

EVERYTHING IS VECCHIA: UNIFYING LOW-RANK AND SPARSE INVERSE CHOLESKY APPROXIMATIONS

EAGAN KAMINETZ* AND ROBERT J. WEBBER†

Abstract. The partial pivoted Cholesky approximation accurately represents matrices that are close to being low-rank. Meanwhile, the Vecchia approximation accurately represents matrices with inverse Cholesky factors that are close to being sparse. What happens if a partial Cholesky approximation is combined with a Vecchia approximation of the residual? This paper shows how the sum can be rewritten as a Vecchia approximation of the original matrix with an augmented sparsity pattern. Thus, the Vecchia approximation is a superset of other factored matrix approximations and it has broad applicability. The paper analyzes the optimality properties of the Vecchia approximation, and it tests this approximation on kernel matrices for high-dimensional machine learning applications.

Key words. Vecchia approximation, partial pivoted Cholesky approximation, kernel matrix, factorized sparse approximate inverse

AMS subject classifications. 65F55, 65C99, 15A23

1. Motivation. The goal of this paper is approximating a large, dense positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ by looking up and processing individual entries $\mathbf{A}(i, j)$. For example, \mathbf{A} might be the kernel matrix for a high-dimensional machine learning data set [24]. In this context, looking up each entry $\mathbf{A}(i, j)$ requires a constant, potentially large number of arithmetic operations. Since a kernel matrix can be very large (e.g., $n \geq 10^5$), kernel computations require an approximation $\hat{\mathbf{A}} \approx \mathbf{A}$ that is generated in *linear* or *sublinear* time. A linear-time algorithm runs in $\mathcal{O}(n^2)$ arithmetic operations, comparable to the cost of looking at each entry of \mathbf{A} once. A sublinear-time algorithm runs in $o(n^2)$ operations while looking at just a fraction of \mathbf{A} 's entries.

Sparse Cholesky approximation is a flexible framework for sublinear- and linear-time matrix approximation. This framework includes the partial pivoted Cholesky approximation [6] and the Vecchia approximation [29], which are traditionally regarded as accurate approximations for different types of matrices. Partial pivoted Cholesky provides an accurate approximation when the target matrix is close to being low-rank [6, Thm. 2.3]. Meanwhile, Vecchia provides an accurate approximation when the inverse Cholesky factor is close to being sparse; see [23, Sec. B.2] and [20]. This paper unifies the two approaches — it shows how the combination of a partial Cholesky approximation and a Vecchia approximation of the residual can be rewritten as a Vecchia approximation with an augmented sparsity pattern.

After positioning the Vecchia approximation as a unifying framework for efficient positive-semidefinite approximations, this paper asks and partially answers, “In what way is the Vecchia approximation optimal?” Optimality theory for the Vecchia approximation was developed in past work [29, 19, 3, 30, 23], but this paper presents a new extension to positive-semidefinite matrices and new error bounds for linear solves and determinant calculations; see [section 3](#) for details. Last, this paper evaluates the Vecchia approximation on kernel matrices for high-dimensional machine learning applications.

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

1.1. Notation. Scalars are written in lower case italics, e.g., m, n, r . Vectors are in lower case boldface, e.g., \mathbf{u}, \mathbf{v} . Matrices are in boldface capital letters, e.g., \mathbf{A}, \mathbf{B} . Index sets are in sans serif font, e.g., R, S . Following standard conventions, $\mathbf{0}$ denotes a vector or matrix of all zeros, \mathbf{I} is an identity matrix, and \mathbf{e}_i denotes a standard basis vector, which is all zeros except for 1 in the i th entry.

We use $\mathbf{u}(i)$ to refer to the i th entry of a vector \mathbf{u} , and we use $\mathbf{A}(i, j)$ to refer to the (i, j) entry of a matrix \mathbf{A} . Given index sets R, S , we use $\mathbf{u}(R)$ to refer to the subvector $(\mathbf{u}(i))_{i \in R}$, and we use $\mathbf{A}(R, S)$ to refer to the submatrix $(\mathbf{A}(i, j))_{i \in R, j \in S}$. Additionally, $\mathbf{A}(i, \cdot)$ and $\mathbf{A}(\cdot, i)$ indicate the i th row and column of \mathbf{A} , and we write $\ell : m$ to refer to $\{\ell, \dots, m\}$ when indexing vectors or matrices.

The complex conjugate of a scalar u is \bar{u} . The conjugate transpose of a vector \mathbf{v} or matrix \mathbf{A} is \mathbf{v}^* or \mathbf{A}^* . The inverse of \mathbf{A} is \mathbf{A}^{-1} , the conjugate transpose inverse is \mathbf{A}^{-*} , and the Moore-Penrose pseudoinverse is \mathbf{A}^+ . Given a positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, the \mathbf{A} -weighted inner product, norm, and distance function are defined as

$$\langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{A}} = \mathbf{u}^* \mathbf{A} \mathbf{v}, \quad \|\mathbf{v}\|_{\mathbf{A}} = \langle \mathbf{v}, \mathbf{v} \rangle_{\mathbf{A}}^{1/2}, \quad \text{and} \quad d_{\mathbf{A}}(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|_{\mathbf{A}}.$$

In the rank-deficient case, $\|\cdot\|_{\mathbf{A}}$ is not technically a norm. It is a seminorm that defines an equivalence relation $\mathbf{u} \sim \mathbf{v}$ when $\mathbf{u} - \mathbf{v} \in \text{null}(\mathbf{A})$. Last, in this paper, the volume of a matrix \mathbf{A} is the product of its nonzero singular values,

$$\text{vol}(\mathbf{A}) = \prod_{\sigma_i(\mathbf{A}) > 0} \sigma_i(\mathbf{A}),$$

which aligns with the determinant if \mathbf{A} is strictly positive definite.

1.2. Organization of paper. The rest of the paper is organized as follows. [Section 2](#) provides background and establishes the main theoretical result. [Section 3](#) derives optimality theory for the Vecchia approximation. [Section 4](#) introduces algorithmic optimization strategies for the Vecchia approximation, and [section 5](#) presents numerical experiments.

2. Background and contributions. This section introduces the framework of sparse Cholesky approximation ([subsection 2.1](#)). Then it proves the paper's main theoretical result ([subsection 2.2](#)) and discusses broader implications ([subsection 2.3](#)).

2.1. Factored approximations based on the Cholesky decomposition. Any positive-semidefinite matrix can be exactly represented through a Cholesky or inverse Cholesky decomposition.

DEFINITION 2.1 (Cholesky and inverse Cholesky decompositions). *The pivoted Cholesky and pivoted inverse Cholesky decompositions for a positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ are defined as*

$$(2.1) \quad \mathbf{A} = \mathbf{P} \mathbf{L} \mathbf{D} \mathbf{L}^* \mathbf{P}^* \quad \text{and} \quad \mathbf{A} = \mathbf{P} \mathbf{C}^{-1} \mathbf{D} \mathbf{C}^{-*} \mathbf{P}^*.$$

Here, $\mathbf{P} \in \{0, 1\}^{n \times n}$ is a permutation matrix, $\mathbf{D} \in \mathbb{R}_+^{n \times n}$ is a nonnegative-valued diagonal matrix, and $\mathbf{L}, \mathbf{C} \in \mathbb{C}^{n \times n}$ are lower triangular matrices with ones on the diagonal.

Cholesky and inverse Cholesky decompositions exist for any positive-semidefinite matrix \mathbf{A} and any permutation matrix \mathbf{P} . See [Figure 1](#) for illustrations. The Cholesky

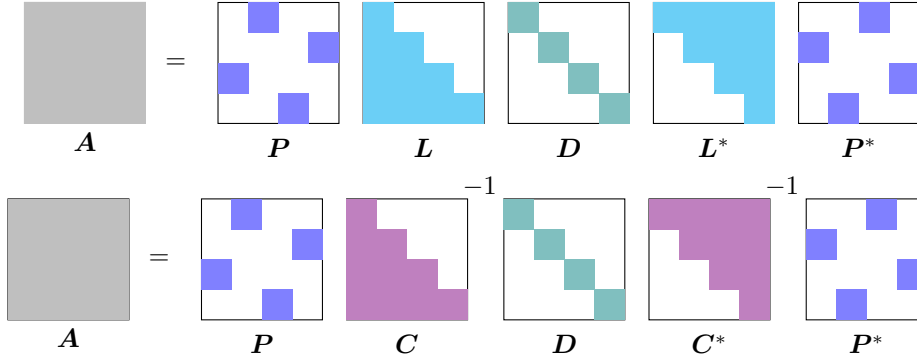


FIG. 1. Cholesky and inverse Cholesky decompositions of a dense matrix \mathbf{A} . Factors $\mathbf{P}, \mathbf{L}, \mathbf{D}$ or $\mathbf{P}, \mathbf{C}, \mathbf{D}$ are stored, and inverses $\mathbf{C}^{-1}, \mathbf{C}^{-*}$ are accessed implicitly. Filled boxes show entries that are allowed to be nonzero.

decomposition can be manipulated into inverse Cholesky form and vice versa through a matrix inversion $\mathbf{C} = \mathbf{L}^{-1}$.

The Cholesky and inverse Cholesky decompositions facilitate fast matrix computations. Once the factors $\mathbf{P}, \mathbf{D}, \mathbf{L}$ or $\mathbf{P}, \mathbf{D}, \mathbf{C}$ have been generated and stored, the matrix \mathbf{A} does not need to be accessed again. Rather, each matrix–vector product $\mathbf{v} \mapsto \mathbf{A}\mathbf{v}$ can be computed by sequentially multiplying the vector with each matrix in the factorization (2.1). The multiplications with \mathbf{C}^{-1} or \mathbf{C}^{-*} can be carried out via efficient triangular solves, without computing the inverse explicitly. Any consistent linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$ can be solved by multiplying the output vector \mathbf{b} by a sequence of matrices,

$$\mathbf{x} = \mathbf{P}\mathbf{L}^{-*}\mathbf{D}^+\mathbf{L}^{-1}\mathbf{P}^*\mathbf{b} \quad \text{or} \quad \mathbf{x} = \mathbf{P}\mathbf{C}^*\mathbf{D}^+\mathbf{C}\mathbf{P}^*\mathbf{b}.$$

Thus, the cost of a linear solve is $\mathcal{O}(n^2)$ arithmetic operations.

Motivated by the Cholesky and inverse Cholesky decompositions, we can generate a sparse Cholesky or sparse inverse Cholesky approximation

$$(2.2) \quad \hat{\mathbf{A}} = \mathbf{P}\hat{\mathbf{L}}\hat{\mathbf{D}}\hat{\mathbf{L}}^*\mathbf{P}^* \quad \text{or} \quad \hat{\mathbf{A}} = \mathbf{P}\hat{\mathbf{C}}^{-1}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-*}\mathbf{P}^*.$$

Here \mathbf{P} is a permutation matrix, $\hat{\mathbf{D}}$ is a nonnegative-valued diagonal matrix, and $\hat{\mathbf{L}}$ or $\hat{\mathbf{C}}$ is a sparse lower triangular matrix with ones on the diagonal. The *sparsity pattern* $\{\mathbf{S}_i\}_{i=1}^n$ is a collection of index sets $\mathbf{S}_i \subseteq \{1, \dots, i-1\}$ describing which off-diagonal entries of $\hat{\mathbf{L}}$ or $\hat{\mathbf{C}}$ are allowed to be nonzero.

The imposition of sparsity leads to three benefits. First, matrix–vector products and linear solves can be computed in $\mathcal{O}(ns)$ arithmetic operations, where s is an upper bound on the cardinality of the sparsity pattern: $|\mathbf{S}_i| \leq s$ for each $i = 1, \dots, n$. Second, the sparse approximation factors can be stored in $\mathcal{O}(sn)$ memory. Third, in many cases generating the approximation $\hat{\mathbf{A}}$ is relatively cheap, as it only requires examining $\mathcal{O}(sn)$ or $\mathcal{O}(s^2n)$ entries of \mathbf{A} .

Subsections 2.1.1 and 2.1.2 describe two specific examples of sparse Cholesky approximations that can be generated in sublinear or linear time.

2.1.1. Partial pivoted Cholesky. The partial pivoted Cholesky approximation is a common rank-revealing factorization for positive-semidefinite matrices [13]. This approximation is generated from r selected columns of the matrix \mathbf{A} , where r is a

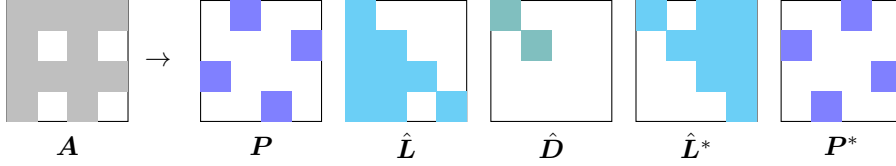


FIG. 2. *Partial pivoted Cholesky* accesses the gray-colored entries of \mathbf{A} . The approximation rank is $r = 2$ and the columns $u_1 = 3$ and $u_2 = 1$ are perfectly replicated.

Algorithm 2.1 Partial pivoted Cholesky approximation

Require: Positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ with entry-wise access; permutation matrix $\mathbf{P} \in \{0, 1\}^{n \times n}$; approximation rank r

Ensure: Sparse Cholesky approximation $\hat{\mathbf{A}} = \mathbf{P} \hat{\mathbf{L}} \hat{\mathbf{D}} \hat{\mathbf{L}}^* \mathbf{P}^*$ in factored form

Initialize sparse matrices $\hat{\mathbf{F}} = \mathbf{P}$ and $\hat{\mathbf{D}} = \mathbf{0} \in \mathbb{C}^{n \times n}$ $\triangleright \hat{\mathbf{F}} = \mathbf{P} \hat{\mathbf{L}}$

for $i = 1, \dots, r$ **do**

Identify pivot $u_i \in \{1, \dots, n\}$ with $\mathbf{P}(u_i, i) = 1$

$\mathbf{v} \leftarrow \mathbf{A}(\cdot, u_i) - \hat{\mathbf{F}} \hat{\mathbf{D}} \hat{\mathbf{F}}(u_i, \cdot)^*$

if $\mathbf{v}(u_i) > 0$ **then**

$\hat{\mathbf{F}}(\cdot, i) \leftarrow \mathbf{v}/\mathbf{v}(u_i)$

$\hat{\mathbf{D}}(i, i) \leftarrow \mathbf{v}(u_i)$

end if

end for

$\hat{\mathbf{L}} = \mathbf{P}^* \hat{\mathbf{F}}$

rank parameter chosen by the user. The approximation exactly replicates the selected columns. See the following definition and see Figure 2 for an illustration.

DEFINITION 2.2 (Partial pivoted Cholesky). *Given a positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, the partial pivoted Cholesky approximation with permutation \mathbf{P} and approximation rank r is a sparse Cholesky approximation $\hat{\mathbf{A}} = \mathbf{P} \hat{\mathbf{L}} \hat{\mathbf{D}} \hat{\mathbf{L}}^* \mathbf{P}^*$ where each row $\hat{\mathbf{L}}(i, \cdot)$ has sparsity pattern $S_i = \{1, \dots, r\} \cap \{1, \dots, i-1\}$ and where*

$$\begin{cases} (\hat{\mathbf{L}} \hat{\mathbf{D}} \hat{\mathbf{L}}^*)(i, j) = (\mathbf{P}^* \mathbf{A} \mathbf{P})(i, j), & \min\{i, j\} \leq r, \\ \hat{\mathbf{D}}(i, i) = 0, & r + 1 \leq i \leq n. \end{cases}$$

In this definition, the first r rows and columns of $\hat{\mathbf{L}} \hat{\mathbf{D}} \hat{\mathbf{L}}^*$ match the first r rows and columns of the permuted matrix $\mathbf{P}^* \mathbf{A} \mathbf{P}$. The last $n - r$ columns of $\hat{\mathbf{L}}$ and $\hat{\mathbf{D}}$ can be trivial, containing just zeros and ones.

Algorithm 2.1 generates a partial Cholesky approximation in $\mathcal{O}(r^2 n)$ arithmetic operations. The algorithm processes the columns of \mathbf{A} identified by the first r “pivots” u_1, \dots, u_r , which are defined by entries $\mathbf{P}(u_i, i) = 1$ in the permutation matrix. The $\mathcal{O}(r^2 n)$ cost arises because the i th step forms linear combinations of the first i selected columns.

A long line of research has investigated optimal pivot selection for the partial pivoted Cholesky approximation. “Randomly pivoted Cholesky” is a randomized selection rule that guarantees near-optimal approximation error in the expected trace norm [6], and other common pivot selection rules are based on greedy selection [13] or farthest point sampling [15]. See subsection 4.1 for more analysis and discussion.

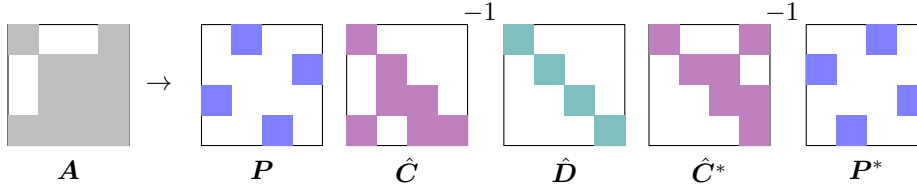


FIG. 3. *Vecchia approximation accesses the gray entries in \mathbf{A} . The pivots are $u_1 = 3$, $u_2 = 1$, $u_3 = 4$, $u_4 = 2$, and the sparsity pattern is $\mathbf{S}_1 = \emptyset$, $\mathbf{S}_2 = \emptyset$, $\mathbf{S}_3 = \{2\}$, $\mathbf{S}_4 = \{1, 3\}$.*

Algorithm 2.2 Conventional Vecchia algorithm

Require: Positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ with entry-wise access; permutation matrix $\mathbf{P} \in \{0, 1\}^{n \times n}$; sparsity pattern $\{\mathbf{S}_i\}_{i=1}^n$

Ensure: Vecchia approximation $\hat{\mathbf{A}} = \mathbf{P}\hat{\mathbf{C}}^{-1}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-*}\mathbf{P}^*$ in factored form

Initialize sparse matrices $\hat{\mathbf{C}} = \mathbf{I} \in \mathbb{C}^{n \times n}$ and $\hat{\mathbf{D}} = \mathbf{0} \in \mathbb{C}^{n \times n}$

for $i = 1, \dots, n$ **do**

▷ Can be executed in parallel

$$\begin{bmatrix} \mathbf{M} & \mathbf{v} \\ \mathbf{v}^* & \alpha \end{bmatrix} \leftarrow \begin{bmatrix} \mathbf{P}(\cdot, \mathbf{S}_i) & \mathbf{P}(\cdot, i) \end{bmatrix}^* \mathbf{A} \begin{bmatrix} \mathbf{P}(\cdot, \mathbf{S}_i) & \mathbf{P}(\cdot, i) \end{bmatrix}$$

Solve positive-semidefinite system $\mathbf{M}\mathbf{x} = -\mathbf{v}$

$$\hat{\mathbf{C}}(i, \mathbf{S}_i) \leftarrow \mathbf{x}^*$$

$$\hat{\mathbf{D}}(i, i) \leftarrow \alpha + \mathbf{x}^* \mathbf{v}$$

end for

2.1.2. Vecchia approximation. The Vecchia approximation [29] is a common approximation of covariance matrices of Gaussian processes. See below for a mathematical definition, and see Figure 3 for an illustration.

DEFINITION 2.3 (Vecchia approximation). *Given a positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, the Vecchia approximation with permutation \mathbf{P} and sparsity pattern $\{\mathbf{S}_i\}_{i=1}^n$ is a sparse inverse Cholesky approximation $\hat{\mathbf{A}} = \mathbf{P}\hat{\mathbf{C}}^{-1}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-*}\mathbf{P}^*$ where each row $\hat{\mathbf{C}}(i, \cdot)$ has sparsity pattern \mathbf{S}_i and it satisfies*

$$(2.3) \quad \begin{cases} (\hat{\mathbf{C}}\tilde{\mathbf{A}})(i, \mathbf{S}_i) = 0, \\ (\hat{\mathbf{C}}\tilde{\mathbf{A}})(i, i) = \hat{\mathbf{D}}(i, i), \end{cases} \quad \text{for } \tilde{\mathbf{A}} = \mathbf{P}^* \mathbf{A} \mathbf{P}.$$

In this definition, the Vecchia approximation is based on a sequence of n linear solves that determine each row vector $\hat{\mathbf{C}}(i, \mathbf{S}_i)$ and associated scalar $\hat{\mathbf{D}}(i, i)$. The Vecchia approximation does not necessarily reproduce the entries of \mathbf{A} that were examined during the algorithm, but it guarantees a different type of approximation quality framed in terms of the Kaporin condition number; see section 3.

Algorithm 2.2 presents a conventional Vecchia implementation that evaluates each linear system separately, creating opportunities for parallelism. The total cost of the Vecchia construction is $\mathcal{O}(s^3 n)$ arithmetic operations, where s is an upper bound on the sparsity: $|\mathbf{S}_i| \leq s$ for $i = 1, \dots, n$. Yet speedups may be possible by solving the linear systems in a different way — this paper investigates a Vecchia construction that has a smaller cost of $\mathcal{O}(s^2 n)$ arithmetic operations when the sparsity pattern takes a particular structured form; see subsection 2.2.

The design choices in the Vecchia approximation are the permutation matrix \mathbf{P} and the sparsity pattern $\{\mathbf{S}_i\}_{i=1}^n$. Historically, researchers chose the sparsity pattern using nearest neighbors in a problem-specific metric [25], and they generated the

permutation matrix by recursively choosing the farthest-away index in the same metric [16]. More recently, Huan et al. [18] proposed optimizing each index set S_i using a greedy heuristic that has been studied in the sparse approximation literature [27].

2.2. Main theoretical result. This paper analyzes a hybrid approximation that combines the partial pivoted Cholesky and Vecchia approximations [31, 5]. We establish the following main result.

THEOREM 2.4 (Partial Cholesky + Vecchia = Vecchia). *Given a target positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, consider the following two-part approximation.*

1. *Generate a partial Cholesky approximation of \mathbf{A} with permutation \mathbf{P} and approximation rank r . Call it $\hat{\mathbf{A}}_{\text{part}}$.*
2. *Generate a Vecchia approximation of the residual $\mathbf{R} = \mathbf{A} - \hat{\mathbf{A}}_{\text{part}}$ with permutation \mathbf{P} and sparsity pattern $\{\mathbf{Q}_i\}_{i=1}^n$. Call it $\hat{\mathbf{A}}_{\text{res}}$.*

Then $\hat{\mathbf{A}}_{\text{part}} + \hat{\mathbf{A}}_{\text{res}}$ can be rewritten as a Vecchia approximation of \mathbf{A} with permutation \mathbf{P} and an augmented sparsity pattern $S_i = (\{1, \dots, r\} \cup \mathbf{Q}_i) \cap \{1, \dots, i-1\}$.

Proof. Introduce the permuted matrices $\tilde{\mathbf{A}} = \mathbf{P}^* \mathbf{A} \mathbf{P}$, $\tilde{\mathbf{R}} = \mathbf{P}^* \mathbf{R} \mathbf{P}$,

$$\begin{aligned} \tilde{\mathbf{A}}_{\text{part}} &= \mathbf{P}^* \hat{\mathbf{A}}_{\text{part}} \mathbf{P} = \begin{bmatrix} \hat{\mathbf{L}}_{11} & \\ \hat{\mathbf{L}}_{21} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{D}}_{11} & \\ & \mathbf{0} \end{bmatrix} \begin{bmatrix} \hat{\mathbf{L}}_{11} & \\ \hat{\mathbf{L}}_{21} & \mathbf{I} \end{bmatrix}^*, \text{ and} \\ \tilde{\mathbf{A}}_{\text{res}} &= \mathbf{P}^* \hat{\mathbf{A}}_{\text{res}} \mathbf{P} = \begin{bmatrix} \mathbf{I} & \\ \mathbf{0} & \hat{\mathbf{C}}_{22} \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{0} & \\ & \hat{\mathbf{D}}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \\ \mathbf{0} & \hat{\mathbf{C}}_{22} \end{bmatrix}^{-*} \end{aligned}$$

where $\tilde{\mathbf{A}}_{\text{part}}$ and $\tilde{\mathbf{A}}_{\text{res}}$ are partitioned into the first r and last $n-r$ entries. We can add the approximations to produce a sparse inverse Cholesky approximation

$$\hat{\mathbf{A}} = \hat{\mathbf{A}}_{\text{part}} + \hat{\mathbf{A}}_{\text{res}} = \mathbf{P} \hat{\mathbf{C}}^{-1} \hat{\mathbf{D}} \hat{\mathbf{C}}^{-*} \mathbf{P}^*,$$

where

$$\hat{\mathbf{C}} = \begin{bmatrix} \hat{\mathbf{L}}_{11}^{-1} & \\ -\hat{\mathbf{C}}_{22} \hat{\mathbf{L}}_{21} \hat{\mathbf{L}}_{11}^{-1} & \hat{\mathbf{C}}_{22} \end{bmatrix} \quad \text{and} \quad \hat{\mathbf{D}} = \begin{bmatrix} \hat{\mathbf{D}}_{11} & \\ & \hat{\mathbf{D}}_{22} \end{bmatrix}.$$

By inspection, $\hat{\mathbf{A}}$ is a sparse Cholesky approximation with permutation \mathbf{P} and sparsity pattern $\{S_i\}_{i=1}^n$. See Figure 4 for an illustration. To show it is a Vecchia approximation, we need to verify the following equalities for $i = 1, \dots, n$:

$$(2.4) \quad (\hat{\mathbf{C}} \tilde{\mathbf{A}})(i, \mathbf{Q}_i) = 0 \quad \text{and} \quad (\hat{\mathbf{C}} \tilde{\mathbf{A}})(i, i) = \hat{\mathbf{D}}(i, i).$$

We will do this separately for $i \leq r$ and $i \geq r+1$.

First consider $i \leq r$. Since $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{A}}_{\text{part}}$ share the same first r rows, we can write

$$\begin{aligned} (\hat{\mathbf{C}} \tilde{\mathbf{A}})(i, \cdot) &= (\hat{\mathbf{C}} \tilde{\mathbf{A}}_{\text{part}})(i, \cdot) \\ &= (\hat{\mathbf{L}}_{11}^{-1} \hat{\mathbf{L}}_{11} \hat{\mathbf{D}}_{11})(i, \cdot) \begin{bmatrix} \hat{\mathbf{L}}_{11}^* & \hat{\mathbf{L}}_{21}^* \end{bmatrix} \\ &= \hat{\mathbf{D}}_{11}(i, i) \begin{bmatrix} \hat{\mathbf{L}}_{11}(\cdot, i)^* & \hat{\mathbf{L}}_{21}(\cdot, i)^* \end{bmatrix}. \end{aligned}$$

Since $\hat{\mathbf{L}}_{11}$ is lower triangular with ones on the diagonal, conclude that $(\hat{\mathbf{C}} \tilde{\mathbf{A}})(i, \cdot)$ is a row vector with the first $i-1$ entries equal to zero and the i th entry equal to $\hat{\mathbf{D}}(i, i) = \hat{\mathbf{D}}_{11}(i, i)$, confirming (2.4).

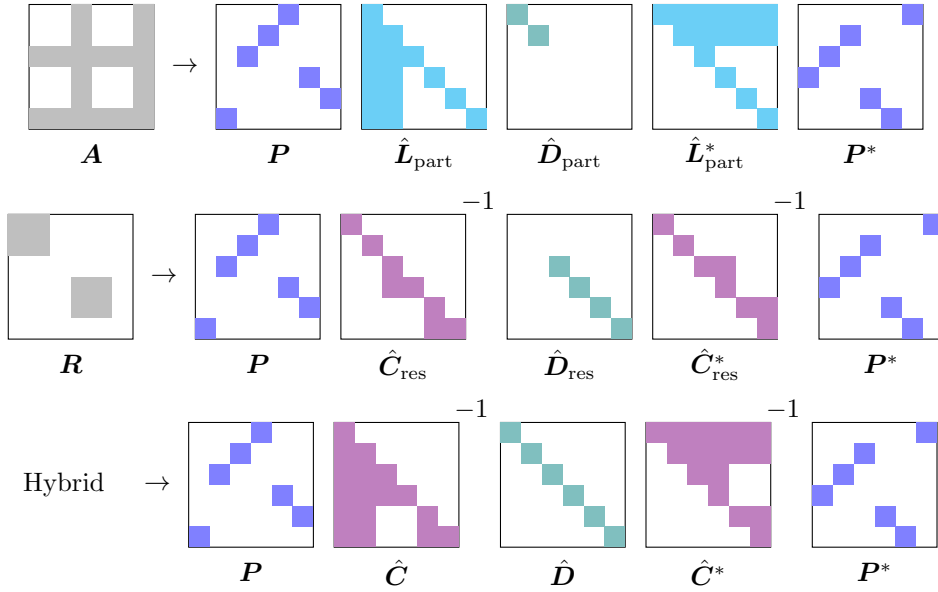


FIG. 4. First row: partial Cholesky accesses the gray entries of \mathbf{A} to generate an approximation $\hat{\mathbf{A}}_{\text{part}} = \mathbf{P}\hat{\mathbf{L}}_{\text{part}}\hat{\mathbf{D}}_{\text{part}}\hat{\mathbf{L}}_{\text{part}}^*\mathbf{P}^*$. Second row: Vecchia accesses the gray entries of $\mathbf{R} = \mathbf{A} - \hat{\mathbf{A}}_{\text{part}}$ to generate an approximation $\hat{\mathbf{A}}_{\text{res}} = \mathbf{P}\hat{\mathbf{C}}_{\text{res}}^{-1}\hat{\mathbf{D}}_{\text{res}}\hat{\mathbf{C}}_{\text{res}}^{*-1}\mathbf{P}^*$. Third row: partial Cholesky + Vecchia yields an improved approximation $\hat{\mathbf{A}} = \mathbf{P}\hat{\mathbf{C}}^{-1}\hat{\mathbf{D}}\hat{\mathbf{C}}^{*-1}\mathbf{P}^*$.

Next consider $i \geq r + 1$. We use the fact that $\tilde{\mathbf{R}}$ is all zeros except for the bottom right $(n - r) \times (n - r)$ submatrix, which we call $\tilde{\mathbf{R}}_{22}$. Therefore,

$$\begin{aligned} (\hat{\mathbf{C}}\tilde{\mathbf{A}})(i, \cdot) &= \hat{\mathbf{C}}(i, \cdot)(\tilde{\mathbf{A}}_{\text{part}} + \tilde{\mathbf{R}}) \\ &= \hat{\mathbf{C}}_{22}(i - r, \cdot) \begin{bmatrix} -\hat{\mathbf{L}}_{21}\hat{\mathbf{L}}_{11}^{-1} & \mathbf{I} \end{bmatrix} \left(\begin{bmatrix} \hat{\mathbf{L}}_{11} \\ \hat{\mathbf{L}}_{21} \end{bmatrix} \hat{\mathbf{D}}_{11} [\hat{\mathbf{L}}_{11}^* \hat{\mathbf{L}}_{21}^*] + \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \tilde{\mathbf{R}}_{22} \end{bmatrix} \right) \\ &= \hat{\mathbf{C}}_{22}(i - r, \cdot) \begin{bmatrix} \mathbf{0} & \tilde{\mathbf{R}}_{22} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{0} & \hat{\mathbf{C}}_{22} \end{bmatrix} (i, \cdot) \tilde{\mathbf{R}}. \end{aligned}$$

Since $\begin{bmatrix} \mathbf{0} & \hat{\mathbf{C}}_{22} \end{bmatrix} (i, \cdot)$ is a row vector generated by a Vecchia approximation of \mathbf{R} , we confirm the stated equalities in (2.4) and complete the proof. \square

2.3. Implications. Theorem 2.4 demonstrates that the partial Cholesky + Vecchia approach which appeared in the previous papers [31, 5] is secretly constructing a Vecchia approximation with the first r indices included in the sparsity pattern. We list two significant implications.

First, the partial Cholesky + Vecchia approach is *computationally efficient*. When $q \geq 1$ is an upper bound on the sparsity level in the residual component: $|\mathbf{Q}_i| \leq q$ for $i = 1, \dots, n$, the construction cost is $\mathcal{O}((r^2 + rq^2 + q^3)n)$ arithmetic operations, If $q = \mathcal{O}(r^{1/2})$, the construction cost is $\mathcal{O}((r+q)^2n)$ operations, which is smaller than the conventional Vecchia construction cost of $\mathcal{O}((r+q)^3n)$ operations using Algorithm 2.2.

Second, the partial Cholesky + Vecchia approach can be very *accurate*. In the simplest case, we set the residual component to have sparsity $\mathbf{Q}_i = \emptyset$ for $i = 1, \dots, n$,

leading to a “partial Cholesky + diagonal” approximation:

$$\hat{\mathbf{A}} = \hat{\mathbf{A}}_{\text{part}} + \hat{\mathbf{A}}_{\text{res}} = \hat{\mathbf{A}}_{\text{part}} + \text{diag}(\mathbf{A} - \hat{\mathbf{A}}_{\text{part}}).$$

In contrast, a previous algorithm for refining a partial Cholesky approximation added a multiple of the identity to the nullspace of $\hat{\mathbf{A}}_{\text{part}}$ [11]. The numerical experiments in [section 5](#) show that partial Cholesky + diagonal leads to $5\times$ faster conjugate gradient convergence for a poorly conditioned target matrix.

Going beyond the minimal sparsity pattern, we also consider a small, carefully chosen sparsity pattern $\{\mathbf{Q}_i\}_{i=1}^n$ for the residual component. Specifically, we investigate a two-part algorithm called “randomly pivoted Cholesky + Vecchia” (RPC+V).

1. Apply randomly pivoted Cholesky [6] to generate a rank- r approximation of a $n \times n$ positive-semidefinite matrix.
2. Generate a Vecchia approximation of the residual matrix, where each component of the sparsity pattern comes from a greedy optimization [18] with c candidate points and s selected nonzero entries.

To ensure a linear construction cost, we use the parameter settings $r = \mathcal{O}(n^{1/2})$, $c = \mathcal{O}(n^{1/4})$, and $s = \mathcal{O}(n^{1/4})$. Numerical experiments in [section 5](#) show that RPC+V leads to $1.5\times$ faster conjugate gradient convergence and $10\times$ more accurate log determinant estimation than randomly pivoted Cholesky + diagonal.

In summary, the paper’s main contribution is a recognition and theoretical justification that the Vecchia approximation is a superset of other factored matrix approximations. This perspective leads to higher accuracy Vecchia approximations that are also generated more quickly.

3. Kaporin optimality theory and applications. In this section, we ask and partially answer, “In what way is the Vecchia approximation optimal?”

Vecchia optimality theory was previously developed in the papers [29, 19, 3, 30, 23]. Here, we review this optimality theory and push it in new directions. Our analysis is based on the following Kaporin condition number [19, 3], which measures the accuracy of a matrix approximation.

DEFINITION 3.1 (Kaporin condition number). *For any positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and positive-semidefinite approximation $\hat{\mathbf{A}} \in \mathbb{C}^{n \times n}$, the Kaporin condition number is*

$$\kappa_{\text{Kap}} = \frac{\left(\frac{1}{r} \text{tr}(\mathbf{A}\hat{\mathbf{A}}^+)\right)^r}{\text{vol}(\mathbf{A}\hat{\mathbf{A}}^+)}, \quad \text{where } r = \text{rank}(\mathbf{A}),$$

if \mathbf{A} and $\hat{\mathbf{A}}$ share the same range. The Kaporin condition number is $\kappa_{\text{Kap}} = \infty$ if \mathbf{A} and $\hat{\mathbf{A}}$ have different ranges.

The Kaporin condition number is the average positive eigenvalue of $\mathbf{A}\hat{\mathbf{A}}^+$ raised to the $\text{rank}(\mathbf{A})$ power, divided by the product of the positive eigenvalues of $\mathbf{A}\hat{\mathbf{A}}^+$. By the arithmetic-geometric mean inequality, the Kaporin condition number is bounded from below by one. Equality holds if and only if all the positive eigenvalues are equal.

Kaporin [19, App. A.3] proved that the Vecchia approximation optimizes κ_{Kap} for strictly positive-definite \mathbf{A} subject to any given sparsity pattern. Here we extend this result to any positive-semidefinite \mathbf{A} . The proof is presented in [Appendix A](#).

THEOREM 3.2 (Optimality of Vecchia). *For any positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, the Vecchia approximation $\hat{\mathbf{A}} = \mathbf{P}\hat{\mathbf{C}}^{-1}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-*}\mathbf{P}^*$ is the inverse Cholesky*

approximation with permutation \mathbf{P} and sparsity pattern $\{\mathbf{S}_i\}_{i=1}^n$ that achieves the smallest possible Kaporin condition number. When $\hat{\mathbf{A}}$ and \mathbf{A} have the same range, the Kaporin condition number is

$$(3.1) \quad \kappa_{\text{Kap}} = \prod_{d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j < i}) > 0} \frac{d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{S}_i})^2}{d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j < i})^2},$$

where $\tilde{\mathbf{A}} = \mathbf{P}^* \mathbf{A} \mathbf{P}$ is the permuted \mathbf{A} matrix.

Several points are worth mentioning. First, the range of the Vecchia approximation always contains the range of the target matrix, $\text{range}(\mathbf{A}) \subseteq \text{range}(\hat{\mathbf{A}})$. The ranges are the same if and only if the sparsity set is sufficiently large. Specifically, when there are “bad” indices i for which

$$d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j < i}) = 0,$$

the sparsity component \mathbf{S}_i must be large enough that

$$d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{S}_i}) = 0.$$

Otherwise, the Kaporin condition number is infinite and there is no sparse Cholesky approximation with the same range as \mathbf{A} .

Second, there is a way to rewrite the Kaporin condition number (3.1) using an inverse Cholesky decomposition $\mathbf{A} = \mathbf{P} \mathbf{C}^{-1} \mathbf{D} \mathbf{C}^{-*} \mathbf{P}^*$.

$$\kappa_{\text{Kap}} = \frac{\text{vol}(\hat{\mathbf{A}})}{\text{vol}(\mathbf{A})} = \frac{\text{vol}(\hat{\mathbf{D}})}{\text{vol}(\mathbf{D})} = \prod_{\mathbf{D}(i,i) > 0} \frac{\hat{\mathbf{D}}(i,i)}{\mathbf{D}(i,i)}.$$

The entries $\hat{\mathbf{D}}(i,i)$ and $\mathbf{D}(i,i)$ are encoding the $\tilde{\mathbf{A}}$ -weighted distances.

$$\begin{cases} \hat{\mathbf{D}}(i,i) = d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{S}_i}), \\ \mathbf{D}(i,i) = d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j < i}). \end{cases}$$

However, the expressions using $\tilde{\mathbf{A}}$ -weighted distances are more intuitive, and they motivate the Vecchia optimization strategies that will be pursued in [section 4](#).

Third, the Kaporin condition number optimality is useful and important, because it controls the error of several linear algebra calculations, especially those involving determinants. For a list of Kaporin condition number bounds, see [Table 1](#). The rest of the section will introduce these bounds in more detail and explain how to use a factored matrix approximation like a Vecchia approximation to speed up linear algebra.

3.1. Direct linear system solves. In order to solve a linear system $\mathbf{A} \mathbf{x} = \mathbf{b}$ or a more general linear least-squares problem $\min_{\mathbf{x}} \|\mathbf{A} \mathbf{x} - \mathbf{b}\|^2$, we can make an initial guess \mathbf{x}_0 and then refine our guess by calculating

$$\hat{\mathbf{x}} = \mathbf{x}_0 + \hat{\mathbf{A}}^+ [\mathbf{b} - \mathbf{A} \mathbf{x}_0].$$

This calculation is very fast. We do not even need to examine the entries of \mathbf{A} once; we only need to perform a linear solve with a factored approximation $\hat{\mathbf{A}}$. The following proposition bounds the error of the resulting direct solve. To our knowledge, the proof is new, and it appears in [Appendix B.1](#).

Method	Error bound	Reference
Linear system, direct solver	$\frac{\ \hat{\mathbf{x}} - \mathbf{x}_\star\ _{\hat{\mathbf{A}}}^2}{\ \mathbf{x}_0 - \mathbf{x}_\star\ _{\mathbf{A}}^2} \leq 2 \text{rank}(\mathbf{A}) \log(\kappa_{\text{Kap}})$	Proposition 3.3
Linear system, iterative solver	$\frac{\ \mathbf{x}_t - \mathbf{x}_\star\ _{\hat{\mathbf{A}}}^2}{\ \mathbf{x}_0 - \mathbf{x}_\star\ _{\mathbf{A}}^2} \leq \left[\frac{3 \log(\kappa_{\text{Kap}})}{t} \right]^t$	[3], Proposition 3.4
Determinant, direct solver	$\log\left(\frac{\det \hat{\mathbf{A}}}{\det \mathbf{A}}\right) = \log(\kappa_{\text{Kap}})$	Proposition 3.5
Determinant, iterative solver	$\mathbb{E} \left \log\left(\frac{e^{st} \det \hat{\mathbf{A}}}{\det \mathbf{A}}\right) \right ^2 \leq \frac{4 \log(\kappa_{\text{Kap}})}{t}$	Proposition 3.6

TABLE 1

Summary of error bounds for direct and iterative solvers in terms of the Kaporin condition number κ_{Kap} . The direct solver bounds require the normalization $\text{tr}(\mathbf{A}\hat{\mathbf{A}}^+) = \text{rank}(\mathbf{A})$ which is automatically satisfied for the Vecchia approximation when $\kappa_{\text{Kap}} < \infty$. The determinant bounds require that \mathbf{A} and $\hat{\mathbf{A}}$ are strictly positive definite.

PROPOSITION 3.3 (Approximate direct solver for linear systems). *If $\hat{\mathbf{A}}$ is normalized so that $\text{tr}(\mathbf{A}\hat{\mathbf{A}}^+) = \text{rank}(\mathbf{A})$, the approximate direct solver satisfies*

$$\frac{\|\hat{\mathbf{x}} - \mathbf{x}_\star\|_{\hat{\mathbf{A}}}^2}{\|\mathbf{x}_0 - \mathbf{x}_\star\|_{\mathbf{A}}^2} \leq 2 \text{rank}(\mathbf{A}) \log(\kappa_{\text{Kap}}).$$

Here, $\mathbf{x}_\star = \mathbf{A}^+\mathbf{b}$ is the minimum-norm solution to $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$.

The Vecchia approximation uses the normalization $\text{tr}(\mathbf{A}\hat{\mathbf{A}}^+) = \text{rank}(\mathbf{A})$, so the guarantees in Proposition 3.3 are valid. Yet the guarantees in Proposition 3.3 are not very strong unless the Kaporin condition number is extremely small, $\log(\kappa_{\text{Kap}}) < 1/\text{rank}(\mathbf{A})$.

3.2. Iterative linear system solves. If the approximate direct solver fails to deliver an accurate solution, we can instead use the more expensive preconditioned conjugate gradient algorithm [14, Sec. 11.5] Starting with an initial guess \mathbf{x}_0 , PCG calculates the residual $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$ and the search direction $\mathbf{d}_1 = \hat{\mathbf{A}}^+\mathbf{r}_0$. At each step $t = 1, 2, \dots$, PCG updates the iterate and residual according to

$$\begin{cases} \mathbf{x}_t = \mathbf{x}_{t-1} + \alpha_t \mathbf{d}_t, \\ \mathbf{r}_t = \mathbf{r}_{t-1} - \alpha_t \mathbf{A} \mathbf{d}_t. \end{cases}$$

The step size and search direction are updated according to

$$\alpha_t = \frac{\mathbf{r}_{t-1}^* \hat{\mathbf{A}}^+ \mathbf{r}_{t-1}}{\mathbf{d}_t^* \mathbf{A} \mathbf{d}_t} \quad \text{and} \quad \mathbf{d}_{t+1} = \hat{\mathbf{A}}^+ \mathbf{r}_t + \frac{\mathbf{r}_t^* \hat{\mathbf{A}}^+ \mathbf{r}_t}{\mathbf{r}_{t-1}^* \hat{\mathbf{A}}^+ \mathbf{r}_{t-1}} \mathbf{d}_t.$$

We can run PCG for any number of iterations, producing better and better estimates of $\mathbf{x}_\star = \mathbf{A}^+\mathbf{b}$ with each iteration. PCG requires just one multiplication with \mathbf{A} and one linear solve with $\hat{\mathbf{A}}$ per iteration.

Axelsson and Kaporin [3, Thm. 4.3] bounded the error of PCG in terms of κ_{Kap} assuming an even number of iterations $t \geq 0$, and Appendix B.2 extends their proof to handle any even or odd number of iterations $t \geq 0$.

PROPOSITION 3.4 (Convergence of PCG [3]). *At each iteration $t \geq 0$, the iterates produced by preconditioned conjugate gradient satisfy*

$$(3.2) \quad \frac{\|\mathbf{x}_t - \mathbf{x}_*\|_{\mathbf{A}}^2}{\|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}^2} \leq \left\lceil \frac{3 \log(\kappa_{\text{Kap}})}{t} \right\rceil^t.$$

Here, $\mathbf{x}_* = \mathbf{A}^+ \mathbf{b}$ is the minimum-norm solution to $\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$.

Proposition 3.4 guarantees *superlinear* convergence, which is even stronger than the standard conjugate gradient linear convergence bounds (e.g., [14, Eq. 11.3.27]). Yet in large-scale applications, the superlinear convergence is mainly observed in the latter iterations when PCG has already achieved high accuracy [3]. As such, Proposition 3.4 provides only partial evidence that optimizing the Kaporin condition number is a good idea for solving linear systems.

3.3. Direct determinant calculation. When \mathbf{A} and $\hat{\mathbf{A}}$ are strictly positive definite, we can use $\det(\hat{\mathbf{A}})$ as an estimator for $\det(\mathbf{A})$. This computation is very fast for the Vecchia approximation $\hat{\mathbf{A}} = \mathbf{P}\hat{\mathbf{C}}^{-1}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-*}\mathbf{P}^*$, because

$$\det(\hat{\mathbf{A}}) = \det(\hat{\mathbf{D}}) = \prod_{i=1}^n \hat{D}(i, i).$$

The next proposition exactly describes the error in the determinant estimation.

PROPOSITION 3.5 (Approximate direct solver for determinants). *If \mathbf{A} and $\hat{\mathbf{A}}$ are strictly positive definite and $\hat{\mathbf{A}}$ is normalized so that $\text{tr}(\mathbf{A}\hat{\mathbf{A}}^{-1}) = n$, then*

$$\log\left(\frac{\det \hat{\mathbf{A}}}{\det \mathbf{A}}\right) = \log(\kappa_{\text{Kap}}).$$

Proof. We rewrite the Kaporin condition number as

$$\kappa_{\text{Kap}} = \frac{\left(\frac{1}{n} \text{tr}(\mathbf{A}\hat{\mathbf{A}}^{-1})\right)^n}{\det(\mathbf{A}\hat{\mathbf{A}}^{-1})} = \frac{1}{\det(\mathbf{A}\hat{\mathbf{A}}^{-1})} = \frac{\det(\hat{\mathbf{A}})}{\det(\mathbf{A})},$$

using the trace normalization. \square

The Vecchia approximation uses the normalization $\text{tr}(\mathbf{A}\hat{\mathbf{A}}^{-1}) = n$, so Proposition 3.5 shows that the Vecchia determinant upper bounds the true determinant of \mathbf{A} . Note, however, the Vecchia approximation does not always provide the sharpest upper bound on the determinant of \mathbf{A} given the revealed entries: this is given by the solution to the *maximum entropy* problem; see [7, pg. 161].

3.4. Iterative determinant calculation. When the approximation $\det \mathbf{A} \approx \det \hat{\mathbf{A}}$ is not sufficiently accurate, we can refine the approximation with a multiplicative correction term. To that end, we rewrite

$$(3.3) \quad \log\left(\frac{\det \mathbf{A}}{\det \hat{\mathbf{A}}}\right) = \log(\det(\hat{\mathbf{A}}^{-1/2} \mathbf{A} \hat{\mathbf{A}}^{-1/2})) = \text{tr}(\log(\hat{\mathbf{A}}^{-1/2} \mathbf{A} \hat{\mathbf{A}}^{-1/2})).$$

Here, the matrix logarithm is defined by taking the logarithm of a matrix's eigenvalues. To approximate the error term (3.3), we sample independent length- \sqrt{n} vectors with uniformly random directions, $\mathbf{z}_1, \dots, \mathbf{z}_m$, and form a *stochastic trace estimator* [28]:

$$s_t = \frac{1}{t} \sum_{i=1}^t \mathbf{z}_i^* \log(\hat{\mathbf{A}}^{-1/2} \mathbf{A} \hat{\mathbf{A}}^{-1/2}) \mathbf{z}_i.$$

This estimator is unbiased and converges as $t \rightarrow \infty$.

The next proposition bounds the mean square error of the stochastic determinant estimate. The proof is in [Appendix B.3](#) and is apparently new.

PROPOSITION 3.6 (Stochastic determinant estimation). *If \mathbf{A} and $\hat{\mathbf{A}}$ are strictly positive definite $n \times n$ matrices with $\log(\kappa_{\text{Kap}}) \leq n$ and the matrix logarithm is computed exactly, the iterates produced by stochastic determinant estimation satisfy*

$$(3.4) \quad \mathbb{E} \left| \log \left(\frac{e^{s_t} \det \hat{\mathbf{A}}}{\det \mathbf{A}} \right) \right|^2 \leq \frac{4 \log(\kappa_{\text{Kap}})}{t}.$$

[Proposition 3.6](#) shows that the stochastic determinant estimator achieves a better mean square error than the direct estimator as soon as the number of iterations is $t \geq 4/\log(\kappa_{\text{Kap}})$. As the one limitation, we cannot really compute the matrix logarithm exactly. In practice, we approximate each term $\log(\hat{\mathbf{A}}^{-1/2} \mathbf{A} \hat{\mathbf{A}}^{-1/2}) \mathbf{z}_i$ by building a Krylov subspace

$$K_q(\mathbf{M}, \mathbf{z}_i) = \text{span}\{\mathbf{z}_i, \mathbf{M}\mathbf{z}_i, \dots, \mathbf{M}^q \mathbf{z}_i\} \quad \text{for } \mathbf{M} = \hat{\mathbf{A}}^{-1/2} \mathbf{A} \hat{\mathbf{A}}^{-1/2}$$

and applying an explicit polynomial approximation [26] or an implicit Krylov–Ritz approximation [12] of the matrix logarithm. Each Krylov subspace calculation requires q matrix multiplications, and we use $q = 100$ in the [section 5](#) experiments.

4. Optimization strategies. This section explores the following open question.

Question 4.1 (Optimal Vecchia approximation). What permutation \mathbf{P} and sparsity pattern $\{\mathbf{S}_i\}_{i=1}^n$ with bounded cardinality $|\mathbf{S}_i| \leq s$ lead to a Vecchia approximation with the lowest possible Kaporin condition number [\(3.1\)](#)?

This question is NP-hard to resolve exactly. Nonetheless, we can suggest several approximation schemes.

4.1. Partial Cholesky + diagonal. First we consider the simpler problem of generating an optimal partial Cholesky + diagonal approximation ([subsection 2.3](#)). When the pivot set is $\mathbf{R} = \{u_1, \dots, u_r\}$, the condition number can be rewritten as

$$(4.1) \quad \kappa_{\text{Kap}} = \prod_{d_{\mathbf{A}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{R}}) > 0} \frac{d_{\mathbf{A}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{R}})^2}{d_{\mathbf{A}}(\mathbf{e}_i, \text{span}(\{\mathbf{e}_j\}_{j < i} \cup \{\mathbf{e}_j\}_{j \in \mathbf{R}}))^2}.$$

The Kaporin condition number [\(4.1\)](#) only depends on the pivot set \mathbf{R} , so we write $\kappa_{\text{Kap}} = \kappa_{\text{Kap}}(\mathbf{R})$. Below we introduce several strategies to find the best pivot set. Only one of these strategies, called “farthest point sampling”, has been considered in the recent literature on Vecchia approximation [23, 31, 5].

4.1.1. Adaptive search. The *adaptive search* algorithm starts with an empty set $\mathbf{R} = \emptyset$ and then adds one pivot at a time by a rule that changes (“adapts”) based on the previous selections. At each stage, this algorithm tries all the pivots and chooses the pivot causing the greatest decrease in the condition number.

$$\mathbf{R} \leftarrow \mathbf{R} \cup \{i\}, \quad \text{where } i \in \underset{1 \leq j \leq n}{\text{argmin}} \kappa_{\text{Kap}}(\mathbf{R} \cup \{j\}) \quad [\text{adaptive search}].$$

The adaptive search is expensive, since each stage requires processing *all* the entries in the matrix. Thus, the cost of producing a cardinality- r pivot set is $\mathcal{O}(n^2 r)$ arithmetic operations.

4.1.2. Adaptive sampling. *Adaptive sampling* is a class of cheaper algorithms that sample new pivots from a distribution that prioritizes distance to the already chosen pivots. Here we describe four variations of adaptive sampling.

- In “randomly pivoted Cholesky” (RPC, [6]), each pivot is sampled as follows.

$$R \leftarrow R \cup \{i\}, \quad \text{with } \text{prob}(i) = \frac{d_{\mathbf{A}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in R})^2}{\sum_{j=1}^n d_{\mathbf{A}}(\mathbf{e}_j, \text{span}\{\mathbf{e}_j\}_{j \in R})^2}. \quad [\text{RPC}]$$

- In “square distance sampling” (SDS, [1]), we select a slightly different set of sampling probabilities.

$$R \leftarrow R \cup \{i\}, \quad \text{with } \text{prob}(i) = \frac{d_{\mathbf{A}}(\mathbf{e}_i, \{\mathbf{e}_j\}_{j \in R})^2}{\sum_{j=1}^n d_{\mathbf{A}}(\mathbf{e}_j, \{\mathbf{e}_j\}_{j \in R})^2}. \quad [\text{SDS}]$$

- In “column pivoted Cholesky” (CPC, [17]), we use a simpler rule for inducting new pivots.

$$R \leftarrow R \cup \{i\}, \quad \text{where } i \in \underset{1 \leq j \leq n}{\text{argmin}} d_{\mathbf{A}}(\mathbf{e}_j, \text{span}\{\mathbf{e}_j\}_{j \in R}). \quad [\text{CPC}]$$

- Last, in “farthest point sampling” (FPS, [15]), we apply the rule

$$R \leftarrow R \cup \{i\}, \quad \text{where } i \in \underset{1 \leq j \leq n}{\text{argmin}} d_{\mathbf{A}}(\mathbf{e}_j, \{\mathbf{e}_j\}_{j \in R}). \quad [\text{FPS}]$$

If there are ties in CPC or FPS, they can be broken arbitrarily. Each of the adaptive sampling algorithms is cheap compared to the overall cost of matrix approximation. In the context of partial Cholesky + diagonal, the cost of the partial Cholesky step is already $\mathcal{O}(nr^2)$ arithmetic operations, and adding an adaptive sampling rule requires an extra $\mathcal{O}(nr)$ operations.

Adaptive sampling theory shows how each algorithm approximately minimizes one of the following distance functionals.

$$(4.2) \quad \begin{aligned} \eta_{\text{RPC}}(R) &= \sum_{i=1}^n d_{\mathbf{A}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in R})^2, \\ \eta_{\text{SDS}}(R) &= \sum_{i=1}^n d_{\mathbf{A}}(\mathbf{e}_i, \{\mathbf{e}_j\}_{j \in R})^2, \\ \eta_{\text{CPC}}(R) &= \max_{1 \leq i \leq n} d_{\mathbf{A}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in R}), \text{ or} \\ \eta_{\text{FPS}}(R) &= \max_{1 \leq i \leq n} d_{\mathbf{A}}(\mathbf{e}_i, \{\mathbf{e}_j\}_{j \in R}). \end{aligned}$$

Each algorithm uses a pivot selection rule naturally related to (4.2), and the approximation accuracy is controlled by the following theorem.

THEOREM 4.2 (Adaptive sampling guarantees). *For any positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ and any cardinality $r \leq \text{rank}(\mathbf{A})$, the adaptive sampling algorithms generate a random or deterministic pivot set $R \subseteq \{1, \dots, n\}$ with cardinality $|R| = r$*

that satisfies

$$\begin{aligned} \frac{\mathbb{E}[\eta_{\text{RPC}}(\mathbf{R})]}{\min_{|\mathbf{S}| \leq r} \eta_{\text{RPC}}(\mathbf{S})} &\leq 2^{r-1}, & [\text{RPC}] \\ \frac{\eta_{\text{CPC}}(\mathbf{R})}{\min_{|\mathbf{S}| \leq r} \eta_{\text{CPC}}(\mathbf{S})} &\leq r!, & [\text{CPC}] \\ \frac{\mathbb{E}[\eta_{\text{SDS}}(\mathbf{R})]}{\min_{|\mathbf{S}| \leq r} \eta_{\text{SDS}}(\mathbf{S})} &\leq 5(\log r + 2), & [\text{SDS}] \\ \frac{\eta_{\text{FPS}}(\mathbf{R})}{\min_{|\mathbf{S}| \leq r} \eta_{\text{FPS}}(\mathbf{S})} &\leq 2. & [\text{FPS}] \end{aligned}$$

The adaptive sampling algorithms approximately minimize the objective functions η_{RPC} , η_{CPC} , η_{SDS} , and η_{FPS} that are defined in (4.2).

Proof. The RPC error bound is [6, Lem. 5.5]. The CPC error bound is implied by [8, Thm. 1] in the limit $p \rightarrow \infty$; see also [10, Thm. 2]. The SDS error bound is [1, Thm. 3.1], with an improved approximation factor due to [21, Lem. 4.1]. The FPS error bound is [15, Thm. 2.2]. \square

Theorem 4.2 provides approximation guarantees that are independent of the dimensionality n and only depend on the cardinality r of the index set. As a limitation, the approximation factors 2^{r-1} and $r!$ in **Theorem 4.2** increase rapidly with r . Yet there is some additional “oversampling” theory that addresses this issue. For example, randomly pivoted Cholesky with slightly more than r pivots controls the error functional nearly as well as the best approximation with r pivots [6, Thm. 2.3].

A second, more major limitation is that the adaptive sampling algorithms optimize the distance functionals η_{RPC} , η_{CPC} , η_{SDS} , and η_{FPS} , which are not the same as the Kaporin condition number κ_{Kap} (4.1). As such, it is no surprise that these algorithms lead to slightly worse approximation quality than adaptive search. However, they run much more quickly; see the experiments in section 5.

4.2. Adding to the sparsity pattern. Next we consider building a partial Cholesky + Vecchia approximation that optimizes the Kaporin condition number

$$(4.3) \quad \kappa_{\text{Kap}} = \prod_{d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j < i}) > 0} \frac{d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{S}_i})^2}{d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j < i})^2}, \quad \text{where } \tilde{\mathbf{A}} = \mathbf{P}^* \mathbf{A} \mathbf{P}.$$

We observe that each numerator term is the square of a distance functional

$$\delta_i(\mathbf{S}_i) = d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{S}_i}).$$

Therefore, we focus on adding q new indices to the index set \mathbf{S}_i that approximately minimize δ_i . Unfortunately, finding the best q indices is NP-hard, by a reduction to the sparse recovery problem studied in [22, Thm. 1]. However, the Vecchia literature suggests a couple heuristic optimization methods [16, 18].

4.2.1. Nearest neighbor search. *Nearest neighbor (NN) search* is the most common method for adding q new indices to a sparsity component [16]. In this method, we recursively find a new index $j \notin \mathbf{S}_i$ that achieves the smallest possible distance $d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \mathbf{e}_j)$ and we add it.

$$\mathbf{S}_i \leftarrow \mathbf{S}_i \cup \{j\}, \quad \text{where } j \in \underset{k \notin \mathbf{S}_i}{\text{argmin}} d_{\tilde{\mathbf{A}}}(\mathbf{e}_i, \mathbf{e}_k) \quad [\text{NN search}].$$

Finding q new indices in this way requires looking up $\mathcal{O}(n)$ entries of the matrix and performing $\mathcal{O}(n)$ arithmetic operations [4]. In special cases we can more efficiently calculate approximate nearest neighbors for all n components, by using advanced data structures [2]. However, even though the NN search can be made relatively efficient, there is no guarantee that it does a good job at optimizing the distance functional δ_i .

4.2.2. Orthogonal matching pursuit. *Orthogonal matching pursuit* (OMP) is an alternative approach for adding q new indices based on the following rule.

$$\mathbf{S}_i \leftarrow \mathbf{S}_i \cup \{j\}, \quad \text{where } j \in \underset{k \notin \mathbf{S}_i}{\operatorname{argmin}} \nu_i(\mathbf{S}_i \cup \{k\}) \quad [\text{OMP}].$$

OMP requires looking up $\mathcal{O}(qn)$ entries of the matrix and performing $\mathcal{O}(q^2n)$ additional arithmetic operations [18, App. C.1]. Since each optimization is expensive, one option is to restrict the search to c candidate points that come from a NN search or uniformly random sampling [18]. Then, the algorithm only requires $\mathcal{O}(qc)$ entry look-ups and $\mathcal{O}(q^2c)$ arithmetic operations where potentially $c \ll n$.

OMP is attractive because it directly targets the distance functional δ_i . However, the only available theoretical guarantees [27, Thms. B and C] require fast off-diagonal decay in the entries of \mathbf{A} , which seldom occurs in practice. Hence, in the worst case, OMP might result in a value δ_i that is incomparably worse than the best choice of q indices.

5. Experiments and analysis. In this section, we study the empirical performance of partial Cholesky + Vecchia approximations as applied to preconditioned conjugate gradient and log determinant estimation on Gaussian kernel matrices in real world kernel ridge regression problems.

5.1. Setup. Specifically, we work with 27 real-world datasets from libsvm and openML. Each dataset consists of a set of training vectors $\{\mathbf{x}_i\}_{i=1}^n \subseteq \mathbb{R}^d$, where we have scaled each coordinate so that the sample variance in each coordinate is 1, as well as additional test vectors whose coordinates we scaled in the same way. We use only the first $n = 20,000$ training vectors on any dataset with more than this due to VRAM constraints. The datasets are downloadable using scripts originally developed for [9], available at [this link](#).

On each dataset, we take $\mathfrak{K}(\mathbf{x}, \mathbf{x}') = \exp(-\frac{(\mathbf{x}-\mathbf{x}')^*(\mathbf{x}-\mathbf{x}')}{2d})$ our Gaussian kernel function, and let $\mathbf{K} \in \mathbb{R}^{n \times n}$ be such that $\mathbf{K}(i, j) = \mathfrak{K}(\mathbf{x}_i, \mathbf{x}_j)$. Finally, we take a nugget parameter and set $\Theta = \mathbf{K} + \mu \mathbf{I}$. For each dataset, Θ is the matrix we attempt to approximate. We test the quality of the approximation on the methods discussed in [section 3](#) for both approximate application of the inverse and log determinant estimation.

5.1.1. Preconditioned Conjugate Gradient. For preconditioned conjugate gradient, our goal is to approximate $\Theta^{-1} \mathbf{z}_{\text{test}}$ for \mathbf{x}_{test} a random test vector withheld from the training vectors (but still normalized according to the sample variance of the training vectors in each coordinate) and $\mathbf{z}_{\text{test}}(i) = \mathfrak{K}(\mathbf{x}_i, \mathbf{x}_{\text{test}})$ using the method described in [subsection 3.2](#). For our experiments, we choose five random \mathbf{x}_{test} from each dataset and run PCG against \mathbf{z}_{test} , starting from the zero vector. For each dataset and preconditioner, we record the number of iterations required so that $\|\mathbf{z}_{\text{test}} - \Theta \hat{\mathbf{w}}^{(k)}\|_2 \leq 10^{-4} \|\mathbf{z}_{\text{test}}\|_2$, where $\hat{\mathbf{w}}^{(k)}$ is PCG's approximation of $\Theta^{-1} \mathbf{z}_{\text{test}}$ at the k th iteration, at which point we consider the problem "solved". We repeat this for nugget values 10^{-10} , 10^{-6} , and 10^{-3} . We then plot, as a function of the number of iterations

k , the percentage of kernel ridge regression problems "solved" by iteration k by each preconditioner.

Note that this is in contrast to the usual procedure for kernel ridge regression, where one forms a vector \mathbf{y} whose i th coordinate $\mathbf{y}(i)$ is the label associated to data point \mathbf{x}_i for all $i \in 1 \dots n$ and computes $\mathbf{y}^* \Theta^{-1}$. Both aim to eventually compute $\mathbf{y}^* \Theta^{-1} \mathbf{z}_{\text{test}}$, but computing $\mathbf{y}^* \Theta^{-1}$ first requires only one application of Θ^{-1} even if there are multiple points $\mathbf{z}_{\text{test}}^{(1)}, \mathbf{z}_{\text{test}}^{(2)}, \dots$ whose labels we want to predict. However, we find empirically that PCG with Vecchia preconditioners requires orders of magnitude fewer iterations when computing $\Theta^{-1} \mathbf{z}_{\text{test}}$ as opposed to $\mathbf{y}^* \Theta^{-1}$ to attain the same 10^{-4} relative error tolerance, especially for small nugget; see Figure 8. While we also observe improvement using partial Cholesky + Vecchia or partial Cholesky + diagonal over existing methods when approximating $\mathbf{y}^* \Theta^{-1}$, we still recommend computing $\Theta^{-1} \mathbf{z}_{\text{test}}$ for each test point as opposed to $\mathbf{y}^* \Theta^{-1}$ once, unless one must estimate labels at several hundred or more test points.

5.1.2. Log-determinant estimation. For log-determinant estimation, we implement both direct and indirect estimators as described in subsection 3.3 and subsection 3.4 respectively. For the indirect estimator, we use a Krylov subspace depth of 100 and estimate the log determinant as the mean of 10 samples. For our figures, we only plot data from the large nugget ($\mu = 10^{-3}$ case, as it was the only case where any tractable preconditioner we tested attained reasonably accurate estimates (though we present data for other nuggets in tables in Table 2, Table 3, and Table 4. After obtaining the direct estimators and the indirect estimates, we obtained the absolute error of the estimates compared to the ground truth log determinant on each dataset (computed with a GPU), then presented the median of this error over the 27 datasets. Note that we are plotting error in the log determinant, and the y-axis is also in log scale, so the y axes are akin to log-of-log scale of absolute error in the determinant.

Importantly, while the indirect estimator is unbiased if the matrix logarithm is computed exactly, this is not a reasonable assumption in practice for all but the best preconditioners, at least with a feasible Krylov depth (in our case, 100). So while Proposition 3.6 focuses mainly on bounding the variance, the error shown in the figures mostly reflects bias.

5.2. Comparison of Pivot Choosers. In this section, we compare the performance of the four different adaptive sampling pivot choosers from subsection 4.1.2 when used for partial Cholesky + diagonal and partial Cholesky + Vecchia preconditioners, both in preconditioned conjugate gradient and log determinant estimation. In figures, we refer to them by their respective three-letter abbreviations from subsection 4.1.2. We also compare them to the adaptive search preconditioner in subsection 4.1.1, which we made tractable for $20,000 \times 20,000$ matrices despite its $O(n^2 r)$ runtime by writing a custom GPU kernel for the algorithm using the package `CUDA.jl` and running it on an RTX 3090. We refer to adaptive search as "AS" in figures.

For the partial Cholesky portion, we always take $r = \lfloor \sqrt{n} \rfloor$ as the rank of the partial Cholesky decomposition. For the Vecchia portion, we take $s = \lfloor n^{1/4} \rfloor$ as the maximum number of off-diagonal nonzeros per column. We always use greedy conditional selection to choose the sparsity pattern for partial Cholesky + Vecchia, and always with $c = 10s$ candidates.

For PCG, we plot the results using the procedure described in subsection 5.1.1 in Figure 5. We observe that among adaptive sampling methods, RPC and SDS, which perform similarly, are significantly better than CPC and FPS. Unsurprisingly,

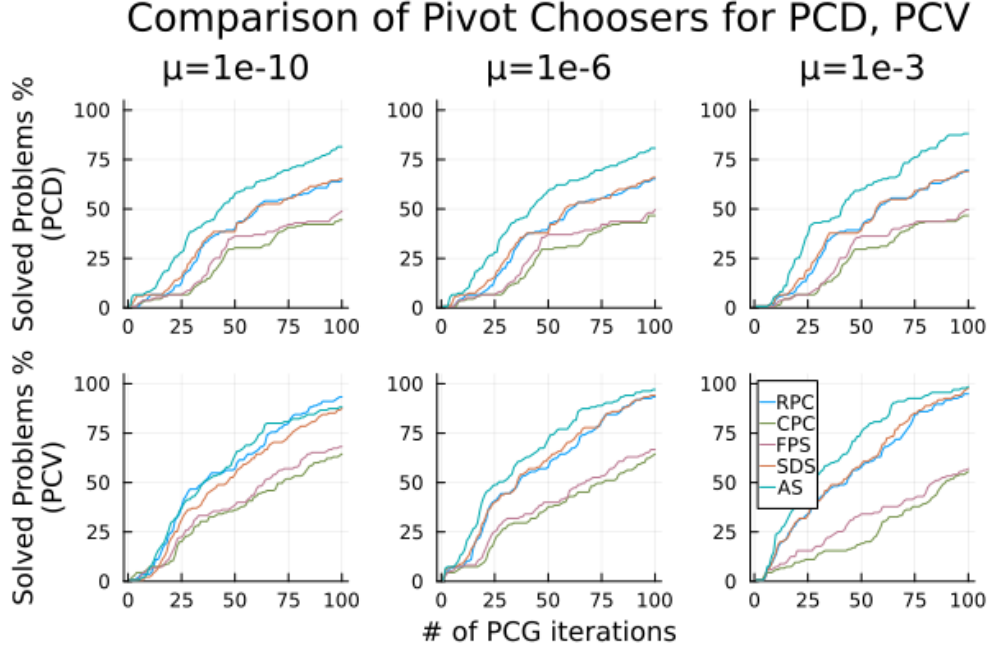


FIG. 5. Comparison of pivot choosers used in partial Cholesky + diagonal and partial Cholesky + Vecchia preconditioners for conjugate gradient

adaptive search is generally better than adaptive sampling methods, and the gap is significant for partial Cholesky + diagonal, suggesting the Kaporin condition number could reasonably be used as a heuristic to improve existing pivot choosing methods. However, the gap is nearly insignificant when applied to partial Cholesky + Vecchia, providing evidence that Kaporin-inspired pivot choosers couldn't hope to improve partial Cholesky + Vecchia approximations much over existing pivot choosers. We hypothesize that this is because Kaporin-oriented greedy conditional selection "picks up the slack" left by adaptive sampling methods when choosing the sparsity pattern. Finally, we observe that relative performance is not significantly affected by varying the nugget parameter.

For log determinant estimation, we give a bar chart using the procedure described in [subsection 5.1.2](#) (in particular, this data is only for $\mu = 10^{-3}$) in [Figure 6](#). We again observe that AS outperforms RPC and SDS, which outperform CPC and FPS. The gap between AS and RPC or SDS, especially for partial Cholesky + Vecchia, is more noticeable than with PCG, suggesting the Kaporin condition number is perhaps a better pivot choosing heuristic when applied to log determinant estimation than PCG..

5.3. How much does Vecchia help?. To study the improvement of partial Cholesky + diagonal and/or partial Cholesky + Vecchia over previous methods, we implement [\[11\]](#)'s preconditioner ("Frangella" in figures), [\[9\]](#)'s preconditioner ("Diaz" in figures), and compare them to Randomly Pivoted Cholesky + diagonal, Randomly Pivoted Cholesky + Vecchia with $s = \lfloor n^{1/4} \rfloor$ nonzero entries per column, and Randomly Pivoted Cholesky + Vecchia with $s = \lfloor n^{1/3} \rfloor$ nonzero entries per column

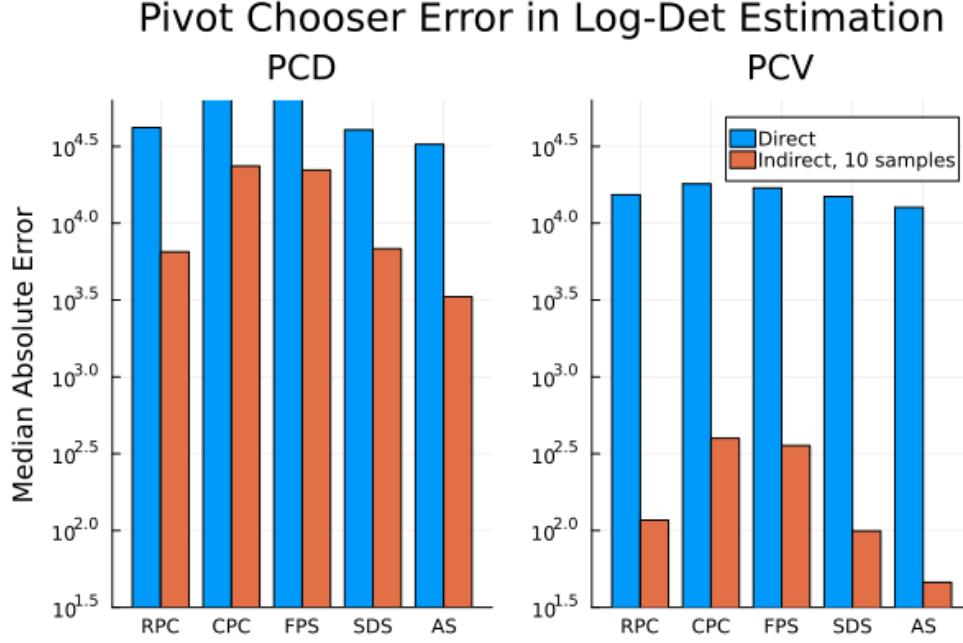


FIG. 6. Comparison of pivot choosers used in partial Cholesky + diagonal and partial Cholesky + Vecchia estimators or preconditioners for log determinant estimation.

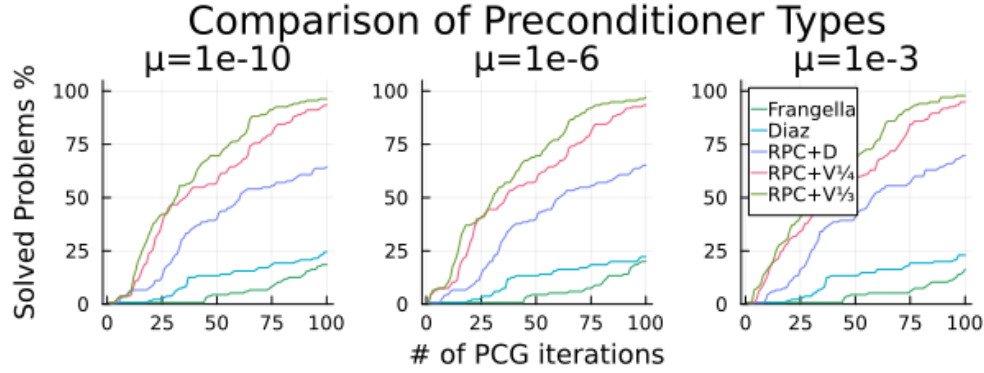


FIG. 7. Comparison of [11]’s and [9]’s RPC-based preconditioners (which use a Kaporin-suboptimal diagonal), RPC plus diagonal, RPC plus Vecchia with $s = n^{1/4}$, and RPC plus Vecchia with $s = n^{1/3}$, as applied to PCG on 5 kernelized test vectors per dataset as described in [subsection 5.1.1](#).

(“RPC+D”, “RPC+V^{1/4}”, and “RPC+V^{1/3}” respectively), each with $c = 10s$ GCS candidates.

For PCG applied to the kernelized test vectors, we observe in [Figure 7](#) that simply moving from [11]’s or [9]’s diagonal heuristics to the Kaporin-optimal choice of diagonal provided by partial Cholesky + Vecchia approximations makes a dramatic difference, more than doubling the number of problems solved within 100 PCG it-

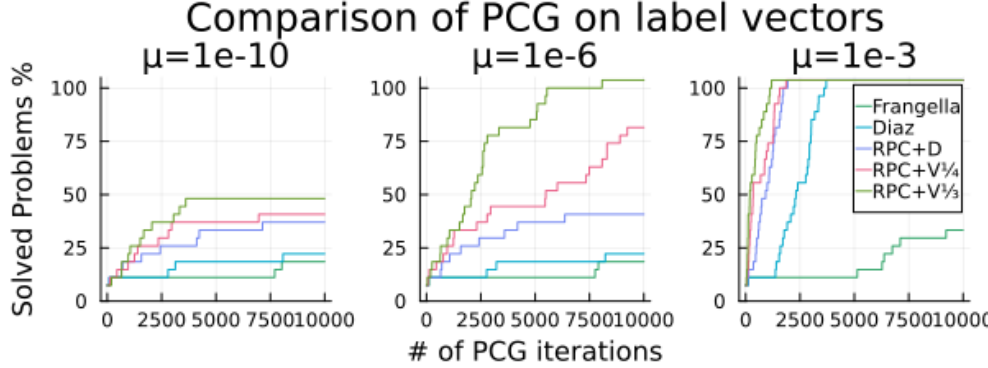


FIG. 8. Comparison of [11]’s and [9]’s RPC-based preconditioners (which use a Kaporin-suboptimal diagonal), RPC plus diagonal, RPC plus Vecchia with $s = n^{1/4}$, and RPC plus Vecchia with $s = n^{1/3}$, as applied to PCG on the vector of labels for each dataset.

erations. As both strategies require the same computational complexity, this result strongly favors partial Cholesky + diagonal approximations. Beyond this, adding $s = n^{1/4}$ entries to the sparsity pattern results in a considerable benefit, solving all but two problems within 100 iterations, but there is only modest benefit from increasing s to $n^{1/3}$. We also observe that changing the nugget parameter does not substantially affect the performance of any of these methods, which means they can be used in a variety of settings.

We also examine PCG applied to the vector of labels in Figure 8, where the results are considerably different. Although the relative order of the methods remains unchanged, performance is worse across the board compared to PCG on the kernelized test vectors. Moreover, the results are highly nugget-dependent. In the extremely small nugget regime, even partial Cholesky + Vecchia with $n = s^{1/3}$ fails to solve half the problems within 10000 iterations. As the nugget increases, all methods perform better, and the significance of choosing a larger sparsity set first grows as more problems are able to be solved, then shrinks as simpler Vecchia preconditioners also become capable of solving the problems.

For log determinant estimation (Figure 9), the relative order remains, except that [9]’s estimator slightly outperforms RPC+D. We suspect this is because [9]’s separation of its approximation of \mathbf{K} from its addition of the nugget is beneficial in large nugget regimes like $\mu = 10^{-3}$. However, noting the logarithmic scale of the figure, the only methods that could reasonably be usable for most tasks depending on log determinant estimates (e.g. approximating likelihood ratios) are partial Cholesky + Vecchia. We also observe that the indirect estimator is essentially necessary to obtain a reasonably small error, even for partial Cholesky + Vecchia.

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

REFERENCES

- [1] D. ARTHUR AND S. VASSILVITSKII, *k-means++: The advantages of careful seeding*, in Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, 2007, <https://dl.acm.org/doi/10.5555/1283383.1283494>.

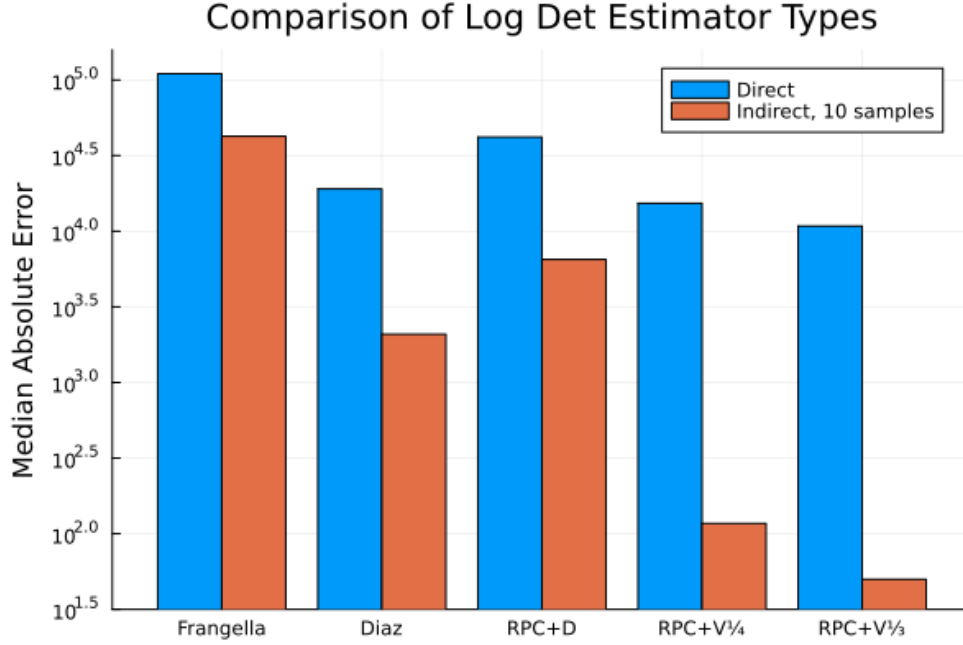


FIG. 9. Comparison of [11]’s and [9]’s RPC-based preconditioners (which use a Kaporin-suboptimal diagonal), RPC plus diagonal, RPC plus Vecchia with $s = n^{1/4}$, and RPC plus Vecchia with $s = n^{1/3}$, as applied to log determinant estimation.

- [2] S. ARYA, D. M. MOUNT, N. S. NETANYAHU, R. SILVERMAN, AND A. Y. WU, *An optimal algorithm for approximate nearest neighbor searching fixed dimensions*, Journal of the ACM, 45 (1998), p. 891–923, <https://doi.org/10.1145/293347.293348>.
- [3] O. AXELSSON AND I. KAPORIN, *On the sublinear and superlinear rate of convergence of conjugate gradient methods*, Numerical Algorithms, 25 (2000), p. 1–22, <https://doi.org/10.1023/a:1016694031362>.
- [4] M. BLUM, R. W. FLOYD, V. PRATT, R. L. RIVEST, AND R. E. TARJAN, *Time bounds for selection*, Journal of Computer and System Sciences, 7 (1973), pp. 448–461, [https://doi.org/10.1016/S0022-0000\(73\)80033-9](https://doi.org/10.1016/S0022-0000(73)80033-9).
- [5] D. CAI, E. CHOW, AND Y. XI, *Posterior covariance structures in gaussian processes*, SIAM Journal on Matrix Analysis and Applications, 46 (2025), pp. 1640–1673, <https://doi.org/10.1137/24M1684918>.
- [6] 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>.
- [7] A. P. DEMPSTER, *Covariance selection*, Biometrics, 28 (1972), pp. 157–175, <https://doi.org/https://doi.org/10.2307/2528966>.
- [8] A. DESHPANDE AND K. VARADARAJAN, *Sampling-based dimension reduction for subspace approximation*, in Proceedings of the Thirty-Ninth Annual ACM Symposium on Theory of Computing, 2007, <https://doi.org/10.1145/1250790.1250884>.
- [9] 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>.
- [10] H. ENGLER, *The behavior of the QR-factorization algorithm with column pivoting*, Applied Mathematics Letters, 10 (1997), pp. 7–11, [https://doi.org/10.1016/S0893-9659\(97\)00098-0](https://doi.org/10.1016/S0893-9659(97)00098-0).
- [11] Z. FRANGELLA, J. A. TROPP, AND M. UDELL, *Randomized nyström preconditioning*, SIAM Journal on Matrix Analysis and Applications, 44 (2023), pp. 718–752, <https://doi.org/10.1137/21M1466244>.

	direct	331.845	-6.24e+04	-1.1253e+05	-1.1754e+05	-3.1688e+04	-3.8907e+04	-7.6956e+04	-6.7110e+04	-1.0900e+05	-1.1035e+05	-1.1344e+05	-1.1280e+05	-1.3681e+05
ACSLincome	mean	-2.3698e+04	-1.0872e+05	-1.2515e+05	-1.2725e+05	-9.0708e+04	-1.1251e+05	-1.0836e+05	-1.2486e+05	-1.2529e+05	-1.2529e+05	-1.2517e+05	-1.2517e+05	-1.2517e+05
	stdev	681.065	631.001	190.929	126.609	657.369	780.049	414.005	430.400	296.647	249.839	242.794	274.351	524.556
AirlinesDepDelayLM (Truth: -1.2863e+05)	direct	369.476	-6.4580e+04	-1.1333e+05	-1.1586e+05	-4.6907e+04	-5.3631e+04	-7.3915e+04	-6.6741e+04	-1.1060e+05	-1.1170e+05	-1.1449e+05	-1.1373e+05	-1.3681e+05
	mean	-1.0232e+04	-1.2912e+05	-1.2824e+05	-1.2860e+05	-1.0671e+05	-1.1120e+05	-1.1177e+05	-1.1346e+05	-1.2824e+05	-1.2828e+05	-1.2828e+05	-1.2828e+05	-1.2828e+05
stdev	1401.108	544.001	248.538	186.685	1046.949	426.923	282.680	699.750	141.401	192.545	346.960	274.493	189.154	
COMET.MC.SAMPLE	direct	544.274	-1.2677e+05	-1.3512e+05	-1.3580e+05	-9.5746e+04	-9.8808e+04	-1.3035e+05	-1.2811e+05	-1.3482e+05	-1.3486e+05	-1.3526e+05	-1.3520e+05	-1.3711e+05
	mean	-1.3159e+05	-1.3631e+05	-1.3654e+05	-1.3652e+05	-1.3610e+05	-1.3604e+05	-1.3618e+05	-1.3617e+05	-1.3653e+05	-1.3653e+05	-1.3654e+05	-1.3655e+05	-1.3643e+05
stdev	307.150	-1.3533e+04	-5.7627e+04	-6.1177e+04	-1.2022e+04	-1.2008e+04	-3.6166e+04	-3.5332e+04	-1.0608e+04	-1.0608e+04	-8.8816e+04	-8.7023e+04	-1.3682e+05	
CreditPredictionsSmall (Truth: -1.3017e+05)	direct	365.986	-8.2072e+04	-1.2416e+05	-1.2674e+05	-3.7946e+04	-3.9718e+04	-8.8757e+04	-8.7895e+04	-1.2192e+05	-1.2210e+05	-1.2344e+05	-1.2344e+05	-1.3687e+05
	mean	1045.022	633.831	616.038	273.716	590.336	654.686	553.356	714.033	511.948	362.210	288.936	590.005	666.023
stdev	1869.905	467.448	190.691	142.414	614.345	130.088	468.703	473.271	218.985	160.347	136.979	148.108	313.733	
HIGGS (Truth: -8.9494e+04)	direct	307.176	-3.5338e+04	-5.7627e+04	-6.1177e+04	-1.2022e+04	-1.2008e+04	-3.9466e+04	-3.5332e+04	-1.0608e+04	-1.0608e+04	-8.8816e+04	-8.7023e+04	-1.3682e+05
	mean	-1.2325e+04	-7.7894e+04	-8.9102e+04	-8.9706e+04	-5.3704e+04	-5.3657e+04	-8.3727e+04	-7.7353e+04	-8.7082e+04	-8.6704e+04	-8.8806e+04	-8.8816e+04	-8.4930e+04
stdev	1045.022	633.831	616.038	273.716	590.336	654.686	553.356	714.033	511.948	362.210	288.936	590.005	666.023	
MNIST (Truth: -1.1774e+05)	direct	334.587	-3.2156e+04	-5.1423e+04	-5.7070e+04	-4.4453e+05	-4655.401	-3.7676e+04	-3.1768e+04	-4.7110e+04	-4.7076e+04	-5.3137e+04	-5.1109e+04	-1.3686e+05
	mean	-4.6163e+04	-6.7671e+04	-7.0694e+04	-7.0737e+04	-4.8851e+04	-4.8573e+04	-7.0307e+04	-6.6813e+04	-7.0309e+04	-7.0301e+04	-7.0606e+04	-7.0500e+04	-6.8281e+04
stdev	494.214	388.971	231.871	216.753	540.955	584.384	361.600	310.046	339.559	386.977	156.285	166.022	707.193	
Medical Appointment (Truth: -1.1774e+05)	direct	344.702	-4.8448e+04	-9.5010e+04	-1.0252e+05	-3.2529e+04	-3.7509e+04	-5.5835e+04	-4.2997e+04	-9.1787e+04	-9.2628e+04	-9.6599e+04	-9.5284e+04	-1.3678e+05
	mean	-2.4468e+04	-9.4161e+04	-1.1734e+05	-1.1765e+05	-6.7118e+04	-7.5690e+04	-1.0070e+05	-9.4446e+04	-1.1638e+05	-1.1684e+05	-1.1751e+05	-1.1733e+05	-1.1087e+05
stdev	735.489	409.146	294.861	264.997	665.033	593.038	565.881	411.805	211.164	216.395	362.523	346.267	602.909	
MiniBooNE (Truth: -1.3493e+05)	direct	434.227	-1.1491e+05	-1.2906e+05	-1.3132e+05	-8.4264e+04	-9.2031e+04	-1.2033e+05	-1.1643e+05	-1.2681e+05	-1.2987e+05	-1.2924e+05	-1.2924e+05	-1.3714e+05
	mean	-1.3090e+05	-1.3439e+05	-1.3494e+05	-1.3492e+05	-1.3394e+05	-1.3421e+05	-1.3460e+05	-1.3447e+05	-1.3497e+05	-1.3488e+05	-1.3493e+05	-1.3487e+05	-1.3479e+05
stdev	1535.619	376.637	177.803	114.152	303.253	331.832	258.444	251.644	136.803	171.195	174.572	214.645	141.837	
YearPredictionMSD (Truth: -7.2344e+04)	direct	300.073	-3.3351e+04	-4.8690e+04	-5.4392e+04	-1.2935e+04	-1.2945e+04	-4.2656e+04	-3.4734e+04	-4.0707e+04	-4.0708e+04	-5.3255e+04	-4.9211e+04	-1.3688e+05
	mean	-3.6243e+04	-6.8782e+04	-7.4172e+04	-7.5143e+04	-4.8014e+04	-4.8183e+04	-7.3891e+04	-7.0087e+04	-7.1299e+04	-7.1232e+04	-7.5190e+04	-7.4431e+04	-7.2954e+04
stdev	464.650	325.411	444.036	294.478	320.207	451.839	245.356	507.840	519.097	238.316	407.982	313.401	1160.069	
a9a	direct	300.617	-6.7051e+04	-6.9396e+04	-7.0535e+04	-1.0181e+04	-1.0282e+04	-4.3540e+04	-3.7320e+04	-4.6430e+04	-6.3690e+04	-7.1366e+04	-6.9154e+04	-1.3800e+05
	mean	-2.6995e+04	-7.7022e+04	-9.0674e+04	-9.0963e+04	-3.8335e+04	-4.0240e+04	-8.6246e+04	-7.8684e+04	-8.8259e+04	-8.8514e+04	-9.0766e+04	-9.0681e+04	-8.6520e+04
stdev	525.440	730.245	436.426	214.374	451.557	367.468	531.301	658.115	408.269	336.619	329.391	271.758	1096.337	
cadata	direct	396.911	-8.8913e+04	-1.0767e+05	-1.0886e+05	-4.8752e+04	-6.0192e+05	-9.3433e+04	-8.8055e+04	-1.0636e+05	-1.0672e+05	-1.0778e+05	-1.0758e+05	-1.1298e+05
	mean	-6.9036e+04	-1.0823e+05	-1.1062e+05	-1.1090e+05	-1.0767e+05	-1.0886e+05	-1.0849e+05	-1.1060e+05	-1.1063e+05	-1.1063e+05	-1.1062e+05	-1.1060e+05	-1.1026e+05
stdev	1015.349	338.949	108.997	82.596	698.542	479.771	191.300	334.357	108.104	142.108	109.371	106.229	46.488	
cod-rna	direct	412.645	-9.3538e+04	-1.2893e+05	-1.3092e+05	-5.6886e+04	-7.0179e+04	-1.0432e+05	-9.7116e+04	-1.2767e+05	-1.2818e+05	-1.2942e+05	-1.2915e+05	-1.3688e+05
	mean	-8.2082e+04	-1.2722e+05	-1.3227e+05	-1.3331e+05	-1.2905e+05	-1.2758e+05	-1.2791e+05	-1.2647e+05	-1.3340e+05	-1.3328e+05	-1.3325e+05	-1.3327e+05	-1.3327e+05
stdev	1693.477	430.720	173.391	140.765	489.192	194.441	210.147	280.583	141.838	168.978	149.338	106.182	110.019	
connect-4 (Truth: -9.6767e+04)	direct	136.692	-3.0312e+04	-7.3114e+04	-8.0005e+04	-2.3851e+04	-2.6701e+04	-1.8105e+04	-3.8151e+04	-6.9031e+04	-6.9733e+04	-7.4614e+04	-7.3558e+04	-1.3680e+05
	mean	-3.0513e+04	-7.9013e+04	-9.6332e+04	-9.6718e+04	-5.7405e+04	-6.1399e+04	-8.7599e+04	-8.1856e+04	-9.5190e+04	-9.5873e+04	-9.6722e+04	-9.6508e+04	-9.2736e+04
stdev	519.676	439.614	456.923	377.811	491.442	551.741	474.423	759.840	444.849	509.929	383.780	299.657	669.897	
covtype.binary	direct	544.444	-7.3366e+04	-1.2538e+05	-1.2765e+05	-4.4486e+04	-5.5897e+04	-8.5953e+04	-7.4830e+04	-1.2136e+05	-1.2477e+05	-1.2858e+05	-1.2541e+05	-1.3682e+05
	mean	-4.2024e+04	-1.1700e+05	-1.2984e+05	-1.2981e+05	-1.0074e+05	-1.1224e+05	-1.1918e+05	-1.1657e+05	-1.2992e+05	-1.2986e+05	-1.2981e+05	-1.2984e+05	-1.2776e+05
stdev	180.154	167.341	164.408	162.286	161.655	161.655	161.655	161.655	161.655	161.655	161.655	161.655	161.655	
creditcard	direct	341.341	-5.7296e+04	-9.0957e+04	-9.6405e+04	-5.475.087	-8025.005	-6.1002e+04	-4.9405e+04	-8.5242e+04	-8.5756e+04	-9.3187e+04	-9.1324e+04	-1.3685e+05
	mean	-3.3638e+04	-9.1483e+04	-1.0921e+05	-1.0942e+05	-3.9206e+04	-4.2008e+04	-1.0021e+05	-9.1323e+04	-1.0755e+05	-1.0755e+05	-1.0938e+05	-1.0938e+05	-9.9553e+04
stdev	919.532	551.257	353.856	201.626	2067.358	942.980	497.463	386.728	230.956	190.024	187.100	257.150	289.465	
diamonds (Truth: -1.3378e+05)	direct	120.947	-9.1323e+04	-1.2907e+05	-1.3113e+05	-5.3691e+04	-6.3989e+04	-1.0174e+05	-9.8005e+04	-1.2727e+05	-1.2831e+05	-1.2927e+05	-1.2912e+05	-1.3601e+05
	mean	-9.1277e+04	-1.2995e+05	-1.3373e+05	-1.3370e+05	-1.2550e+05	-1.2790e+05	-1.3038e+05	-1.2905e+05	-1.3373e+05	-1.3358e+05	-1.3370e+05	-1.3370e+05	-1.3330e+05
stdev	1523.201	459.712	147.119	95.742	855.481	786.498	320.502	374.067	202.342	163.998	96.625	95.143	143.682	
hlsfm.lhc_jets_hlf (Truth: -1.3455e+05)	direct	458.793	-1.0550e+05	-1.2988e+05	-1.3207e+05	-8.2491e+04	-9.0911e+04	-1.1175e+05	-1.0578e+05	-1.2966e+05	-1.2941e+05	-1.3066e+05	-1.2982e+05	-1.3700e+05
	mean	-1.2962e+05	-1.3324e+05	-1.3458e+05	-1.3454e+05	-1.3289e+05	-1.3309e+05	-1.3330e+05	-1.3270e+05	-1.3453e+05	-1.3452e+05	-1.3460e+05	-1.3457e+05	-1.3430e+05
stdev	1642.532	309.136	63.948	122.235	651.027	431.501	194.544	267.382	211.056	136.712	126.157	114.718	111.022	
ijcnn1 (Truth: -1.2838e+05)	direct	383.848	-6.1634e+04	-1.1393e+05	-1.1832e+05	-2.9214e+04	-3.6695e+04	-7.7972e+04	-6.2816e+04	-1.1238e+05	-1.1273e+05	-1.1420e+05	-1.1389e+05	-1.3680e+05
	mean	-3.8573e+04	-1.0553e+05	-1.2836e+05	-1.2706e+05	-7.1030e+04	-8.2484e+04	-1.1340e+05	-1.0493e+05	-1.2371e+05	-1.2373e+05	-1.2387e+05	-1.2382e+05	-1.1926e+05
stdev	643.282	679.864	262.984	162.727	962.454	528.583	220.525	614.284	242.612	112.757	114.162	232.451	291.496	
jannis	direct	296.254	-2.7846e+04	-4.1473e+04	-4.5695e+04	-1.5077e+04	-1.6102e+04	-3.1524e+04	-2.7863e+04	-3.6655e+04	-3.7092e+04	-4.2590e+04	-4.1157e+04	-1.3684e+05
	mean	-5.3977e+04	-6.5018e+04	-6.7421e+04	-6.7706e+04	-5.8763e+04	-5.9286e+04	-6.6624e+04	-6.5532e+04	-6.6918e+04	-6.6830e+04	-6.7887e+04	-6.7549e+04	-6.7112e+04
stdev	180.737	167.864	162.984	162.527	162.527	162.527	162.527	162.527	162.527	162.527	162.527	162.527	162.527	
phishing (Truth: -2.0079e+04)	direct	151.107	-7.411.038	-1.5058e+04	-1.6369e+04	-4.627.919	-4866.678	-8934.265	-7484.812	-1.4249e+04	-1.4398e+04	-1.5353e+04	-1.5978e+04	-3.4623e+04
	mean	-8713.223	-1.8508e+04	-2.0118e+04	-2.0123e+04	-1.2770e+04	-1.2991e+04	-1.9654e+04	-1.8633e+04	-2.0051e+04	-1.9908e+04	-2.0042e+04	-2.0042e+04	-1.9272e+04
stdev	141.840	366.												

name	quantity	Frangella	RPC+D	RPC+V $\frac{1}{4}$	RPC+V $\frac{1}{2}$	CPC+D	FPS+D	GKL+D	SDS+D	CPC+V $\frac{1}{4}$	FPS+V $\frac{1}{4}$	GKL+V $\frac{1}{4}$	SDS+V $\frac{1}{4}$	Diss	
ACIncome	direct	334.925	-6.699e+04	-1.442e+05	-1.632e+05	-3.181e+04	-8.319e+04	-6.781e+04	-1.351e+05	-1.368e+05	-1.515e+05	-1.448e+05	-2.739e+05		
	mean	326.924	-7.199e+04	-1.735e+05	-1.975e+05	-3.461e+04	-9.251e+04	-7.219e+04	-1.594e+05	-1.638e+05	-1.849e+05	-1.721e+05	-2.947e+05		
	stdv	317.645	-1.849e8	-3.664e8	-444.975	-248.157	-183.212	-270.306	-299.461	-317.979	-458.527	-357.319	-479.340	-948.544	
AirlinesDepDelayLM	direct	369.589	-6.674e+04	-1.309e+05	-1.484e+05	-7.729e+04	-5.609e+04	-7.793e+04	-9.960e+04	-1.211e+05	-1.262e+05	-1.392e+05	-1.320e+05	-2.739e+05	
	mean	-212.74	-6.972e+04	-1.557e+05	-1.899e+05	-4.949e+04	-5.609e+04	-8.151e+04	-1.027e+05	-1.145e+05	-1.175e+05	-1.268e+05	-1.580e+05	-2.426e+05	
	stdv	128.847	-177.521	-230.649	-260.602	-224.668	-270.464	-257.910	-174.131	-227.897	-395.783	-511.491	-495.315	-902.888	
COMET.MC.SAMPLE	direct	547.855	-2.034e+05	-2.659e+05	-2.686e+05	-1.029e+05	-2.502e+05	-2.242e+05	-2.642e+05	-2.649e+05	-2.670e+05	-2.663e+05	-2.742e+05		
	mean	-8.320e+04	-2.451e+05	-2.707e+05	-2.710e+05	-1.058e+05	-1.124e+05	-2.595e+05	-2.523e+05	-2.702e+05	-2.765e+05	-2.708e+05	-2.707e+05	-2.287e+05	
	stdv	2.0270e-04	404.252	170.426	97.918	1205.659	1883.910	171.093	149.794	128.086	158.429	148.999	252.683	954.794	
Click.prediction.small	direct	366.155	-8.801e+04	-1.958e+05	-2.118e+05	-3.828e+04	-4.000e+04	-1.258e+05	-1.000e+05	-1.899e+05	-1.850e+05	-2.022e+05	-1.983e+05	-2.740e+05	
	mean	-3053.860	-9.008e+04	-2.210e+05	-2.350e+05	-3.959e+04	-4.137e+04	-1.377e+05	-1.034e+05	-2.076e+05	-2.089e+05	-2.275e+05	-2.238e+05	-2.510e+05	
	stdv	278.677	165.415	317.801	539.764	127.587	74.074	265.943	183.803	276.564	357.236	384.792	350.907	580.354	
HIGGS	direct	307.249	-3.593e+04	-5.902e+04	-6.389e+04	-1.224e+04	-1.395e+04	-3.982e+04	-3.772e+04	-1.110e+04	-3.106e+04	-3.960e+04	-8.847e+04	-2.700e+05	
	mean	-3.2567e+04	-6.855e+04	-9.3617e+04	-9.6118e+04	-4.3267e+04	-4.2890e+04	-7.745e+04	-6.952e+04	-8.7790e+04	-8.8048e+04	-9.4790e+04	-9.4025e+04	-1.0906e+05	
	stdv	678.839	513.307	323.742	498.868	757.878	619.386	297.815	251.946	682.111	425.105	349.786	611.421	1692.906	
MNIST	direct	332.968	-3.173e+04	-5.212e+04	-5.690e+04	-4.475e+04	-4.685.436	-3.799e+04	-3.212e+04	-4.803e+04	-4.803e+04	-5.115e+04	-5.210e+04	-2.704e+05	
	mean	-4.2797e+04	-6.404e+04	-7.3831e+04	-7.4447e+04	-4.551e+04	-4.5305e+04	-7.0893e+04	-6.4710e+04	-7.3053e+04	-7.3109e+04	-7.4554e+04	-7.3906e+04	-8.2306e+04	
	stdv	391.498	457.497	308.236	260.081	683.829	833.240	504.378	506.625	254.668	301.245	409.419	319.472	207.463	
Medical.Appointment	direct	344.779	-4.838e+04	-1.038e+05	-1.1648e+05	-3.2677e+04	-3.7706e+04	-5.6973e+04	-4.8880e+04	-9.9462e+04	-1.0651e+05	-1.0676e+05	-1.0405e+05	-2.7397e+05	
	mean	-7862.124	-5.8412e+04	-1.3380e+05	-1.5332e+05	-4.0462e+04	-4.5949e+04	-6.9628e+04	-5.8998e+04	-1.2652e+05	-1.2790e+05	-1.4007e+05	-1.3443e+05	-2.0457e+05	
	stdv	347.785	342.926	298.101	649.913	284.693	156.877	268.785	323.524	622.214	369.176	397.115	343.315	913.384	
MiniBooNE	direct	436.213	-1.3317e+05	-1.7701e+05	-1.9326e+05	-8.6621e+04	-9.5535e+04	-1.5398e+05	-1.4180e+05	-1.6236e+05	-1.6352e+05	-1.8781e+05	-1.8183e+05	-2.7433e+05	
	mean	-1.0970e+04	-2.885.110	303.743	375.604	634.117	319.867	274.811	348.341	458.840	471.799	328.092	487.580	682.160	
	stdv	285.510	303.743	375.604	436.117	319.867	274.811	348.341	458.840	471.799	328.092	487.580	682.160	478.726	
YearPredictionMSD	direct	300.172	-3.2956e+04	-4.8930e+04	-5.4265e+04	-1.2999e+04	-1.3009e+04	-4.3430e+04	-3.4469e+04	-4.1146e+04	-4.1148e+04	-5.4490e+04	-4.9791e+04	-2.7406e+05	
	mean	-3.0988e+04	-6.2187e+04	-7.6926e+04	-7.9996e+04	-4.1818e+04	-4.1986e+04	-7.5328e+04	-6.631e+04	-7.1072e+04	-7.0857e+04	-8.0555e+04	-7.7224e+04	-8.9637e+04	
	stdv	409.526	381.288	575.546	458.809	404.891	379.834	682.798	400.727	378.022	469.249	448.637	371.422	1555.579	
a9a	direct	307.887	-3.2920e+04	-7.1621e+04	-7.943.999	-1.942.672	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186
	mean	-2.225e+04	-5.9850e+04	-1.005e+05	-1.0412e+05	-3.113e+04	-3.3183e+04	-7.3280e+04	-6.0063e+04	-9.3385e+04	-9.3594e+04	-1.0229e+05	-1.0089e+05	-1.3405e+05	
	stdv	385.449	255.444	513.662	395.157	440.785	361.225	463.507	411.840	460.235	496.559	419.359	368.032	1883.872	
cadata	direct	-1.8143e+04	-9.8200e+04	-1.9711e+05	-2.0810e+05	-5.4139e+04	-6.2524e+04	-1.3997e+05	-1.1411e+05	-1.8295e+05	-1.8618e+05	-2.0342e+05	-2.0006e+05	-2.5111e+05	
	mean	4174.340	234.610	567.249	323.322	825.211	336.002	303.583	485.779	387.897	361.191	362.794	280.979	277.213	
	stdv	412.913	-1.0154e+05	-2.0426e+05	-2.2395e+05	-5.7488e+04	-7.1388e+04	-1.3327e+05	-1.1230e+05	-1.9223e+05	-1.9539e+05	-2.1448e+05	-2.0823e+05	-2.7407e+05	
cod-rnn	direct	-215.36	-1.0287e+05	-2.1103e+05	-2.4745e+05	-5.8141e+04	-7.1681e+04	-1.4665e+05	-1.2924e+05	-2.1770e+05	-2.2163e+05	-2.4070e+05	-2.3516e+05	-2.5974e+05	
	mean	429.133	299.086	359.058	420.114	52.705	96.386	300.256	235.355	367.632	452.731	349.753	407.336	219.823	
	stdv	436.867	-3.7532e+04	-7.4977e+04	-8.2746e+04	-2.3940e+04	-2.6807e+04	-4.2144e+04	-3.7768e+04	-7.0527e+04	-7.1240e+04	-7.6531e+04	-7.5026e+04	-2.7398e+05	
connect-4	direct	-2.3698e+04	-6.4899e+04	-1.0610e+05	-1.0840e+05	-4.6986e+04	-5.0457e+04	-7.4027e+04	-6.5703e+04	-1.0101e+05	-1.0174e+05	-1.0698e+05	-1.0611e+05	-1.2762e+05	
	mean	586.906	421.356	400.485	369.317	593.864	497.894	452.092	295.092	578.151	367.587	507.872	345.889	1400.016	
	stdv	544.685	-7.5471e+04	-1.3300e+05	-2.0548e+05	-4.4821e+04	-5.6541e+04	-9.9056e+04	-7.9148e+04	-1.7718e+05	-1.7942e+05	-1.9256e+05	-1.8470e+05	-2.7406e+05	
covtype.binary	direct	-1573.697	-7.7408e+04	-2.1320e+05	-2.2940e+05	-4.5848e+04	-5.7678e+04	-1.0744e+05	-1.8358e+04	-2.0682e+05	-2.0887e+05	-2.2227e+05	-2.1466e+05	-2.5115e+05	
	mean	165.233	291.190	311.374	346.723	352.154	352.154	352.154	352.154	352.154	352.154	352.154	352.154	352.154	
	stdv	341.659	-8.8310e+04	-1.1253e+05	-1.2348e+05	-5.505.813	-8.065.583	-6.6485e+04	-5.1014e+04	-1.0355e+05	-1.0372e+05	-1.1730e+05	-1.1328e+05	-2.7403e+05	
creditcard	direct	-1.1073e+04	-6.0350e+04	-1.1429e+05	-1.5376e+05	-1.5881e+04	-1.7998e+04	-8.4080e+04	-6.3407e+04	-1.3001e+05	-1.3049e+05	-1.4893e+05	-1.4219e+05	-1.8814e+05	
	mean	434.569	403.245	311.751	410.784	332.249	270.735	303.787	312.231	619.893	524.833	422.812	490.900	993.545	
	stdv	421.273	-1.0389e+05	-2.0277e+05	-2.2146e+05	-5.1184e+04	-7.1190e+04	-1.2397e+05	-1.1286e+05	-1.9334e+05	-1.9502e+05	-2.0948e+05	-2.0520e+05	-2.7408e+05	
diamonds	direct	-5727.127	-1.0447e+05	-2.3077e+05	-2.4728e+05	-5.4827e+04	-7.1434e+04	-1.4079e+05	-1.1623e+05	-2.1795e+05	-2.2186e+05	-2.3681e+05	-2.3335e+05	-2.5090e+05	
	mean	1516.542	135.963	521.442	349.360	74.216	143.409	406.337	203.733	279.037	416.968	369.196	529.177	519.439	
	stdv	459.655	-1.2606e+05	-1.9345e+05	-2.1402e+05	-8.5792e+04	-9.4616e+04	-1.4038e+05	-1.2498e+05	-1.8244e+05	-1.8454e+05	-2.0046e+05	-1.9544e+05	-2.7418e+05	
hlsfm.lhc.jets.hlf	direct	-1.4217e+04	-1.2444e+05	-2.2878e+05	-2.4669e+05	-8.5846e+04	-9.5549e+04	-1.5501e+05	-1.3128e+05	-2.1627e+05	-2.1944e+05	-2.3418e+05	-2.3048e+05	-2.5695e+05	
	mean	5180.807	168.923	428.815	521.642	181.025	286.255	348.394	187.203	395.036	490.716	591.908	509.586	557.662	
	stdv	383.970	-6.3410e+04	-1.1837e+05	-1.6726e+05	-2.9368e+04	-3.6906e+04	-8.6161e+04	-6.5058e+04	-1.4354e+05	-1.4448e+05	-1.5183e+05	-1.4858e+05	-2.7398e+05	
ijcnn1	direct	-4477.291	-6.8540e+04	-1.7942e+05	-1.9834e+05	-3.3066e+04	-4.0752e+04	-9.5012e+04	-7.6010e+04	-1.7317e+05	-1.7441e+05	-1.8380e+05	-1.7979e+05	-2.2965e+05	
	mean	375.887	-3.2920e+04	-7.1621e+04	-7.943.999	-1.942.672	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186	-1.08.186
	stdv	275.887	331.717	313.774	524.992	224.672	198.186	200.138	226.092	547.534	508.348	341.781	429.496	987.895	
jannis	direct	-296.325	-2.7786e+04	-4.1640e+04	-4.5929e+04	-1.5126e+04	-1.6156e+04	-3.1709e+04	-2.8248e+04	-3.6832e+04	-3.7273e+04	-4.2850e+04	-4.1805e+04	-7.4026e+05	
	mean	5.2013e+04	-6.5125e+04	-6.8943e+04	-6.9009e+04	-5.7267e+04	-5.7942e+04	-6.7747e+04	-6.5652e+04	-6.8008e+04	-6.8185e+04	-6.8912e+04	-6.9032e+04	-7.0206e+04	
	stdv	389.299	309.218	346.723	352.154	352.154	352.154	352.154	352.154	352.154	352.154	352.154	352.154	352.154	
phishing	direct	-1.0225e+04	-1.4283e+04	-2.1815e+04	-2.1966e+04	-1.2081e+04	-1.2248e+04	-1.7629e+04	-1.3705e+04	-2.1023e+04	-2.1340e+04	-2.2070e+04	-2.1750e+04	-2.7897e+04	
	mean	271.929	165.892	336.294	124.941	173.470	161.635	324.259	221.586	168.240	255.960	270.588	270.082	1218.773	
	stdv	173.836	-1.5656e+04	-1.7230e+04	-1.9799e+04	-1.2864e+04	-1.4444e+04	-1.3384e+04	-1.4444e+04	-1.7824e+04	-1.7824e+04	-2.0519e+04	-2.0519e+04	-2.0519e+04	
santander	direct	-305.306	-3.0530e+04	-3.0590e+04</											

name	quantity	Frangella	RPC+D	RPC+V $\frac{1}{4}$	RPC+V $\frac{1}{2}$	CPC+D	FPS+D	GKL+D	SDS+D	CPC+V $\frac{1}{4}$	FPS+V $\frac{1}{4}$	GKL+V $\frac{1}{4}$	SDS+V $\frac{1}{4}$	Dvar
ACIncome	direct	334.926	-6.699e+04	-1.1523e+05	-1.6638e+05	-3.1831e+04	-3.3126e+04	-6.7833e+04	-1.3577e+05	-1.3742e+05	-1.3314e+05	-1.4341e+05	-1.3698e+05	
	mean	3339.010	-7.1076e+04	-1.5813e+05	-1.8456e+05	-3.4846e+04	-4.2211e+04	-9.0304e+04	-2.7106e+04	-1.4658e+05	-1.4878e+05	-1.5829e+05	-1.3756e+05	
	stdev	177.103	225.613	334.239	370.668	315.055	226.360	221.193	161.919	291.461	418.628	402.563	285.111	176.066
AchillesDepDelayLM	direct	369.589	-6.6714e+04	-1.3696e+05	-1.4857e+05	-4.7225e+04	-5.4100e+04	-7.7599e+04	-6.9067e+04	-1.2418e+05	-1.2631e+05	-1.3908e+05	-1.3225e+05	-1.5000e+05
	mean	3339.010	-6.5984e+04	-1.4296e+05	-1.6670e+05	-5.6526e+04	-8.1076e+04	-7.2061e+04	-1.3461e+05	-1.3728e+05	-1.4476e+05	-1.5829e+05	-1.3937e+05	
	stdev	157.861	189.842	320.845	417.094	424.034	168.595	186.575	244.076	288.810	330.053	429.548	319.384	1040.703
COMET_MC_SAMPLE	direct	547.858	-1.9501e+05	-4.2275e+05	-4.2967e+05	-1.0851e+05	-3.9891e+05	-3.7083e+05	-4.1893e+05	-4.1994e+05	-4.0409e+05	-3.6210e+05	-4.5720e+05	
	mean	-8.8374e+04	-2.1578e+05	-4.2558e+05	-4.3465e+05	-1.6477e+05	-1.5557e+05	-4.0923e+05	-3.9051e+05	-4.2090e+05	-4.2212e+05	-4.1510e+05	-3.7666e+05	-4.8824e+05
	stdev	2.2968e+04	3445.151	188.708	326.777	1.0672e+04	91.46189	252.291	463.163	101.984	134.929	374.004	258.163	1453.342
ClickPredictionSmall	direct	366.155	-8.8884e+04	-2.1700e+05	-2.1897e+05	-3.8238e+04	-4.0002e+04	-1.2878e+05	-1.0010e+05	-2.0278e+05	-2.0377e+05	-2.2686e+05	-2.2049e+05	-1.5696e+05
	mean	-3235.334	-9.0198e+04	-2.2572e+05	-2.5900e+05	-3.9596e+04	-4.1348e+04	-1.3150e+05	-1.0154e+05	-2.0966e+05	-2.1066e+05	-2.3686e+05	-2.3007e+05	-4.1120e+05
	stdev	653.732	176.713	286.028	346.910	139.272	99.851	163.150	151.423	304.025	380.108	339.326	282.470	907.353
HIGGS	direct	307.249	-3.5030e+04	-5.9030e+04	-6.3896e+04	-1.1254e+04	-1.3596e+04	-3.9828e+04	-3.7727e+04	-1.1104e+04	-3.1077e+04	-5.9605e+04	-8.8767e+04	-1.5601e+05
	mean	-3.2232e+04	-6.8214e+04	-9.3872e+04	-9.6006e+04	-4.3265e+04	-4.2915e+04	-7.7337e+04	-6.9536e+04	-8.7908e+04	-8.7871e+04	-9.4975e+04	-9.3767e+04	-1.3850e+05
	stdev	502.085	535.646	301.267	226.850	559.836	605.775	510.127	567.065	224.794	460.271	478.466	697.309	2760.891
MNIST	direct	334.968	-3.1735e+04	-5.2112e+04	-5.9801e+04	-4.475.069	-4.685.466	-3.9959e+04	-3.2118e+04	-4.8036e+04	-4.8034e+04	-5.4116e+04	-5.2120e+04	-1.5694e+05
	mean	-4.2762e+04	-6.4110e+04	-7.3891e+04	-7.4531e+04	-4.5171e+04	-4.5151e+04	-7.0901e+04	-6.4913e+04	-7.2903e+04	-7.3068e+04	-7.4276e+04	-7.4009e+04	-9.6383e+04
	stdev	557.866	356.698	358.959	279.699	645.415	454.630	204.412	490.663	394.362	370.578	336.795	238.108	3240.331
MedicalAppointment	direct	344.779	-4.8381e+04	-1.0385e+05	-1.1602e+05	-3.2677e+04	-3.7707e+04	-5.6975e+04	-4.8880e+04	-9.9519e+04	-1.0657e+05	-1.0684e+05	-1.0412e+05	-1.5667e+05
	mean	7749.806	-5.8306e+04	-1.3082e+05	-1.4390e+05	-4.4540e+04	-4.5905e+04	-6.9541e+04	-5.8939e+04	-1.2377e+05	-1.2550e+05	-1.3570e+05	-1.3100e+05	-1.1756e+05
	stdev	358.068	353.878	620.635	495.720	295.938	267.868	324.816	273.675	407.336	358.542	464.972	613.920	2255.523
MiniBooneNE	direct	436.215	-1.3320e+05	-1.7755e+05	-1.9434e+05	-8.6624e+04	-9.5537e+04	-1.5447e+05	-1.4190e+05	-1.6252e+05	-1.6369e+05	-1.8882e+05	-1.8234e+05	-4.5723e+05
	mean	9048.227	-1.3624e+05	-1.8613e+05	-2.0862e+05	-8.8510e+04	-9.7397e+04	-1.5905e+05	-1.4534e+05	-1.6853e+05	-1.6977e+05	-2.0050e+05	-1.9216e+05	-1.1286e+05
	stdev	1452.024	134.491	367.541	384.217	113.615	301.321	219.088	285.084	220.001	259.763	352.72	211.822	852.002
YearPredictionMSD	direct	300.172	-3.2956e+04	-4.8943e+04	-5.1269e+04	-1.2999e+04	-1.3009e+04	-4.3460e+04	-3.4469e+04	-4.1149e+04	-4.1151e+04	-5.4494e+04	-4.9795e+04	-1.5697e+05
	mean	-3.1253e+04	-6.2062e+04	-7.6970e+04	-7.9810e+04	-4.1773e+04	-4.1698e+04	-7.5364e+04	-6.4397e+04	-7.1107e+04	-7.1185e+04	-8.0519e+04	-7.7458e+04	-1.0914e+05
	stdev	480.687	578.636	360.989	353.342	370.798	472.019	259.345	324.944	353.531	315.151	397.085	256.223	3129.016
a9a	direct	340.686	-3.0101e+04	-7.1521e+04	-7.9374e+04	-1.0522e+04	-1.2477e+04	-4.6854e+04	-3.6689e+04	-6.5247e+04	-6.5092e+04	-7.9024e+04	-7.7669e+04	-1.8454e+05
	mean	-2.2477e+04	-5.9566e+04	-1.0027e+05	-1.0451e+05	-3.1107e+04	-3.2823e+04	-7.3046e+04	-6.0129e+04	-9.3371e+04	-9.3642e+04	-1.0268e+05	-1.0060e+05	-1.3691e+05
	stdev	409.895	484.529	519.729	258.743	458.462	413.731	464.814	193.123	485.078	602.851	265.019	403.412	3923.541
cadata	direct	397.433	-9.6354e+04	-1.7490e+05	-1.9781e+05	-4.9533e+04	-6.1198e+04	-4.3621e+05	-1.0528e+05	-1.5029e+05	-1.5028e+05	-1.5180e+05	-1.4145e+05	-1.7778e+05
	mean	-2.1145e+04	-9.4450e+04	-1.7684e+05	-2.0203e+05	-5.4495e+04	-6.3696e+04	-1.2655e+05	-1.0603e+05	-1.6058e+05	-1.6558e+05	-1.9500e+05	-1.8354e+05	-3.5621e+05
	stdev	6594.422	95.269	208.381	201.983	615.619	525.068	167.193	154.528	155.223	245.662	218.657	198.393	718.995
cod-rnn	direct	412.913	-9.8517e+04	-2.1388e+05	-2.4249e+05	-5.7484e+04	-7.1390e+04	-1.3440e+05	-1.1233e+05	-2.0075e+05	-2.0429e+05	-2.3121e+05	-2.2031e+05	-1.5697e+05
	mean	-2170.166	-9.8774e+04	-2.1810e+05	-2.5110e+05	-5.8305e+04	-7.7194e+04	-1.3111e+05	-1.2274e+05	-2.0303e+05	-2.0775e+05	-2.3604e+05	-2.2506e+05	-1.2908e+05
	stdev	257.941	155.649	122.478	328.796	152.917	107.675	140.256	132.764	134.915	170.341	382.938	267.321	744.061
connect-4	direct	436.867	-3.7352e+04	-7.4978e+04	-8.2749e+04	-2.3940e+04	-2.6807e+04	-4.2145e+04	-3.7768e+04	-7.0529e+04	-7.1250e+04	-7.6533e+04	-7.5028e+04	-1.5698e+05
	mean	-2.3712e+04	-6.4730e+04	-1.0593e+05	-1.0839e+05	-4.6918e+04	-5.0381e+04	-7.4903e+04	-6.5320e+04	-1.0065e+05	-1.0168e+05	-1.0692e+05	-1.0575e+05	-1.6898e+05
	stdev	460.085	346.412	460.874	512.058	369.489	479.955	519.089	318.451	378.653	512.281	361.080	334.645	3180.772
covtype.binary	direct	544.685	-7.5457e+04	-1.8452e+05	-2.1177e+05	-4.4821e+04	-9.9192e+04	-7.9153e+04	-1.7828e+05	-1.8066e+05	-1.9643e+05	-1.8642e+05	-1.5091e+05	
	mean	-1515.947	-7.7275e+04	-1.9290e+05	-2.2375e+05	-4.5772e+04	-5.7766e+04	-1.0202e+05	-8.1048e+04	-1.5852e+05	-1.5852e+05	-2.0700e+05	-1.9478e+05	-1.1187e+05
	stdev	106.424	175.829	140.092	138.760	190.007	215.546	339.521	414.902	215.546	339.521	414.902	215.546	875.215
creditcard	direct	341.659	-4.8310e+04	-1.1903e+05	-1.3233e+05	-5.505.874	-8.065.619	-6.7241e+04	-5.1016e+04	-1.1077e+05	-1.1088e+05	-1.2440e+05	-1.1964e+05	-1.5693e+05
	mean	-1.3142e+04	-6.0106e+04	-1.4371e+05	-1.6122e+05	-1.6027e+04	-1.8136e+04	-8.2936e+04	-6.3388e+04	-1.3286e+05	-1.3315e+05	-1.4961e+05	-1.4473e+05	-2.9186e+05
	stdev	343.457	348.329	448.763	672.172	292.922	330.411	351.531	321.906	510.583	289.976	393.057	448.742	1990.670
diamonds	direct	421.274	-1.0308e+05	-2.1157e+05	-2.1148e+05	-5.1184e+04	-7.1191e+04	-1.2991e+05	-1.1396e+05	-2.0108e+05	-2.0438e+05	-2.2164e+05	-2.1453e+05	-1.5696e+05
	mean	-5266.104	-1.0321e+05	-2.1585e+05	-2.5042e+05	-5.4970e+04	-7.1427e+04	-1.3048e+05	-1.1329e+05	-2.0502e+05	-2.0816e+05	-2.2727e+05	-2.1907e+05	-4.2920e+05
	stdev	1752.397	97.815	237.677	281.260	141.189	80.403	138.477	114.769	232.663	249.523	239.662	237.197	574.169
hlsfm_lhc_lets_hlf	direct	459.656	-1.1882e+05	-1.9447e+05	-2.1939e+05	-8.5794e+04	-9.4629e+04	-1.4113e+05	-1.2505e+05	-1.8330e+05	-1.8545e+05	-2.0471e+05	-1.9761e+05	-1.5709e+05
	mean	-1.5747e+04	-1.1970e+05	-1.9946e+05	-2.2821e+05	-8.6371e+04	-9.2545e+04	-1.4243e+05	-1.2594e+05	-1.8728e+05	-1.8966e+05	-2.1089e+05	-2.0208e+05	-4.2178e+05
	stdev	5824.100	150.941	294.245	295.252	146.448	217.592	147.496	131.183	249.787	267.296	266.668	239.378	1066.903
ijcnn1	direct	383.971	-6.5152e+04	-1.1894e+05	-1.6942e+05	-2.9368e+04	-3.6990e+04	-8.6202e+04	-6.4701e+04	-1.4400e+05	-1.4496e+05	-1.5275e+05	-1.4926e+05	-1.5698e+05
	mean	-4422.527	-7.0536e+04	-1.4448e+05	-1.9035e+05	-3.3657e+04	-4.0638e+04	-9.3186e+04	-6.9805e+04	-1.5508e+05	-1.5976e+05	-1.6974e+05	-1.6510e+05	-3.6750e+05
	stdev	128.634	178.829	342.235	387.451	192.999	190.007	447.766	271.532	527.452	370.624	370.624	370.624	185.446
jannis	direct	296.325	-2.7786e+04	-4.1640e+04	-4.5900e+04	-1.5126e+04	-1.6157e+04	-3.1709e+04	-2.8248e+04	-3.6832e+04	-3.7273e+04	-4.2850e+04	-4.1805e+04	-1.5693e+05
	mean	-5.2033e+04	-6.5210e+04	-6.8856e+04	-6.8964e+04	-5.7427e+04	-5.7999e+04	-6.7092e+04	-6.5661e+04	-6.8255e+04	-6.7953e+04	-6.9049e+04	-6.9016e+04	-1.0502e+05
	stdev	151.154	-7.456.314	-1.5399e+04	-1.6919e+04	-4.644.952	-4.884.723	-9031.730	-7203.437	-1.4539e+04	-1.4607e+04	-1.5949e+04	-1.5278e+04	-1.1578e+05
phishing	direct	-1.0322e+04	-1.4124e+04	-2.1814e+04	-2.2140e+04	-1.2047e+04	-1.2236e+04	-1.7579e+04	-1.3675e+04	-2.1163e+04	-2.1158e+04	-2.1934e+04	-2.1679e+04	-3.6837e+04
	mean	231.541	338.905	300.222	250.855	234.340	233.184	271.942	323.619	169.899	213.493	142.974	234.861	234.356
	stdev	176.830	136.644	1.7803e+04	1									
santander														

a few proofs deferred from the main text.

A.1. Sparse approximation characterization. This section provides a linear algebraic characterization of the Vecchia approximation. We are not aware of this characterization appearing elsewhere.

PROPOSITION A.1 (Vecchia characterization). *For any positive-semidefinite matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$, the Vecchia approximation $\hat{\mathbf{A}} = \hat{\mathbf{C}}^{-1} \hat{\mathbf{D}} \hat{\mathbf{C}}^{-*}$ with sparsity pattern $\{\mathcal{S}_i\}_{i=1}^n$ is characterized by*

$$(A.1) \quad \begin{cases} \mathbf{e}_i - \hat{\mathbf{C}}(i, \cdot)^* \in \underset{\mathbf{v} \in \text{span}\{\mathbf{e}_j\}_{j \in \mathcal{S}_i}}{\text{argmin}} d_{\mathbf{A}}(\mathbf{e}_i, \mathbf{v}), & 1 \leq i \leq n, \\ \hat{\mathbf{D}}(i, i) = \min_{\mathbf{v} \in \text{span}\{\mathbf{e}_j\}_{j \in \mathcal{S}_i}} d_{\mathbf{A}}(\mathbf{e}_i, \mathbf{v}), & 1 \leq i \leq n. \end{cases}$$

In contrast, the inverse Cholesky decomposition $\mathbf{A} = \mathbf{C}^{-1} \mathbf{D} \mathbf{C}^{-*}$ is characterized by

$$(A.2) \quad \begin{cases} \mathbf{e}_i - \mathbf{C}(i, \cdot)^* \in \underset{\mathbf{v} \in \text{span}\{\mathbf{e}_j\}_{j < i}}{\text{argmin}} d_{\mathbf{A}}(\mathbf{e}_i, \mathbf{v}), & 1 \leq i \leq n, \\ \mathbf{D}(i, i) = \min_{\mathbf{v} \in \text{span}\{\mathbf{e}_j\}_{j < i}} d_{\mathbf{A}}(\mathbf{e}_i, \mathbf{v}), & 1 \leq i \leq n. \end{cases}$$

Proof. For any $\mathbf{J} \subseteq \{1, \dots, i-1\}$ and $\mathbf{v} \in \text{span}\{\mathbf{e}_j\}_{j \in \mathbf{J}}$, we can write

$$\begin{aligned} d_{\mathbf{A}}(\mathbf{e}_i, \mathbf{v})^2 &= \begin{bmatrix} -\mathbf{v}(\mathbf{J}) \\ 1 \end{bmatrix}^* \begin{bmatrix} \mathbf{A}(\mathbf{J}, \mathbf{J}) & \mathbf{A}(\mathbf{J}, i) \\ \mathbf{A}(i, \mathbf{J}) & \mathbf{A}(i, i) \end{bmatrix} \begin{bmatrix} -\mathbf{v}(\mathbf{J}) \\ 1 \end{bmatrix} \\ &= \mathbf{A}(i, i) - \mathbf{A}(i, \mathbf{J}) \mathbf{A}(\mathbf{J}, \mathbf{J})^+ \mathbf{A}(\mathbf{J}, i) \\ &\quad + [\mathbf{v}(\mathbf{J})^* \mathbf{A}(\mathbf{J}, \mathbf{J}) - \mathbf{A}(i, \mathbf{J})] \mathbf{A}(\mathbf{J}, \mathbf{J})^+ [\mathbf{A}(\mathbf{J}, \mathbf{J}) \mathbf{v}(\mathbf{J}) - \mathbf{A}(\mathbf{J}, i)]. \end{aligned}$$

Here we are using the fact that $\mathbf{A}(\mathbf{J}, i) \in \text{range } \mathbf{A}(\mathbf{J}, \mathbf{J})$ and consequently

$$\mathbf{A}(\mathbf{J}, \mathbf{J}) \mathbf{A}(\mathbf{J}, \mathbf{J})^+ \mathbf{A}(\mathbf{J}, i) = \mathbf{A}(\mathbf{J}, i)$$

because \mathbf{A} is positive-semidefinite. The minimizers of $d(\mathbf{e}_i, \mathbf{v})^2$ are characterized by

$$(A.3) \quad \mathbf{v}(\mathbf{J})^* \mathbf{A}(\mathbf{J}, \mathbf{J}) = \mathbf{A}(i, \mathbf{J}),$$

and the minimum value of $d(\mathbf{e}_i, \mathbf{v})^2$ can be written using (A.3) as

$$(A.4) \quad \|\mathbf{e}_i - \mathbf{v}\|_{\mathbf{A}}^2 = \mathbf{A}(i, i) - \mathbf{v}(\mathbf{J})^* \mathbf{A}(\mathbf{J}, i).$$

Thus, we have characterized the minimizers and minimum value of the distance.

Next we check that formulas (A.3) and (A.4) apply to the Vecchia approximation and inverse Cholesky decomposition. Indeed, the vector $\mathbf{v} = \mathbf{e}_i - \hat{\mathbf{C}}(i, \cdot)^*$ and diagonal entry $\hat{\mathbf{D}}(i, i)$ in the Vecchia approximation perfectly match the characterization (A.3)-(A.4) with $\mathbf{J} = \mathcal{S}_i$. Now observe the matrix $\mathbf{D} \mathbf{C}^{-*}$ in the inverse Cholesky decomposition is upper triangular, so

$$(A.5) \quad \begin{aligned} \mathbf{0} &= (\mathbf{D} \mathbf{C}^{-*})(i, 1:i-1) \\ &= (\mathbf{C} \mathbf{A})(i, 1:i-1) \\ &= \mathbf{C}(i, 1:i-1) \mathbf{A}(1:i-1, 1:i-1) + \mathbf{A}(i, 1:i-1). \end{aligned}$$

Using (A.5), each diagonal entry satisfies

$$\mathbf{D}(i, i) = \mathbf{C}(i, \cdot) \mathbf{A} \mathbf{C}(i, \cdot)^* = \mathbf{A}(i, i) + \mathbf{C}(i, 1:i-1) \mathbf{A}(1:i-1, i)$$

Hence, the vector $\mathbf{v} = \mathbf{e}_i - \mathbf{C}(i, \cdot)^*$ and diagonal entry $\mathbf{D}(i, i)$ in the inverse Cholesky decomposition match the characterization (A.3) and (A.4) with $\mathbf{J} = \{1, \dots, i-1\}$.

Last, if \mathbf{D} and \mathbf{C} are generated according to (A.2), then (A.3) shows that

$$\mathbf{C}(i, \cdot) \mathbf{A}(\cdot, 1:i-1) = \mathbf{0}, \quad \text{for each } i = 1, \dots, n,$$

so \mathbf{CA} must be upper triangular. By multiplication with another upper triangular matrix, \mathbf{CAC}^* is upper triangular also. By symmetry, \mathbf{CAC}^* must be diagonal, and by (A.2) the diagonal entries are given by

$$\mathbf{D}(i, i) = \mathbf{C}(i, \cdot) \mathbf{AC}(i, \cdot)^*, \quad \text{for each } i = 1, \dots, n,$$

We conclude that $\mathbf{A} = \mathbf{C}^{-*} \mathbf{DC}^{-1}$ is an inverse Cholesky decomposition. \square

Proposition A.1 gives a sparse variational characterization of the Vecchia approximation. Identifying each row $\mathbf{C}(i, \cdot)$ and diagonal entry $\mathbf{D}(i, i)$ in the inverse Cholesky decomposition via (A.2) is possible but expensive — it requires examining and processing a large submatrix $\mathbf{A}(1:i, 1:i)$. It is cheaper to restrict the optimization to a sparsity set $\mathbf{S}_i \subseteq \{1, \dots, i-1\}$ so the solution only depends on a submatrix $\mathbf{A}(\{i\} \cup \mathbf{S}_i, \{i\} \cup \mathbf{S}_i)$. This sparse optimization results in the Vecchia approximation.

A.2. Volume and trace formulas. The inverse Cholesky decomposition $\mathbf{A} = \mathbf{C}^{-*} \mathbf{DC}^{-1}$ leads to pseudoinverse and volume formulas given in the next lemma.

LEMMA A.2 (Pseudoinverse and volume formulas). *Fix a positive-semidefinite matrix with inverse Cholesky decomposition $\mathbf{A} = \mathbf{C}^{-*} \mathbf{DC}^{-1}$. Let \mathbf{Q} be an orthonormal basis for $\text{range}(\mathbf{A})$ and set $\mathbf{R} = \mathbf{I}(\cdot, \mathbf{N})$, where $\mathbf{N} \subseteq \{1, \dots, n\}$ picks out the nonzero elements of \mathbf{D} . Then*

$$(A.6) \quad \mathbf{A}^+ = \mathbf{QQ}^* \mathbf{CD}^+ \mathbf{C}^* \mathbf{QQ}^* \quad \text{and} \quad \text{vol}(\mathbf{A}) = \frac{\text{vol}(\mathbf{D})}{\det(\mathbf{R}^* \mathbf{Q})^2}.$$

Proof. The pseudoinverse \mathbf{A}^+ is characterized as the unique positive-semidefinite matrix that satisfies $\mathbf{A}^+ \mathbf{x} = \mathbf{0}$ if $\mathbf{x} \perp \text{range}(\mathbf{Q})$ and $\mathbf{AA}^+ \mathbf{x} = \mathbf{x}$ if $\mathbf{x} \in \text{range}(\mathbf{Q})$. We can check that these conditions are satisfied for $\mathbf{A}^+ = \mathbf{QQ}^* \mathbf{CD}^+ \mathbf{C}^* \mathbf{QQ}^*$. If $\mathbf{x} \perp \text{range}(\mathbf{Q})$, we verify

$$\mathbf{QQ}^* \mathbf{CD}^+ \mathbf{C}^* \mathbf{QQ}^* \mathbf{x} = \mathbf{0}.$$

If $\mathbf{x} \in \text{range}(\mathbf{Q})$, then $\mathbf{x} = \mathbf{C}^{-*} \mathbf{Dy}$ for some $\mathbf{y} \in \mathbb{C}^n$ and we verify

$$\mathbf{AQQ}^* \mathbf{CD}^+ \mathbf{C}^* \mathbf{QQ}^* \mathbf{x} = \mathbf{C}^{-*} \mathbf{DD}^+ \mathbf{C}^* \mathbf{x} = \mathbf{C}^{-*} \mathbf{DD}^+ \mathbf{C}^* \mathbf{C}^{-*} \mathbf{Dy} = \mathbf{C}^{-*} \mathbf{Dy} = \mathbf{x}.$$

This confirms the pseudoinverse formula in (A.6).

Next observe that $\mathbf{R}^* \mathbf{C}^{-*} \mathbf{R}$ is an upper triangular matrix with ones on the diagonal; hence, it has determinant one. Using the definition of \mathbf{Q} as an orthonormal basis for $\text{range}(\mathbf{A}) = \text{range}(\mathbf{C}^{-*} \mathbf{R})$, it follows

$$1 = \det(\mathbf{R}^* \mathbf{C}^{-*} \mathbf{R}) = \det(\mathbf{R}^* \mathbf{QQ}^* \mathbf{C}^{-*} \mathbf{R}) = \det(\mathbf{R}^* \mathbf{Q}) \det(\mathbf{Q}^* \mathbf{C}^{-*} \mathbf{R}).$$

The volume of \mathbf{A} is given by

$$\begin{aligned} \text{vol}(\mathbf{A}) &= \text{vol}(\mathbf{QQ}^* \mathbf{C}^{-*} \mathbf{RR}^* \mathbf{DRR}^* \mathbf{C}^{-1} \mathbf{QQ}^*) \\ &= \det(\mathbf{Q}^* \mathbf{C}^{-*} \mathbf{RR}^* \mathbf{DRR}^* \mathbf{C}^{-1} \mathbf{Q}) = \det(\mathbf{Q}^* \mathbf{C}^{-*} \mathbf{R})^2 \det(\mathbf{R}^* \mathbf{DR}) \\ &= \text{vol}(\mathbf{D}) / \det(\mathbf{R}^* \mathbf{Q})^2. \end{aligned}$$

This confirms the volume formula in (A.6) and completes the proof. \square

The next lemma gives explicit formulas for the trace and volume of $\mathbf{A}\hat{\mathbf{A}}^+$, where $\hat{\mathbf{A}}$ is an approximation with the same range as \mathbf{A} . In particular, when $\hat{\mathbf{A}}$ is the Vecchia approximation with the same range as \mathbf{A} , this guarantees $\text{tr}(\hat{\mathbf{A}}^+\mathbf{A}) = \text{rank}(\mathbf{A})$.

LEMMA A.3 (Trace and volume formulas). *When two inverse Cholesky decompositions $\mathbf{A} = \mathbf{C}^{-1}\mathbf{D}\mathbf{C}^{-*}$ and $\hat{\mathbf{A}} = \hat{\mathbf{C}}^{-1}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-*}$ share the same range, it follows*

$$\text{vol}(\hat{\mathbf{A}}^+\mathbf{A}) = \frac{\text{vol}(\mathbf{A})}{\text{vol}(\hat{\mathbf{A}})} = \frac{\text{vol}(\mathbf{D})}{\text{vol}(\hat{\mathbf{D}})} \quad \text{and} \quad \text{tr}(\hat{\mathbf{A}}^+\mathbf{A}) = \text{tr}(\hat{\mathbf{D}}^+\hat{\mathbf{C}}^*\mathbf{A}\hat{\mathbf{C}}).$$

Proof. First we argue that $\text{range}(\hat{\mathbf{D}}) = \text{range}(\mathbf{D})$. To that end, introduce the affine subspace

$$\mathbf{A}_i = \{(\mathbf{a}(1), \dots, \mathbf{a}(i-1), 1, 0, \dots, 0) \mid \mathbf{a} \in \mathbb{C}^{i-1}\}$$

for $i = 1, \dots, n$. Recall that \mathbf{C}^{-*} is upper triangular, and fix an index $i \in \{1, \dots, n\}$ for which $\mathbf{D}(i, i) > 0$. Then the range of $\hat{\mathbf{A}}$ contains the vector $\mathbf{C}^{-*}(\cdot, i) \in \mathbf{A}_i$ and it is also spanned by vectors $\hat{\mathbf{C}}^{-*}(\cdot, j) \in \mathbf{A}_j$ for which $\hat{\mathbf{D}}(j, j) > 0$. Any set of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ with $\mathbf{x}_j \in \mathbf{A}_j$ is linearly independent, so it must be the case that $\hat{\mathbf{D}}(i, i) > 0$. This argument shows $\text{range}(\mathbf{D}) \subseteq \text{range}(\hat{\mathbf{D}})$ and reversing the argument shows $\text{range}(\hat{\mathbf{D}}) \subseteq \text{range}(\mathbf{D})$ also.

At this point, take an orthonormal basis \mathbf{Q} for $\text{range}(\hat{\mathbf{A}}) = \text{range}(\mathbf{A})$ and an orthonormal basis \mathbf{R} for $\text{range}(\hat{\mathbf{D}}) = \text{range}(\mathbf{D})$. To calculate $\text{vol}(\hat{\mathbf{A}}^+\mathbf{A})$, observe

$$\begin{aligned} \text{vol}(\hat{\mathbf{A}}^+\mathbf{A}) &= \text{vol}(\mathbf{Q}\mathbf{Q}^*\hat{\mathbf{A}}^+\mathbf{Q}\mathbf{Q}^*\mathbf{A}\mathbf{Q}\mathbf{Q}^*) \\ &= \det(\mathbf{Q}^*\hat{\mathbf{A}}^+\mathbf{Q}\mathbf{Q}^*\mathbf{A}\mathbf{Q}) = \det(\mathbf{Q}^*\hat{\mathbf{A}}^+\mathbf{Q}) \det(\mathbf{Q}^*\mathbf{A}\mathbf{Q}) \\ &= \text{vol}(\mathbf{Q}\mathbf{Q}^*\hat{\mathbf{A}}^+\mathbf{Q}\mathbf{Q}^*) \text{vol}(\mathbf{Q}\mathbf{Q}^*\mathbf{A}\mathbf{Q}\mathbf{Q}^*) = \text{vol}(\hat{\mathbf{A}}^+) \text{vol}(\mathbf{A}) = \frac{\text{vol}(\mathbf{A})}{\text{vol}(\hat{\mathbf{A}})}. \end{aligned}$$

The last line uses the fact that $\text{vol}(\hat{\mathbf{A}}^+)$ is the inverse of $\text{vol}(\hat{\mathbf{A}})$. Then apply the exact formula for $\text{vol}(\mathbf{A})$ and $\text{vol}(\hat{\mathbf{A}})$ given in Lemma A.2:

$$\frac{\text{vol}(\mathbf{A})}{\text{vol}(\hat{\mathbf{A}})} = \frac{\text{vol}(\mathbf{D}) / \det(\mathbf{R}^*\mathbf{Q})^2}{\text{vol}(\hat{\mathbf{D}}) / \det(\mathbf{R}^*\mathbf{Q})^2} = \frac{\text{vol}(\mathbf{D})}{\text{vol}(\hat{\mathbf{D}})}.$$

Similarly, apply the exact formula for the pseudoinverse $\hat{\mathbf{A}}$ in Lemma A.2 to calculate $\text{tr}(\hat{\mathbf{A}}^+\mathbf{A})$ as follows:

$$\begin{aligned} \text{tr}(\hat{\mathbf{A}}^+\mathbf{A}) &= \text{tr}(\mathbf{Q}\mathbf{Q}^*\hat{\mathbf{C}}\hat{\mathbf{D}}^+\hat{\mathbf{C}}^*\mathbf{Q}\mathbf{Q}^*\mathbf{A}) \\ &= \text{tr}(\hat{\mathbf{D}}^+\hat{\mathbf{C}}^*\mathbf{Q}\mathbf{Q}^*\mathbf{A}\mathbf{Q}\mathbf{Q}^*\hat{\mathbf{C}}) = \text{tr}(\hat{\mathbf{D}}^+\hat{\mathbf{C}}^*\mathbf{A}\hat{\mathbf{C}}) \end{aligned}$$

This completes the proof. \square

A.3. Proof of Theorem 3.2. We introduce an inverse Cholesky decomposition $\mathbf{A} = \mathbf{C}^{-*}\mathbf{D}\mathbf{C}^{-1}$, and we will derive from first principles an inverse Cholesky approximation $\hat{\mathbf{A}} = \hat{\mathbf{C}}^{-*}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-1}$ that minimizes κ_{Kap} , assuming $\kappa_{\text{Kap}} < \infty$.

To that end, we partition the indices $\{1, \dots, n\}$ into a set of “good” indices $\mathbf{G} = \{i \mid \mathbf{D}(i, i) > 0\}$ and a set of “bad” indices $\mathbf{B} = \{i \mid \mathbf{D}(i, i) = 0\}$. The number of

good indices is $r = \text{rank}(\mathbf{A})$ and the number of bad indices is $n - r$. We first consider what happens to the bad indices $i \in \mathbf{B}$. If \mathbf{A} and $\hat{\mathbf{A}}$ have the same nullspace, we establish the following chain of implications:

- (i) $\mathbf{D}(i, i) = 0$;
- (ii) $\hat{\mathbf{D}}(i, i) = 0$;
- (iii) $\mathbf{A}\hat{\mathbf{C}}(\cdot, i) = \mathbf{0}$.

We can check the implications as follows:

- (i) \implies (ii). If $\mathbf{D}(i, i) = 0$, the vector $\mathbf{C}(\cdot, i)$ is in the shared nullspace of $\mathbf{A} = \mathbf{C}^{-*}\mathbf{D}\mathbf{C}^{-1}$ and $\hat{\mathbf{A}} = \hat{\mathbf{C}}^{-*}\hat{\mathbf{D}}\hat{\mathbf{C}}^{-1}$, so consequently the vector $\hat{\mathbf{C}}^{-1}\mathbf{C}(\cdot, i)$ is in the nullspace of $\hat{\mathbf{D}}$. Now recall that a lower triangular matrix with ones on the diagonal is a mapping from \mathcal{V}_i into \mathcal{V}_i , where

$$\mathcal{V}_i = \{\mathbf{v} \in \mathbb{C}^n \mid \mathbf{v}(i) = 1, \mathbf{v}(j) = 0 \text{ for } j < i\}.$$

Since $\mathbf{C}(\cdot, i) \in \mathcal{V}_i$ and $\hat{\mathbf{C}}^{-1}$ is lower triangular, it follows that $\hat{\mathbf{C}}^{-1}\mathbf{C}(\cdot, i) \in \mathcal{V}_i$. Since $\hat{\mathbf{C}}^{-1}\mathbf{C}(\cdot, i) \in \mathcal{V}_i$ lies in the nullspace of the nonnegative-valued diagonal matrix $\hat{\mathbf{D}}$, it follows that $\hat{\mathbf{D}}(i, i) = 0$.

- (ii) \implies (iii). If $\hat{\mathbf{D}}(i, i) = 0$, the vector $\hat{\mathbf{C}}(\cdot, i)$ is in the shared nullspace of \mathbf{A} and $\hat{\mathbf{A}}$.

Statements (ii) and (iii) imply that we must choose $\hat{\mathbf{C}}(\cdot, i)$ so that $\hat{\mathbf{C}}(\cdot, i)^*\mathbf{A}\hat{\mathbf{C}}(\cdot, i) = 0$ and we must set $\hat{\mathbf{D}}(i, i) = 0$. By the variational characterization in [Proposition A.1](#), $\hat{\mathbf{C}}(\cdot, i)$ and $\hat{\mathbf{D}}(i, i)$ are constructed in exactly this way in the Vecchia approximation.

Now we focus on the good indices $i \in \mathbf{G}$. To that end, we rewrite κ_{Kap} using the volume and trace formulas from [Lemma A.3](#) in the appendix.

$$\kappa_{\text{Kap}} = \frac{\left(\frac{1}{r} \text{tr}(\mathbf{A}\hat{\mathbf{A}}^+)\right)^r}{\text{vol}(\mathbf{A}\hat{\mathbf{A}}^+)} = \frac{\text{vol}(\hat{\mathbf{D}}) \text{tr}(\hat{\mathbf{D}}^+\hat{\mathbf{C}}^*\mathbf{A}\hat{\mathbf{C}})^r}{r^r \text{vol}(\mathbf{D})}.$$

When we optimize $\hat{\mathbf{C}}(\cdot, i)$ to make $(\hat{\mathbf{C}}^*\mathbf{A}\hat{\mathbf{C}})(i, i)$ as small as possible, [Proposition A.1](#) shows that we are reproducing the column $\hat{\mathbf{C}}(\cdot, i)$ in the Vecchia approximation. Next consider the impact of varying $\hat{\mathbf{D}}(i, i)$ after $\hat{\mathbf{C}}(\cdot, i)$ is selected.

$$\partial_{\hat{\mathbf{D}}(i, i)} \log(\kappa_{\text{Kap}}) = \frac{1}{\hat{\mathbf{D}}(i, i)} - \frac{r(\hat{\mathbf{C}}^*\mathbf{A}\hat{\mathbf{C}})(i, i)}{\hat{\mathbf{D}}(i, i)^2 \text{tr}(\hat{\mathbf{D}}^+\hat{\mathbf{C}}^*\mathbf{A}\hat{\mathbf{C}})}.$$

The logarithmic derivative is negative for small $\hat{\mathbf{D}}(i, i)$ and positive for large $\hat{\mathbf{D}}(i, i)$, and it achieves a zero value when $\hat{\mathbf{D}}(i, i) = (\hat{\mathbf{C}}^*\mathbf{A}\hat{\mathbf{C}})(i, i)$. This matches the Vecchia definition as well, due to [Proposition A.1](#).

With the Vecchia parameter settings, it follows that $\kappa_{\text{Kap}} = \text{vol}(\hat{\mathbf{D}})/\text{vol}(\mathbf{D})$. We can use the distance formulas in [Proposition A.1](#) to rewrite κ_{Kap} as given in [\(3.1\)](#) and complete the proof.

Appendix B. Applications of the Kaporin condition number. The section proves several upper bounds for linear algebra calculations in terms of the Kaporin condition number.

B.1. Proof of Proposition 3.3. We can use the fact that $\mathbf{b} - \mathbf{A}\mathbf{x}_\star$ is orthogonal to $\text{range}(\mathbf{A}) = \text{range}(\hat{\mathbf{A}})$ to calculate

$$\begin{aligned} \|\hat{\mathbf{x}} - \mathbf{x}_\star\|_{\mathbf{A}} &= \|\mathbf{x}_0 + \hat{\mathbf{A}}^+[\mathbf{b} - \mathbf{A}\mathbf{x}_0] - \mathbf{x}_\star\|_{\mathbf{A}} \\ &= \|[\mathbf{I} - \hat{\mathbf{A}}^+\mathbf{A}][\mathbf{x}_0 - \mathbf{x}_\star]\|_{\mathbf{A}} \\ &= \|[\mathbf{I} - \mathbf{A}^{1/2}\hat{\mathbf{A}}^+\mathbf{A}^{1/2}]\mathbf{A}^{1/2}[\mathbf{x}_0 - \mathbf{x}_\star]\| \\ &\leq \|\mathbf{I} - \mathbf{A}^{1/2}\hat{\mathbf{A}}^+\mathbf{A}^{1/2}\| \|\hat{\mathbf{x}} - \mathbf{x}_0\|_{\mathbf{A}} \end{aligned}$$

The eigenvalues of $\mathbf{A}^{1/2}\hat{\mathbf{A}}^+\mathbf{A}^{1/2}$ are the eigenvalues of $\mathbf{A}\hat{\mathbf{A}}^+$, which we write as $\lambda_1 \geq \dots \geq \lambda_r > 0$. Therefore,

$$\frac{\|\hat{\mathbf{x}} - \mathbf{x}_\star\|_{\mathbf{A}}}{\|\hat{\mathbf{x}} - \mathbf{x}_0\|_{\mathbf{A}}} \leq \|\mathbf{I} - \mathbf{A}^{1/2}\hat{\mathbf{A}}^+\mathbf{A}^{1/2}\| = \max_{1 \leq i \leq r} |1 - \lambda_i|,$$

which gives a sharp upper bound.

By the concavity of $x \mapsto \log x$, $\sum_{j=1}^r \log(\lambda_j)$ is maximized when the eigenvalues λ_j for $j \neq i$ are equal. By the normalization $\sum_{j=1}^r \lambda_j = r$, we obtain

$$-\log(\kappa_{\text{Kap}}) = \sum_{j=1}^r \log(\lambda_j) \leq \log(\lambda_i) + (r-1) \log\left(\frac{r - \lambda_i}{r-1}\right).$$

Over the interval $(0, r]$, we can bound $\log(x)$ from above by a concave quadratic that passes through $(1, 0)$ and $(r, \frac{r}{2} - \frac{1}{2r})$, so

$$\log(\lambda_i) \leq -1 + \lambda_i - \frac{1}{2r}(1 - \lambda_i)^2.$$

Also, $\log(1+x) \leq x$ holds globally, so

$$(r-1) \log\left(\frac{r - \lambda_i}{r-1}\right) = (r-1) \log\left(1 + \frac{1 - \lambda_i}{r-1}\right) \leq 1 - \lambda_i.$$

It follows that

$$-\log(\kappa_{\text{Kap}}) \leq \log(\lambda_i) + (r-1) \log\left(\frac{r - \lambda_i}{r-1}\right) \leq -\frac{1}{2r}(1 - \lambda_i)^2.$$

We conclude that $(1 - \lambda_i)^2 \leq 2r \log(\kappa_{\text{Kap}})$ for each $i = 1, \dots, r$, which completes the proof.

B.2. Proof of Proposition 3.4. The starting point is a generic error bound for PCG iterates [3, eq. (3.4)].

$$\frac{\|\mathbf{x}_t - \mathbf{x}_\star\|_{\mathbf{A}}}{\|\mathbf{x}_0 - \mathbf{x}_\star\|_{\mathbf{A}}} \leq \max_{1 \leq i \leq r} |p_t(\lambda_i)|.$$

Here, p_t is any degree- t polynomial satisfying $p_t(0) = 1$, and $\lambda_1 \geq \dots \geq \lambda_r$ are the sorted positive eigenvalues of $\mathbf{A}\hat{\mathbf{A}}^+$. In the case of even t , we construct

$$p_t(\lambda) = \prod_{i=1}^{t/2} \left(1 - \frac{\lambda}{\lambda_i}\right) \left(1 - \frac{\lambda}{\lambda_{r+1-i}}\right).$$

In the case of odd t , we set $t_* = (t+1)/2$ and construct

$$p_t(\lambda) = \left(1 - \frac{2\lambda}{\lambda_{t_*} + \lambda_{r+1-t_*}}\right) \prod_{i=1}^{(t-1)/2} \left(1 - \frac{\lambda}{\lambda_i}\right) \left(1 - \frac{\lambda}{\lambda_{r+1-i}}\right).$$

Now observe that

$$\max_{1 \leq i \leq r} |p_t(\lambda_i)| = \max_{\lfloor t/2 \rfloor \leq i \leq r+1-\lfloor t/2 \rfloor} |p_t(\lambda_i)|,$$

and the right-hand side is bounded by the product of terms

$$\max_{\lambda_i \leq \lambda \leq \lambda_{r+1-i}} \left(\frac{\lambda}{\lambda_i} - 1\right) \left(1 - \frac{\lambda}{\lambda_{r+1-i}}\right) = \frac{(\lambda_i - \lambda_{r+1-i})^2}{4\lambda_i \lambda_{r+1-i}}$$

and potentially a term

$$\max_{\lambda_{t_*} \leq \lambda \leq \lambda_{r+1-t_*}} \left|1 - \frac{2\lambda}{\lambda_{t_*} + \lambda_{r+1-t_*}}\right| = \frac{\lambda_{t_*} - \lambda_{r+1-t_*}}{\lambda_{t_*} + \lambda_{r+1-t_*}} \leq \frac{\lambda_{t_*} - \lambda_{r+1-t_*}}{2(\lambda_{t_*} \lambda_{r+1-t_*})^{1/2}},$$

where the last line follows from the arithmetic-geometric mean inequality

$$(\lambda_{t_*} \lambda_{r+1-t_*})^{1/2} \leq \frac{\lambda_{t_*} + \lambda_{r+1-t_*}}{2}.$$

In the case of even t , we can use the arithmetic-geometric mean inequality twice to write

$$\begin{aligned} & \prod_{i=1}^{t/2} \left(\frac{4\lambda_i \lambda_{r+1-i}}{(\lambda_i + \lambda_{r+1-i})^2}\right)^{2/t} + \prod_{i=1}^{t/2} \left(\frac{(\lambda_i - \lambda_{r+1-i})^2}{(\lambda_i + \lambda_{r+1-i})^2}\right)^{2/t} \\ & \leq \frac{2}{t} \sum_{i=1}^{t/2} \frac{4\lambda_i \lambda_{r+1-i} + (\lambda_i - \lambda_{r+1-i})^2}{(\lambda_i + \lambda_{r+1-i})^2} = 1, \end{aligned}$$

and consequently

$$1 + \left[\frac{\|\mathbf{x}_t - \mathbf{x}_*\|_{\mathbf{A}}}{\|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}} \right]^{2/t} \leq 1 + \prod_{i=1}^{t/2} \left(\frac{(\lambda_i - \lambda_{r+1-i})^2}{4\lambda_i \lambda_{r+1-i}}\right)^{2/t} \leq \prod_{i=1}^{t/2} \left(\frac{(\lambda_i + \lambda_{r+1-i})^2}{4\lambda_i \lambda_{r+1-i}}\right)^{2/t}.$$

In the case of odd t , we can use the generalized arithmetic-geometric mean inequality with weights $1/t, 2/t, \dots, 2/t$ twice to write

$$\begin{aligned} & \left(\frac{4\lambda_{t_*} \lambda_{r+1-t_*}}{(\lambda_{t_*} + \lambda_{r+1-t_*})^2}\right)^{1/t} \prod_{i=1}^{(t-1)/2} \left(\frac{4\lambda_i \lambda_{r+1-i}}{(\lambda_i + \lambda_{r+1-i})^2}\right)^{2/t} \\ & + \left(\frac{(\lambda_{t_*} - \lambda_{r+1-t_*})^2}{(\lambda_{t_*} + \lambda_{r+1-t_*})^2}\right)^{1/t} \prod_{i=1}^{(t-1)/2} \left(\frac{(\lambda_i - \lambda_{r+1-i})^2}{(\lambda_i + \lambda_{r+1-i})^2}\right)^{2/t} \\ & \leq \frac{1}{t} \frac{4\lambda_{t_*} \lambda_{r+1-t_*} + (\lambda_{t_*} - \lambda_{r+1-t_*})^2}{(\lambda_{t_*} + \lambda_{r+1-t_*})^2} + \frac{2}{t} \sum_{i=1}^{(t-1)/2} \frac{4\lambda_i \lambda_{r+1-i} + (\lambda_i - \lambda_{r+1-i})^2}{(\lambda_i + \lambda_{r+1-i})^2} = 1, \end{aligned}$$

and consequently

$$\begin{aligned}
& 1 + \left[\frac{\|\mathbf{x}_t - \mathbf{x}_*\|_{\mathbf{A}}}{\|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}} \right]^{2/t} \\
& \leq 1 + \left(\frac{(\lambda_{t_*} - \lambda_{r+1-t_*})^2}{4\lambda_{t_*}\lambda_{r+1-t_*}} \right)^{1/t} \prod_{i=1}^{(t-1)/2} \left(\frac{(\lambda_i - \lambda_{r+1-i})^2}{4\lambda_i\lambda_{r+1-i}} \right)^{2/t} \\
& \leq \left(\frac{(\lambda_{t_*} + \lambda_{r+1-t_*})^2}{4\lambda_{t_*}\lambda_{r+1-t_*}} \right)^{1/t} \prod_{i=1}^{t/2} \left(\frac{(\lambda_i + \lambda_{r+1-i})^2}{4\lambda_i\lambda_{r+1-i}} \right)^{2/t}.
\end{aligned}$$

Last, when r is even, we can use the arithmetic-geometric mean inequality two more times to write

$$\left(1 + \left[\frac{\|\mathbf{x}_t - \mathbf{x}_*\|_{\mathbf{A}}}{\|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}} \right]^{2/t} \right)^{t/2} \leq \frac{1}{\prod_{i=1}^r \lambda_i} \prod_{i=1}^{r/2} \left(\frac{\lambda_i + \lambda_{r+1-i}}{2} \right)^2 \leq \kappa_{\text{Kap}}.$$

Similarly when r is odd, we can use the generalized arithmetic-geometric mean inequality two more times to obtain

$$\left(1 + \left[\frac{\|\mathbf{x}_t - \mathbf{x}_*\|_{\mathbf{A}}}{\|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}} \right]^{2/t} \right)^{t/2} \leq \frac{1}{\prod_{i=1}^r \lambda_i} \lambda_{(r+1)/2} \prod_{i=1}^{(r-1)/2} \left(\frac{\lambda_i + \lambda_{r+1-i}}{2} \right)^2 \leq \kappa_{\text{Kap}}.$$

We conclude by recalling that

$$\begin{aligned}
2x &= \int_0^x 2dy \leq \int_0^x \frac{2dy}{1-y^2} = \int_0^x \left[\frac{1}{1-y} + \frac{1}{1+y} \right] dy \\
&= \left[-\log(1-y) + \log(1+y) \right]_{y=0}^{y=x} = \log\left(\frac{1+x}{1-x} \right)
\end{aligned}$$

and therefore $e^{2x} - 1 \leq \frac{2x}{1-x}$. Therefore, we make the calculation

$$\begin{aligned}
\frac{\|\mathbf{x}_t - \mathbf{x}_*\|_{\mathbf{A}}}{\|\mathbf{x}_0 - \mathbf{x}_*\|_{\mathbf{A}}} &\leq (\kappa_{\text{Kap}}^{2/t} - 1)^{t/2} \leq \left(\frac{2\log(\kappa_{\text{Kap}})/t}{1 - \log(\kappa_{\text{Kap}})/t} \right)^{t/2} \\
&= \left(\frac{2\log(\kappa_{\text{Kap}})}{2t/3 + [t/3 - \log(\kappa_{\text{Kap}})]} \right)^{t/2} \leq \left(\frac{3\log(\kappa_{\text{Kap}})}{t} \right)^{t/2}
\end{aligned}$$

on the event that $t/3 \geq \log(\kappa_{\text{Kap}})$. On the other hand, if $t/3 < \log(\kappa_{\text{Kap}})$, the result (3.2) holds vacuously. This completes the proof.

B.3. Proof of Proposition 3.6. Consider the matrix $\mathbf{M} = \hat{\mathbf{A}}^{-1/2} \mathbf{A} \hat{\mathbf{A}}^{-1/2}$ which has eigendecomposition

$$\mathbf{M} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^* \quad \text{where} \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n).$$

We need to prove the mean square error bound

$$\mathbb{E} \left| \text{tr}(\log \mathbf{M}) - \frac{1}{t} \sum_{i=1}^t \mathbf{z}_i^* (\log \mathbf{M}) \mathbf{z}_i \right|^2 \leq \frac{4\log(\kappa_{\text{Kap}})}{t}.$$

To that end, we introduce a complex Gaussian vector $\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and set $\mathbf{z} = \sqrt{n} \boldsymbol{\omega} / \|\boldsymbol{\omega}\|$. It suffices to show

$$(B.1) \quad \mathbb{E}[\mathbf{z}^*(\log \mathbf{M})\mathbf{z}] = \text{tr}(\log \mathbf{M}) \quad \text{and} \quad \text{Var}[\mathbf{z}^*(\log \mathbf{M})\mathbf{z}] \leq 4 \log(\kappa_{\text{Kap}})$$

when the Kaporin condition number is $\log(\kappa_{\text{Kap}}) \leq n$.

We first observe that rotational invariance implies $\mathbf{z} \stackrel{\mathcal{D}}{=} \mathbf{Q}^* \mathbf{z}$ so

$$\mathbf{z}^*(\log \mathbf{M})\mathbf{z} \stackrel{\mathcal{D}}{=} \mathbf{z}^*(\log \mathbf{\Lambda})\mathbf{z} = \sum_{i=1}^n \log(\lambda_i) |\mathbf{z}(i)|^2.$$

Then we check the stochastic trace estimator is unbiased.

$$(B.2) \quad \mathbb{E} \left[\sum_{i=1}^n \log(\lambda_i) |\mathbf{z}(i)|^2 \right] = \sum_{i=1}^n \log(\lambda_i) \mathbb{E} |\mathbf{z}(i)|^2 = \sum_{i=1}^n \log(\lambda_i) = \text{tr}(\log \mathbf{M}).$$

This confirms the first part of (B.2).

To bound the variance, we make a calculation using the Gaussian vector $\boldsymbol{\omega} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^n \log(\lambda_i) |\boldsymbol{\omega}(i)|^2 \right]^2 - \left[\sum_{i=1}^n \log(\lambda_i) \right]^2 = \text{Var} \left[\sum_{i=1}^n \log(\lambda_i) |\boldsymbol{\omega}(i)|^2 \right] \\ &= \sum_{i=1}^n (\log \lambda_i)^2 \text{Var} [|\boldsymbol{\omega}(i)|^2] = \sum_{i=1}^n (\log \lambda_i)^2. \end{aligned}$$

Here we have used the independence of $\boldsymbol{\omega}(i)$ variables and the identity $\text{Var} [|\boldsymbol{\omega}(i)|^2] = 1$, which holds for complex Gaussians. Since $\boldsymbol{\omega}$ has independent length and direction, we calculate

$$\mathbb{E} \left[\sum_{i=1}^n \log(\lambda_i) |\boldsymbol{\omega}(i)|^2 \right]^2 = \frac{\mathbb{E} \|\boldsymbol{\omega}\|^4}{n^2} \mathbb{E} \left[\sum_{i=1}^n \log(\lambda_i) |\mathbf{z}(i)|^2 \right]^2 = \frac{n^2 + n}{n^2} \mathbb{E} \left[\sum_{i=1}^n \log(\lambda_i) |\mathbf{z}(i)|^2 \right]^2,$$

where again we have used the fact that $\text{Var} [|\boldsymbol{\omega}(i)|^2] = 1$. By rearrangement, it follows

$$\mathbb{E} \left[\sum_{i=1}^n \log(\lambda_i) |\mathbf{z}(i)|^2 \right]^2 = \frac{n}{n+1} \left[\sum_{i=1}^n (\log \lambda_i)^2 + \left(\sum_{i=1}^n \log(\lambda_i) \right)^2 \right].$$

Subtracting the square mean (B.2) shows that

$$\begin{aligned} & \text{Var} \left[\sum_{i=1}^n \log(\lambda_i) |\mathbf{z}(i)|^2 \right] \\ & \leq \frac{n}{n+1} \left[\frac{1}{n} \sum_{i=1}^n (\log \lambda_i)^2 - \left(\sum_{i=1}^n \log(\lambda_i) \right)^2 \right] = \frac{1}{n+1} \sum_{i=1}^n (\log \lambda'_i)^2, \end{aligned}$$

where we have introduced $\lambda'_i = \lambda_i / \prod_{j=1}^n \lambda_j$ and we observe that $\prod_{i=1}^n \log(\lambda'_i) = 0$ and $\frac{1}{n} \sum_{i=1}^n \lambda'_i = \kappa_{\text{Kap}}^{1/n}$. Then recall the standard identity that

$$1 + t + \frac{1}{2} t^2 \leq e^t \quad \text{and therefore} \quad (\log u)^2 \leq 2(u - 1 - \log u).$$

We obtain the variance bound

$$\mathrm{Var} \left[\sum_{i=1}^n \log(\lambda_i) |\mathbf{z}(i)|^2 \right] \leq \frac{2n}{n+1} \sum_{i=1}^n (\lambda'_i - 1 - \log \lambda'_i) = \frac{2n^2}{n+1} (\kappa_{\mathrm{Kap}}^{1/n} - 1).$$

Last observe that $t \mapsto (e^t - 1)/t$ is strictly increasing for $t > 0$ and therefore

$$\frac{\kappa_{\mathrm{Kap}}^{1/n} - 1}{\log(\kappa_{\mathrm{Kap}})/n} \leq 2 \quad \text{if} \quad \log(\kappa_{\mathrm{Kap}}) \leq 1.42n.$$

This completes the proof.