

EXPERIMENTAL DESIGN FOR FAST LINEAR ALGEBRA

STEPHEN HUAN*, JOE GUINNESS†, MATTHIAS KATZFUSS‡, HOUMAN OWHADI§,
AND FLORIAN SCHÄFER¶

Abstract. Dense kernel matrices resulting from pointwise evaluations of a kernel function arise naturally in machine learning and statistics. Previous work in constructing sparse transport maps or sparse approximate inverse Cholesky factors for such matrices by minimizing Kullback-Leibler divergence recovers the Vecchia approximation for Gaussian processes. However, this method for Cholesky factorization relies only on geometry to construct the sparsity pattern, ignoring the conditional effect of adding an entry. In this work, we construct the sparsity pattern by leveraging a greedy selection algorithm which selects points that maximize mutual information with target points, conditional on all points selected previously. For selecting k points out of N , the naive time complexity is $\mathcal{O}(Nk^4)$, but by maintaining a partial Cholesky factor we reduce this to $\mathcal{O}(Nk^2)$. Furthermore, for multiple (m) targets we achieve a time complexity of $\mathcal{O}(Nk^2 + Nm^2 + m^3)$ which is maintained in the setting of Cholesky factorization where a selected point need not condition every target. We directly apply the selection algorithm to image classification and recovery of sparse Cholesky factors, improving upon k -th nearest neighbors in every case. Through Kullback-Leibler minimization we apply the algorithm to Cholesky factorization and Gaussian process regression, improving in high dimensional geometries as well as when preconditioning with the conjugate gradient.

Key words.

to do

AMS subject classifications.

1. Introduction. [4] test citation

2. Greedy selection for directed inference.

2.1. Conditional k -th nearest neighbors. Consider the simple regression algorithm k th-nearest neighbors (k -NN). Given a training set $X_{\text{Tr}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ and corresponding labels $\mathbf{y}_{\text{Tr}} = \{y_1, \dots, y_n\}$, the goal is to estimate the unknown label y_{Pr} of some unseen prediction point \mathbf{x}_{Pr} . Stated informally, the k -NN approach is to select the k points in X_{Tr} *most informative* about \mathbf{x}_{Pr} and combine their results.

Algorithm 2.1 Idealized k -NN regression

Given $(X_{\text{Tr}}, \mathbf{y}_{\text{Tr}}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ and \mathbf{x}_{Pr}

1. Select the k points in X_{Tr} most informative about \mathbf{x}_{Pr}
 2. Combine the labels of the selected points to generate a prediction
-

One specific approach is intuitively, points close to \mathbf{x}_{Pr} should be similar to it. So we select the k closest points in X_{Tr} to \mathbf{x}_{Pr} and pool their labels (e.g., by averaging).

Algorithm 2.2 k -NN regression

1. Select the k points $\{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}\} \subseteq X_{\text{Tr}}$ closest to \mathbf{x}_{Pr}
 2. Compute \mathbf{y}_{Pr} by $\frac{1}{k} \sum_{j=1}^k y_{i_j}$
-

*Georgia Institute of Technology

†Joe

‡Matthias

§Houman

¶Georgia Institute of Technology, S1317 CODA, 756 W Peachtree St Atlanta, GA 30332,
florian.schaefer@cc.gatech.edu,
Corresponding Author

Florian: "specific" and "intuitively" are fluff, meaning that they don't really add information. In order to achieve crisp, high-quality academic writing, it is important to try to those fluff words as much as possible. It's normal to add them out of reflex initially, so it requires active postprocessing.

However, we can generalize the notion of “closest” with the *kernel trick*, by using an arbitrary kernel function to measure similarity. For example, commonly used kernels like the Gaussian kernel and Matérn family of covariance functions are *isotropic*; they depend only on the distance between the two vectors. If such isotropic kernels monotonically decrease with distance, then selecting points based on the largest kernel similarity recovers k -NN. However, kernels need not be isotropic in general — they just need to capture some sense of “similarity”, motivating kernel k -NN.

Algorithm 2.3 Kernel k -NN regression

Given kernel function $K(\mathbf{x}, \mathbf{y})$

1. Select the k points $\{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}\} \subseteq X_{\text{Tr}}$ most similar to \mathbf{x}_{Pr}
 2. Compute \mathbf{y}_{Pr} by an average weighted by similarity
-

Stephen: not sure whether “stationary” or “isotropic” are the right word(s) to use here

Although the kernel k -NN approach is more general than its normed counterpart, it still suffers from a fundamental issue. Suppose the closest point to \mathbf{x}_{Pr} has many duplicates in the training set. Then the algorithm will select the same point multiple times, even though in some sense the duplicate point has stopped giving additional information about the prediction point. In order to fix this issue, we should be selecting new points *conditional* on the points we’ve already selected. This preserves the idealized algorithm of selecting points based on the information they tell us about the prediction point — once we’ve selected a point, conditioning on it reduces the information similar points tells us, encouraging diverse point selection.

Algorithm 2.4 Conditional kernel k -NN regression

1. Select the k points $\{\mathbf{x}_{i_1}, \dots, \mathbf{x}_{i_k}\} \subseteq X_{\text{Tr}}$ most **informative** to \mathbf{x}_{Pr}
after conditioning on all points selected beforehand
 2. Compute \mathbf{y}_{Pr} by an average weighted by **information**
-

In order to make the notions of conditioning and information precise, we need a specific framework. Kernel methods lead naturally to Gaussian processes, whose covariance matrices naturally result from kernel functions and allows us to use the rigorous statistical and information-theoretic notions of conditioning and information.

2.2. Sparse Gaussian process regression. A *Gaussian process* is a prior distribution over functions, such that for any finite set of points, the corresponding function over the points is distributed according to a multivariate Gaussian. In order to generate such a distribution over an uncountable number of points consistently, a Gaussian process is specified by a *mean function* $\mu(\mathbf{x})$ and *covariance function* or *kernel function* $K(\mathbf{x}, \mathbf{y})$. For any finite set of points $X = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, $f(X) \sim \mathcal{N}(\boldsymbol{\mu}, \Theta)$, where $\boldsymbol{\mu}_i = \mu(\mathbf{x}_i)$ and $\Theta_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$.

In order to compute a prediction at \mathbf{x}_{Pr} , we can simply condition the desired prediction \mathbf{y}_{Pr} on the observed outputs and compute the conditional expectation. We can also find the conditional variance, which will quantify the uncertainty of our prediction. If we block our covariance matrix $\Theta = \begin{pmatrix} \Theta_{\text{Tr}, \text{Tr}} & \Theta_{\text{Tr}, \text{Pr}} \\ \Theta_{\text{Pr}, \text{Tr}} & \Theta_{\text{Pr}, \text{Pr}} \end{pmatrix}$ where $\Theta_{\text{Tr}, \text{Tr}}, \Theta_{\text{Pr}, \text{Pr}}, \Theta_{\text{Tr}, \text{Pr}}, \Theta_{\text{Pr}, \text{Tr}}$ are the covariance matrices of the training data, testing data, and training and test data respectively, then the conditional expectation and

Stephen: mention sensor placement/spatial statistics perspective/literature

covariance are:

$$(2.1) \quad \mathbb{E}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}] = \boldsymbol{\mu}_{\text{Pr}} + \Theta_{\text{Pr},\text{Tr}} \Theta_{\text{Tr},\text{Tr}}^{-1} (\mathbf{y}_{\text{Tr}} - \boldsymbol{\mu}_{\text{Tr}})$$

$$(2.2) \quad \text{Cov}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}] = \Theta_{\text{Pr},\text{Pr}} - \Theta_{\text{Pr},\text{Tr}} \Theta_{\text{Tr},\text{Tr}}^{-1} \Theta_{\text{Tr},\text{Pr}}$$

For brevity of notation, we will often denote the conditional covariance matrix as

$$(2.3) \quad \Theta_{\text{Pr},\text{Pr}|\text{Tr}} := \Theta_{\text{Pr},\text{Pr}} - \Theta_{\text{Pr},\text{Tr}} \Theta_{\text{Tr},\text{Tr}}^{-1} \Theta_{\text{Tr},\text{Pr}}$$

When conditioning on multiple sets, the sets are given in order of computation. Although the resulting covariance matrix is the same, a different order of conditioning means different intermediate results in repeated application of (2.2). In general,

$$(2.4) \quad \Theta_{I,J|K_1,K_2,\dots,K_n} := \text{Cov}[\mathbf{y}_I, \mathbf{y}_J \mid \mathbf{y}_{K_1 \cup K_2 \cup \dots \cup K_n}]$$

denotes the covariance between the variables in index sets I and J , conditional on the variables in K_1, K_2, \dots, K_n . Note that by the quotient rule of Schur complementation:

$$(2.5) \quad \Theta_{I,J|K_1\dots n} = \Theta_{I,J|K_1\dots n-1} - \Theta_{I,K_n|K_1\dots n-1} \Theta_{K_n,K_n|K_1\dots n-1}^{-1} \Theta_{K_n,J|K_1\dots n-1}$$

Note that calculating the posterior mean and variance requires inverting the training covariance matrix $\Theta_{\text{Tr},\text{Tr}}$, which costs $\mathcal{O}(N^3)$, where N is the number of training points. This scaling is prohibitive for large datasets, so many *sparse* Gaussian process regression techniques have been proposed. These methods often focus on selecting a subset of the training data that is most informative about the prediction points, which naturally aligns with our k -NN perspective. If s points are selected out of the N , then the inversion will cost $\mathcal{O}(s^3)$, which could be substantially cheaper if s is significantly smaller than N . The question is then how to select as few points as possible while maintaining predictive accuracy.

2.3. Problem: optimal selection. The natural criterion justified from the k -NN perspective is to maximize the *mutual information* between the selected points and the target point for prediction. Such information-theoretic objectives have seen success in the spatial statistics community [1], who use such criteria to determine the best locations to place sensors in a Gaussian process regression context. The mutual information, or *information gain* is defined as

$$(2.6) \quad \text{MI}[\mathbf{y}_{\text{Pr}}; \mathbf{y}_{\text{Tr}}] = \mathbb{H}[\mathbf{y}_{\text{Pr}}] - \mathbb{H}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}]$$

We can use the fact that the entropy of a multivariate Gaussian is monotonically increasing with the log determinant of its covariance matrix to efficiently compute these entropies. Because the entropy of \mathbf{y}_{Pr} is constant, maximizing the mutual information is equivalent to minimizing the conditional entropy. From (2.2) we see that minimizing the conditional entropy is equivalent to minimizing the log determinant of the posterior covariance matrix. Note that for a single predictive point, this is monotonic with its variance. So another justification is that we are reducing the *conditional variance* of the desired point as much as possible. In particular, because our estimator is the conditional expectation (2.1), it is unbiased because $\mathbb{E}[\mathbb{E}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}]] = \mathbb{E}[\mathbf{y}_{\text{Pr}}]$. Because it is unbiased, its expected mean squared error is simply the conditional variance since $\mathbb{E}[(\mathbf{y}_{\text{Pr}} - \mathbb{E}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}])^2 \mid \mathbf{y}_{\text{Tr}}] = \text{Var}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}]$ where the expectation is taken under conditioning because of the assumption that \mathbf{y}_{Pr} is distributed according to the Gaussian process. So maximizing the mutual information is equivalent to minimizing the conditional variance which is in turn equivalent to minimizing the expected mean

Stephen: I is overloaded to be identity, an index set, and mutual information/information gain. I should refer to identity and the other two must be disambiguated.

Stephen: cite sparse Gaussian regression papers

squared error of the prediction. Another perspective on the objective can be derived from comparing the mutual information to the EV-VE identity, which states

$$\begin{aligned} H[\mathbf{y}_{\text{Pr}}] &= H[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}] + \text{MI}[\mathbf{y}_{\text{Pr}}; \mathbf{y}_{\text{Tr}}] \\ \text{Var}[\mathbf{y}_{\text{Pr}}] &= \text{E}[\text{Var}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}]] + \text{Var}[\text{E}[\mathbf{y}_{\text{Pr}} \mid \mathbf{y}_{\text{Tr}}]] \end{aligned}$$

On the left hand side, entropy is monotone with variance. On the right hand side, the expectation of the conditional variance can be interpreted to be the fluctuation of the prediction point after conditioning, and is monotone with the conditional entropy. Because the expectation of conditional variance and variance of conditional expectation add to a constant, minimizing the expectation of the conditional variance is equivalent to maximizing the variance of conditional expectation, which we see corresponds to the mutual information term. Supposing \mathbf{y}_{Pr} was independent of \mathbf{y}_{Tr} , then the conditional expectation becomes simply the expectation, whose variance is 0. Thus, the variance of the conditional expectation can be interpreted to be the “information” shared between \mathbf{y}_{Pr} and \mathbf{y}_{Tr} , as the larger it is, the more the prediction for \mathbf{y}_{Pr} (the conditional expectation) depends on the observed results of \mathbf{y}_{Tr} .

2.4. A greedy approach. We now consider how to efficiently minimize the conditional variance objective using a greedy approach. At each iteration, we pick the training point which most reduces the conditional variance of the prediction point. Let $I = \{i_1, i_2, \dots, i_j\}$ be the set of indexes of training points selected already. Let the prediction point have index $n+1$, the last index. For a candidate index k , we update the covariance matrix after conditioning on y_k , in addition to the indices already selected according to (2.2):

$$\begin{aligned} \Theta_{:,|I,k} &= \Theta_{:,|I} - \Theta_{:,k|I} \Theta_{k,k|I}^{-1} \Theta_{k,|I} \\ (2.7) \quad &= \Theta_{:,|I} - \frac{\Theta_{:,k|I} \Theta_{:,k|I}^\top}{\Theta_{k,k|I}} \end{aligned}$$

We see that conditioning on a new entry is a rank-one update on the current covariance matrix $\Theta_{|I}$, given by the vector $\mathbf{u} = \frac{\Theta_{:,k|I}}{\sqrt{\Theta_{k,k|I}}}$. Thus, the amount that the variance of \mathbf{y}_{Pr} will decrease after selecting k is given by u_{n+1}^2 , or

$$(2.8) \quad \frac{\text{Cov}[\mathbf{y}_{\text{Tr}}[k], \mathbf{y}_{\text{Pr}} \mid I]^2}{\text{Var}[\mathbf{y}_{\text{Tr}}[k] \mid I]} = \frac{\Theta_{k,n+1|I}^2}{\Theta_{k,k|I}} = \text{Var}[\mathbf{y}_{\text{Pr}} \mid I] \text{Corr}[\mathbf{y}_{\text{Tr}}[k], \mathbf{y}_{\text{Pr}} \mid I]^2$$

For each candidate k , we need to keep track of its conditional variance and conditional covariance with the prediction point after conditioning on the points already selected to compute (2.8). We then simply choose the candidate with the largest decrease in predictive variance. To keep track of the conditional variance and covariance, we can simply start with the initial values given by Θ_{kk} and Θ_{nk} and update after selecting an index j . We compute \mathbf{u} for j directly according to (2.2) and update k 's conditional variance by subtracting u_k^2 and update its conditional covariance by subtracting $u_k u_{n+1}$.

In order to efficiently compute \mathbf{u} , we rely on two main strategies. The direct method is to keep track of $\Theta_{I,I}^{-1}$, or the precision of the selected entries, and update the precision every time a new index is added to I . This can be done efficiently in $\mathcal{O}(s^2)$, see Appendix A.1. Once $\Theta_{I,I}^{-1}$ has been computed, \mathbf{u} is computed trivially

according to (2.2). For each of the s rounds of selection, it takes s^2 to update the precision, and costs Ns to compute \mathbf{u} , costing $\mathcal{O}(Ns^2 + s^3) = \mathcal{O}(Ns^2)$ overall.

The second approach is to take advantage of the quotient rule of Schur complementation. Stated statistically, the quotient rule states that conditioning on I and then conditioning on J is the same as conditioning on $I \cup J$. We then remind ourselves that Cholesky factorization can be viewed as iterative conditioning:

Re-writing the joint covariance matrix,

$$(2.9) \quad \begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \Theta_{2,1} & \Theta_{2,2} \end{pmatrix} = \begin{pmatrix} I & 0 \\ \Theta_{2,1}\Theta_{1,1}^{-1} & I \end{pmatrix} \begin{pmatrix} \Theta_{1,1} & 0 \\ 0 & \Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2} \end{pmatrix} \begin{pmatrix} I & \Theta_{1,1}^{-1}\Theta_{1,2} \\ 0 & I \end{pmatrix}$$

so we see that the Cholesky factorization of the joint covariance Θ is

$$(2.10) \quad \begin{aligned} \text{chol}(\Theta) &= \begin{pmatrix} I & 0 \\ \Theta_{2,1}\Theta_{1,1}^{-1} & I \end{pmatrix} \begin{pmatrix} \text{chol}(\Theta_{1,1}) & 0 \\ 0 & \text{chol}(\Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2}) \end{pmatrix} \\ &= \begin{pmatrix} \text{chol}(\Theta_{1,1}) & 0 \\ \Theta_{2,1}\text{chol}(\Theta_{1,1})^{-\top} & \text{chol}(\Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2}) \end{pmatrix} \end{aligned}$$

Here $\Theta_{2,1}\Theta_{1,1}^{-1}$ corresponds to the conditional expectation in (2.1) and $\Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2}$ corresponds to the conditional covariance in (2.2). Thus, we see that Cholesky factorization is iteratively conditioning the Gaussian process. From the iterative conditioning perspective, the i th column of the Cholesky factor corresponds precisely to the corresponding \mathbf{u} for i since a iterative sequence of conditioning on $i_1, i_2 \dots$ is equivalent to conditioning on I by the quotient rule.

The Cholesky factorization can be efficiently computed without excess dependence on N with left-looking, so the conditioning only happens when we need it. For each of the s rounds of selection, it costs $\mathcal{O}(Ns)$ to compute the next column of the Cholesky factorization, for a total cost of $\mathcal{O}(Ns^2)$, matching the time complexity of the explicit precision approach.

Algorithm 2.5 Point selection by explicit precision

Input: $\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Pr}}, K(\cdot, \cdot), s$ **Output:** I

```

1:  $n \leftarrow |\mathbf{x}_{\text{Tr}}|$ 
2:  $\mathbf{x} \leftarrow \begin{pmatrix} \mathbf{x}_{\text{Tr}} \\ \mathbf{x}_{\text{Pr}} \end{pmatrix}$ 
3:  $I \leftarrow \emptyset$ 
4:  $-I \leftarrow \{1, 2, \dots, n\}$ 
5:  $\Theta_{I,I}^{-1} \leftarrow \mathbb{R}^{0 \times 0}$ 
6:  $\Theta_{\text{Tr}, \text{Pr}|I} \leftarrow K(\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Pr}})$ 
7:  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I}) \leftarrow \text{diag}(K(\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Tr}}))$ 
8: while  $|-I| > 0$  and  $|I| < s$  do
9:    $k \leftarrow \max_{j \in -I} \frac{\Theta_{j, \text{Pr}|I}^2}{\Theta_{jj|I}}$ 
10:   $I \leftarrow I \cup \{k\}$ 
11:   $-I \leftarrow -I - \{k\}$ 
12:   $\mathbf{v} \leftarrow \Theta_{I,I}^{-1} K(\mathbf{x}_{\text{Tr}}[I - \{k\}], \mathbf{x}_{\text{Tr}}[k])$ 
13:   $\Theta_{I,I}^{-1} \leftarrow \begin{pmatrix} \Theta_{I,I}^{-1} + \frac{\mathbf{v}\mathbf{v}^\top}{\Theta_{kk|I}} & \frac{-\mathbf{v}}{\Theta_{kk|I}} \\ \frac{-\mathbf{v}^\top}{\Theta_{kk|I}} & \frac{1}{\Theta_{kk|I}} \end{pmatrix}$ 
14:   $\Theta_{:,k|I} \leftarrow K(\mathbf{x}, \mathbf{x}_k) - K(\mathbf{x}, \mathbf{x}_{I-\{k\}})\mathbf{v}$ 
15:   $\mathbf{u} \leftarrow \frac{\Theta_{:,k|I}}{\sqrt{\Theta_{kk|I}}}$ 
16:  for  $j \in -I$  do
17:     $\Theta_{jj|I} \leftarrow \Theta_{jj|I} - \mathbf{u}_j^2$ 
18:     $\Theta_{j, \text{Pr}|I} \leftarrow \Theta_{j, \text{Pr}|I} - \mathbf{u}_j \mathbf{u}_{n+1}$ 
19:  end for
20: end while
21: return  $I$ 

```

Algorithm 2.6 Point selection by Cholesky factorization

Input: $\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Pr}}, K(\cdot, \cdot), s$ **Output:** I

```

1:  $n \leftarrow |\mathbf{x}_{\text{Tr}}|$ 
2:  $\mathbf{x} \leftarrow \begin{pmatrix} \mathbf{x}_{\text{Tr}} \\ \mathbf{x}_{\text{Pr}} \end{pmatrix}$ 
3:  $I \leftarrow \emptyset$ 
4:  $-I \leftarrow \{1, 2, \dots, n\}$ 
5:  $L \leftarrow \mathbf{0}^{(n+1) \times s}$ 
6:  $\Theta_{\text{Tr}, \text{Pr}|I} \leftarrow K(\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Pr}})$ 
7:  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I}) \leftarrow \text{diag}(K(\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Tr}}))$ 
8: while  $|-I| > 0$  and  $|I| < s$  do
9:    $k \leftarrow \max_{j \in -I} \frac{\Theta_{j, \text{Pr}|I}^2}{\Theta_{jj|I}}$ 
10:   $I \leftarrow I \cup \{k\}$ 
11:   $-I \leftarrow -I - \{k\}$ 
12:   $i \leftarrow |I|$ 
13:   $L_{:,i} \leftarrow K(\mathbf{x}, \mathbf{x}_k) - L_{:,i-1} L_{k,i-1}^\top$ 
14:   $L_{:,i} \leftarrow \frac{L_{:,i}}{\sqrt{L_{k,i}}}$ 
15:  for  $j \in -I$  do
16:     $\Theta_{jj|I} \leftarrow \Theta_{jj|I} - L_{j,i}^2$ 
17:     $\Theta_{j, \text{Pr}|I} \leftarrow \Theta_{j, \text{Pr}|I} - L_{j,i} L_{n,i}$ 
18:  end for
19: end while
20: return  $I$ 

```

While both approaches have the same time complexity, the explicit precision algorithm uses $\mathcal{O}(s^2)$ space to store the precision while the Cholesky factorization uses $\mathcal{O}(Ns)$ to store the first s columns of the Cholesky factorization of Θ , which is always more memory than the precision ($N > s$). Both algorithms use an additional $\mathcal{O}(N)$ space to store the conditional variances and covariances. Although the precision algorithm uses asymptotically less memory than the Cholesky algorithm, the Cholesky algorithm is preferred for better constant factor speed and ease of implementation.

Once the indices have been computed according to [Algorithm 2.5](#) or [Algorithm 2.6](#), inferring the conditional mean and covariance of the unknown data can be done directly according to (2.1) and (2.2) in time $\mathcal{O}(s^3)$ using [Algorithm 2.7](#).

This algorithm is in fact the covariance equivalent of the signal recovery algorithm orthogonal matching pursuit (OMP) [5], a connection elaborated in [Appendix A.9](#).

2.5. Supernodes and blocked selection. We now consider how to efficiently deal with multiple prediction points. The first question is how to generalize the previous objective for a single point (2.8) to multiple points. Following the same mutual information justification as before, a natural criterion is to minimize the log determinant of the prediction points' covariance matrix after conditioning on the

Stephen: what is a good way to say non-asymptotically faster?

Algorithm 2.7 Gaussian process inference by selection**Input:** $\mathbf{x}_{\text{Tr}}, \mathbf{y}_{\text{Tr}}, \mathbf{x}_{\text{Pr}}, K(\cdot, \cdot), s$ **Output:** $\mathbb{E}[\mathbf{y}_{\text{Pr}} | \mathbf{y}_{\text{Tr}}], \text{Cov}[\mathbf{y}_{\text{Pr}} | \mathbf{y}_{\text{Tr}}]$

- 1: Compute I using [Algorithm 2.5](#) or [Algorithm 2.6](#)
- 2: $\Theta_{\text{Tr}, \text{Tr}} \leftarrow K(\mathbf{x}_{\text{Tr}}[I], \mathbf{x}_{\text{Tr}}[I])$
- 3: $\Theta_{\text{Pr}, \text{Pr}} \leftarrow K(\mathbf{x}_{\text{Pr}}, \mathbf{x}_{\text{Pr}})$
- 4: $\Theta_{\text{Tr}, \text{Pr}} \leftarrow K(\mathbf{x}_{\text{Tr}}[I], \mathbf{x}_{\text{Pr}})$
- 5: $\mathbb{E}[\mathbf{y}_{\text{Pr}} | \mathbf{y}_{\text{Tr}}] \leftarrow \Theta_{\text{Pr}, \text{Tr}} \Theta_{\text{Tr}, \text{Tr}}^{-1} \mathbf{y}_{\text{Tr}}[I]$
- 6: $\text{Cov}[\mathbf{y}_{\text{Pr}} | \mathbf{y}_{\text{Tr}}] \leftarrow \Theta_{\text{Pr}, \text{Pr}} - \Theta_{\text{Tr}, \text{Pr}}^{\top} \Theta_{\text{Tr}, \text{Tr}}^{-1} \Theta_{\text{Tr}, \text{Pr}}$
- 7: **return** $\mathbb{E}[\mathbf{y}_{\text{Pr}} | \mathbf{y}_{\text{Tr}}], \text{Cov}[\mathbf{y}_{\text{Pr}} | \mathbf{y}_{\text{Tr}}]$

selected points, or $\log\det(\Theta_{\text{Pr}, \text{Pr}|I})$. This objective, known as D-optimal [1], has many intuitive interpretations — for example, as the volume of the region of uncertainty or as the scaling factor in the density function for the Gaussian process.

We now need to be able to efficiently compute the effect of selecting an index k on the log determinant. From (2.7), we know that selecting an index is a rank-one update on the prediction covariance. Using the matrix determinant lemma,

(2.11)

$$\begin{aligned} \log\det(\Theta_{\text{Pr}, \text{Pr}|I \cup \{k\}}) &= \log\det\left(\Theta_{\text{Pr}, \text{Pr}|I} - \frac{\Theta_{\text{Pr}, k|I} \Theta_{\text{Pr}, k|I}^{\top}}{\Theta_{kk|I}}\right) \\ &= \log\det(\Theta_{\text{Pr}, \text{Pr}|I}) + \log\left(1 - \frac{\Theta_{\text{Pr}, k|I}^{\top} \Theta_{\text{Pr}, \text{Pr}|I}^{-1} \Theta_{\text{Pr}, k|I}}{\Theta_{kk|I}}\right) \end{aligned}$$

Focusing on the second term, we can turn the quadratic form into conditioning:

$$= \log\det(\Theta_{\text{Pr}, \text{Pr}|I}) + \log\left(\frac{\Theta_{kk|I} - \Theta_{k, \text{Pr}|I} \Theta_{\text{Pr}, \text{Pr}|I}^{-1} \Theta_{\text{Pr}, k|I}}{\Theta_{kk|I}}\right)$$

By the quotient rule, we combine the conditioning:

$$(2.12) \quad = \log\det(\Theta_{\text{Pr}, \text{Pr}|I}) + \log\left(\frac{\Theta_{kk|I, \text{Pr}}}{\Theta_{kk|I}}\right)$$

The greedy objective (2.12) tells us that to minimize the log determinant, we can simply select the index k with the smallest ratio between the conditional variance after conditioning on the previously selected points as well as the prediction points, and the conditional variance after just conditioning on the selected points. Intuitively this tells us that we can place sensors *backwards*, where we imagine placing sensors at the *prediction points* instead of the candidates. We then measure the conditional variance at a candidate, and pick the candidate whose conditional variance decreases the most (relative from what it started out as). Intuitively, these candidates are likely to give information about the prediction points, because the prediction points give information about the candidate.

Re-writing the objective in this way also gives an efficient algorithm to compute the necessary quantities. We condition on the prediction points essentially the same as described in [subsection 2.4](#), by simply maintaining two structures instead of one, one for the conditional variance after conditioning on the previously selected points, and the other for the conditional variance after also conditioning on the prediction points. By the quotient rule, the order of conditioning does not matter as long as the order is

Stephen: move to appendix

Stephen: this maybe should be changed to information gain by flipping the sign

consistent. For the second structure, we therefore condition on the prediction points *first* before any points have been selected. We again have two strategies, one which explicitly maintains precisions and the other which relies on Cholesky factorization.

For the precision algorithm, using (2.2) directly, for m prediction points it costs $\mathcal{O}(m^3)$ to compute $\Theta_{\text{Pr},\text{Pr}}^{-1}$ and then $\mathcal{O}(Nm^2)$ to compute $\Theta_{kk|\text{Pr}}$ for the N candidates k . For each of the s rounds of selecting candidates, it costs s^2 and m^2 to update the precisions $\Theta_{I,I}^{-1}$ and $\Theta_{\text{Pr},\text{Pr}}^{-1}$ respectively, where the details of efficiently updating $\Theta_{\text{Pr},\text{Pr}}^{-1}$ after the rank-one update in (2.11) are given in Appendix A.3. Given the precisions, $\mathbf{u} = \frac{\Theta_{:,k|I}}{\sqrt{\Theta_{kk|I}}}$ and $\mathbf{u}_{\text{Pr}} = \frac{\Theta_{:,k|I,\text{Pr}}}{\sqrt{\Theta_{kk|I,\text{Pr}}}}$ are computed as usual according to (2.2) in time Ns and Nm . Finally, for each candidate j , the conditional variance $\Theta_{jj|I}$ is updated by subtracting u_j^2 , the conditional covariance $\Theta_{\text{Pr},k|I}$ is updated for each prediction point index c each by subtracting $u_j u_c$, and the conditional variance $\Theta_{jj|I,\text{Pr}}$ is updated by subtracting $u_{\text{Pr},j}^2$. The total time complexity after simplification is $\mathcal{O}(Ns^2 + Nm^2 + m^3)$.

For the Cholesky algorithm, two Cholesky factorizations are stored. We first compute the Cholesky factorization after selecting each prediction point, for a cost of $(N+m)m$ for each of the m columns. We then begin selecting candidates, which requires updating both Cholesky factors in time $(N+m)(m+s)$ which is dominated by updating the preconditioned Cholesky factor. The columns of the Cholesky factors correspond precisely to \mathbf{u} and \mathbf{u}_{Pr} and both conditional variances $\Theta_{jj|I}$ and $\Theta_{jj|I,\text{Pr}}$ can be computed as above. The conditional covariances do not need to be computed. Over s rounds the total time complexity is $\mathcal{O}((N+m)m^2 + s(N+m)(m+s))$ which simplifies to $\mathcal{O}(Ns^2 + Nm^2 + m^3)$.

Although both approaches have the same time complexity, like the single point case they differ in memory usage. The explicit precision requires $\mathcal{O}(s^2 + m^2)$ memory to store both precisions, as well as $\mathcal{O}(Nm)$ memory to store the conditional covariances. The Cholesky algorithm, on the other hand, requires $\mathcal{O}((n+m)(m+s))$ to store the first $m+s$ columns of the Cholesky factorization of the joint covariance matrix between training and prediction points, which simplifies to $\mathcal{O}(Ns + Nm + m^2)$. Comparing the memory usage, they are the same except for s^2 versus Ns , so the Cholesky algorithm again uses more memory than the explicit precision algorithm.

Algorithm 2.8 Multiple prediction point selection by explicit precision

Input: $\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Pr}}, K(\cdot, \cdot), s$ **Output:** I

```

1:  $n \leftarrow |\mathbf{x}_{\text{Tr}}|$ 
2:  $m \leftarrow |\mathbf{x}_{\text{Pr}}|$ 
3:  $\mathbf{x} \leftarrow \begin{pmatrix} \mathbf{x}_{\text{Tr}} \\ \mathbf{x}_{\text{Pr}} \end{pmatrix}$ 
4:  $I \leftarrow \emptyset$ 
5:  $-I \leftarrow \{1, 2, \dots, n\}$ 
6:  $\Theta_{I,I}^{-1} \leftarrow \mathbb{R}^{0 \times 0}$ 
7:  $\Theta_{\text{Pr}, \text{Pr}|I}^{-1} \leftarrow K(\mathbf{x}_{\text{Pr}}, \mathbf{x}_{\text{Pr}})$ 
8:  $\Theta_{\text{Tr}, \text{Pr}|I} \leftarrow K(\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Pr}})$ 
9:  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I}) \leftarrow \text{diag}(K(\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Tr}}))$ 
10:  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I, \text{Pr}}) \leftarrow \text{diag}(\Theta_{\text{Tr}, \text{Tr}|I})$ 
     $\quad - \text{diag}(\Theta_{\text{Tr}, \text{Pr}|I} \Theta_{\text{Pr}, \text{Pr}|I}^{-1} \Theta_{\text{Pr}, \text{Tr}|I})$ 
11: while  $|-I| > 0$  and  $|I| < s$  do
12:    $k \leftarrow \min_{j \in -I} \frac{\Theta_{jj|I, \text{Pr}}}{\Theta_{jj|I}}$ 
13:    $I \leftarrow I \cup \{k\}$ 
14:    $-I \leftarrow -I - \{k\}$ 
15:    $\mathbf{v} \leftarrow \Theta_{I,I}^{-1} K(\mathbf{x}_{\text{Tr}}[I - \{k\}], \mathbf{x}_{\text{Tr}}[k])$ 
16:    $\Theta_{I,I}^{-1} \leftarrow \begin{pmatrix} \Theta_{I,I}^{-1} + \frac{\mathbf{v}\mathbf{v}^\top}{\Theta_{kk|I}} & \frac{-\mathbf{v}}{\Theta_{kk|I}} \\ \frac{-\mathbf{v}^\top}{\Theta_{kk|I}} & \frac{1}{\Theta_{kk|I}} \end{pmatrix}$ 
17:    $\mathbf{w} \leftarrow \Theta_{\text{Pr}, \text{Pr}|I}^{-1} \Theta_{k, \text{Pr}|I}^\top$ 
18:    $\Theta_{\text{Pr}, \text{Pr}|I}^{-1} \leftarrow \Theta_{\text{Pr}, \text{Pr}|I}^{-1} + \frac{\mathbf{w}\mathbf{w}^\top}{\Theta_{kk|I, \text{Pr}}}$ 
19:    $\Theta_{:,k|I} \leftarrow K(\mathbf{x}, \mathbf{x}_k) - K(\mathbf{x}, \mathbf{x}_{I - \{k\}}) \mathbf{v}$ 
20:    $\Theta_{:,k|I, \text{Pr}} \leftarrow \Theta_{:,k|I} - \Theta_{:, \text{Pr}|I} \mathbf{w}$ 
21:    $\mathbf{u} \leftarrow \frac{\Theta_{:,k|I}}{\sqrt{\Theta_{kk|I}}}$ 
22:    $\mathbf{u}_{\text{Pr}} \leftarrow \frac{\Theta_{:,k|I, \text{Pr}}}{\sqrt{\Theta_{kk|I, \text{Pr}}}}$ 
23:   for  $j \in -I$  do
24:      $\Theta_{jj|I} \leftarrow \Theta_{jj|I} - \mathbf{u}_j^2$ 
25:      $\Theta_{jj|I, \text{Pr}} \leftarrow \Theta_{jj|I, \text{Pr}} - (\mathbf{u}_{\text{Pr}})_j^2$ 
26:     for  $c \in \{1, 2, \dots, m\}$  do
27:        $\Theta_{j, \text{Pr}[c]|I} \leftarrow \Theta_{j, \text{Pr}[c]|I} - \mathbf{u}_j \mathbf{u}_{n+c}$ 
28:     end for
29:   end for
30: end while
31: return  $I$ 

```

Algorithm 2.9 Multiple prediction point selection by Cholesky factorization

Input: $\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Pr}}, K(\cdot, \cdot), s$ **Output:** I

```

1:  $n \leftarrow |\mathbf{x}_{\text{Tr}}|$ 
2:  $m \leftarrow |\mathbf{x}_{\text{Pr}}|$ 
3:  $\mathbf{x} \leftarrow \begin{pmatrix} \mathbf{x}_{\text{Tr}} \\ \mathbf{x}_{\text{Pr}} \end{pmatrix}$ 
4:  $I \leftarrow \emptyset$ 
5:  $-I \leftarrow \{1, 2, \dots, n\}$ 
6:  $L \leftarrow \mathbf{0}^{(n+m) \times s}$ 
7:  $L_{\text{Pr}} \leftarrow \mathbf{0}^{(n+m) \times (s+m)}$ 
8:  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I}) \leftarrow \text{diag}(K(\mathbf{x}_{\text{Tr}}, \mathbf{x}_{\text{Tr}}))$ 
9:  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I, \text{Pr}}) \leftarrow \text{diag}(\Theta_{\text{Tr}, \text{Tr}|I})$ 
10: for  $i \in \{1, 2, \dots, m\}$  do
11:   Update  $L_{\text{Pr}}$  and  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I, \text{Pr}})$ 
    with  $k = n + i$  by Algorithm 2.10.
12: end for
13: while  $|-I| > 0$  and  $|I| < s$  do
14:    $k \leftarrow \max_{j \in -I} \frac{\Theta_{j, \text{Pr}|I}}{\Theta_{jj|I}}$ 
15:    $I \leftarrow I \cup \{k\}$ 
16:    $-I \leftarrow -I - \{k\}$ 
17:    $i \leftarrow |I|$ 
18:   Update  $L$  and  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I})$  by
    Algorithm 2.10.
19:   Update  $L_{\text{Pr}}$  and  $\text{diag}(\Theta_{\text{Tr}, \text{Tr}|I, \text{Pr}})$ 
    with  $i = i + m$  by Algorithm 2.10.
20: end while
21: return  $I$ 

```

Algorithm 2.10 Update Cholesky factor

Input: $\mathbf{x}, K(\cdot, \cdot), i, k, L, \text{diag}(\Theta)$ **Output:** $L_{:,i}, \text{diag}(\Theta_{|k})$

```

 $n \leftarrow |\text{diag}(\Theta)|$ 
 $L_{:,i} \leftarrow K(\mathbf{x}, \mathbf{x}_k) - L_{:,i-1} L_{k,i-1}^\top$ 
 $L_{:,i} \leftarrow \frac{L_{:,i}}{\sqrt{L_{k,i}}}$ 
for  $j \in \{1, 2, \dots, n\}$  do
   $\Theta_{jj} \leftarrow \Theta_{jj} - L_{j,i}^2$ 
end for

```

2.6. Near optimality by empirical submodularity. If the objective satisfies the property of submodularity, then it is guaranteed that the greedy algorithm produces an objective within a $1 - \frac{1}{e}$ of the optimal objective. Unfortunately the information gain objective is not submodular in general, see [Appendix A.10](#). However,

it is possible to empirically observe for a particular point set and choice of kernel function whether it is submodular. We note that one-dimensional points with a Matérn kernel function seems to be submodular, but not for any higher dimension.

3. Greedy selection for *global* approximation by KL-minimization. We have a covariance matrix Θ and wish to compute the Cholesky factorization of Θ into a lower triangular factor L such that $\Theta = LL^\top$. This can be done in $\mathcal{O}(N^3)$ with standard algorithms, which is often prohibitive. Recall the problem of inference in Gaussian process regression as described in [subsection 2.2](#) also took $\mathcal{O}(N^3)$ to invert the covariance matrix Θ . Thus, similar to Gaussian process regression, we will use *sparsity* to mitigate the computational cost. In fact, we will be able to reuse our previous algorithms [Algorithms 2.5](#) and [2.8](#) on each column of the Cholesky factorization.

Stephen: justify importance/downstream applications of Cholesky factorization

We will first compute the Cholesky factorization of Θ^{-1} , also known as the *precision matrix*, and use the resulting sparse factorization to efficiently compute an approximation for Θ . Because the precision matrix encodes the distribution of the full conditionals, the (i, j) th entry of the precision matrix is 0 if and only if the variables x_i and x_j are conditionally independent, conditional on the rest of the variables. Thus, the precision matrix Θ^{-1} can be sparse as a result of conditional independence even if the original covariance matrix Θ is dense. It therefore makes sense to attempt to approximately “sparsify” Θ^{-1} instead of Θ with iterated conditioning.

Because of sparsity, we can only get an approximate Cholesky factor L , \hat{L} belonging to a pre-specified sparsity pattern S — a set of (row, column) indices that are allowed to be nonzero. In order to measure the performance of the estimator, we treat the matrices as covariance matrices of centered Gaussian processes (mean $\mathbf{0}$). In order to compare the resulting distributions, we use the *KL divergence* according to [\[3\]](#), or the expected difference in log-densities:

$$(3.1) \quad L := \operatorname{argmin}_{\hat{L} \in S} \mathbb{D}_{\text{KL}} \left(\mathcal{N}(\mathbf{0}, \Theta) \parallel \mathcal{N}(\mathbf{0}, (\hat{L}\hat{L}^\top)^{-1}) \right)$$

Note that here we are computing the Cholesky factorization of Θ^{-1} . Surprisingly enough, it is possible to exactly compute L . First, we re-write the KL-divergence:

$$(3.2) \quad 2\mathbb{D}_{\text{KL}} \left(\mathcal{N}(\mathbf{0}, \Theta_1) \parallel \mathcal{N}(\mathbf{0}, \Theta_2) \right) = \operatorname{trace}(\Theta_2^{-1}\Theta_1) + \log\det(\Theta_2) - \log\det(\Theta_1) - N$$

where Θ_1 and Θ_2 are both of size $N \times N$. See [Appendix B.1](#) for details.

THEOREM 3.1. [\[3\]](#). *The non-zero entries of the i th column of L in (3.1) are:*

$$(3.3) \quad L_{s_i, i} = \frac{\Theta_{s_i, s_i}^{-1} \mathbf{e}_1}{\sqrt{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}}$$

Plugging the optimal L [\(3.3\)](#) back into the KL divergence [\(3.2\)](#), we obtain:

$$(3.4) \quad 2\mathbb{D}_{\text{KL}} \left(\mathcal{N}(\mathbf{0}, \Theta) \parallel \mathcal{N}(\mathbf{0}, (LL^\top)^{-1}) \right) = \sum_{i=1}^N [\log((\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1)^{-1})] - \log\det(\Theta)$$

See [Appendix B.2](#) for details. In particular, it is important which direction the KL divergence is or else cancellation of the $\operatorname{trace}(\Theta_2^{-1}\Theta_1)$ term may not occur.

In order to maximize [\(3.4\)](#), we can ignore $\log\det(\Theta)$ since it does not depend on L and maximize over each column independently, since each term in the sum

only depends on a single column. We want to minimize $(\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1)^{-1}$, the term corresponding to the diagonal entry in the inverse of the submatrix of Θ corresponding to the entries we've taken. We can give this value statistical interpretation by using the fact that marginalization in covariance is conditioning in precision.

$$(3.5) \quad \Theta_{1,1|2} = ((\Theta^{-1})_{1,1})^{-1}$$

where Θ is blocked according to

$$(3.6) \quad \Theta = \begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \Theta_{2,1} & \Theta_{2,2} \end{pmatrix}$$

Thus, we see that

$$(3.7) \quad (\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1)^{-1} = ((\Theta_{s_i, s_i}^{-1})_{11})^{-1} = \Theta_{ii|s_i - \{i\}}$$

So our objective on each column is to minimize the conditional variance of the i th variable, conditional on the entries we've selected — s_i contains i to begin with, so $s_i - \{i\}$ is the selected entries. We can therefore use algorithm [Algorithm 2.5](#) directly on each column, where the prediction point is the i variable and the number of points selected is the number of nonzero entries per column. The only difference is that the candidates is limited to indices lower than i , that is, candidate indices k such that $k > i$ to maintain the lower triangularity of L . Once s_i has been computed for each i , L can be constructed according to [Theorem 3.1](#). Each column costs $\mathcal{O}(s^3)$ to compute Θ_{s_i, s_i}^{-1} for a total cost of $\mathcal{O}(Ns^3)$ for the N columns of L .

3.1. Aggregated sparsity pattern. We can also use the Gaussian process regression viewpoint to efficiently aggregate multiple columns, that is, to use the same sparsity pattern for multiple columns. We denote aggregating the column indices i_1, \dots, i_m into the same group as $\tilde{i} = \{i_1, i_2, \dots, i_m\}$, letting $s_{\tilde{i}} = \bigcup_{i \in \tilde{i}} s_i$ be the aggregated sparsity pattern, and letting $\tilde{s} = s_{\tilde{i}} - \tilde{i}$ be the set of selected entries excluding the diagonal entries. Each $s_i = \tilde{s} \cup \{j \in \tilde{i} \mid j \geq i\}$, that is, the sparsity pattern of the i column is the selected entries plus all the diagonal entries lower than it. We will enforce that all the selected entries, excluding the indices of the diagonals of the columns themselves, are below the lowest index so that indices are not selected “partially” — that is, an index could be above some indices in the aggregated columns, and therefore invalid to add to their column, but below others. That is, we restrict the candidate indices $k > \max \tilde{i}$ so that the selected index can be added to each column in \tilde{i} without violating the lower triangularity of L . It is in fact possible to properly account for these partial updates, but the reasoning and eventual algorithm becomes more complicated. We defer a detailed discussion of the partial update case to [Appendix A.6](#).

We now show that the KL-minimization objective on the aggregated indices corresponds precisely to [\(2.11\)](#), the objective multiple point Gaussian regression with the chain rule of log determinant through conditioning.

$$(3.8) \quad \log \det(\Theta) = \log \det(\Theta_{1,1|2}) + \log \det(\Theta_{2,2})$$

where Θ is blocked according to [\(3.6\)](#). The KL divergence objective for \tilde{i} is:

$$(3.9) \quad \begin{aligned} \sum_{i \in \tilde{i}} \log(\Theta_{ii|s_i - \{i\}}) &= \log(\Theta_{i_m i_m | \tilde{s}}) + \log(\Theta_{i_{m-1} i_{m-1} | \tilde{s} \cup \{i_m\}}) + \dots \\ &= \log \det(\Theta_{\{i_m, i_{m-1}\}, \{i_m, i_{m-1}\} | \tilde{s}}) + \log(\Theta_{i_{m-2} i_{m-2} | \tilde{s} \cup \{i_m, i_{m-1}\}}) + \dots \\ &= \log \det(\Theta_{\tilde{i}, \tilde{i} | \tilde{s}}) \end{aligned}$$

We see that the objective (3.9) is equivalent to the objective (2.11), that is, to minimize the log determinant of the conditional covariance matrix corresponding to a set of prediction points, conditional on the selected entries. We can therefore directly use Algorithm 2.8 on the aggregated columns, where the prediction points correspond to indices in the aggregation and where we restrict the candidates k to those below each column in the aggregation, $k > \max \tilde{i}$.

Hence the sparse Cholesky factorization motivated by KL divergence can be viewed as sparse Gaussian process selection over each column, where entries are selected to maximize mutual information with the entry on the diagonal of the current column. In the aggregated case, the multiple columns in the aggregated group correspond directly to predicting for multiple prediction points, where entries are again selected to maximize mutual information with each diagonal entry in the aggregation. This viewpoint leads directly to Algorithm 3.1.

Algorithm 3.1 Cholesky factorization by selection

Input: $\mathbf{x}, K(\cdot, \cdot), s, g = \{\tilde{i}_1, \dots, \tilde{i}_{N/m}\}$

Output: L such that $(LL^\top)^{-1} \approx K(\mathbf{x}, \mathbf{x})$

```

1:  $n \leftarrow |\mathbf{x}|$ 
2: for  $\tilde{i} \in g$  do
3:    $J \leftarrow \{\max(\tilde{i}) + 1, \max(\tilde{i}) + 2, \dots, n\}$ 
4:   Compute  $I$  using Algorithm 2.8 or Algorithm 2.9
     where  $\mathbf{x}_{\text{Tr}} = \mathbf{x}[J], \mathbf{x}_{\text{Pr}} = \mathbf{x}[\tilde{i}], s = s - |\tilde{i}|$ 
5:    $\tilde{s} \leftarrow J[I]$ 
6:   for  $i \in \text{reversed}(\text{sorted}(\tilde{i}))$  do
7:      $\tilde{s} \leftarrow \tilde{s} \cup \{i\}$ 
8:    $s_i \leftarrow \text{reversed}(\tilde{s})$ 
9:   end for
10: end for
11: return  $L$  computed with Algorithm 3.3
```

Algorithm 3.2 Computing L without aggregation

Input: $\mathbf{x}, K(\cdot, \cdot), s_i$

Output: $L_{s_i, i}$

```

1:  $\Theta_{s_i, s_i}^{-1} \leftarrow K(\mathbf{x}[s_i], \mathbf{x}[s_i])^{-1}$ 
2:  $L_{s_i, i} \leftarrow \frac{\Theta_{s_i, s_i}^{-1} \mathbf{e}_1}{\sqrt{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}}$ 
3: return  $L_{s_i, i}$ 
```

Algorithm 3.3 Computing L with aggregation

Input: $\mathbf{x}, K(\cdot, \cdot), \tilde{s}, \tilde{i}$

Output: $L_{s_i, i}$ for all $i \in \tilde{i}$

```

1:  $s \leftarrow \tilde{i} \cup \tilde{s}$ 
2:  $U \leftarrow P^\dagger \text{chol}(P^\dagger \Theta_{s, s} P^\dagger) P^\dagger$ 
3: for  $i \in \tilde{i}$  do
4:    $k \leftarrow \text{index of } i \text{ in } \tilde{i}$ 
5:    $L_{s_i, i} \leftarrow U^{-\top} \mathbf{e}_k$ 
6: end for
7: return  $L$ 
```

Once the sparsity pattern has been determined, we need to compute each column of L according to Theorem 3.1. Because the sparsity pattern for each column in the same group are subsets of each other, we can efficiently compute all their columns at once. The observation is that the smallest index in the group (corresponding to the entry highest in the matrix) will have the largest sparsity pattern, the next index will have one less entry (lacking the entry above it, which would violate lower

triangularity), and so on. We need to compute $\Theta_{s_i, s_i}^{-1} \mathbf{e}_1$ for each $i \in \tilde{i}$, or the precision of the marginalized covariance corresponding to the selected entries. By (3.5), we can turn marginalization in covariance into conditioning in precision:

$$\begin{aligned} L_{s_i, i} &= \frac{\Theta_{s_i, s_i}^{-1} \mathbf{e}_1}{\sqrt{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}} \\ (3.10) \quad &= \frac{(\Theta_{s, s})_{k:, k: | : k-1}^{-1} \mathbf{e}_1}{\sqrt{\mathbf{e}_1^\top (\Theta_{s, s})_{k:, k: | : k-1}^{-1} \mathbf{e}_1}} \end{aligned}$$

where $s = \tilde{i} \cup \tilde{s}$ and k is i 's index in \tilde{i} . So we want the k th column of the precision of the marginalized covariance, conditional on all the entries before it. From (2.10), this can be directly read off the Cholesky factorization. Thus, we can simply compute:

$$(3.11) \quad L = \text{chol}(\Theta_{s, s}^{-1})$$

and read off the k th column to compute (3.10) for each $i \in \tilde{i}$. However, instead of computing a lower triangular factor for the precision, we can compute an *upper* triangular factor the covariance whose inverse transpose will be a *lower* triangular factor for the original matrix. In particular, we see that

$$(3.12) \quad U = P^\dagger \text{chol}(P^\dagger \Theta_{s, s} P^\dagger) P^\dagger$$

satisfies $UU^\top = \Theta_{s, s}$ where P^\dagger is the order-reversing permutation. Thus,

$$\Theta_{s, s}^{-1} = U^{-\top} U^{-1}$$

where $U^{-\top}$ is an *lower* triangular factor for $\Theta_{s, s}^{-1}$ equal to (3.3) because the Cholesky factorization is unique. Computing $U^{-\top}$ leads directly to Algorithm 3.3.

Recall that the complexity of selecting s out of N total training points for m prediction points using Algorithm 2.8 or Algorithm 2.9 was $\mathcal{O}(Ns^2 + Nm^2 + m^3)$. In the context of Cholesky factorization, N is the size of the matrix, m is the number of columns to aggregate, and s is the number of nonzero entries in each column of L . We therefore need to do $\frac{N}{m}$ selections, one for each aggregated group, where we only need to select $s - m$ entries (since the m prediction points are automatically added). We then need to actually construct each column of L after determining the sparsity pattern, with Algorithm 3.3. This costs $\mathcal{O}(s^3)$ for each aggregated group to compute the Cholesky factor of the submatrix, which dominates the time to compute each column of L for the m columns in the group, $\mathcal{O}(ms^2)$ ($N > s > m$). Thus, the overall complexity is $\mathcal{O}(\frac{N}{m}(N(s-m)^2 + Nm^2 + m^3 + s^3))$, which simplifies to $\mathcal{O}(\frac{N^2 s^2}{m})$ by making use of the bound that $(s-m)^2 = \mathcal{O}(s^2 + m^2)$.

Note that the non-aggregated factorization is equivalent to $m = 1$, which yields $\mathcal{O}(N^2 s^2)$ (using the non-aggregated algorithms Algorithms 2.5 and 3.2, but one can also use the aggregated versions Algorithms 2.8 and 3.3 with $m = 1$ and achieve equivalent complexity). Thus, we see that the aggregated version is m times faster than its non-aggregated counterpart, at the cost that the resulting sparsity pattern will be lower quality (since the algorithm is forced to select the same entry for *all* columns in the group).

Unlike the geometric algorithms of [3, 4] which rely on the pairwise distance between points, and whose covariance matrix is implicitly determined by a list of points and kernel function, this algorithm relies only on the entries of the covariance matrix Θ . Thus, it can factor arbitrary symmetric positive definite matrices without access to points or an explicit kernel function.

4. Numerical experiments. All experiments were run on the Partnership for an Advanced Computing Environment (PACE) Phoenix cluster at the Georgia Institute of Technology, with 8 cores of a Intel(R) Xeon(R) Gold 6226 CPU @ 2.70GHz and 6 GB of RAM per core. Python code for all numerical experiments can be found at <https://github.com/stephen-huan/conditional-knn>.

4.1. k th-nearest neighbors selection. We justify that diverse point selection based on conditional information can lead to better performance than simply selecting the nearest neighbor in a toy example on the MNIST dataset. We compare k th-nearest neighbors (KNN) directly to conditional k th-nearest neighbors (CKNN) in the following experiment. We randomly select 1000 images to form the training set and 100 to form the testing set. For each image in the testing set, we select the k “closest” training points with either KNN or CKNN. For KNN we use the standard Euclidean distance and for CKNN we use Matérn kernel with smoothness $\nu = 1.5$ and length scale $l = 2^{10}$. Finally, we predict the label of the test point by taking the most frequently occurring label in the k selected points.

Stephen: cite python libraries that want citations (e.g. scikit-learn)

Stephen: cite mnist dataset

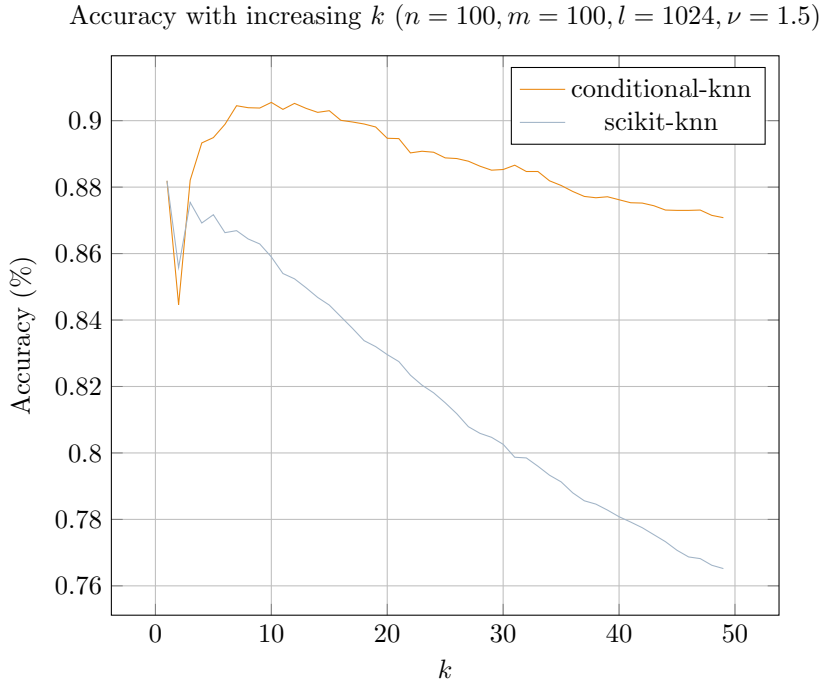


FIG. 1. Accuracy with increasing k

As k increases, KNN degrades near-linearly in accuracy. We hypothesize that nearby images are more likely to have the same label as a given test image. By forcing the algorithm to select more points, it increases the likelihood that the algorithm becomes confused by differently labeled images. However, CKNN is more accurate than KNN for nearly every k , suggesting that conditional selection is able to take advantage of selecting more points. We emphasize that the difference in accuracy is solely a result of conditional selection — because the Matérn kernel degrades monotonically with distance, sorting by covariance is identical to sorting by distance. In addition, we use the mode to aggregate the labels of the selected points, rather than

performing Gaussian process classification. The difference in accuracy can therefore be attributed to precisely the difference in which points were selected.

4.2. Recovery of sparse Cholesky factors. As noted in [Appendix A.9](#), the selection algorithm can be viewed as orthogonal matching pursuit [5] in feature space. We experiment with the sparse recovery properties of the selection algorithm by randomly generating a sparse Cholesky factor L . We prescribe a fixed number of nonzeros per column s over N columns. For each column, we uniformly randomly pick s entry that satisfies lower triangularity to make nonzero. We randomly generate values according to i.i.d. standard normal $\mathcal{N}(0, 1)$. Finally, we fill the diagonal with a “large” positive value (10) to almost guarantee that the resulting matrix $\Theta = LL^\top$ is positive definite. The selection algorithms are then given Θ and s and are asked to reconstruct L . The strategies are as follows: “cknn” uses conditional selection on each column to minimize the conditional variance of the diagonal entry, “knn” selects entries with the largest covariance with the diagonal entry, “corr” selects entries with the highest correlation objective (2.8) without accounting for conditioning, and “random” simply randomly samples entries uniformly. The strategies are given either the covariance Θ or the precision Θ^{-1} depending on which results in higher accuracy, in particular, the “cknn” strategy is given the precision while the rest of the methods are given the covariance. Accuracy is measured by taking the cardinality of the intersection of the recovered sparsity set with the ground truth sparsity set over the cardinality of their union, intersection over union (IOU).

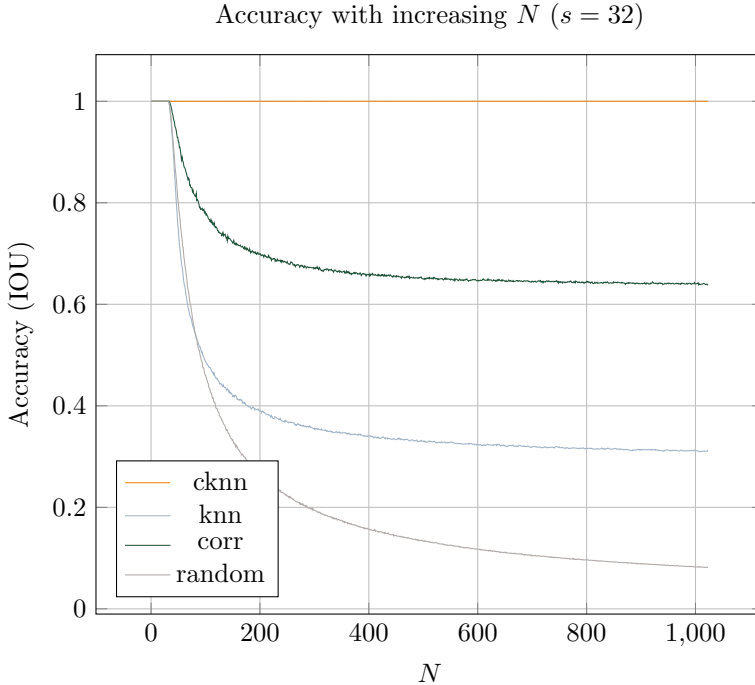
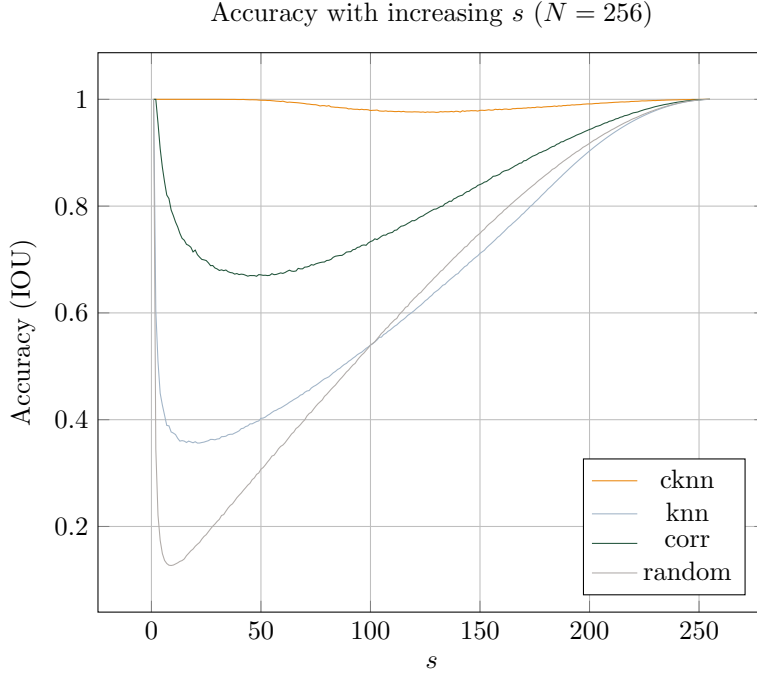


FIG. 2. Accuracy with increasing N

As the number of nonzero entries per column is fixed and the number of rows and columns is increased, the “cknn” retains high accuracy near perfect recovery, and the rest of the methods quickly degrade and asymptote to their final accuracies.

FIG. 3. Accuracy with increasing s

If the number of rows and columns is fixed while the number of nonzero entries per column is increased, all methods drop in accuracy with increasing density into a tipping point where the problem starts to become easier. Accuracy then increases until the Cholesky factor becomes fully dense, in which case perfect recovery is trivial. The “cknn” strategy exhibits the same behavior, but maintains much higher accuracy than the rest of the strategies.

Finally, we experiment with the addition of noise. Noise sampled i.i.d from $\mathcal{N}(0, \sigma^2)$ is added to each entry of Θ symmetrically (i.e. Θ_{ij} receives the same noise as Θ_{ji}) to preserve the symmetry of Θ . As expected, accuracy degrades with increasing noise, but the algorithm is fairly robust to low levels of noise. At higher levels of noise, Θ can lose positive definiteness, which causes the algorithm to break down.

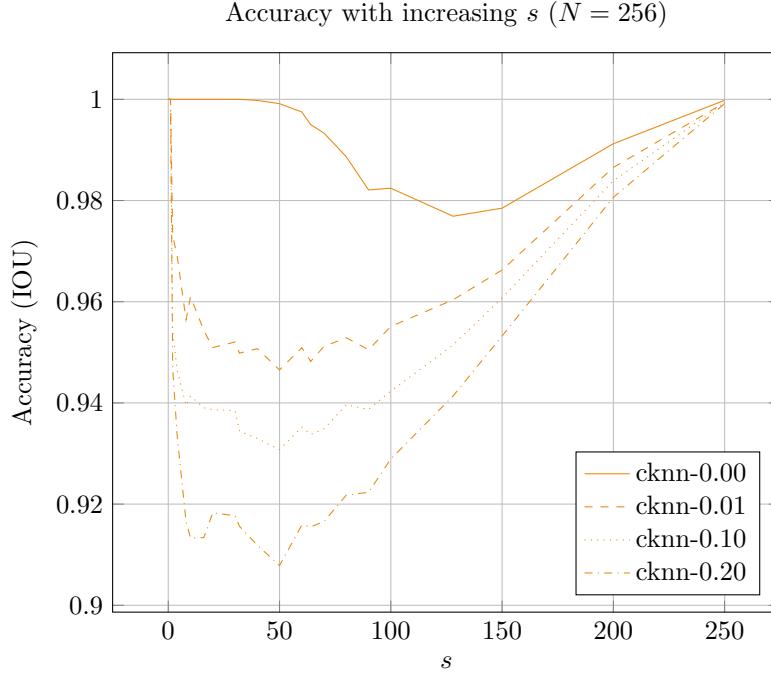
4.3. Cholesky factorization.

4.4. Gaussian process regression.

4.5. Preconditioning for conjugate gradient.

Acknowledgments. This research was supported in part through research cyberinfrastructure resources and services provided by the Partnership for an Advanced Computing Environment (PACE) at the Georgia Institute of Technology, Atlanta, Georgia, USA.

 add more funding
information

FIG. 4. Accuracy with increasing s

- [1] A. KRAUSE, A. SINGH, AND C. GUESTRIN, *Near-Optimal Sensor Placements in Gaussian Processes: Theory, Efficient Algorithms and Empirical Studies*, The Journal of Machine Learning Research, 9 (2008), pp. 235–284.
- [2] O. KRAUSE AND C. IGEL, *A More Efficient Rank-one Covariance Matrix Update for Evolution Strategies*, in Proceedings of the 2015 ACM Conference on Foundations of Genetic Algorithms XIII, FOGA '15, New York, NY, USA, Jan. 2015, Association for Computing Machinery, pp. 129–136, <https://doi.org/10.1145/2725494.2725496>.
- [3] F. SCHÄFER, M. KATZFUSS, AND H. OWHADI, *Sparse Cholesky factorization by Kullback-Leibler minimization*, arXiv:2004.14455 [cs, math, stat], (2021), <https://arxiv.org/abs/2004.14455>.
- [4] F. SCHÄFER, T. J. SULLIVAN, AND H. OWHADI, *Compression, inversion, and approximate PCA of dense kernel matrices at near-linear computational complexity*, arXiv:1706.02205 [cs, math], (2020), <https://arxiv.org/abs/1706.02205>.
- [5] J. A. TROPP AND A. C. GILBERT, *Signal Recovery From Random Measurements Via Orthogonal Matching Pursuit*, IEEE Transactions on Information Theory, 53 (2007), pp. 4655–4666, <https://doi.org/10.1109/TIT.2007.909108>.

add proofs, if any,
in appendix

Appendix A. Computation in sparse Gaussian process selection.

A.1. Updating precision after insertion. We have the matrix $\Theta_{I,I}^{-1}$ corresponding to the precision of the selected entries, and wish to take into account the addition of a new entry k into I . That is, we wish to compute $\Theta_{I',I'}^{-1}$ for $I' = I \cup \{k\}$, which in effect adds a new row and column to $\Theta_{I,I}^{-1}$. In order to invert the new matrix efficiently, we can block the matrix to separate the new and old information.

(A.1)

$$\begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \Theta_{2,1} & \Theta_{2,2} \end{pmatrix} = \begin{pmatrix} I & 0 \\ \Theta_{2,1}\Theta_{1,1}^{-1} & I \end{pmatrix} \begin{pmatrix} \Theta_{1,1} & 0 \\ 0 & \Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2} \end{pmatrix} \begin{pmatrix} I & \Theta_{1,1}^{-1}\Theta_{1,2} \\ 0 & I \end{pmatrix}$$

For notational convenience, we denote the Schur complement $\Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2}$ as $\Theta_{2,2|1}$. Inverting both sides of the equation,

$$(A.2) \quad \Theta^{-1} = \begin{pmatrix} I & -\Theta_{1,1}^{-1}\Theta_{1,2} \\ 0 & I \end{pmatrix} \begin{pmatrix} \Theta_{1,1}^{-1} & 0 \\ 0 & \Theta_{2,2|1}^{-1} \end{pmatrix} \begin{pmatrix} I & 0 \\ -\Theta_{2,1}\Theta_{1,1}^{-1} & I \end{pmatrix}$$

$$(A.3) \quad = \begin{pmatrix} \Theta_{1,1}^{-1} + (\Theta_{1,1}^{-1}\Theta_{1,2})\Theta_{2,2|1}^{-1}(\Theta_{2,1}\Theta_{1,1}^{-1}) & -(\Theta_{1,1}^{-1}\Theta_{1,2})\Theta_{2,2|1}^{-1} \\ -\Theta_{2,2|1}^{-1}(\Theta_{2,1}\Theta_{1,1}^{-1}) & \Theta_{2,2|1}^{-1} \end{pmatrix}$$

In the context of adding a new entry to the matrix, $\Theta_{1,1} = \Theta_{I,I}$, $\Theta_{1,2} = \Theta_{I,k}$, and $\Theta_{2,2} = \Theta_{kk}$. Also note that $\Theta_{kk|I}^{-1}$ is the inverse of the variance of k conditional on the entries in I , which has already been computed in [Algorithm 2.5](#). If we let $\mathbf{v} = \Theta_{I,I}^{-1}\Theta_{I,k}$, then we can write the update as:

$$(A.4) \quad = \begin{pmatrix} \Theta_{I,I}^{-1} + \Theta_{kk|I}^{-1}\mathbf{v}\mathbf{v}^T & -\Theta_{kk|I}^{-1}\mathbf{v} \\ -\Theta_{kk|I}^{-1}\mathbf{v}^T & \Theta_{kk|I}^{-1} \end{pmatrix}$$

which is precisely the update in line 13 of [Algorithm 2.5](#). Note that the update is a rank-one update to $\Theta_{1,1}^{-1}$, which can be computed in $\mathcal{O}(|I|^2) = \mathcal{O}(s^2)$.

A.2. Updating precision after marginalization. Suppose we have the precision Θ^{-1} and wish to compute the precision of the marginalized covariance after ignoring an index k . That is, we wish to compute the inverse of a matrix after deleting a row and column, given the inverse of the original matrix. We could use the result in [Appendix A.1](#) by “reading” the update backwards. That is, we could identify $\Theta_{2,2|1}^{-1}$ from $(\Theta^{-1})_{kk}$ and $\mathbf{v} = \Theta_{1,1}^{-1}\Theta_{1,2}$ from $-\frac{(\Theta^{-1})_{-k,k}}{\Theta_{2,2|1}^{-1}}$ where $-k$ denotes all rows excluding the k th row. We can then revert the rank-one update by subtracting out the update, computing $\Theta_{-k,-k}^{-1} = (\Theta^{-1})_{-k,-k} - \Theta_{kk|I}^{-1}\mathbf{v}\mathbf{v}^T$. However, a more intuitive derivation relies on the fact that marginalization in covariance is conditioning in precision. Using [\(3.5\)](#), we see that $\Theta_{-k,-k}^{-1} = (\Theta^{-1})_{-k,-k|k}$, or the precision conditional on the deleted entry. By [\(2.2\)](#), we immediately obtain the equivalent update

$$(A.5) \quad (\Theta^{-1})_{-k,-k|k} = \Theta_{-k,-k}^{-1} - \frac{(\Theta^{-1})_{-k,k}(\Theta^{-1})_{-k,k}^T}{(\Theta^{-1})_{kk}}$$

Since this is a rank-one update to the precision Θ^{-1} , this can be computed in $\mathcal{O}(\# \text{ rows}(\Theta^{-1}))^2$.

Stephen: this is not used in the paper but is nice to know + used in the sensor placement

A.3. Updating precision after conditioning. We have the matrix $\Theta_{\text{Pr}, \text{Pr}|I}^{-1}$, or the precision of the prediction points, conditional on the selected entries. We want to take into account selecting an entry k , or to compute $\Theta_{\text{Pr}, \text{Pr}|I \cup \{k\}}^{-1}$ which is a rank-one update to the original matrix from (2.11). We can directly apply the Sherman–Morrison–Woodbury formula which states that:

$$(A.6) \quad \Theta_{1,1|2}^{-1} = \Theta_{1,1}^{-1} + (\Theta_{1,1}^{-1} \Theta_{1,2}) \Theta_{2,2|1}^{-1} (\Theta_{2,1} \Theta_{1,1}^{-1})$$

Expanding the conditioning by definition,

$$(A.7) \quad (\Theta_{1,1} - \Theta_{1,2} \Theta_{2,2}^{-1} \Theta_{2,1})^{-1} = \Theta_{1,1}^{-1} + (\Theta_{1,1}^{-1} \Theta_{1,2}) \Theta_{2,2|1}^{-1} (\Theta_{2,1} \Theta_{1,1}^{-1})$$

Letting $\mathbf{u} = \Theta_{1,2}$ and $\mathbf{v} = \Theta_{1,1}^{-1} \Theta_{1,2} = \Theta_{1,1}^{-1} \mathbf{u}$,

$$(A.8) \quad (\Theta_{1,1} - \Theta_{2,2}^{-1} \mathbf{u} \mathbf{u}^\top)^{-1} = \Theta_{1,1}^{-1} + \Theta_{2,2|1}^{-1} \mathbf{v} \mathbf{v}^\top$$

So we see that a rank-one update to $\Theta_{1,1}$ then inverting is a rank-one update to $\Theta_{1,1}^{-1}$. In our context, $\Theta_{1,1} = \Theta_{\text{Pr}, \text{Pr}|I}$, $\mathbf{u} = \Theta_{\text{Pr}, k|I}$, $\Theta_{2,2} = \Theta_{kk|I}$ so $\Theta_{2,2|1}^{-1} = \Theta_{kk| \text{Pr}, I}^{-1}$ (this can be rigorously shown by expanding the Schur complement and taking advantage of the quotient rule as in (2.12)). \mathbf{v} can be computed according to definition as $\Theta_{\text{Pr}, \text{Pr}|I}^{-1} \mathbf{u}$. Thus, we can write the update as

$$(A.9) \quad \left(\Theta_{\text{Pr}, \text{Pr}|I} - \frac{\Theta_{\text{Pr}, k|I} \Theta_{\text{Pr}, k|I}^\top}{\Theta_{kk|I}} \right)^{-1} = \Theta_{1,1}^{-1} + \Theta_{kk| \text{Pr}, I}^{-1} \mathbf{v} \mathbf{v}^\top$$

which is the update in line 18 of Algorithm 2.8. Since the update is a rank-one update, it can be computed in $\mathcal{O}(|\text{Pr}|^2) = \mathcal{O}(m^2)$.

A.4. Updating Cholesky factorization after rank-one downdate.

We use the approach from Lemma 1 of [2], slightly adapted to use in-place operations and to make no assumption on the particular row ordering of the Cholesky factor. Let L be a Cholesky factorization of Θ , that is, $L = \text{chol}(\Theta)$. We wish to compute the updated Cholesky factor $L' = \text{chol}(\Theta')$ where $\Theta' = \Theta - \mathbf{u} \mathbf{u}^\top$. To do so, assume L and L' are blocked according to the same block structure:

$$(A.10) \quad L = \begin{pmatrix} r_1 & \mathbf{0} \\ \mathbf{r}_2 & L_2 \end{pmatrix}, L' = \begin{pmatrix} r'_1 & \mathbf{0} \\ \mathbf{r}'_2 & L'_2 \end{pmatrix}$$

Multiplying, we find

$$(A.11) \quad LL^\top = \Theta = \begin{pmatrix} r_1^2 & r_1 \mathbf{r}_2^\top \\ r_1 \mathbf{r}_2 & L_2 L_2^\top + \mathbf{r}_2 \mathbf{r}_2^\top \end{pmatrix}$$

$$(A.12) \quad L' L'^\top = \Theta' = \begin{pmatrix} r'^2_1 & r'_1 \mathbf{r}'^\top_2 \\ r'_1 \mathbf{r}'_2 & L'_2 L'^\top_2 + \mathbf{r}'_2 \mathbf{r}'^\top_2 \end{pmatrix}$$

Stephen: remove because unnecessary, describe insertion instead

From here, we solve for r'_1 , \mathbf{r}' , and L'_2

$$(A.13) \quad r_1'^2 = \Theta_{11}' = \Theta_{11} - u_1^2$$

$$(A.14) \quad = r_1^2 - u_1^2$$

$$(A.15) \quad r'_1 = \sqrt{r_1^2 - u_1^2}$$

$$(A.16) \quad r'_1 \mathbf{r}'_2 = \Theta_{2:,1}' = \Theta_{2:,1} - u_1 \mathbf{u}_2$$

$$(A.17) \quad = r_1 \mathbf{r}_2 - u_1 \mathbf{u}_2$$

$$(A.18) \quad \mathbf{r}'_2 = \frac{1}{r'_1} (r_1 \mathbf{r}_2 - u_1 \mathbf{u}_2)$$

$$(A.19) \quad L'_2 L_2'^\top = L_2 L_2^\top + \mathbf{r}_2 \mathbf{r}_2^\top - \mathbf{u}_2 \mathbf{u}_2^\top - \mathbf{r}'_2 \mathbf{r}'_2{}^\top$$

Plugging in the expression for \mathbf{r}'_2 ,

$$(A.20) \quad = L_2 L_2^\top + \mathbf{r}_2 \mathbf{r}_2^\top - \mathbf{u}_2 \mathbf{u}_2^\top - \left(\frac{r_1}{r'_1} \mathbf{r}_2 - \frac{u_1}{r'_1} \mathbf{u}_2 \right) \left(\frac{r_1}{r'_1} \mathbf{r}_2 - \frac{u_1}{r'_1} \mathbf{u}_2 \right)^\top$$

$$(A.21) \quad = L_2 L_2^\top + \left(1 - \frac{r_1^2}{r_1'^2} \right) \mathbf{r}_2 \mathbf{r}_2^\top + \frac{r_1 u_1}{r_1'^2} \mathbf{r}_2 \mathbf{u}_2^\top + \frac{u_1 r_1}{r_1'^2} \mathbf{u}_2 \mathbf{r}_2^\top - \left(1 + \frac{u_1^2}{r_1'^2} \right) \mathbf{u}_2 \mathbf{u}_2^\top$$

Using $r'_1 = \sqrt{r_1^2 - u_1^2}$,

$$(A.22) \quad = L_2 L_2^\top - \frac{u_1^2}{r_1'^2} \mathbf{r}_2 \mathbf{r}_2^\top + \frac{r_1 u_1}{r_1'^2} \mathbf{r}_2 \mathbf{u}_2^\top + \frac{u_1 r_1}{r_1'^2} \mathbf{u}_2 \mathbf{r}_2^\top - \frac{r_1^2}{r_1'^2} \mathbf{u}_2 \mathbf{u}_2^\top$$

After factoring we find

$$(A.23) \quad L'_2 L_2'^\top = L_2 L_2^\top - \left(\frac{r_1}{r'_1} \mathbf{u}_2 - \frac{u_1}{r'_1} \mathbf{r}_2 \right) \left(\frac{r_1}{r'_1} \mathbf{u}_2 - \frac{u_1}{r'_1} \mathbf{r}_2 \right)^\top$$

which is a rank-one downdate to the subfactor L_2 . Recursively updating L_2 yields a $\mathcal{O}(N^2)$ algorithm. We now re-write the algorithm to be in-place to take advantage of BLAS routines. The updates can be summarized as:

$$(A.24) \quad r'_1 = \sqrt{r_1^2 - u_1^2}$$

$$(A.25) \quad \mathbf{r}' = \frac{r_1}{r'_1} \mathbf{r} - \frac{u_1}{r'_1} \mathbf{u}$$

$$(A.26) \quad \mathbf{u}' = \frac{r_1}{r'_1} \mathbf{u} - \frac{u_1}{r'_1} \mathbf{r}$$

Note that we drop the subscripting on \mathbf{r} and \mathbf{u} . By updating the entire vector on each iteration, we can avoid keeping track of the lower triangular structure of L . We will first update \mathbf{r}' and then use it to update \mathbf{u} . Solving for \mathbf{r} in terms of \mathbf{r}' ,

$$(A.27) \quad \mathbf{r} = \frac{r'_1}{r_1} \mathbf{r}' + \frac{u_1}{r_1} \mathbf{u}$$

Plugging the expression for \mathbf{r} into the update for \mathbf{u}' ,

$$(A.28) \quad \mathbf{u}' = \frac{r_1}{r'_1} \mathbf{u} - \frac{u_1}{r'_1} \left(\frac{r'_1}{r_1} \mathbf{r}' + \frac{u_1}{r_1} \mathbf{u} \right)$$

$$(A.29) \quad = \frac{r_1^2 - u_1^2}{r_1 r'_1} \mathbf{u} - \frac{u_1}{r_1} \mathbf{r}'$$

$$(A.30) \quad = \frac{r'_1}{r_1} \mathbf{u} - \frac{u_1}{r_1} \mathbf{r}'$$

Thus, the updates proceed sequentially as follows:

$$(A.31) \quad \gamma \leftarrow \sqrt{r_1^2 - u_1^2}$$

$$(A.32) \quad \alpha \leftarrow \frac{r_1}{\gamma}$$

$$(A.33) \quad \beta \leftarrow \frac{u_1}{\gamma}$$

$$(A.34) \quad \mathbf{r} \leftarrow \alpha \mathbf{r}' - \beta \mathbf{u}$$

$$(A.35) \quad \mathbf{u} \leftarrow \frac{1}{\alpha} \mathbf{u} - \frac{\beta}{\alpha} \mathbf{r}$$

These can be efficiently performed in-place by BLAS as level-one **daxpy** operations.

A.5. Updating Cholesky factor after insertion. Suppose we have a Cholesky factor L of Θ and we insert a new point into Θ . We wish to update the Cholesky L to account for this insertion. Using the recursive conditioning interpretation of Cholesky factorization in (2.10), we see that the columns of L before the insertion will remain unchanged, the column at the insertion point is a new column given by the conditional covariance of the new point with the rest of the points, conditional on the points before it, which can be computed with standard left-looking, and the columns of L after the insertion correspond to the Cholesky factor of the conditional covariance, conditional on the newly inserted point in addition to the previous points. From (2.7) we know that conditioning on an additional point is a rank-one update of the covariance, so we can use rank-one downdating from Appendix A.4 to update L for the columns after the insertion point, where the vector in the rank-one downdate is the newly inserted column. This update touches every value in the Cholesky factor exactly once, so its complexity is $\mathcal{O}(N^2)$ as opposed to the $\mathcal{O}(N^3)$ cost of completely regenerating the Cholesky factor from scratch.

A.6. Partial updates in the selection algorithm. In the context of the selection algorithm, we have M prediction points and wish to minimize the log determinant of the resulting covariance matrix of the prediction points, conditional on the points we've selected from the training data. In the specific context of Cholesky factorization, it is possible to add a training point and have it apply *partially* on the prediction points. If nonadjacent columns indices are aggregated, an entry selected between two indices can be higher than one column, but lower than another. Adding the entry to the sparsity pattern would therefore only add to some, but not all, columns in the aggregation. We will model this as partially conditioning the variables of interest. In particular, if we have prediction variables y_1, y_2, \dots, y_M , a partial condition ignoring the first j variables on the selected index k would result in $y_1, y_2, \dots, y_j, y_{j+1|k}, \dots, y_{M|k}$.

The first question is to compute the resulting covariance matrix. We know $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \Theta)$ and $\mathbf{y}_{|k}$ has conditional distribution according to (2.2), $\mathbf{y}_{|k} \sim \mathcal{N}(\mu, \Theta -$

$\Theta_{:,k} \Theta_{k,k}^{-1} \Theta_{k,:}$). Taking the Cholesky factorization of both covariance matrices, let $L = \text{chol}(\Theta)$ and $L_{|k} = \text{chol}(\Theta_{|k})$. We can then view \mathbf{y} as $L\mathbf{z}$, where \mathbf{z} is distributed according to $\mathcal{N}(\mathbf{0}, I)$. Similarly, $\mathbf{y}_{|k} = L_{|k}\mathbf{z} + \boldsymbol{\mu}$. For unconditioned y_i and y_j , the covariance between them is defined to be Θ_{ij} . Similarly, for conditioned y_i and y_j , the covariance is $\Theta_{ij|k}$. The only question is what the covariance between unconditioned y_i and conditioned y_j is. By definition,

$$(A.36) \quad \text{Cov}[y_i, y_j] = \mathbb{E}[(y_i - \mathbb{E}[y_i])(y_j - \mathbb{E}[y_j])]$$

$$(A.37) \quad = \mathbb{E}[(L_i \mathbf{z})(L_{|k} \mathbf{z})]$$

$$(A.38) \quad = \mathbb{E}[(L_{1,i} z_1 + \dots + L_{N,i} z_N)(L_{1,j|k} z_1 + \dots + L_{N,j|k} z_N)]$$

For $i \neq j$, $\mathbb{E}[z_i z_j] = \mathbb{E}[z_i] \mathbb{E}[z_j] = 0$ since z_i is independent of z_j and has mean 0.

$$(A.39) \quad = \mathbb{E}[L_{1,i} L_{1,j|k} z_1^2 + \dots + L_{N,i} L_{N,j|k} z_N^2]$$

$$(A.40) \quad = L_{1,i} L_{1,j|k} \mathbb{E}[z_1^2] + \dots + L_{N,i} L_{N,j|k} \mathbb{E}[z_N^2]$$

For any i , $\mathbb{E}[z_i^2] = \text{Var}[z_i] + \mathbb{E}[z_i]^2 = 1 + 0 = 1$

$$(A.41) \quad = L_{1,i} L_{1,j|k} + \dots + L_{N,i} L_{N,j|k}$$

$$(A.42) \quad = L_i \cdot L_{j|k}$$

Thus, the new covariance matrix can be written as:

$$(A.43) \quad \begin{pmatrix} L_{:,j} L_{:,j}^\top & L_{:,j} L_{j:,|k}^\top \\ L_{j:,|k} L_{:,j}^\top & L_{j:,|k} L_{j:,|k}^\top \end{pmatrix} = \begin{pmatrix} L_{:,j} \\ L_{j:,|k} \end{pmatrix} \begin{pmatrix} L_{:,j} \\ L_{j:,|k} \end{pmatrix}^\top$$

We will denote a partially conditioned matrix as

$$(A.44) \quad \Theta_{:, :, |k}$$



FIG. 5. Illustration of the Cholesky factorization of a partially conditioned matrix.

We can now connect minimization of the log determinant of the partially updated covariance matrix to the KL divergence objective of Cholesky factorization. Computing the log determinant of the partially updated covariance matrix, we make use of (A.43) and make use of the fact that the determinant of a triangular matrix is the product of its diagonal entries:

$$(A.45) \quad \frac{1}{2} \log \det(\Theta_{:, :, |k}) = \underbrace{\log(L_{11}) + \dots + \log(L_{jj})}_{\text{the same}} + \underbrace{\log(L_{j+1,j+1|k}) + \dots + \log(L_{M,M|k})}_{\text{conditioned}}$$

Comparing to the KL divergence (3.4), $\mathbb{D}_{\text{KL}} \left(\mathcal{N}(\mathbf{0}, \Theta) \parallel \mathcal{N}(\mathbf{0}, (LL^\top)^{-1}) \right)$ which is equivalent to maximizing

$$(A.46) \quad = \sum_{i=1}^M \log(\Theta_{ii|s_i-\{i\}})$$

Recalling that k is added partially to some s_i , only those $i > j$

$$(A.47) \quad = \underbrace{\log(\Theta_{11|s_1-\{1\}}) + \cdots + \log(\Theta_{jj|s_j-\{j\}})}_{\text{the same}} + \underbrace{\log(\Theta_{j+1,j+1|s_{j+1}-\{j+1\}}) + \cdots + \log(\Theta_{MM|s_M-\{M\}})}_{\text{conditioned}}$$

Since L_{ii} is the square root of the variance of the i th variable conditional on each entry before it in the ordering, we have

$$(A.48) \quad 2\log(L_{ii}) = \log(\Theta_{ii|s_i-\{i\}})$$

So minimizing the log determinant of the partially conditioned covariance matrix (A.45) is the same as minimizing the KL divergence (3.4).

A.7. Algorithm for partial updates. We now need an efficient algorithm to keep track of partial updates. The key idea is to implicitly maintain the prediction matrix with selected points inserted to maintain proper ordering, and keep track of the log determinant throughout selection. We first give how this different perspective affects the interpretation of the multiple point selection algorithm. In the example, let x and y be selected points and 1 and 2 be prediction points.

$$(A.49) \quad \Theta = \begin{pmatrix} \Theta_{xx} & \Theta_{xy} & \Theta_{x1} & \Theta_{x2} \\ \Theta_{yx} & \Theta_{yy} & \Theta_{y1} & \Theta_{y2} \\ \Theta_{1x} & \Theta_{1y} & \Theta_{11} & \Theta_{12} \\ \Theta_{2x} & \Theta_{2y} & \Theta_{21} & \Theta_{22} \end{pmatrix}$$

Computing the log determinant by chain rule,

$$(A.50) \quad \log\det(\Theta) = \log(\Theta_{xx}) + \log(\Theta_{yy|x}) + \log(\Theta_{11|x,y}) + \log(\Theta_{22|x,y,1})$$

Isolating the objective — the variances of the prediction points

$$(A.51) \quad \log(\Theta_{11|x,y}) + \log(\Theta_{22|x,y,1}) = \log\det(\Theta) - \log(\Theta_{xx}) - \log(\Theta_{yy|x})$$

Now consider how inserting y changed the objective from when it was just x .

$$(A.52) \quad \log(\Theta_{11|x}) + \log(\Theta_{22|x,1}) = \log\det(\Theta_{-y,-y}) - \log(\Theta_{xx})$$

$$(A.53) \quad \Delta = \log\det(\Theta) - \log(\Theta_{yy|x}) - \log\det(\Theta_{-y,-y})$$

But from (2.12) we know

$$(A.54) \quad \Delta = \log\left(\frac{\Theta_{yy|x,1,2}}{\Theta_{yy|x}}\right)$$

Substituting,

$$(A.55) \quad \log(\Theta_{yy|x,1,2}) - \log(\Theta_{yy|x}) = \log\det(\Theta) - \log(\Theta_{yy|x}) - \log\det(\Theta_{-y,-y})$$

$$(A.56) \quad \log(\Theta_{yy|x,1,2}) = \log\det(\Theta) - \log\det(\Theta_{-y,-y})$$

In general,

$$(A.57) \quad \log(\Theta_{kk|I,\text{Pr}}) = \log\det(\Theta) - \log\det(\Theta_{-k,-k})$$

Another way to arrive at the same result is to note that if we inserted y at the *end* of $\Theta_{-y,-y}$, to compute the log determinant of the new, bigger matrix Θ we would add the variance of y conditional on every entry in the matrix to the old determinant by chain rule. Since the determinant is invariant to symmetric permutation, the matrix inserting y at the end has the same determinant as inserting y where it should be.

So we see that the conditional variance of a candidate point conditional on everything else in the matrix is the difference in log determinant between the matrix with the candidate inserted and the original matrix. The multiple prediction point algorithm can therefore be interpreted as we insert the candidate *after* all the previously selected points (so it is conditional on all the previous points) and *before* the prediction points (which conditions all of them). We then compute $\log(\Theta_{kk|I,\text{Pr}})$ for some candidate k which represents the difference in log determinant and then subtract $\log(\Theta_{kk|I})$ which is the spurious variance introduced by inserting k into the matrix. We do not need to subtract the spurious variances from the previously selected points because k does not affect them, and we select candidates by *relative* score.

We now apply this result to partial selection. In the example, let 1 and 2 be prediction points while x and y are both a selected points below 2 but above 1, where x has already been selected and y is a candidate.

$$(A.58) \quad \Theta = \begin{pmatrix} \Theta_{11} & \Theta_{1y} & \Theta_{1x} & \Theta_{12} \\ \Theta_{y1} & \Theta_{yy} & \Theta_{yx} & \Theta_{y2} \\ \Theta_{x1} & \Theta_{xy} & \Theta_{xx} & \Theta_{x2} \\ \Theta_{21} & \Theta_{2y} & \Theta_{2x} & \Theta_{22} \end{pmatrix}$$

Computing the log determinant by chain rule,

$$(A.59) \quad \log\det(\Theta) = \log(\Theta_{11}) + \log(\Theta_{yy|1}) + \log(\Theta_{xx|1,y}) + \log(\Theta_{22|1,y,x})$$

We see that y conditions 2 but not 1, precisely what we want to encode. However, we have introduced a spurious term $\log(\Theta_{yy|1})$ and changed the variance of x , both of which must be subtracted out.

$$(A.60) \quad \log(\Theta_{11}) + \log(\Theta_{22|1,y,x}) = \log\det(\Theta) - \log(\Theta_{yy|1}) - \log(\Theta_{xx|1,y})$$

We can substitute $\log(\Theta_{yy|1,x,2})$ for $\log\det(\Theta)$ by (A.57). Although it differs by a constant, this does not change the objective.

$$(A.61) \quad = \log(\Theta_{yy|1,x,2}) - \log(\Theta_{yy|1}) - \log(\Theta_{xx|1,y})$$

As long as we can compute conditional variances of our candidate on each *prefix* of the current ordering of prediction points interspersed with selected points, we can use the conditional variances to compute the updated conditional variances of the selected points by using their conditional covariances with the candidate. We are then able to compute every term in the objective. To do so, we maintain a partial Cholesky factor whose ordering is given by the current ordering. When we select a new point, we insert it in its appropriate place in the Cholesky factor. To update the Cholesky factor after an insertion efficiently, we left-look to get the column of its insertion position, and then update all columns right of the column by a rank-one downdate as described in [Appendix A.5](#) which touches every entry in the Cholesky factor, $\mathcal{O}(N(m+s))$ per update for a total cost of $\mathcal{O}(N(m+s)(s))$ over s selections. In addition, the algorithm can be considerably simplified by simply adding the conditional variances of the prediction points, instead of starting with a proxy for the log determinant of the entire matrix and subtracting out the spurious interactions from the training points.

By inspecting the Cholesky factor, we get the covariance of a selected point with a candidate, conditional on all the points prior to the selected point in the ordering. The conditional variance of the selected point is the diagonal entry. We can then compute the new conditional variance given the variance of the candidate, conditional on all points prior to the selected point. Suppose we are at index i and the candidate is index j , the updates are as follows:

$$(A.62) \quad \Theta_{ii|i-1} = (L_{ii})^2$$

$$(A.63) \quad \Theta_{ij|i-1} = L_{ij} \cdot L_{ii}$$

$$(A.64) \quad \Theta_{ii|i-1,j} = \Theta_{ii|i-1} - \frac{\Theta_{ij|i-1}^2}{\Theta_{jj|i-1}}$$

$$(A.65) \quad \Theta_{jj|i-1,i} = \Theta_{jj|i-1} - \frac{\Theta_{ij|i-1}^2}{\Theta_{ii|i-1}}$$

$$(A.66) \quad = \Theta_{jj|i}$$

Of course, the base case Θ_{jj} is simply $K(\mathbf{x}_j, \mathbf{x}_j)$, the variance of the j th point.

For each of the N candidates, it requires $m + s$ operations from the above updates to compute the objective. Over s selections, the total time is the same as the cost to update the Cholesky factor, matching the complexity of the non-partial multiple point algorithm. However, the asymptotic work in the non-partial algorithm can be implemented as BLAS level-2 calls, while the partial algorithm relies heavily on vector (level-1) calls, affecting the constant-factor performance of the algorithm.

A.8. Global greedy selection. Although each column is essentially independent from the perspective of selection, if there is a prescribed budget for the number of nonzeros then there is the problem of distributing the nonzeros over the columns. A natural method is to distribute as evenly as possible, this is efficient and practically useful. However, one principled way of allocating nonzeros is to maintain a *global* priority queue over all columns, and selecting from this queue determines not only which entry out of the candidate set is added as a nonzero, but also which column to select from. This allows the algorithm to greedily select the next entry which will decrease the global objective (3.4) as much as possible. The main change is that within a column, any monotonic transformation of the objective will preserve the relative ranking of candidates, for example adding a constant, multiplying by a constant, taking the log, etc. However, from the global perspective, if one column adds a different constant to their objectives than another column, the relative ranking of candidates between columns is skewed. Thus, each column must maintain an objective that corresponds directly to minimizing the global objective (3.4). Here we describe the modifications that must be made to the selection algorithms to support global comparison.

A.8.1. Single column selection. For a single prediction point, the objective is $\frac{\Theta_{k,Pr|I}^2}{\Theta_{kk|I}}$ (2.8) which is exactly the amount the variance of the prediction point is decreased if the k th candidate is selected, that is, $\Theta_{Pr,Pr|I} - \Theta_{Pr,Pr|I,k} = \frac{\Theta_{k,Pr|I}^2}{\Theta_{kk|I}}$. From the global perspective, all other prediction points are untouched, so the amount

the sum of the log variances of all the prediction points changes is

$$(A.67) \quad \min \Delta = \min [\log(\Theta_{Pr,Pr|I,k}) - \log(\Theta_{Pr,Pr|I})]$$

$$(A.68) \quad = \min \frac{\Theta_{Pr,Pr|I,k}}{\Theta_{Pr,Pr|I}}$$

$$(A.69) \quad = \min \frac{\Theta_{Pr,Pr|I} - \frac{\Theta_{k,Pr|I}^2}{\Theta_{kk|I}}}{\Theta_{Pr,Pr|I}}$$

$$(A.70) \quad = \min \left[1 - \frac{\Theta_{k,Pr|I}^2}{\Theta_{kk|I} \Theta_{Pr,Pr|I}} \right]$$

$$(A.71) \quad = \max \frac{\Theta_{k,Pr|I}^2}{\Theta_{kk|I} \Theta_{Pr,Pr|I}}$$

where (A.71) can be interpreted as the *percentage* of the decrease in variance from selecting the k th point to the variance before selecting the point.

A.8.2. Aggregated selection. Since the nonadjacent algorithm directly computes the sum of the log of the conditional variances of the prediction points, few modifications have to be made. One heuristical improvement is to take into account for “bang-for-buck”, that is, to account for the fact that different candidates cost a different amount of nonzero entries. Selecting a candidate can add between 1 and the number of columns in its aggregated group, depending on its relative index. Thus, candidates with larger groups will appear to decrease the global variance more, even if they are not as efficient as a candidate with a single group. In practice, it is better to use the objective $\frac{\Delta}{n}$ where Δ is the change in variance after selecting the candidate and n is the number of nonzero entries selecting the candidate adds.

A.9. Equivalence of selection and orthogonal matching pursuit. We show that the single-point selection algorithm described in Algorithm 2.5 is the covariance space equivalent to the feature space orthogonal matching pursuit (OMP) algorithm described in [5]. The equivalence comes from the fact that Cholesky factorization is Gram-Schmitt in feature space.

Let Θ be a symmetric positive definite matrix such that

$$(A.72) \quad \Theta = F^\top F$$

For some matrix F whose columns are vectors in feature space,

$$(A.73) \quad F = (\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_N)$$

Immediately we have

$$(A.74) \quad \Theta_{ij} = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

where $\langle \cdot, \cdot \rangle$ denotes the ordinary inner product on \mathbb{R}^N .

It suffices to see a single step of Cholesky factorization. Selecting \mathbf{x}_1 ,

$$(A.75) \quad \Theta' = \Theta - \frac{\mathbf{x}_1 \mathbf{x}_1^\top}{\Theta_{11}}$$

$$(A.76) \quad \Theta'_{ij} = \Theta_{ij} - \frac{\Theta_{i1} \Theta_{j1}}{\Theta_{11}}$$

Switching to the feature space perspective, if we select \mathbf{x}_1 we force the rest of the feature vectors to be orthogonal to \mathbf{x}_1 ,

$$(A.77) \quad \mathbf{x}'_i = \mathbf{x}_i - \frac{\langle \mathbf{x}_i, \mathbf{x}_1 \rangle}{\langle \mathbf{x}_1, \mathbf{x}_1 \rangle} \mathbf{x}_1$$

$$(A.78) \quad \langle \mathbf{x}'_i, \mathbf{x}'_j \rangle = \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \frac{\langle \mathbf{x}_i, \mathbf{x}_1 \rangle \langle \mathbf{x}_j, \mathbf{x}_1 \rangle}{\langle \mathbf{x}_1, \mathbf{x}_1 \rangle}$$

Comparing (A.76) and (A.78), we see that they are the same as expected. As a corollary, the objective of selecting the point \mathbf{x}_k that minimizes the residual of some target point \mathbf{x}_{Pr} can be written as

$$(A.79) \quad \|\mathbf{x}_{Pr} - \text{proj}_{\mathbf{x}_k} \mathbf{x}_{Pr}\|^2 = \langle \mathbf{x}_{Pr}, \mathbf{x}_{Pr} \rangle - \frac{\langle \mathbf{x}_{Pr}, \mathbf{x}_k \rangle^2}{\langle \mathbf{x}_k, \mathbf{x}_k \rangle}$$

Which is precisely the squared covariance of the candidate with the prediction over the variance of the candidate, as in (2.8).

$$(A.80)$$

This shows the equivalence as the objective is the same.

A.10. Checking submodularity. Our objective is the mutual information between the prediction and training points (2.6). A natural question is whether this objective is submodular with respect to the training set. The answer is no in general, see [1], section 8.3 for a simple counterexample. However, we can empirically check submodularity for particular geometries and choices of kernel function. If Pr is the set of prediction points, I is a set of indexes, and x_1, x_2 are additional indices not in I , then the set function is submodular if and only if

$$MI(Pr, I \cup \{x_2\}) - MI(Pr, I) \stackrel{?}{\geq} MI(Pr, I \cup \{x_1, x_2\}) - MI(Pr, I \cup \{x_1\})$$

Expanding from the definition of mutual information (2.6),

$$H[Pr | I] - H[Pr | I \cup \{x_2\}] \stackrel{?}{\geq} H[Pr | I \cup \{x_1\}] - H[Pr | I \cup \{x_1, x_2\}]$$

Since this is the objective (2.11) (with additional log),

$$\frac{\Theta_{x_2, x_2 | I}}{\Theta_{x_2, x_2 | I, Pr}} \stackrel{?}{\geq} \frac{\Theta_{x_2, x_2 | I} - \frac{\Theta_{x_1, x_2 | I}^2}{\Theta_{x_1, x_1 | I}}}{\Theta_{x_2, x_2 | I, Pr} - \frac{\Theta_{x_1, x_2 | I, Pr}^2}{\Theta_{x_1, x_1 | I, Pr}}}$$

From the fact that $\frac{a}{b} \geq \frac{a-c}{b-d}$ if and only if $\frac{a}{b} \leq \frac{c}{d}$,

$$\frac{\Theta_{x_2, x_2 | I}}{\Theta_{x_2, x_2 | I, Pr}} \leq \frac{\frac{\Theta_{x_1, x_2 | I}^2}{\Theta_{x_1, x_1 | I}}}{\frac{\Theta_{x_1, x_2 | I, Pr}^2}{\Theta_{x_1, x_1 | I, Pr}}}$$

Multiplying by $\frac{\Theta_{x_1, x_1 | I}}{\Theta_{x_1, x_1 | I, Pr}}$ on both sides,

$$\frac{\Theta_{x_1, x_1 | I} \Theta_{x_2, x_2 | I}}{\Theta_{x_1, x_1 | I, Pr} \Theta_{x_2, x_2 | I, Pr}} \stackrel{?}{\leq} \frac{\Theta_{x_1, x_2 | I}^2}{\Theta_{x_1, x_2 | I, Pr}^2}$$

Multiplying by $\frac{\Theta_{x_1, x_2 | I, Pr}^2}{\Theta_{x_1, x_1 | I} \Theta_{x_2, x_2 | I}}$ on both sides,

$$\frac{\Theta_{x_1, x_2 | I, Pr}^2}{\Theta_{x_1, x_1 | I, Pr} \Theta_{x_2, x_2 | I, Pr}} \stackrel{?}{\leq} \frac{\Theta_{x_1, x_2 | I}^2}{\Theta_{x_1, x_1 | I} \Theta_{x_2, x_2 | I}}$$

By definition, this is

$$\text{Corr}[x_1, x_2 \mid I, \text{Pr}] \stackrel{?}{\leq} \text{Corr}[x_1, x_2 \mid I]$$

so the mutual information objective is submodular if and only if conditioning on additional point(s) decreases the correlation between every pair of points. Intuitively, this corresponds to the screening effect observed in spatial statistics literature — conditioning on a nearby point decreases the correlation for far away points.

Appendix B. Derivations in KL-minimization.

B.1. Linear-algebraic formulation of objective. We want to show that the KL divergence between two multivariate Gaussians centered at $\mathbf{0}$ with covariance matrices Θ_1 and Θ_2 can be written as

(B.1)

$$2\mathbb{D}_{\text{KL}}(\mathcal{N}(\mathbf{0}, \Theta_1) \parallel \mathcal{N}(\mathbf{0}, \Theta_2)) = \text{trace}(\Theta_2^{-1}\Theta_1) + \log\det(\Theta_2) - \log\det(\Theta_1) - N$$

where Θ_1 and Θ_2 are both of size $N \times N$. Recall that the log density $\log \pi(\mathbf{x})$ for $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Theta)$ is

$$(B.2) \quad \log \pi(\mathbf{x}) = -\frac{1}{2}(N \log(2\pi) + \log\det(\Theta) + \mathbf{x}^\top \Theta^{-1} \mathbf{x})$$

By the definition of KL divergence,

(B.3)

$$2\mathbb{D}_{\text{KL}}(\mathcal{N}(\mathbf{0}, \Theta_1) \parallel \mathcal{N}(\mathbf{0}, \Theta_2)) = 2\mathbb{E}_P[\log P - \log Q]$$

where P and Q are the corresponding densities for Θ_1 and Θ_2 respectively, and \mathbb{E}_P denotes expectation under P .

$$(B.4) \quad \begin{aligned} &= 2\mathbb{E}_P\left[-\frac{1}{2}(N \log(2\pi) + \log\det(\Theta_1) + \mathbf{x}^\top \Theta_1^{-1} \mathbf{x})\right. \\ &\quad \left. + \frac{1}{2}(N \log(2\pi) + \log\det(\Theta_2) + \mathbf{x}^\top \Theta_2^{-1} \mathbf{x})\right] \end{aligned}$$

$$(B.5) \quad = \mathbb{E}_P[\mathbf{x}^\top \Theta_2^{-1} \mathbf{x} - \mathbf{x}^\top \Theta_1^{-1} \mathbf{x}] + \log\det(\Theta_2) - \log\det(\Theta_1)$$

(B.6)

$$\mathbb{E}_P[\mathbf{x}^\top \Theta_2^{-1} \mathbf{x} - \mathbf{x}^\top \Theta_1^{-1} \mathbf{x}] = \mathbb{E}_P[\text{trace}(\mathbf{x}^\top \Theta_2^{-1} \mathbf{x}) - \text{trace}(\mathbf{x}^\top \Theta_1^{-1} \mathbf{x})]$$

because the trace of a scalar is a scalar, and the linearity of trace.

$$(B.7) \quad = \mathbb{E}_P[\text{trace}(\Theta_2^{-1} \mathbf{x} \mathbf{x}^\top) - \text{trace}(\Theta_1^{-1} \mathbf{x} \mathbf{x}^\top)] \quad \text{cyclic property of trace}$$

$$(B.8) \quad = \mathbb{E}_P[\text{trace}(\Theta_2^{-1} \mathbf{x} \mathbf{x}^\top - \Theta_1^{-1} \mathbf{x} \mathbf{x}^\top)] \quad \text{linearity of trace}$$

$$(B.9) \quad = \mathbb{E}_P[\text{trace}((\Theta_2^{-1} - \Theta_1^{-1}) \mathbf{x} \mathbf{x}^\top)] \quad \text{factoring}$$

$$(B.10) \quad = \text{trace}(\mathbb{E}_P[(\Theta_2^{-1} - \Theta_1^{-1}) \mathbf{x} \mathbf{x}^\top]) \quad \text{swapping trace and expectation}$$

$$(B.11) \quad = \text{trace}((\Theta_2^{-1} - \Theta_1^{-1}) \mathbb{E}_P[\mathbf{x} \mathbf{x}^\top]) \quad \text{linearity of expectation}$$

$$(B.12) \quad = \text{trace}((\Theta_2^{-1} - \Theta_1^{-1}) \Theta_1) \quad \Theta_1 = \mathbb{E}_P[\mathbf{x} \mathbf{x}^\top]$$

$$(B.13) \quad = \text{trace}(\Theta_2^{-1} \Theta_1 - I) \quad \text{multiplying}$$

$$(B.14) \quad = \text{trace}(\Theta_2^{-1} \Theta_1) - \text{trace}(I) \quad \text{linearity of trace}$$

$$(B.15) \quad = \text{trace}(\Theta_2^{-1} \Theta_1) - N \quad \text{trace of } N \times N \text{ identity } N$$

Combining (B.15) with (B.5), we obtain

$$2\mathbb{D}_{\text{KL}}\left(\mathcal{N}(\mathbf{0}, \Theta_1) \parallel \mathcal{N}(\mathbf{0}, \Theta_2)\right) = \text{trace}(\Theta_2^{-1}\Theta_1) + \log\det(\Theta_2) - \log\det(\Theta_1) - N$$

as desired.

B.2. Reduction for optimal factor. We wish to compute the KL divergence between Θ and the Cholesky factor L computed according to Theorem 3.1. From (3.2),

$$(B.16)$$

$$2\mathbb{D}_{\text{KL}}\left(\mathcal{N}(\mathbf{0}, \Theta) \parallel \mathcal{N}(\mathbf{0}, (LL^\top)^{-1})\right) = \text{trace}(LL^\top\Theta) - \log\det(LL^\top) - \log\det(\Theta) - N$$

Ignoring terms not depending on L ,

$$(B.17) \quad = \text{trace}(LL^\top\Theta) - \log\det(LL^\top)$$

By the cyclic property of trace,

$$(B.18) \quad = \text{trace}(L\Theta L^\top) - \log\det(LL^\top)$$

Focusing on $\text{trace}(L\Theta L^\top)$ and expanding on the columns of L ,

$$(B.19) \quad \text{trace}(L\Theta L^\top) = \sum_{i=1}^N L_{s_i, i}^\top \Theta_{s_i, s_i} L_{s_i, i}$$

Plugging in $L_{s_i, i}$ from Theorem 3.1,

$$(B.20) \quad = \sum_{i=1}^N \left(\frac{(\Theta_{s_i, s_i}^{-1} \mathbf{e}_1)^\top}{\sqrt{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}} \right) \Theta_{s_i, s_i} \left(\frac{\Theta_{s_i, s_i}^{-1} \mathbf{e}_1}{\sqrt{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}} \right)$$

$$(B.21) \quad = \sum_{i=1}^N \frac{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \Theta_{s_i, s_i} \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}$$

$$(B.22) \quad = \sum_{i=1}^N 1 = N$$

Using N for $\text{trace}(LL^\top\Theta)$ in (B.16),

$$(B.23) \quad 2\mathbb{D}_{\text{KL}}\left(\mathcal{N}(\mathbf{0}, \Theta) \parallel \mathcal{N}(\mathbf{0}, (LL^\top)^{-1})\right) = -\log\det(LL^\top) - \log\det(\Theta)$$

L^\top has the same log determinant as L , and because L is lower triangular, its log determinant is just the sum of its diagonal entries:

$$(B.24) \quad = -2 \sum_{i=1}^N [\log(L_{ii})] - \log\det(\Theta)$$

Plugging (3.3) for the diagonal entries,

$$(B.25) \quad = - \sum_{i=1}^N [\log(\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1)] - \log\det(\Theta)$$

Bringing the negative inside,

$$(B.26) \quad = \sum_{i=1}^N [\log((\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1)^{-1})] - \log\det(\Theta)$$

So minimizing the KL divergence (given optimal L) corresponds to minimizing the sum of the inverse of the diagonal entries. We can give an intuitive view of this result by making use of (3.7) and expanding the log determinant by the chain rule (3.8).

(B.27)

$$\sum_{i=1}^N [\log ((\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1)^{-1})] - \log \det(\Theta) = \sum_{i=1}^N [\log (\Theta_{ii|s_i - \{i\}})] - \log \det(\Theta)$$

(B.28)

$$= \sum_{i=1}^N \log (\Theta_{ii|s_i - \{i\}}) - \sum_{i=1}^N \log (\Theta_{ii|i+1:})$$

(B.29)

$$= \sum_{i=1}^N [\log (\Theta_{ii|s_i - \{i\}}) - \log (\Theta_{ii|i+1:})]$$

We can view this sum as the accumulated *difference* in prediction error for a series of prediction problems, where each prediction problem is to predict the value of the i th variable given variables $i+1, i+2, \dots, N$. The left term $\log (\Theta_{ii|s_i - \{i\}})$ is restricted to only using those variables in the sparsity pattern s_i , while the right term $\log (\Theta_{ii|i+1:})$ is able to use every variable after i . The left term will necessarily have greater variance than the right, and the goal is to minimize the accumulated deviation. Thus, the KL divergence gives a natural way to measure the quality of a sparsity pattern as a good sparsity pattern should maintain predictive accuracy while subject to the constraint that some variables have no interaction with others.

Because the KL divergence is not symmetric, it matters which way the KL divergence is taken as well as whether both matrices have been inverted or not. This seems to imply that there are four possible ways to compare two covariance matrices. However, note that

$$(B.30) \quad \mathbb{D}_{\text{KL}} \left(\mathcal{N}(\mathbf{0}, \Theta) \parallel \mathcal{N}(\mathbf{0}, (LL^\top)^{-1}) \right) = \mathbb{D}_{\text{KL}} \left(\mathcal{N}(\mathbf{0}, LL^\top) \parallel \mathcal{N}(\mathbf{0}, \Theta^{-1}) \right)$$

from (3.2) and the cyclic property of trace, so inverting both matrices implicitly reverses the order of the KL divergence. There are therefore only two possible ways to compare the two, which depends on the order of the arguments. A statistical interpretation comes from the fact that the KL divergence can be interpreted as the likelihood-ratio test, so the non-symmetry of the order of the arguments corresponds to the asymmetry between the null and alternative hypotheses.