(\mathbf{A}_{S_i, S_i})^{-1}\mathbf{e}_1}{\sqrt{\mathbf{e}_1^*(\mathbf{A}_{S_i, S_i})^{-1}\mathbf{e}_1}} \quad \text{and} \quad \hat{\mathbf{C}}_{S_i^c, i} = \mathbf{0}, \quad \text{for each } i = 1, \dots, n.$$

The approximation (2.6) is known as the Vecchia approximation in spatial statistics [22].

The Vecchia approximation has several optimality properties among approximations $\hat{\mathbf{A}}^{-1} = \mathbf{C}\mathbf{C}^*$ with a fixed sparsity pattern:

1. In 1990, Kaporin [12] discovered that the Vecchia approximation optimizes the ratio of the trace to the normalized determinant

$$(2.7) \quad \frac{\text{tr}(\mathbf{C}\mathbf{A}\mathbf{C}^*)}{\det(\mathbf{C}\mathbf{A}\mathbf{C}^*)^{1/n}}.$$

2. Later, Yereimin et al. [23] found that the Vecchia approximation minimizes the Frobenius norm of the residual

$$(2.8) \quad \|\mathbf{I} - \mathbf{C}\mathbf{A}\mathbf{C}^*\|_{\text{F}}$$

subject to the constraint $\text{diag}(\mathbf{C}\mathbf{A}\mathbf{C}^T) = \mathbf{1}$.

3. Most recently, Schäfer et al. [18] observed that the Vecchia approximation minimizes the KL divergence

$$(2.9) \quad D_{\text{KL}}(\mathcal{N}(0, \mathbf{A}) \| \mathcal{N}(0, (\mathbf{C}\mathbf{C}^*)^{-1})).$$

We also note two desirable computational properties of the Vecchia formula (2.6). First, this formula does not use all entries of \mathbf{A} , only blocks of entries corresponding to each column's nonzero set. Second, as a result, we can compute the nonzero entries of $\hat{\mathbf{C}}$ in only $\mathcal{O}(nr^3)$ operations if $|S_i| \leq r$ for each index set i .

Finally, we notice that the stipulation $S_i \subseteq \{i, i+1 \dots n\}$ depends heavily on the ordering of the indices. Therefore, we consider a permutation matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ that puts the important indices last, and then we apply Vecchia approximation of the permuted matrix

$$\mathbf{M} = \mathbf{P}^* \mathbf{A} \mathbf{P}.$$

We obtain an approximation $\hat{\mathbf{M}}^{-1} = \hat{\mathbf{C}}\hat{\mathbf{C}}^*$, which transfers to an approximation of the original matrix via

$$(2.10) \quad \hat{\mathbf{A}}^{-1} = \mathbf{P}\hat{\mathbf{M}}^{-1}\mathbf{P}^* = (\mathbf{P}\hat{\mathbf{C}}\mathbf{P}^*)(\mathbf{P}\hat{\mathbf{C}}\mathbf{P}^*)^*.$$

We call the approximation (2.10) a generalized Vecchia approximation.

2.4. Column Nyström plus Vecchia. Last, we consider an approach that combines low-rank and sparse structure when approximating a positive definite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$.

In this approach, we start by generating a column Nyström approximation using an index set $R = \{k_1, \dots, k_r\}$:

$$\mathbf{A}\langle R \rangle = \mathbf{A}_{:,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R,:}.$$

Following the exposition in subsection 2.2, the structure of the Nyström approximation become clearer when we introduce a permutation matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ that puts the R indices last: $\mathbf{P}\mathbf{e}_{n-r+i} = \mathbf{e}_{k_i}$ for each $i = 1, \dots, r$. The Nyström approximation can then be written as

$$\mathbf{P}^* \mathbf{A}\langle R \rangle \mathbf{P} = \begin{bmatrix} \mathbf{A}_{R^c,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R,R^c} & \mathbf{A}_{R^c,R} \\ \mathbf{A}_{R,R^c} & \mathbf{A}_{R,R} \end{bmatrix}.$$

and the Schur complement is

$$\mathbf{P}^*(\mathbf{A}/R)\mathbf{P} = \begin{bmatrix} \mathbf{A}_{R^c,R^c} - \mathbf{A}_{R^c,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R,R^c} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

The residual from this approximation is a positive-definite residual matrix $\mathbf{B} = \mathbf{A}_{R^c,R^c} - \mathbf{A}_{R^c,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R,R^c}$, which we can systematically improve.

To improve our approximation, we apply sparse Vecchia approximation to \mathbf{B} :

$$\hat{\mathbf{B}}^{-1} = \mathbf{F}\mathbf{F}^*.$$

Combining the Nyström and sparse approximations then gives

$$(2.11) \quad \mathbf{P}^* \hat{\mathbf{A}} \mathbf{P} = \begin{bmatrix} \mathbf{A}_{R^c,R}(\mathbf{A}_{R,R})^{-1}\mathbf{A}_{R,R^c} + \mathbf{F}^{-*}\mathbf{F}^{-1} & \mathbf{A}_{R^c,R} \\ \mathbf{A}_{R,R^c} & \mathbf{A}_{R,R} \end{bmatrix}.$$

This approximation was discovered by [24], who refer to their implementation as “adaptive Factorized Nyström” due to an adaptive determination of the rank of the Nyström part. However, for the rest of paper, we refer to the general class of approximations based on (2.11) as “column Nyström plus Vecchia” (CNV) approximations.

Next, we make a new observation that the *inverse* of a CNV approximation can be written explicitly as

$$(2.12) \quad \mathbf{P}^* \hat{\mathbf{A}}^{-1} \mathbf{P} = \begin{bmatrix} \mathbf{F}\mathbf{F}^* & -\mathbf{F}\mathbf{F}^* \mathbf{A}_{R^c, R} (\mathbf{A}_{R, R})^{-1} \\ -(\mathbf{A}_{R, R}^{-1}) \mathbf{A}_{R, R^c} \mathbf{F}\mathbf{F}^* & (\mathbf{A}_{R, R})^{-1} + (\mathbf{A}_{R, R})^{-1} \mathbf{A}_{R, R^c} \mathbf{F}\mathbf{F}^* \mathbf{A}_{R^c, R} (\mathbf{A}_{R, R})^{-1} \end{bmatrix}$$

This formula is a consequence of the block matrix inversion formula [3, p420], together with the fact that

$$(\hat{\mathbf{A}}_{R^c, R^c} - \hat{\mathbf{A}}_{R^c, R} (\hat{\mathbf{A}}_{R, R})^{-1} \hat{\mathbf{A}}_{R, R^c})^{-1} = \mathbf{F}\mathbf{F}^*.$$

Equation (2.12) can be simplified further, by writing the approximation in terms of its Cholesky decomposition:

$$(2.13) \quad \mathbf{P}^* \hat{\mathbf{A}}^{-1} \mathbf{P} = \begin{bmatrix} \mathbf{F} & \mathbf{0} \\ -(\mathbf{A}_{R, R})^{-1} \mathbf{A}_{R, R^c} \mathbf{F} & \mathbf{C} \end{bmatrix} \begin{bmatrix} \mathbf{F} & \mathbf{0} \\ -(\mathbf{A}_{R, R})^{-1} \mathbf{A}_{R, R^c} \mathbf{F} & \mathbf{C} \end{bmatrix}^*,$$

where $(\mathbf{A}_{R, R})^{-1} = \mathbf{F}\mathbf{F}^*$.

Equation (2.13) gives an intuitive formula for the approximate Cholesky factor that can be evaluated in $\mathcal{O}(r^3 + r\|\mathbf{C}\|_0)$ operations, where $\|\mathbf{C}\|_0$ denotes the number of nonzero entries in \mathbf{C} .

3. Column Nyström plus Vecchia approximations. This section extends the theory of Vecchia approximation to show that column-Nyström-plus-Vecchia approximations are themselves Vecchia approximations in disguise.

Our main theoretical result is the following:

THEOREM 3.1 (CNV = Vecchia). *Let $\mathbf{A} \in \mathbb{C}^{n \times n}$ be a positive definite matrix and partition the indices $\{1, \dots, n\}$ into sets $U = \{1, \dots, n-r\}$ and $V = \{n-r+1, \dots, n\}$. Consider the column Nyström plus Vecchia approximation*

$$(3.1) \quad \hat{\mathbf{A}}^{-1} = \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ -(\mathbf{A}_{V, V})^{-1} \mathbf{A}_{V, U} \mathbf{C} & \mathbf{F} \end{bmatrix} \begin{bmatrix} \mathbf{C} & \mathbf{0} \\ -(\mathbf{A}_{V, V})^{-1} \mathbf{A}_{V, U} \mathbf{C} & \mathbf{F} \end{bmatrix}^*,$$

that incorporates the following components:

1. $\mathbf{A}_{V, V}^{-1} = \mathbf{F}\mathbf{F}^*$ is a dense Cholesky factorization of $\mathbf{A}_{V, V}^{-1}$.
2. $\hat{\mathbf{B}}^{-1} = \mathbf{C}\mathbf{C}^*$ is a Vecchia approximation of $\mathbf{B} = \mathbf{A}_{U, U} - \mathbf{A}_{U, V} (\mathbf{A}_{V, V})^{-1} \mathbf{A}_{V, U}$ with sparsity pattern $\{S_i\}_{i=1}^n$.

Then $\hat{\mathbf{A}}$ is the Vecchia approximation of \mathbf{A} with sparsity pattern $T_i = S_i \cup V$ for $i = 1, \dots, n-r$ and $T_i = \{i, i+1, \dots, n\}$ for $i = n-r+1, \dots, n$.

Proof. Columns $i = n-r+1, \dots, n$ of the approximate Cholesky factor in (3.1) are generated via dense Cholesky factorization, which is the same as a Vecchia approximation with a dense sparsity pattern $T_i = \{i, i+1, \dots, n\}$.

Next consider columns $i = 1, \dots, n-r$ of the approximate Cholesky factor in (3.1). Using the definition (2.6) of the sparse Vecchia approximation $\hat{\mathbf{B}}^{-1} = \mathbf{C}\mathbf{C}^*$,

each column i satisfies

$$\begin{bmatrix} \mathbf{I} \\ -(\mathbf{A}_{V,V})^{-1}\mathbf{A}_{V,U} \end{bmatrix} \mathbf{C}_{:,i}, \quad \text{where } \mathbf{C} = \frac{(\mathbf{B}_{S_i,S_i})^{-1}\mathbf{e}_1}{\mathbf{e}_1^*(\mathbf{B}_{S_i,S_i})^{-1}\mathbf{e}_1}$$

Next, using the block matrix inversion formula [3, p420], we observe

$$\begin{bmatrix} \mathbf{I} \\ -(\mathbf{A}_{V,V})^{-1}\mathbf{A}_{V,U} \end{bmatrix} \frac{(\mathbf{B}_{S_i,S_i})^{-1}\mathbf{e}_1}{\mathbf{e}_1^*(\mathbf{B}_{S_i,S_i})^{-1}\mathbf{e}_1} = \frac{(\mathbf{A}_{S_i \cup V, S_i \cup V})^{-1}\mathbf{e}_1}{\sqrt{\mathbf{e}_1(\mathbf{A}_{S_i \cup V, S_i \cup V})^{-1}\mathbf{e}_1}}.$$

This is precisely the formula (2.6) for the i th column in a Vecchia approximation of \mathbf{A} with sparsity pattern $T_i = S_i \cup V$. \square

Theorem 3.1 has several ramifications. Most immediately, the Kaporin, KL, and Frobenius optimality results from [12, 18, 23] carry over to CNV approximations. Also, the theorem implies that for fixed approximation rank r , the choice of column Nyström approximation is equivalent to the choice of the last r indices in the ordering. Two existing methods with the same index ordering are “adaptive factorized Nyström” approximation [24] and “sparse Cholesky approximation” [18]: these techniques produce the same approximation when the bottom r rows are included to the sparsity pattern.

4. Application: Gaussian Processes. The results in section 3 apply to all positive-definite matrices, but one application which merits further elaboration is covariance matrices in Gaussian processes (GPs). A Gaussian process is a multivariate Gaussian distribution whose dimensions are \mathbb{R}^d , i.e., with mean $\mathbf{m} \in \mathbb{R}^{\mathbb{R}^d}$ and covariance matrix $\bar{\mathbf{K}} \in \mathbb{R}^{\mathbb{R}^d \times \mathbb{R}^d}$. Typically, $\bar{\mathbf{K}}$ is specified by setting $\bar{\mathbf{K}}_{\mathbf{x},\mathbf{y}} = \kappa(\mathbf{x},\mathbf{y})$ or $\mathbf{K}_{\mathbf{x},\mathbf{y}} = \kappa(\mathbf{x},\mathbf{y}) + \mu\delta_{\mathbf{x},\mathbf{y}}$, for κ a *kernel function*, or a function with $\kappa(\mathbf{x},\mathbf{x}) = 1$ for which $\kappa(\mathbf{x},\mathbf{y})$ decreases in $\|\mathbf{x} - \mathbf{y}\|$ for some norm, and $\mu \in \mathbb{R}_{\geq 0}$ a *nugget parameter*. A random variable conforming to such a distribution can be thought of as a random function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ such that the covariance of $f(\mathbf{x})$ and $f(\mathbf{y})$ is given by $\bar{\mathbf{K}}_{\mathbf{x},\mathbf{y}}$. Since estimating and subtracting off the mean is usually not a problem in applications, we will take $\mathbf{m} = \mathbf{0}$ in this paper.

When working with GPs, we often work in terms of a realization of a GP at some finite number n points $\{\mathbf{x}^{(i)}\}_{i=1}^n \subseteq \mathbb{R}^d$. When we’re only concerned with these points, $\bar{\mathbf{K}}$ behaves like a matrix in $\mathbb{R}^{n \times n}$; we thus often define some matrix \mathbf{K} such that $\mathbf{K}_{i,j}$ refers to $\bar{\mathbf{K}}_{\mathbf{x}^{(i)},\mathbf{x}^{(j)}}$. If n is large, practitioners must often work with approximations of \mathbf{K} over these n points with advantageous computational properties for improved computational speed at the cost of some approximation error. In this section, we apply CNV to this task. Specifically, we discuss the compatibility of existing algorithms for choosing pivot sets and sparsity patterns with CNV in the context of GPs. We also propose our own pivot-choosing algorithm in subsection 4.1.3.

4.1. Matérn Kernels. We begin by examining existing methods for finding sparse Vecchia approximations of covariance matrices corresponding to a particular class of kernels, the Matérn kernels. Later in this section, we will show that these methods are compatible with CNV approximations and adapt their theoretical bounds to the CNV case, as well as propose heuristics based on these for choosing pivots during the Nyström part of CNV-type approximations.

The Matérn kernel functions take the form

$$\kappa_{\nu,\rho,\sigma}(\mathbf{x}, \mathbf{x}') = \sigma^2 \frac{2^{1-\nu}}{\Gamma(\nu)} \left(\frac{\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\rho} \right)^\nu B_\nu \left(\frac{\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\rho} \right)$$

where B_ν is the modified Bessel function of the second kind, Γ is the gamma function, ν is called a smoothness parameter, and $\rho, \sigma > 0$ are other parameters. For $\nu \in \{\frac{1}{2}, \frac{3}{2}, \frac{5}{2}, \dots\} = \mathbb{N} + \frac{1}{2}$, this can be written

$$\begin{aligned} & \kappa_{\nu,\rho,\sigma}(\mathbf{x}, \mathbf{x}') = \\ & \sigma^2 \exp \left(-\frac{\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\rho} \right) \frac{(\nu - \frac{1}{2})!}{(2\nu - 1)!} \sum_{i=0}^{\nu - \frac{1}{2}} \frac{(\nu - \frac{1}{2} + i)!}{i!(\nu - \frac{1}{2} - i)!} \left(\frac{2\sqrt{2\nu} \|\mathbf{x} - \mathbf{x}'\|_2}{\rho} \right)^{\nu - \frac{1}{2} - i}. \end{aligned}$$

4.1.1. Reverse Maximin Ordering and Theoretical Sparse Vecchia Error Bounds. Let κ be a Matérn kernel function on \mathbb{R}^d , let $\bar{\mathbf{K}} \in \mathbb{R}^{\mathbb{R}^d \times \mathbb{R}^d}$ be the corresponding covariance matrix, and let \mathbf{K} be the covariance matrix it induces given $X = \{\mathbf{x}^{(i)}\}_{i=1}^n$ according to $\mathbf{K}_{i,j} = \kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$. Additionally, let $Z \subseteq \mathbb{R}^d$ (possibly empty), which in this context can be thought of as an index set. Define

$$\mathbf{K}/Z := (\bar{\mathbf{K}}/Z)_{XX} = \bar{\mathbf{K}}_{XX} - \bar{\mathbf{K}}_{XZ} \bar{\mathbf{K}}_{ZZ}^{-1} \bar{\mathbf{K}}_{ZX}.$$

Later in this subsection, we will state a result of [18] regarding the sparsity of the Cholesky factor of $(\mathbf{K}/Z)^{-1}$ that gives us reason to believe the Cholesky factor of \mathbf{K}^{-1} is sparse, too. But before we do, we recall from subsection 2.3 that if we are to find a sparse Vecchia approximation of \mathbf{K} , the ordering of the points $\{\mathbf{x}^{(i)}\}_{i=1}^n$ matters significantly to the quality of our approximation. When $d = 1$, ordering the points according to the usual ordering on \mathbb{R} , or its reverse, tend to work well in practice. However, methods of this type (e.g. lexicographically coordinate-by-coordinate in some basis) tend to perform poorly in \mathbb{R}^d ([10]).

An alternative which [10] shows performs well empirically, and which [18] require a slightly weaker version of in their aforementioned result, is known as the reverse-maximin ordering. It follows a procedure known as Farthest Point Sampling (FPS) to choose the next-to-last point in the ordering at any given step. Mathematically, assuming $\{k_j\}_{j=i+1}^n$ in the reverse maximin ordering have been determined, and for a given Z as above, one finds

$$(4.1) \quad k_i = \underset{k \notin \{k_j\}_{j=i+1}^n}{\operatorname{argmax}} \min_{\mathbf{x} \in \{\mathbf{x}^{(i+1)} \dots \mathbf{x}^{(n)}\} \cup Z} \|\mathbf{x}^{(k)} - \mathbf{x}^{(k_j)}\|_2$$

for the next point in the ordering. Procedurally, a naïve algorithm for construction such an ordering is given in Algorithm 4.1.

This algorithm takes $O(n(n + |Z|)d)$ operations; whether this would be a computational bottleneck generally depends on d and one's application. If a faster algorithm is needed and d is not too large, [19]'s Algorithm 4.1 yields a reverse-maximin ordering in $O(2^d dn \log^2 n + n|Z|d)$ operations.

We can now state a slightly weaker version of the result of [18]:

THEOREM 4.1 ([18], Theorem 3.4). *Let $\Omega \subseteq \mathbb{R}^d$ be a bounded domain with boundary $\partial\Omega$, and let $\{\mathbf{x}^{(i)}\}_{i=1}^n \subseteq \Omega$ be ordered according to the reverse-maximin ordering. Also, let*

$$(4.2) \quad \delta = \frac{\min_{i \in [n]} \min_{\mathbf{x} \in \{\mathbf{x}^{(j)}\}_{j=1}^n \cup \partial\Omega \setminus \mathbf{x}^{(i)}} \|\mathbf{x}^{(i)} - \mathbf{x}\|_2}{\max_{\mathbf{x} \in \Omega} \min_{\mathbf{x}' \in \{\mathbf{x}^{(i)}\}_{i=1}^n \cup \partial\Omega} \|\mathbf{x} - \mathbf{x}'\|_2},$$

Algorithm 4.1 Construction of Reverse-Maximin Ordering

```

for  $z \in Z$  do
  for  $i \in [n]$  do
     $\mathbf{d}_i \leftarrow \min(\mathbf{d}_i, \|\mathbf{x}^{(i)} - z\|_2)$ 
  end for
end for
for  $i = n \dots 1$  do
   $k_i \leftarrow \arg \max_k \mathbf{d}_k$ 
  for  $j = 1 \dots n$  do
     $\mathbf{d}_j \leftarrow \min(\mathbf{d}_j, \|\mathbf{x}^{(j)} - \mathbf{x}^{(k_i)}\|_2)$ 
  end for
end for

```

and suppose \mathbf{K} is the Matérn covariance matrix at $\{\mathbf{x}^{(i)}\}_{i=1}^n$ for $\kappa_{\nu, \rho, \sigma}$ with $\nu \in \mathbb{N} + \frac{1}{2}$. Further suppose S_i is given by the indices of the $O((\log(\frac{n}{\varepsilon}))^d)$ -nearest neighbors of $\mathbf{x}^{(i)}$ which satisfy lower triangularity along with $\mathbf{x}^{(i)}$ itself, with implicit constant depending only on $\delta, \nu, \rho, \sigma, \Omega$. Then if \mathbf{C} is the sparse Vecchia approximation of $\mathbf{K}/\partial\Omega$ with sparsity pattern $\{S_i\}_{i=1}^n$ as given by (2.6), we have

$$(4.3) \quad D_{KL}(\mathcal{N}(\mathbf{0}, \mathbf{K}/\partial\Omega) \parallel \mathcal{N}(\mathbf{0}, (\mathbf{C}\mathbf{C}^*)^{-1})) + \|(\mathbf{K}/\partial\Omega) - (\mathbf{C}\mathbf{C}^*)^{-1}\|_{FRO} < \varepsilon$$

We remark that while Theorem 4.1 requires conditioning on the boundary of a bounded domain containing $\{\mathbf{x}^{(i)}\}_{i=1}^n$, [18] observed that this did not substantially affect empirical results in settings that did not condition on the boundary. We also remark that this result does not hold with nonzero nugget, and indeed various authors ([13, 14]) observe that the sparsity of \mathbf{C} decreases as μ the nugget parameter increases.

4.1.2. Adapting Nearest Neighbors to CNV case. Now, we adapt Theorem 4.1 from the sparse Vecchia case to the CNV case.

Before proceeding more generally, we first consider the case of CNV approximations whose Nyström pivot set is simply the last r indices in the reverse-maximin ordering. Let $\{\mathbf{x}^{(i)}\}_{i=1}^n$, \mathbf{K} , Ω , and the sparsity pattern $\{S_i\}_{i=1}^n$ as in Theorem 4.1, the latter for error tolerance $\varepsilon > 0$. Here, Theorem 3.1 implies that the sparse Vecchia approximation with sparsity pattern $S_i \cup \{\max(i, n - r + 1) \dots n\}$ equals the CNV approximation with r pivots and sparsity pattern $\{S_i\}_{i=1}^n$ because the index ordering is the same. By the optimality of the Vecchia approximation in the KL divergence loss (2.9) and the Frobenius loss (2.8)¹, we cannot make our approximation worse in these objectives by adding entries to the sparsity pattern. So by Theorem 4.1, which we crucially can only apply because the ordering is still reverse-maximin, this CNV approximation has error no worse than (4.3).

Though [24] did not know of Theorem 3.1 nor this result, this provides considerable justification for their proposal to use farthest point sampling (FPS) to select Nyström pivots. However, the literature on choosing Nyström pivots is very extensive, and we might hope to use other algorithms than FPS for this purpose. Motivated by this, we extend the bounds of Theorem 4.1 to other choices for pivot sets:

THEOREM 4.2. *Let $\{\mathbf{x}^{(i)}\}$, Ω , \mathbf{K} , ε as in the setting of Theorem 4.1, let $R \subseteq [n]$ of size r . Let $\{k_i\}_{i=1}^n$ an index reordering such that $\{k_i\}_{i=n-r+1}^n = R$ and $\{k_i\}_{i=1}^{n-r}$*

¹Actually, the Frobenius losses in (2.8) and (4.3) are not exactly of the same form, though the bound in Theorem 4.1 remains valid due to [19]’s Lemma B.8.

contains the rest of the indices according to the positions of their respective $\mathbf{x}^{(k_i)}$ s in the reverse maximin ordering initialized with $\partial\Omega \cup \{\mathbf{x}^{(R)}\}$. Also, let \mathbf{P} be the $n \times n$ permutation matrix such that $\mathbf{P}\mathbf{e}_1 = \mathbf{e}_{k_1}$. Finally, suppose that for $i \notin R$, we have that S_i is given by the indices of the $O((\log(\frac{n}{\varepsilon}))^d)$ -nearest neighbors of $\mathbf{x}^{(i)}$ among $\{\mathbf{x}^{R^C \cap \{k_j\}_{j=1}^n}\}$, along with $\mathbf{x}^{(i)}$ itself, with implicit constant depending only on $\delta, \nu, \rho, \sigma, \Omega$, as well as $\{\mathbf{x}^{(R)}\}$ the choice of pivots. Then for $(\mathbf{C}\mathbf{C}^*)$ the CNV approximation of $\mathbf{P}^*(\mathbf{K}/\partial\Omega)\mathbf{P}$, we have

$$(4.4) \quad D_{\text{KL}}(\mathcal{N}(\mathbf{0}, \mathbf{P}^*(\mathbf{K}/\partial\Omega)\mathbf{P}) \parallel \mathcal{N}(\mathbf{0}, (\mathbf{C}\mathbf{C}^*)^{-1})) + \|\mathbf{P}^*(\mathbf{K}/\partial\Omega)\mathbf{P} - (\mathbf{C}\mathbf{C}^*)^{-1}\|_{\text{FRO}} < \varepsilon$$

Proof. We'll start by showing that $((\mathbf{K}/\partial\Omega)/R)^{-1}$, the inverse column Nyström residual of $\mathbf{K}/\partial\Omega$ for some index set $R \subseteq [n]$, is also well-approximated by a matrix with $\mathcal{O}((\log(\frac{n}{\varepsilon}))^d)$ entries per column in its Cholesky factor.

We first observe that

$$\begin{aligned} (\mathbf{K}/\partial\Omega)/R &= (\bar{\mathbf{K}}_{X \cup \partial\Omega, X \cup \partial\Omega} - \bar{\mathbf{K}}_{X \cup \partial\Omega, \partial\Omega} \bar{\mathbf{K}}_{\partial\Omega, \partial\Omega}^{-1} \bar{\mathbf{K}}_{\partial\Omega, X \cup \partial\Omega})_{X, X} / R \\ &= ((\bar{\mathbf{K}}_{X \cup \partial\Omega, X \cup \partial\Omega} - \bar{\mathbf{K}}_{X \cup \partial\Omega, \partial\Omega} \bar{\mathbf{K}}_{\partial\Omega, \partial\Omega}^{-1} \bar{\mathbf{K}}_{\partial\Omega, X \cup \partial\Omega}) / \{\mathbf{x}^{(R)}\})_{X, X} = (\mathbf{K}/\partial\Omega \cup \{\mathbf{x}^{(R)}\}). \end{aligned}$$

Motivated by this, define $\tilde{\Omega} = \Omega \setminus \{\mathbf{x}^{(R)}\}$ for any index set $R \subseteq [n]$, which is a bounded domain containing $\{\mathbf{x}^{([n] \setminus R)}\}$. Then $\partial\tilde{\Omega} = \partial\Omega \cup \{\mathbf{x}^{(R)}\}$, and we may apply Theorem 4.1 to $\{\mathbf{x}^{[n] \setminus R}\}$ on $\tilde{\Omega}$.

As such, let $\{k_i\}_{i=1}^n$ and \mathbf{P} be the index reordering and permutation matrix in the proposition, and let $\tilde{\mathbf{P}} = \mathbf{P}_{1:n-r, k_{1:n-r}}$. Then by Theorem 4.1, the Frobenius and KL divergence errors of the Vecchia approximation of $\tilde{\mathbf{P}}^*(\mathbf{K}/\partial\tilde{\Omega})\tilde{\mathbf{P}}$ decay exponentially in the number of nonzeros per column.

Finally, we demonstrate that this propagates to the entire CNV approximation. Letting $\mathbf{G}\mathbf{G}^* \approx \tilde{\mathbf{P}}^*(\mathbf{K}/\partial\tilde{\Omega})\tilde{\mathbf{P}} = \tilde{\mathbf{P}}^*(\mathbf{K}/\partial\Omega \cup \{\mathbf{x}^{(R)}\})\tilde{\mathbf{P}}$ the sparse Vecchia part from the last paragraph, we note that by (2.11) the residual of the entire CNV approximation of $\mathbf{P}^*(\mathbf{K}/\partial\Omega)\mathbf{P}$ with this sparse component takes form

$$\begin{aligned} &\begin{bmatrix} (\mathbf{K}/\partial\Omega)_{R^C, R^C} & (\mathbf{K}/\partial\Omega)_{R^C, R} \\ (\mathbf{K}/\partial\Omega)_{R, R^C} & (\mathbf{K}/\partial\Omega)_{R, R} \end{bmatrix} \\ &\quad - \begin{bmatrix} (\mathbf{K}/\partial\Omega)_{R^C, R} (\mathbf{K}/\partial\Omega)_{R, R}^{-1} (\mathbf{K}/\partial\Omega)_{R, R^C} + \mathbf{G}^* \mathbf{G}^{-1} & (\mathbf{K}/\partial\Omega)_{R^C, R} \\ (\mathbf{K}/\partial\Omega)_{R, R^C} & (\mathbf{K}/\partial\Omega)_{R, R} \end{bmatrix} \\ &= \begin{bmatrix} \tilde{\mathbf{P}}^*(\mathbf{K}/\partial\Omega \cup \{\mathbf{x}^{(R)}\})\tilde{\mathbf{P}} + \mathbf{G}^* \mathbf{G}^{-1} & 0 \\ 0 & 0 \end{bmatrix}. \end{aligned}$$

We showed in the last paragraph that the Frobenius error of the upper right block decays exponentially, so it follows that the Frobenius error of the entire RHS does, too. For the KL divergence, this can also be shown directly, but we refer the reader to Lemma B.8 in [18] for brevity. \square

4.1.3. Alternative pivot chooser: Approximate greedy KL minimization. Motivated by Theorem 4.2, we now propose an alternative to [24]'s FPS for choosing the indices in a column Nyström approximation. As theoretical motivation, Theorem 3.1 shows how the column Nyström part of a CNV approximation is equivalent to a Vecchia approximation. Moreover, a Vecchia approximation optimizes the KL divergence

$$D_{\text{KL}}(\mathcal{N}(0, \mathbf{K}) \parallel \mathcal{N}(0, (\mathbf{C}\mathbf{C}^*)^{-1})).$$

given a specific sparsity pattern. Therefore, it is natural to search for the column Nyström approximation that leads to the smallest KL divergence.

The KL divergence can be expanded according to

$$(4.5) \quad 2D_{\text{KL}}(\mathcal{N}(0, \mathbf{K}), \mathcal{N}(0, (\mathbf{C}\mathbf{C}^*)^{-1}))$$

$$(4.6) \quad = \text{Tr}(\mathbf{C}^* \mathbf{K} \mathbf{C}) - n + \log \det(\mathbf{C}\mathbf{C}^*)^{-1} - \log \det \mathbf{K}.$$

Given a sparsity pattern S_i , the corresponding Vecchia approximation (2.6) that optimizes the KL divergence is given by

$$\mathbf{C}_{S_i, i} = \frac{(\mathbf{A}_{S_i, S_i})^{-1} \mathbf{e}_1}{\sqrt{\mathbf{e}_1^* (\mathbf{A}_{S_i, S_i})^{-1} \mathbf{e}_1}}.$$

This form of Cholesky factor \mathbf{C} ensures that $\mathbf{C}^* \mathbf{K} \mathbf{C}$ has a diagonal of all ones. Hence, (4.6) reduces to

$$(4.7) \quad = \log \det(\mathbf{C}\mathbf{C}^*)^{-1} - \log \det \mathbf{K},$$

which can be shown (see [11], appendix A.1) to equal

$$(4.8) \quad \sum_{i=1}^n \left[\log((\mathbf{K}/S_i \setminus \{i\})_{ii}) - \log((\mathbf{K}/\{i+1:n\})_{ii}) \right].$$

Next, we try to minimize the KL divergence with a “low rank plus diagonal” sparsity pattern, defined by the diagonal entries together with the bottom rows that indicate the low-rank component. Let k_1, \dots, k_n be an arbitrary index reordering, let \mathbf{P} be the permutation matrix with $\mathbf{P}\mathbf{e}_i = \mathbf{e}_{k_i}$, and let $\mathbf{M} = \mathbf{P}^* \mathbf{K} \mathbf{P}$. In particular, the k_i th row of \mathbf{K} exactly equals the i th row of \mathbf{M} , so $(\mathbf{M}/J)_{i,i} = (\mathbf{K}/\{k_J\})_{k_i, k_i}$ for any $J \subseteq [n]$.

Suppose we are approximating \mathbf{M} with a sparse Vecchia approximation, and that that the sparsity pattern currently consists of rows $n-j+1 \dots n$ together with the diagonal. If we added the $n-j$ th row to this sparsity pattern, the change in KL divergence would be given by

$$\begin{aligned} & \sum_{i=1}^{n-j} \log((\mathbf{M}/\{n \dots n-j\} \setminus \{i\})_{i,i}) - \sum_{i=1}^{n-j} \log((\mathbf{M}/\{n \dots n-j+1\})_{i,i}) \\ &= \left(\log \det((\mathbf{M}/\{n \dots n-j\})_{1:n-j-1, 1:n-j-1}) \right. \\ & \quad \left. + \log((\mathbf{M}/\{n \dots n-j+1\})_{n-j, n-j}) \right) \\ & \quad - \left(\log \det((\mathbf{M}/\{n \dots n-j+1\})_{1:n-j, 1:n-j}) \right) \\ &= \log \det((\mathbf{M}/\{n \dots n-j\})_{1:n-j-1, 1:n-j-1}) \\ & \quad - \log \det((\mathbf{M}/\{n \dots n-j+1\})_{1:n-j-1, 1:n-j-1}) \end{aligned}$$

where the higher terms of the sums on the first line cancel. For the second and third lines, again see [11]’s appendix A.1. The above expression can be simplified by an application of the matrix determinant lemma to the first term according to [11], appendix B.5:

$$(4.9) \quad \begin{aligned} & \log((\mathbf{M}/[n] \setminus n-j)_{n-j, n-j}) - \log((\mathbf{M}/\{n-j+1:n\})_{n-j, n-j}) \\ &= \log((\mathbf{K}/[n] \setminus k_{n-j})_{k_{n-j}, k_{n-j}}) - \log((\mathbf{K}/\{k_{n-j+1:n}\})_{k_{n-j}, k_{n-j}}) \end{aligned}$$

This simplified expression leads to the consequence that the greedy optimization strategy can be highly successful.

THEOREM 4.3 (Greedy optimization strategy for Nyström approximation). *Fix a positive definite matrix $\mathbf{K} \in \mathbb{C}^{n \times n}$. For each index set I with arbitrary ordering \prec , let $(\mathbf{F}^{(I)})^{-*}(\mathbf{F}^{(I)})^{-1}$ denote the generalized Vecchia approximation of \mathbf{K} with sparsity pattern $S_j = \{i \in I \mid i \prec j\} \cup \{j\}$ for each $j = 1, \dots, n$. Also, introduce the objective function*

$$f(I) := \frac{1}{2} D_{KL}[\mathcal{N}(0, \mathbf{K}) \parallel \mathcal{N}(0, (\mathbf{F}^{(I)})^{-*}(\mathbf{F}^{(I)})^{-1})].$$

Then the greedy approximation strategy

$$I_0 \leftarrow \emptyset, \quad I_t \leftarrow I_{t-1} \cup \{i_t\} \quad \text{for } i_t \in \underset{1 \leq i \leq n}{\operatorname{argmin}} f(I_{t-1} \cup \{i\}),$$

results in exponential convergence to the best cardinality- m objective function value:

$$f(I_t) \leq \min_{|J|=k} f(J) + e^{-t/m} \left[f(I_0) - \min_{|J|=m} f(J) \right], \quad \text{for each } t = 0, 1, \dots$$

Proof. The proof technique is based on the classic idea of supermodular maximization, first developed by Nemhauser [16]. We will show f is supermodular, i.e.,

$$(4.10) \quad f(I \cup \{\ell\}) - f(I) \leq f(I \cup J \cup \{\ell\}) - f(I \cup J),$$

for any disjoint index sets $I, J \subseteq \{1, \dots, n\}$ and any $\ell \in \{1, \dots, n\}$. Using (4.9), we have

$$\begin{aligned} f(I \cup \{\ell\}) - f(I) &= \log((\mathbf{K}/\{-\ell\})_{\ell\ell}) - \log((\mathbf{K}/I)_{\ell\ell}), \\ f(I \cup J \cup \{\ell\}) - f(I \cup J) &= \log((\mathbf{K}/\{-\ell\})_{\ell\ell}) - \log((\mathbf{K}/I \cup J)_{\ell\ell}). \end{aligned}$$

Subtracting the second line from the first line, we see

$$\begin{aligned} & [f(I \cup \{\ell\}) - f(I)] - [f(I \cup J \cup \{\ell\}) - f(I \cup J)] \\ &= \log((\mathbf{K}/I \cup J)_{\ell\ell}) - \log((\mathbf{K}/I)_{\ell\ell}). \end{aligned}$$

By taking the Schur complement with respect to the index set I first and with respect to the index set J second,

$$(\mathbf{K}/I \cup J)_{\ell\ell} = (\mathbf{K}/I)_{\ell\ell} - (\mathbf{K}/I)_{\ell J} [(\mathbf{K}/I)_{JJ}]^{-1} (\mathbf{K}/I)_{J\ell} \leq (\mathbf{K}/I)_{\ell\ell},$$

where the inequality comes from the fact that $[(\mathbf{K}/I)_{JJ}]^{-1}$ is positive definite and $(\mathbf{K}/I)_{J\ell}$ is a vector. This establishes the supermodularity property (4.10).

Let $J^* = \{j_1, \dots, j_m\}$ be an optimal index set that solves $\min_{|J|=m} f(J)$. Then for each time index $t = 0, 1, \dots$, calculate

$$(4.11) \quad f(J^*) \geq f(J^* \cup I_t)$$

$$(4.12) \quad = f(I_t) + \sum_{\ell=1}^m f(I_j \cup \{j_1, \dots, j_\ell\}) - f(I_t \cup \{j_1, \dots, j_{\ell-1}\})$$

$$(4.13) \quad \geq f(I_t) + \sum_{\ell=1}^m f(I_t \cup \{j_\ell\}) - f(I_t)$$

$$(4.14) \quad \geq f(I_t) + m(f(I_{t+1}) - f(I_t)),$$

In this expression, (4.11) comes from supermodularity, (4.12) is a telescoping sum, (4.13) comes from another application of supermodularity, and (4.14) comes from the definition of the greedy optimization approach. By rearrangement,

$$f(I_{t+1}) - f(J^*) \leq \left(1 - \frac{1}{m}\right) [f(I_t) - f(J^*)], \quad \text{for each } t = 0, 1, \dots$$

By induction, it follows

$$f(I_t) - f(J^*) \leq \left(1 - \frac{1}{m}\right)^t [f(I_0) - f(J^*)] \leq e^{-t/m} [f(I_0) - f(J^*)]$$

which can be rearranged into the desired result. \square

To our knowledge, [Theorem 4.3](#) is new. Krause et al. [15] studied a mutual information objective function, which leads to the same greedy optimization strategy based on (4.7). Yet the Kullback-Leibler divergence and mutual information are different functions. Our result establishes rigorous guarantees for optimization of the the Kullback-Leibler divergence, which is directly relevant to obtaining a good matrix approximation.

Unfortunately, there are barriers to a practical implementation of the greedy optimization strategy based on minimizing (4.9)

$$\log((\mathbf{K}/[n] \setminus i)_{i,i}) - \log((\mathbf{K}/\{k_{n-j+1:n}\})_{i,i}) = -\log(\mathbf{K}_{ii}^{-1}) - \log((\mathbf{K}/\{k_{n-j+1:n}\})_{i,i})$$

We can keep track of the second term, which is the logarithm of the diagonal entry in the residual of our partial Nyström approximation. However, keeping track of the first term is not so easy. Our idea is to approximate $\text{diag}(\mathbf{K}^{-1})$ or each $(\mathbf{K}/[n] \setminus i)_{i,i}$ and use the approximation in the optimization. We call this approach Approximate Greedy KL Minimization (GKL). For a particular approximation \mathbf{h} with $\mathbf{h}_i \approx \mathbf{K}_{ii}^{-1} = \frac{1}{(\mathbf{K}/[n] \setminus i)_{i,i}}$, [Algorithm 4.2](#) codifies such a routine in pseudocode.

Algorithm 4.2 Column Nyström via Greedy KL Minimization

```

F ← 0 ∈ ℝn×r
d ← diag(K) (= 1)                                ▷ Diagonal of Nyström residual
for  $i \in [r]$  do
   $k_{n-i+1} \leftarrow \arg \max_l \mathbf{d}_l \mathbf{h}_l$                 ▷ =  $\arg \min_l -\log(\mathbf{h}_l) - \log(\mathbf{d}_l)$ 
   $\mathbf{F}_{:,i} \leftarrow \mathbf{K}_{:,i}$ 
   $\mathbf{F}_{:,i} \leftarrow \mathbf{F}_{:,i} - \mathbf{F}_{:,1:i-1} (\mathbf{F}_{k_{n-i+1},1:i-1})^*$ 
   $\mathbf{F}_{:,i} \leftarrow \frac{\mathbf{F}_{:,i}}{\mathbf{F}_{k_{n-i+1},i}}$ 
  for  $j \in [n]$  do
     $\mathbf{d}_j \leftarrow \mathbf{d}_j - |\mathbf{F}_{j,i}|^2$ 
  end for
end for
return  $\{k_{n-r+1} \dots k_n\}, \mathbf{F}$   ▷ Pivot set, Nyström approximation's Cholesky factor

```

To approximate \mathbf{K}_{ii}^{-1} for each i , we observe that if $\mathbf{C}\mathbf{C}^* \approx \mathbf{K}$, then $\mathbf{K}_{i,i}^{-1} \approx \mathbf{C}_{i,:} \mathbf{C}_{:,i}^* = \|\mathbf{C}_{i,:}\|_2^2$. For Matérn covariance matrices specifically, we already know from the last two sections that the pure-sparse Vecchia approximation based on order $\log n$ nearest neighbors per column is of high quality. Thus, we can produce such an

Algorithm 4.3 Inverse Diagonal Approximation by Vecchia Row Norm

```

{ $S_i\}_{i=1}^n \leftarrow \text{kNNSparsityPattern}(\mathbf{K}, s)$   ▷ Or any other suitable sparsity pattern
 $\mathbf{C} \leftarrow \text{SparseVecchia}(\mathbf{K}, \{S_i\}_{i=1}^n)$ 
for  $i \in [n]$  do
     $\mathbf{h}_i \leftarrow \|\mathbf{C}_{i,:}\|_2^2$ 
end for
return  $\mathbf{h}$ 

```

approximation in order $n \log^{3d} n$ operations and use the square 2-norm of its i th row as an approximation of $K_{i,i}^{-1}$. See [Algorithm 4.3](#) for pseudocode.

Alternatively, we could approximate each $(\mathbf{K}/[n] \setminus \{i\})_{ii}$ by $(\mathbf{K}/S_i \setminus \{i\})_{ii}$ for any suitable set S_i (not necessarily respecting lower triangularity). This also has an interpretation in terms of [\(2.6\)](#) and Vecchia approximations. By lower triangularity, $\mathbf{K}_{1,1}^{-1} \approx \mathbf{C}_{1,:} \mathbf{C}_{:,1}^* = \mathbf{C}_{1,1}^2$. If we had reordered the indices such that point $\mathbf{x}^{(i)}$ were first instead of $\mathbf{x}^{(1)}$ and adjusted the sparsity pattern to reflect that now *all* indices of nearest neighbors of $\mathbf{x}^{(i)}$ satisfy lower triangularity, we could approximate \mathbf{K} as the square of just the entry corresponding to index i of the vector [\(2.6\)](#) if S_i is allowed to ignore triangularity. We provide [Algorithm 4.4](#) to illustrate this connection, although in practice it is equivalent to inverting \mathbf{K}_{S_i, S_i} and taking the first diagonal entry (which should be used for practical implementations).

Algorithm 4.4 Inverse Diagonal Approximation by Columnwise Vecchia Diagonal

```

for  $i \in [n]$  do
     $S_i \leftarrow \text{kNN}(\mathbf{x}^{(i)}, \{\mathbf{x}^{(1)} \dots \mathbf{x}^{(n)}\}, s)$   ▷ Or any suitable set, ignoring triangularity
     $\mathbf{v} \leftarrow \text{SparseVecchiaColumn}(\mathbf{K}_{:,i}, S_i)$ 
     $\mathbf{h}_i \leftarrow \mathbf{v}_i^2$ 
end for
return  $\mathbf{h}$ 

```

4.1.4. Alternative Pivot Choosers: Randomly Pivoted Cholesky. Instead of approximating first term of [\(4.9\)](#), we could also ignore it entirely. Surprisingly, this recovers a diagonal residual heuristic common to an already well-known class of column Nyström algorithms for general positive-definite matrices ([\[7, 4, 20\]](#)). Since the relative influence of the first term in [\(4.9\)](#) on the objective function diminishes as r increases, such pre-existing column Nyström algorithms might be appropriate when r is taken to be large. These methods also benefit from already-known error bounds when applied to generic positive definite matrices.

This "diagonal residual heuristic" class of algorithms computes Nyström approximations via the procedure in [Algorithm 4.5](#), where `ChoosePivot` is a subroutine which chooses the next pivot based on the information in the diagonal of the residual.

Several choices of `ChoosePivot` are used in practice. One such choice simply ignores \mathbf{d} and returns a uniformly random index among those not previously chosen. This approach is well-motivated: especially in the case of kernel matrices \mathbf{K} where columns do not differ much in scale, invariant subspaces with respect to \mathbf{K} must align closely with many columns of \mathbf{K} in order to be large. By sampling the columns uniformly, we thus have a good chance of choosing such columns ([\[8\]](#)). We refer to this strategy as uniform pivoting.

Another, which we'll refer to as the greedy pivoting ([\[2, 7\]](#)), simply chooses

Algorithm 4.5 Column Nyström via Diagonal Residual Heuristic

```

F  $\leftarrow$  0  $\in \mathbb{R}^{n \times r}$ 
d  $\leftarrow$  diag(K) (= 1) ▷ Diagonal of Nyström residual
for  $i \in [r]$  do
   $k_{n-i+1} = \text{ChoosePivot}(\mathbf{d})$ 
   $\mathbf{F}_{:,i} \leftarrow \mathbf{K}_{:,i}$ 
   $\mathbf{F}_{:,i} \leftarrow \mathbf{F}_{:,i} - \mathbf{F}_{:,1:i-1}(\mathbf{F}_{k_{n-i+1},1:i-1})^*$ 
   $\mathbf{F}_{:,i} \leftarrow \frac{\mathbf{F}_{:,i}}{\mathbf{F}_{k_{n-i+1},i}}$ 
  for  $j \in [n]$  do
     $\mathbf{d}_j \leftarrow \mathbf{d}_j - |\mathbf{F}_{j,i}|^2$ 
  end for
end for

```

$\text{ChoosePivot}(\mathbf{d}) = \arg \max_i \mathbf{d}_i$. Historically, this been motivated by the fact that for a generic positive definite matrix \mathbf{A} , we have $|\mathbf{A}_{i,j}| < \sqrt{\mathbf{A}_{i,i}\mathbf{A}_{j,j}}$; as the residual of a Nyström approximation is positive definite, controlling the magnitude of the diagonal entries of its residual controls the magnitude of the off-diagonal entries, too. The analysis in [subsection 4.1.3](#) provides additional motivation when \mathbf{A} is a covariance matrix.

Both uniform and greedy pivoting have important failure cases associated with them ([4]). In the context of GPs, if the data $\{\mathbf{x}^{(i)}\}_{i=1}^n$ is grouped into order r small clusters, uniform pivoting will often fail to select a pivot corresponding to each cluster, yielding poor approximations for columns corresponding to those missed clusters. On the other hand, greedy pivoting is often derailed in the presence of order r points of $\{\mathbf{x}^{(i)}\}_{i=1}^n$ which are far away from each other and the rest, as it often chooses columns corresponding to these outliers that yield poor approximations for the rest of the columns.

One method which mitigates both these failure modes is known as Randomly Pivoted Cholesky (RPCholesky). To find $\text{ChoosePivot}(\mathbf{d})$, this method samples index i with likelihood \mathbf{d}_i . This is a middle ground between uniform and greedy sampling: [20] observes that all three methods can be viewed as sampling with likelihood \mathbf{d}_i^β , with $\beta = 0$ for uniform, $\beta \rightarrow \infty$ for greedy, and $\beta = 1$ for RPCholesky. For general positive definite matrices, [4] prove the following bound on RPCholesky error:

PROPOSITION 4.4. *Suppose \mathbf{A} is positive-definite, \mathbf{X}^* is such that $\mathbf{A}\langle\mathbf{X}^*\rangle$ is the best rank- l Nyström approximation of \mathbf{A} in the Frobenius and spectral norms with positive semi-definite residual, and*

$$(4.15) \quad r \geq \frac{l}{\varepsilon} + l \log \left(\frac{\varepsilon \text{Tr}(\mathbf{A})}{\text{Tr}(\mathbf{A} - \mathbf{A}\langle\mathbf{X}^*\rangle)} \right).$$

Then for $\hat{\mathbf{A}}$ a rank- r RPCholesky approximation of \mathbf{A} , we have

$$(4.16) \quad \mathbb{E} \text{Tr}(\mathbf{A} - \hat{\mathbf{A}}) \leq (1 + \varepsilon) \text{Tr}(\mathbf{A} - \mathbf{A}\langle\mathbf{X}^*\rangle).$$

Note that this bound is in terms of the trace rather than the Frobenius or spectral norms. [5] also prove bounds on the spectral norm of the residual when \mathbf{A} is a kernel matrix with nonzero nugget parameter, but they are not useful in the general case with no nugget.

In our case where diagonal heuristic methods are motivated by the form of (4.9), RPCholesky can be justified over greedy pivoting because the second term's magnitude

will approach the first as r increases, but the terms have different signs. In other words, optimizing only the second term of (4.9) may correlate with choosing an index with worse first term, which the randomization attempts to adjust for. Together with the above theoretical bounds, there is a strong case for using RPCholesky for the Nyström part; we analyze its performance in section 5.

4.2. Greedy Conditional Selection. We now explore an alternative sparsity pattern chooser to the nearest neighbors type method in subsection 4.1.1 and subsection 4.1.2. While Theorem 4.2 is encouraging for said method, ultimately, no purely Euclidean nearest-neighbors based algorithm can incorporate any information gained from the Nyström part. This seems at odds with the idea of CNV approximations, where the column Nyström pivots and sparsity pattern should ideally complement each other to improve approximation quality. Is there any way our sparsity pattern chooser can be made to take into account column Nyström information?

4.2.1. GCS for FSAI Sparsity Pattern Selection. [11]’s Greedy Conditional Selection (GCS) is one possible idea, which we first develop in a purely-sparse-Vecchia context.

Let \mathbf{K} be the covariance matrix induced on some $\{\mathbf{x}^{(i)}\}_{i=1}^n \subseteq \mathbb{R}^d$ by some kernel function κ . In subsection 4.1.3, we motivated an algorithm which uses sparse approximation to estimate the change in KL loss function (4.5) from adding a pivot to the column Nyström approximation of \mathbf{K} , which is then greedily minimized to choose pivots.

One might similarly hope to greedily minimize (4.5) to choose the sparsity patterns $\{S_i\}_{i=1}^n$ of individual columns. Conveniently, (4.8)

$$\begin{aligned} 2D_{\text{KL}}(\mathcal{N}(0, \mathbf{K}), \mathcal{N}(0, (\mathbf{C}\mathbf{C}^*)^{-1})) \\ = \sum_{i=1}^n \left[\log((\mathbf{K}/S_i \setminus \{i\})_{ii}) - \log((\mathbf{K}/\{i+1:n\})_{ii}) \right] \end{aligned}$$

has the form of a sum in which terms correspond to individual columns’ sparsity patterns and are constant in changes to the sparsity patterns of other columns. In particular, the change in objective function from adding index k to the sparsity pattern in the i th column S_i is given by

$$\begin{aligned} (4.17) \quad & \log((\mathbf{K}/S_i \cup \{k\} \setminus \{i\})_{i,i}) - \log((\mathbf{K}/S_i \setminus \{i\})_{ii}) \\ & = \log\left((\mathbf{K}/S_i \setminus \{i\})_{ii} - \frac{(\mathbf{K}/S_i \setminus \{i\})_{ik}^2}{(\mathbf{K}/S_i \setminus \{i\})_{kk}} \right) - \log((\mathbf{K}/S_i \setminus \{i\})_{ii}) \\ & = \log\left(1 - \frac{(\mathbf{K}/S_i \setminus \{i\})_{ik}^2}{(\mathbf{K}/S_i \setminus \{i\})_{ii}(\mathbf{K}/S_i \setminus \{i\})_{kk}} \right) \end{aligned}$$

which is minimized by maximizing

$$(4.18) \quad \frac{(\mathbf{K}/S_i \setminus \{i\})_{ik}^2}{(\mathbf{K}/S_i \setminus \{i\})_{ii}(\mathbf{K}/S_i \setminus \{i\})_{kk}} = \text{Corr}(\mathbf{y}_i, \mathbf{y}_k \mid \mathbf{y}_{S_i \setminus \{i\}}) \text{ for } \mathbf{y} \sim \mathcal{N}(0, \mathbf{K})$$

over k . Since $(\mathbf{K}/S_i \setminus \{i\})_{ii}$ does not depend on k , we may omit division by it from the optimization if we wish.

Unfortunately, computing (4.18) for the order n points per column satisfying lower triangularity, for each of the n columns, even once would give such an algorithm a $\Omega(n^2)$ computational complexity. To make an algorithm based on (4.18) tractable,

we must only maximize (4.18) over a *candidate set* $C_i \subseteq \{i + 1 \dots n\}$ for each i th column, with $\max_i |C_i| = c$ much less than n . In other words, we must constrain our optimization space by enforcing $S_i \subseteq C_i$. Motivated by Theorem 4.1, [11] recommends choosing each C_i to correspond to nearest neighbors of $\mathbf{x}^{(i)}$, especially in the case of Matérn kernels. We note that c is taken to be greater than s , or else this simply recovers the s -nearest-neighbors-based sparsity pattern of subsection 4.1.1.

When restricted to optimizing over C_i instead of $[n]$ for each i , [11] gives two different algorithms to iteratively construct S_i by computing and greedily maximizing (4.18). One is based on Sherman-Morrison updates to $(\mathbf{K}_{S_i \setminus \{i\}, S_i \setminus \{i\}})^{-1}$, while the other updates select entries of a column Nyström residual based on S_i by maintaining a Cholesky factor. They are just as fast as each other; we give pseudocode for and analyze computational complexity of the latter in Appendix A.1, and the former is Algorithm C.1 in [11]. Both require $O(ncd + ncs^2)$ operations.

4.2.2. GCS for Adaptive Factorized Nyström Approximations. Now, we'll adapt GCS to tractably select the sparsity pattern when computing the sparse Vecchia approximation of the residual during CNV approximations.

Like any method for choosing a sparsity pattern, GCS will work off-the-shelf when applied directly to the Nyström residual. However, this naïve implementation requires explicitly computing the entire column Nyström residual, which takes order $n^2 r$ operations for explicit computation of the Nyström approximation plus n^2 times the operations needed for a kernel evaluation to compute \mathbf{K} explicitly (which could be even larger than r).

Fortunately, the interpretation discussed in section 3 that CNV approximations are sparse Vecchia approximations with the bottom r rows filled in allows us to avoid this issue. In particular, as each successive pivot k_j is added to the pivot set R , we can update $(\mathbf{K}_{R \setminus \{i\}, R \setminus \{i\}})^{-1}$ as required in [11]'s Algorithm C.1 or the pivoted Cholesky matrix required in Algorithm A.1 by treating k_j as if it were the index that GCS had chosen to add to S_i at that step. This is necessary for each $i \in [n] \setminus R$, but $(\mathbf{K}_{R \setminus \{i\}, R \setminus \{i\}})^{-1}$ does not even depend on i for $i \notin R$, and the pivoted Cholesky matrices required in Algorithm A.1 for each i at the step where $S_i = R$ are all submatrices of a common matrix which only needs to be computed once. After that, we can simply continue GCS as normal on a column-by-column basis. We provide pseudocode for such a modified version of the latter in A.1. The CNV version requires $\mathcal{O}(nr^2 + ns(c+r)(r+s) + n(c+r)d) = \mathcal{O}(nr^2 + ns^2c + nscr + ncd)$ operations, as shown in Appendix A.1 (note that we always have $c \geq s$).

5. Experiments and analysis. In this section, we compare the performance of algorithms discussed in subsection 4.1 and subsection 4.2 on Matérn kernel matrices. Specifically, we examine the various methods' recovery in the Frobenius norm

$$(5.1) \quad \|\mathbf{L}^* \mathbf{K} \mathbf{L} - I\|_{\text{FRO}},$$

and how well they control the condition number

$$(5.2) \quad \mathfrak{K}(\mathbf{L} \mathbf{K} \mathbf{L}^T) = \frac{\lambda_{\max}(\mathbf{L} \mathbf{K} \mathbf{L}^T)}{\lambda_{\min}(\mathbf{L} \mathbf{K} \mathbf{L}^T)}.$$

The Frobenius error is indicative of recovery of individual entries in the matrix; we use $\|\mathbf{L} \mathbf{K} \mathbf{L}^T - I\|_{\text{FRO}}$ rather than $\|\mathbf{K} - (\mathbf{L} \mathbf{L}^*)\|_{\text{FRO}}$ due to the Vecchia approximation's optimality in the former ([23]) and the latter depending significantly on $\lambda_{\min}(\mathbf{K})$ which could cause higher variance when testing multiple samplings of synthetic data.

The condition number gives a worst-case bound on the number of iterations of preconditioned conjugate gradient (see [subsection 2.1](#)). We also measure the actual performance of CNV preconditioners on PCG by running PCG with these preconditioners on a set of withheld data points and averaging the number of iterations required before the vector 2-norm error is below a threshold (10^{-4} times the error at initialization).

5.1. Setup. We implemented various methods for both the column Nyström pivot choosing and sparsity pattern choosing steps of the CNV approximation framework in the Julia programming language. Specifically, for the Nyström part, we implemented [\[24\]](#)’s Farthest Point Sampling (FPS), [\[4\]](#)’s Randomly Pivoted Cholesky (RPC), and our own approximate greedy KL minimization (GKL). We distinguish algorithms [Algorithm 4.3](#) and [Algorithm 4.4](#) by referring to the first as GKLR and the second as GKLC. For the sparsity pattern choosing part, we implemented [\[18\]](#)’s Nearest Neighbors (kNN) and [\[11\]](#)’s Greedy Conditional Selection (GCS). GCS requires us to specify a number of candidate points c to consider; in this section, we generally set this to a scalar multiple of s (nonzeros per column in the sparse part) and refer to GCS with a specific scalar by appending the scalar (e.g. GCS10 refers to GCS with candidate set size $c = 10s$). We specifically note that GCS with $c = s$ is equivalent to kNN, though we always write kNN rather than GCS1.

For data, we always take $n = 5000$. We work with synthetic data, namely points uniformly sampled from d -dimensional hypercubes. Working with synthetic data in this manner allows us to isolate the effects of changes in specific parameters through experimentation, and we can also sample many synthetic datasets and appropriately average performance metrics across these samples for robustness. In particular, all quantities reported are arithmetic means (Frobenius error, PCG iterations) or geometric means (condition numbers, max/min eigenvalues) over 3 different synthetic datasets sampled in the same manner. We work only with Matérn kernels in this section, and take $\nu = 2.5$ unless otherwise specified. Finally, we set bandwidths of our Matérn kernels to equal the trace of the $d \times d$ sample covariance matrix of the 5000 data points ([\[1\]](#)). This is useful, among other reasons, to compare performance data apples-to-apples as d is varied.

5.2. Results. We first attempt to determine the best value of c , which because kNN is the special case of GCS with $c = s$ has choosing between kNN and GCS as a subproblem. We plot, for two choices $r = 20$ and $r = 60$ of r the rank in the column Nyström step, the performance of the resulting preconditioners for several values of c and s in [Figure 1](#).

Unsurprisingly, $c = 10s$ universally performs the best. Recalling from [subsection 4.2.2](#) that the runtime of GCS is at most linear in c , and that all of our preconditioner-finding algorithms were designed to be $o(n^2)$, we believe the few PCG iterations saved by GCS10 over GCS5 is worth the increased complexity of the GCS step. Moreover, in cases where d is higher, GCS performs even better. In [Figure 2](#), GCS10 improves PCG convergence rate by an order of magnitude compared to kNN or GCS2 for $d = 5$ and $d = 10$. For d even larger than this, taking $c > 10$ may even be advisable.

We also observe that in [Figure 2](#), the condition number and mean PCG iterations behave counterintuitively in that they initially increase with increasing s . Note that this does not contradict the optimality properties of CNV approximations in [section 3](#), as the condition number is not one of the objectives which Vecchia approximations are optimal in, and the Frobenius error, which was such an objective, did still decrease

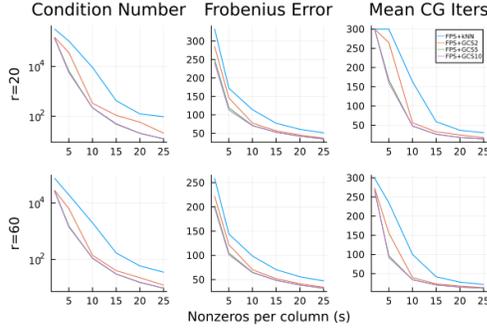


FIG. 1. Performance comparison of GCS for different values of c as r and s vary. Pivots are chosen by FPS.

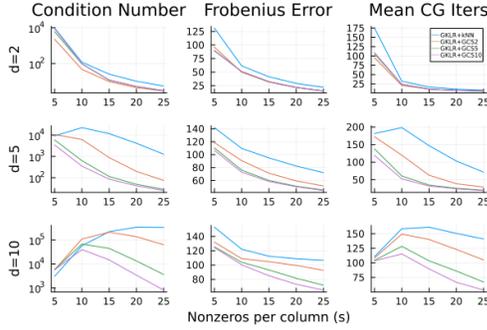


FIG. 2. Performance comparison of GCS for different values of c as d and s vary. $r = 60$ is fixed. Pivots are chosen by GKLR.

monotonically. We lack a complete explanation for this, but empirical observations indicate this behavior is due more to the largest eigenvalue increasing than the smallest eigenvalue decreasing as s increases.

Now that we have suitably chosen c , we turn our attention to comparing algorithms for the column Nyström step. We first compare GKLR and GKLC, which we plot in Figure 3 alongside the (completely intractable) algorithm which directly inverts the kernel matrix to recover (4.14) exactly. We refer to the latter as GKLE (as it is "exact"). Remarkably, the performance of each of these three algorithms is nearly indistinguishable, implying GKLR and GKLC both approximate (4.14) well even for single-digit values of s . We choose GKLR to represent GKL going forward, more or less arbitrarily.

Now, we compare [24]’s FPS, [4]’s RPC, and our GKLR, in Figure 4. We observe that GKLR is usually best, especially for large r and small s . GKLR takes as long as kNN, which takes $\Omega(ns^3)$, while GCS10 takes $\Omega(ns^2c)$ which for $c = 10s$ means approximately 10 times as long. This means that if we’re using GCS to choose the sparsity pattern, using GKLR as opposed to FPS or RPC requires only about 10% more operations for the creation of the preconditioner altogether.

As the preconditioning methods in this paper generally take order n^2 or less operations, which is less than one iteration of PCG, this is well worth saving even a few PCG iterations as it usually does for s relatively small. However, for relatively large s , the choice of column Nyström pivot chooser matters much less in terms

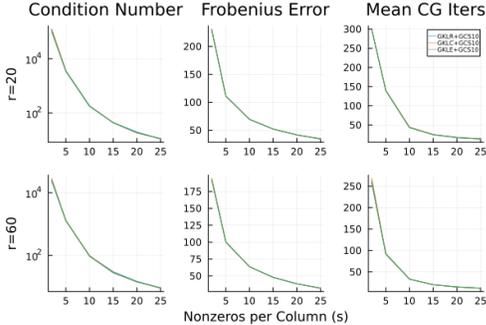


FIG. 3. Performance comparison of GKLR, GKLC, and GKLE for different values of r and s .

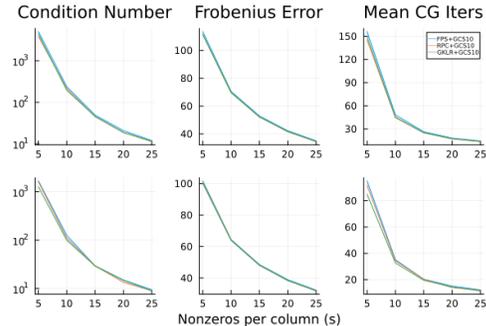


FIG. 4. Comparison of methods for column Nystrom step for different values of r and s . $r = 20$ in the top row and $r = 60$ in the bottom row.

of iterations for PCG convergence while GKLR (and preconditioner construction in general) takes comparatively longer, meaning using FPS or RPC instead and slightly increasing s might be a more efficient use of computational resources.

However, we remark that as ν increases, the gap between GKLR and FPS, and to a lesser extent RPC, widens. In Figure 5, we fix $s = 15$ and compare the algorithms for different values of ν and r . For $\nu = \frac{7}{2}$, one notices a significant gap between FPS and the other methods, though RPC seems to keep pace with GKLR relatively well.

Finally, we fix $s = 15$ and examine the effects of varying r in Figure 6, in attempt to answer how much of a performance improvement the low-rank component of CNV really contributes. As can be seen, increasing r results in rapid improvements while r is still small, though the marginal benefit decreases substantially for $r > 5$ or so. Thus, CNV results in dramatic improvement over purely sparse Vecchia, here a three to four times improvement for the condition number and number of CG iterations, even for small r . Nevertheless, given that (in the case of GCS) the complexity of the preconditioning step is order $nr^2 + ns^2c$, and given that the improvement from increasing r past 5 or so is still significant, it seems that taking r somewhere between s and $s\sqrt{c}$ is worthwhile.

Acknowledgments. I would like to thank my advisor, Robert Webber, for his mathematical expertise as well as all he has taught me about the process of mathematical research. This project would not have come together if it were not for his help directing it, and I feel I have improved significantly as a mathematician thanks to his advice throughout the past two quarters.

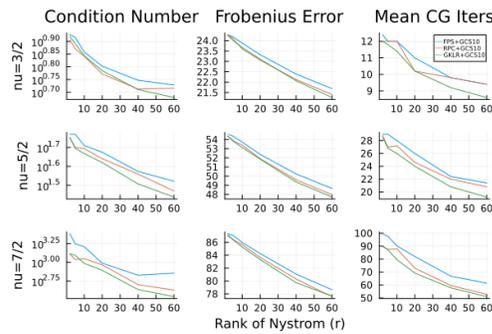


FIG. 5. Performance comparison of methods for column Nyström step for different values of ν and r . $s = 15$ is fixed.

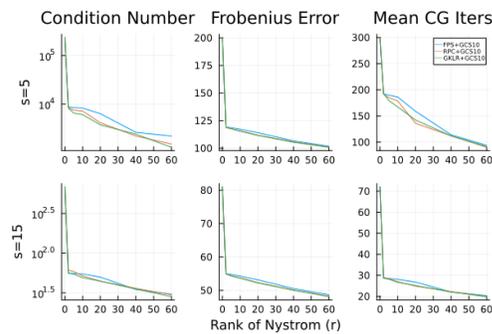


FIG. 6. Plotting objective functions as r increases from 0, to determine the impact of the column Nyström part in CNV approximations.

I would also like to thank Christopher J. Geoga, Chris Camaño, Ethan N. Epperly, and Florian Schäfer for helpful discussions.

REFERENCES

- [1] D. ARISTOFF, M. JOHNSON, G. SIMPSON, AND R. J. WEBBER, *The fast committor machine: Interpretable prediction with kernels*, The Journal of Chemical Physics, 161 (2024), p. 084113, <https://doi.org/10.1063/5.0222798>, <https://doi.org/10.1063/5.0222798>, https://arxiv.org/abs/https://pubs.aip.org/aip/jcp/article-pdf/doi/10.1063/5.0222798/20131959/084113_1_5.0222798.pdf.
- [2] F. R. BACH AND M. I. JORDAN, *Predictive low-rank decomposition for kernel methods*, in Proceedings of the 22nd International Conference on Machine Learning, ICML '05, New York, NY, USA, 2005, Association for Computing Machinery, p. 33–40, <https://doi.org/10.1145/1102351.1102356>.
- [3] D. BERNSTEIN, *Matrix Mathematics: Theory, Facts, and Formulas with Application to Linear Systems Theory*, Princeton University Press, 2005, <https://books.google.com/books?id=pmNRPwOFHKoC>.
- [4] Y. CHEN, E. N. EPPERLY, J. A. TROPP, AND R. J. WEBBER, *Randomly pivoted Cholesky: Practical approximation of a kernel matrix with few entry evaluations*, Communications on Pure and Applied Mathematics, 78 (2025), pp. 995–1041, <https://doi.org/10.1002/cpa.22234>.
- [5] M. DÍAZ, E. N. EPPERLY, Z. FRANGELLA, J. A. TROPP, AND R. J. WEBBER, *Robust, randomized preconditioning for kernel ridge regression*, 2024, <https://arxiv.org/abs/2304.12465>.
- [6] E. N. EPPERLY, J. A. TROPP, AND R. J. WEBBER, *Embrace rejection: Kernel matrix approximation by accelerated randomly pivoted Cholesky*, 2024, <https://arxiv.org/abs/2410.03969>.

- [7] S. FINE AND K. SCHEINBERG, *Efficient svm training using low-rank kernel representations*, J. Mach. Learn. Res., 2 (2002), p. 243–264.
- [8] A. GITTENS, *The spectral norm error of the naive nystrom extension*, arXiv preprint arXiv:1110.5305, (2011).
- [9] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations - 4th Edition*, Johns Hopkins University Press, Philadelphia, PA, 4 ed., 2013, <https://doi.org/10.1137/1.9781421407944>.
- [10] J. GUINNESS, *Permutation and grouping methods for sharpening gaussian process approximations*, Technometrics, 60 (2018), pp. 415–429.
- [11] S. HUAN, J. GUINNESS, M. KATZFUSS, H. OWHADI, AND F. SCHÄFER, *Sparse Cholesky factorization by greedy conditional selection*, 2023, <https://arxiv.org/abs/2307.11648>.
- [12] I. KAPORIN, *An alternative approach to estimating the convergence rate of the cg method*, Numerical Methods and Software, Yu. A. Kuznetsov, ed., Dept. of Numerical Mathematics, USSR Academy of Sciences, Moscow, (1990), pp. 55–72.
- [13] M. KATZFUSS AND J. GUINNESS, *A general framework for vecchia approximations of gaussian processes*, Statistical Science, 36 (2021), pp. 124–141.
- [14] M. KATZFUSS, J. GUINNESS, W. GONG, AND D. ZILBER, *Vecchia approximations of gaussian-process predictions*, Journal of Agricultural, Biological and Environmental Statistics, 25 (2020), pp. 383–414.
- [15] A. KRAUSE, A. SINGH, AND C. GUESTRIN, *Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies.*, Journal of Machine Learning Research, 9 (2008).
- [16] G. L. NEMHAUSER, L. A. WOLSEY, AND M. L. FISHER, *An analysis of approximations for maximizing submodular set functions—i*, Mathematical programming, 14 (1978), pp. 265–294.
- [17] A. J. ROTHMAN, E. LEVINA, AND J. ZHU, *Sparse multivariate regression with covariance estimation*, Journal of Computational and Graphical Statistics, 19 (2010), pp. 947–962, <http://www.jstor.org/stable/25765382> (accessed 2025-04-25).
- [18] F. SCHÄFER, M. KATZFUSS, AND H. OWHADI, *Sparse cholesky factorization by kullback–leibler minimization*, SIAM Journal on Scientific Computing, 43 (2021), pp. A2019–A2046, <https://doi.org/10.1137/20M1336254>.
- [19] F. SCHÄFER, T. J. SULLIVAN, AND H. OWHADI, *Compression, inversion, and approximate pca of dense kernel matrices at near-linear computational complexity*, arXiv preprint arXiv:1706.02205, (2017).
- [20] S. STEINERBERGER, *Randomly pivoted partial cholesky: Random how?*, arXiv preprint arXiv:2404.11487, (2024).
- [21] J. A. TROPP AND R. J. WEBBER, *Randomized algorithms for low-rank matrix approximation: Design, analysis, and applications*, arXiv preprint arXiv:2306.12418, (2023).
- [22] A. V. VECCHIA, *Estimation and model identification for continuous spatial processes*, Journal of the Royal Statistical Society Series B: Statistical Methodology, 50 (1988), pp. 297–312.
- [23] A. Y. YEREMIN, L. Y. KOLOTILINA, AND A. NIKISHIN, *Factorized sparse approximate inverse preconditionings. III. Iterative construction of preconditioners*, Journal of Mathematical Sciences, 101 (2000), pp. 3237–3254.
- [24] S. ZHAO, T. XU, H. HUANG, E. CHOW, AND Y. XI, *An adaptive factorized Nyström preconditioner for regularized kernel matrices*, SIAM Journal on Scientific Computing, 46 (2024), pp. A2351–A2376, <https://doi.org/10.1137/23M1565139>.

Appendix A. Postponed Pseudocode and Complexity Analysis.

A.1. Greedy Conditional Selection. In this section, we provide pseudocode for and examine the computational complexity of the two algorithms for GCS (one for purely sparse Vecchia and one for CNV) proposed in [subsection 4.2](#).

First, we discuss Algorithm [A.1](#). We write $\ker(x^{(\cdot)}, x^{(\cdot)})$ as opposed to \mathbf{K}_{\cdot} , to emphasize that the entire kernel matrix need not be explicitly computed. Examining the pseudocode, we require order ncs kernel evaluations which takes order ncd operations for most commonly-used kernels, plus ns instances of up to $c \times s$ matrix by $s \times 1$ vector multiplication which takes order ncs^2 operations, plus nsc scalar updates in the innermost loop, for a complexity of $O(ncd + ncs^2)$. Correctness of the updates to the diagonal and respective column is more rigorously proven in [\[11\]](#).

We now discuss our modification to work with CNV, given in Algorithm [A.1](#). Compared to Algorithm [A.1](#), it adds $O(nr^2)$ operations up front (which would have

Algorithm A.1 GCS via Truncated Cholesky

```

for  $i \in [n]$  do
   $\{k_l\}_{l=1}^c \leftarrow C_i$  ▷ Order arbitrary
   $k_{c+1} \leftarrow i$ 
   $\mathbf{d} \leftarrow \kappa(\mathbf{x}^{\{\{k_l\}_{l=1}^c\}}, \mathbf{x}^{\{\{k_l\}_{l=1}^c\}})$  ( $= \mathbf{1}$  for kernels) ▷ Nyström residual diagonal
   $\mathbf{b} \leftarrow \kappa(\mathbf{x}^{\{\{k_l\}_{l=1}^c\}}, \mathbf{x}^{(i)})$  ▷ Nyström residual  $i$ th column
   $\mathbf{L} = \mathbf{0} \in \mathbb{R}^{c+1 \times c+1}$ 
  for  $j \in [s]$  do
     $m = \arg \max_{l \in k_{[c]}} \frac{b_l^2}{d_l}$ 
     $S_i \leftarrow S_i \cup \{m\}$ 
     $\mathbf{L}_{:,j} \leftarrow \kappa(\mathbf{x}^{\{\{k_l\}_{l=1}^c\}}, \mathbf{x}^{(m)})$ 
     $\mathbf{L}_{:,j} \leftarrow \mathbf{L}_{:,j} - \mathbf{L}_{:,1:j-1}(\mathbf{L}_{m,1:j-1})^*$ 
     $\mathbf{L}_{:,j} \leftarrow \frac{\mathbf{L}_{:,j}}{\sqrt{\mathbf{L}_{m,j}}}$ 
    for  $l \in [c]$  do
       $d_l \leftarrow d_l - \mathbf{L}_{k_l,j}^2$ 
       $b_l \leftarrow b_l - \mathbf{L}_{k_l,j} \mathbf{L}_{c+1,j}$ 
    end for
  end for
end for

```

been needed for truncated Cholesky-based column Nyström anyway), then adds $O(ncr)$ operations across all columns for updating Nyström residuals' diagonals and corresponding columns (this can be optimized to $O(nr)$ by computing them all at once, but is not a bottleneck). It also increases the complexity of the matrix-vector multiplication in the second block of loops to order $(c+r)(r+s)$ operations, which happen order ns times for a bound of $O(ns(c+r)(r+s))$ overall, as well as requiring order rn more kernel evaluations which usually takes $O(rnd)$ more operations. So the CNV version requires $O(nr^2 + ns(c+r)(r+s) + n(c+r)d) = O(nr^2 + ns^2c + nscr + ncd)$ operations (note that we always have $c \geq s$).

Algorithm A.2 CNV using GCS via Truncated Cholesky

```

 $\bar{\mathbf{L}} \leftarrow \mathbf{0} \in \mathbb{R}^{n \times r+s}$ 
for  $j \in [r]$  do
   $\bar{\mathbf{L}}_{:,j} \leftarrow \kappa(\mathbf{x}^{(n)}, \mathbf{x}^{(k_j)})$  ▷  $\{k_j\}_{j=1}^r$  should be the Nyström pivots
   $\bar{\mathbf{L}}_{:,j} \leftarrow \mathbf{L}_{:,j} - \mathbf{L}_{:,1:j-1}(\mathbf{L}_{k_j,1:j-1})^*$ 
   $\mathbf{L}_{:,j} \leftarrow \frac{\bar{\mathbf{L}}_{:,j}}{\sqrt{\mathbf{L}_{k_j,j}}}$ 
end for
for  $i \in [n]$  do
   $\{k_j\}_{j=r+1}^{c+r} \leftarrow C_i$  ▷ Order arbitrary
   $k_{c+r+1} \leftarrow i$ 
   $\mathbf{d} \leftarrow \kappa(\mathbf{x}^{\{\{k_l\}_{l=r+1}^{c+r}\}}, \mathbf{x}^{\{\{k_l\}_{l=r+1}^{c+r}\}})$  (= 1 for kernels) ▷ Nyström residual diagonal
   $\mathbf{b} \leftarrow \kappa(\mathbf{x}^{\{\{k_l\}_{l=r+1}^{c+r}\}}, \mathbf{x}^{(i)})$  ▷ Nyström residual  $i$ th column
  for  $j \in [r]$  do
    for  $l \in [c]$  do
       $\mathbf{d}_l \leftarrow \mathbf{d}_l - \mathbf{L}_{l,j}^2$ 
       $\mathbf{b}_l \leftarrow \mathbf{b}_l - \mathbf{L}_{l,j} \mathbf{L}_{c+1,j}$ 
    end for
  end for
   $\mathbf{L} = \bar{\mathbf{L}}_{R \cup C, :}$ 
  for  $j \in [s] + r$  do
     $m = \arg \max_{l \in k_{[c]+r}} \frac{\mathbf{b}_l^2}{\mathbf{d}_l}$ 
     $S_i \leftarrow S_i \cup \{m\}$ 
     $\mathbf{L}_{:,j} \leftarrow \kappa(\mathbf{x}^{\{\{k_l\}_{l=1}^c\}}, \mathbf{x}^{(m)})$ 
     $\mathbf{L}_{:,j} \leftarrow \mathbf{L}_{:,j} - \mathbf{L}_{:,1:j-1}(\mathbf{L}_{m,1:j-1})^*$ 
     $\mathbf{L}_{:,j} \leftarrow \frac{\mathbf{L}_{:,j}}{\sqrt{\mathbf{L}_{m,j}}}$ 
    for  $l \in [c] + r$  do
       $\mathbf{d}_l \leftarrow \mathbf{d}_l - \mathbf{L}_{l,j}^2$ 
       $\mathbf{b}_l \leftarrow \mathbf{b}_l - \mathbf{L}_{l,j} \mathbf{L}_{c+1,j}$ 
    end for
  end for
end for

```
