# Sparse Cholesky Factorization by Greedy Conditional Selection

**Georgia Institute of Technology**

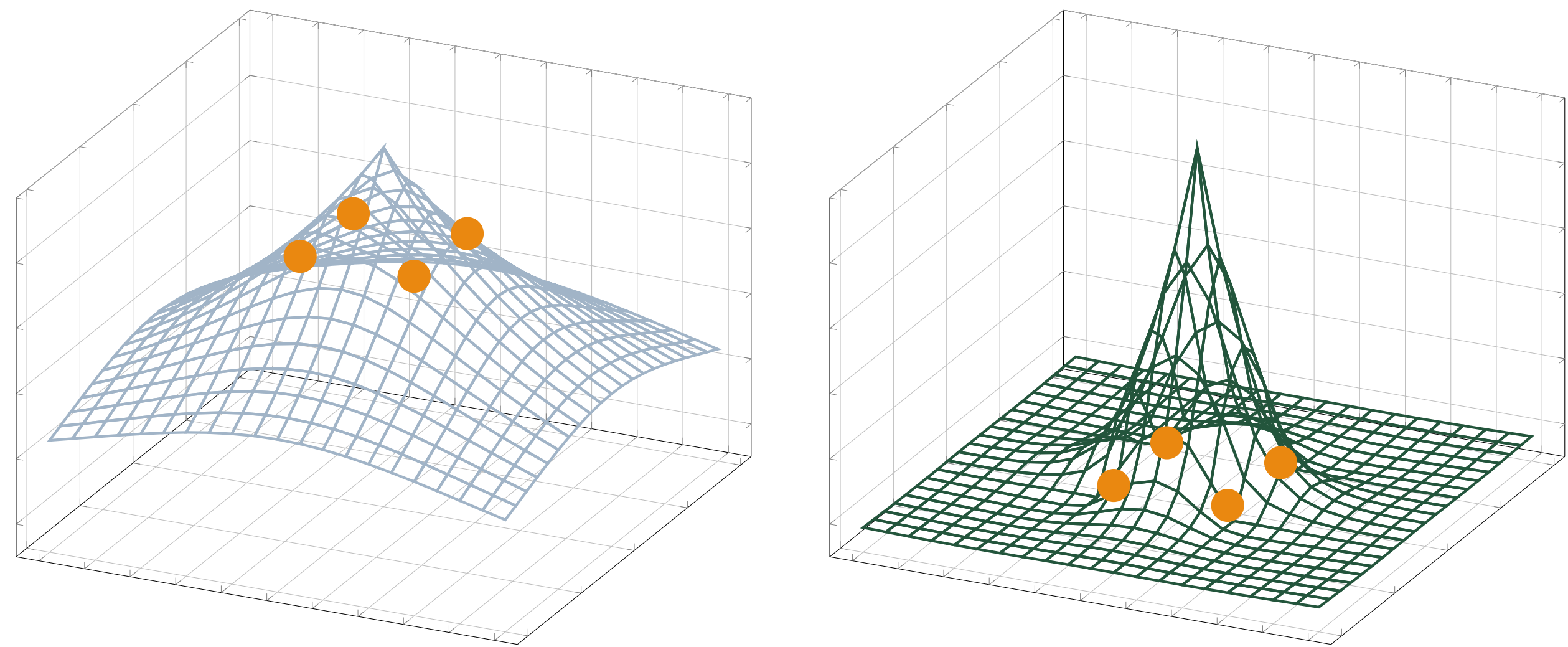Stephen Huan and Florian Schäfer

## The problem: Gaussian process regression

Given measurements $\boldsymbol{y}_{\mathrm{Tr}}$ at $N$ points $X_{\mathrm{Tr}}$, we wish to estimate unseen data $\boldsymbol{y}_{\mathrm{Pr}}$ at $X_{\mathrm{Pr}}$. Estimation of $\boldsymbol{y}_{\mathrm{Pr}}$ can be done by conditioning on $\boldsymbol{y}_{\mathrm{Tr}}$:

$$\mathbb{E}[\boldsymbol{y}_{\mathrm{Pr}} \mid \boldsymbol{y}_{\mathrm{Tr}}] = \boldsymbol{\mu}_{\mathrm{Pr}} + \Theta_{\mathrm{Pr},\mathrm{Tr}}\Theta_{\mathrm{Tr},\mathrm{Tr}}^{-1}(\boldsymbol{y}_{\mathrm{Tr}} - \boldsymbol{\mu}_{\mathrm{Tr}})$$
$$\mathbb{C}\mathrm{ov}[\boldsymbol{y}_{\mathrm{Pr}} \mid \boldsymbol{y}_{\mathrm{Tr}}] = \Theta_{\mathrm{Pr},\mathrm{Pr}} - \Theta_{\mathrm{Pr},\mathrm{Tr}}\Theta_{\mathrm{Tr},\mathrm{Tr}}^{-1}\Theta_{\mathrm{Tr},\mathrm{Pr}} := \Theta_{\mathrm{Pr},\mathrm{Pr}|\mathrm{Tr}}$$
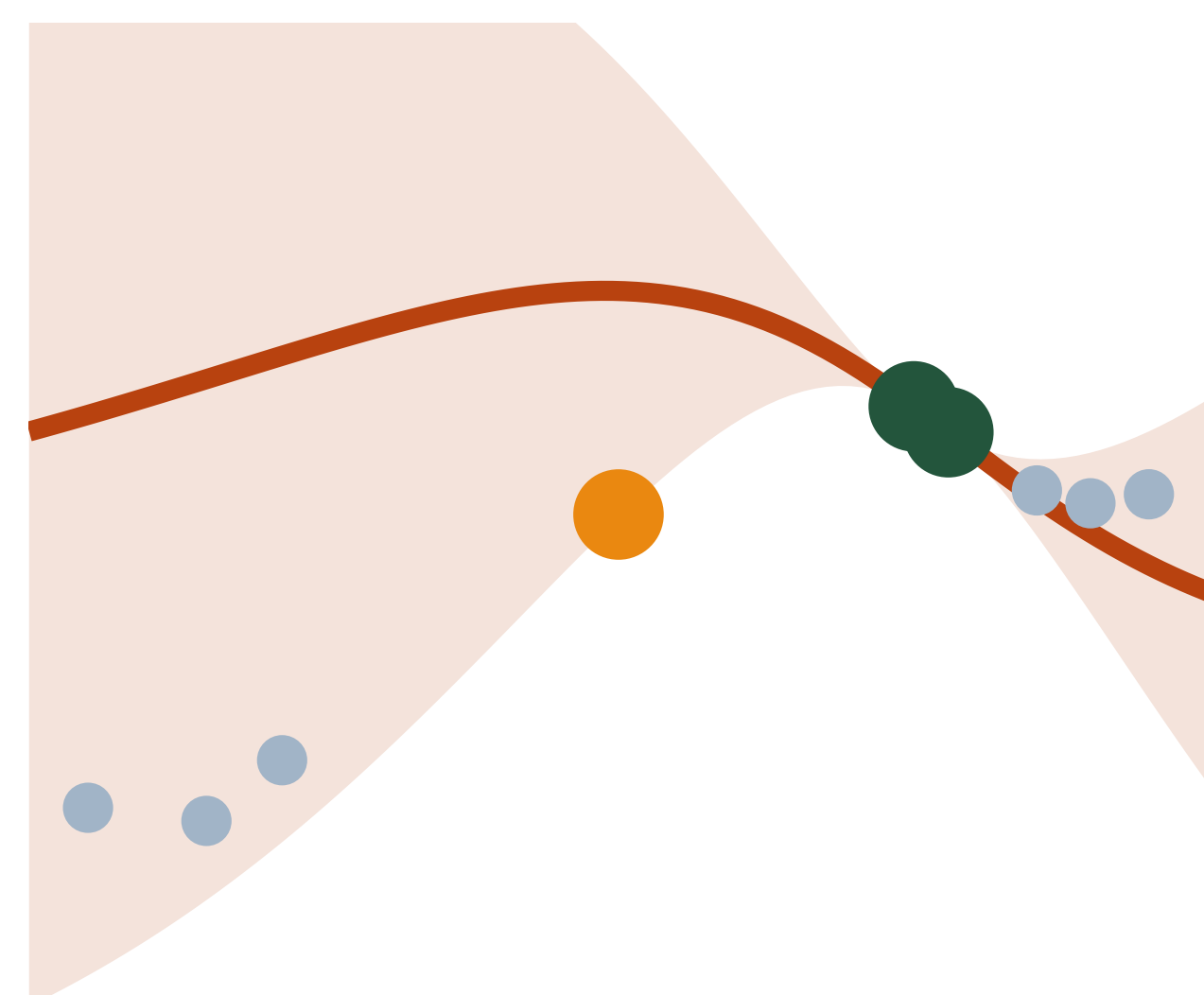
## Cubic bottleneck and the screening effect

Computing the conditional distribution has computational cost $\mathcal{O}(N^3)$, which is infeasible for many points. Instead, exploit the *screening effect*: conditional on nearby points, far away points have little correlation.

## $k$-nearest neighbors?

The screening effect suggests that one should simply pick the $k$ closest points, recovering the $k$-nearest neighbors ($k$-NN) algorithm.

Here, the blue points are the candidates, the orange point is the unknown point, and the green points are the $k$ selected points (in this example, $k = 2$).
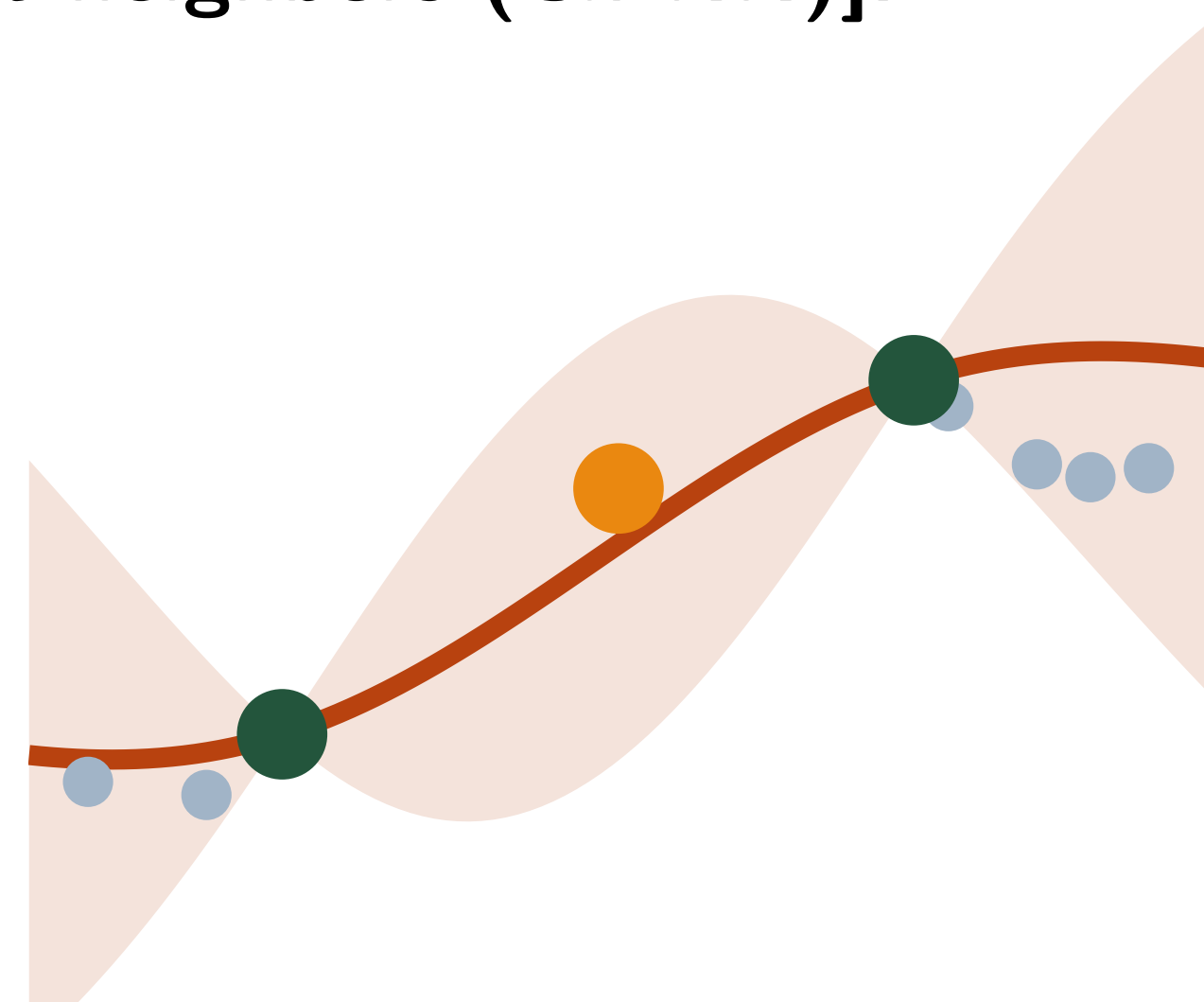
## $k$-NN is myopic, account for conditioning!

**Algorithm [conditional $k$-nearest neighbors (C$k$-NN)]:**

Selecting the closest point every iteration leads to redundancy.

Instead, select points *conditional* on points already selected. Selecting points by *information* instead of by distance motivates conditional $k$-th nearest neighbors (C$k$-NN).

## Greedy mutual information maximization

Greedily select the next training point with highest mutual information with the target point. If $I$ is the set of selected indices, select by:

$$\underset{j \notin I}{\mathrm{argmax}} \; \mathbb{C}\mathrm{orr}[y_{\mathrm{Pr}}, y_j \mid I]^2$$

## Efficient computation from Cholesky factor

A direct computation of the objective takes $\mathcal{O}(Nk^4)$ to select $k$ points. This computational cost can be reduced to $\mathcal{O}(Nk^2)$ by storing a partial Cholesky factor, since each column is conditional on everything before it:

$$\mathrm{chol}(\Theta) = \begin{pmatrix} \mathrm{Id} & 0 \\ \Theta_{2,1}\Theta_{1,1}^{-1} & \mathrm{Id} \end{pmatrix} \begin{pmatrix} \mathrm{chol}(\Theta_{1,1}) & 0 \\ 0 & \mathrm{chol}(\Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2}) \end{pmatrix}$$

## Generalization to multiple prediction points

For multiple prediction points, the objective becomes to minimize the log determinant of the conditional covariance matrix of prediction points. By making use of the matrix determinant lemma, one can show that:

$$\mathrm{logdet}(\Theta_{\mathrm{Pr},\mathrm{Pr}|I,k}) - \mathrm{logdet}(\Theta_{\mathrm{Pr},\mathrm{Pr}|I}) = \log(\Theta_{k,k|I,\mathrm{Pr}}) - \log(\Theta_{k,k|I})$$

We can efficiently compute the objective by storing *two* Cholesky factors, yielding a complexity of $\mathcal{O}(Nk^2 + Nm^2 + m^3)$ for $m$ prediction points.

## Global approximation by KL-minimization

Approximate a Gaussian process by a sparse approximate Cholesky factor of its precision. Measure the resulting approximation accuracy by the KL divergence between the corresponding centered Gaussian processes:

$$L := \underset{\hat{L} \in \mathcal{S}}{\mathrm{argmin}} \; \mathbb{D}_{\mathrm{KL}}\Big(\mathcal{N}(\boldsymbol{0}, \Theta) \; \Big\| \; \mathcal{N}(\boldsymbol{0}, (\hat{L}\hat{L}^{\top})^{-1})\Big)$$

Using the optimal unique minimizer $L$ from closed-form computation:

$$L_{s_i,i} = \frac{\Theta_{s_i,s_i}^{-1}\boldsymbol{e}_1}{\sqrt{\boldsymbol{e}_1^{\top}\Theta_{s_i,s_i}^{-1}\boldsymbol{e}_1}}$$
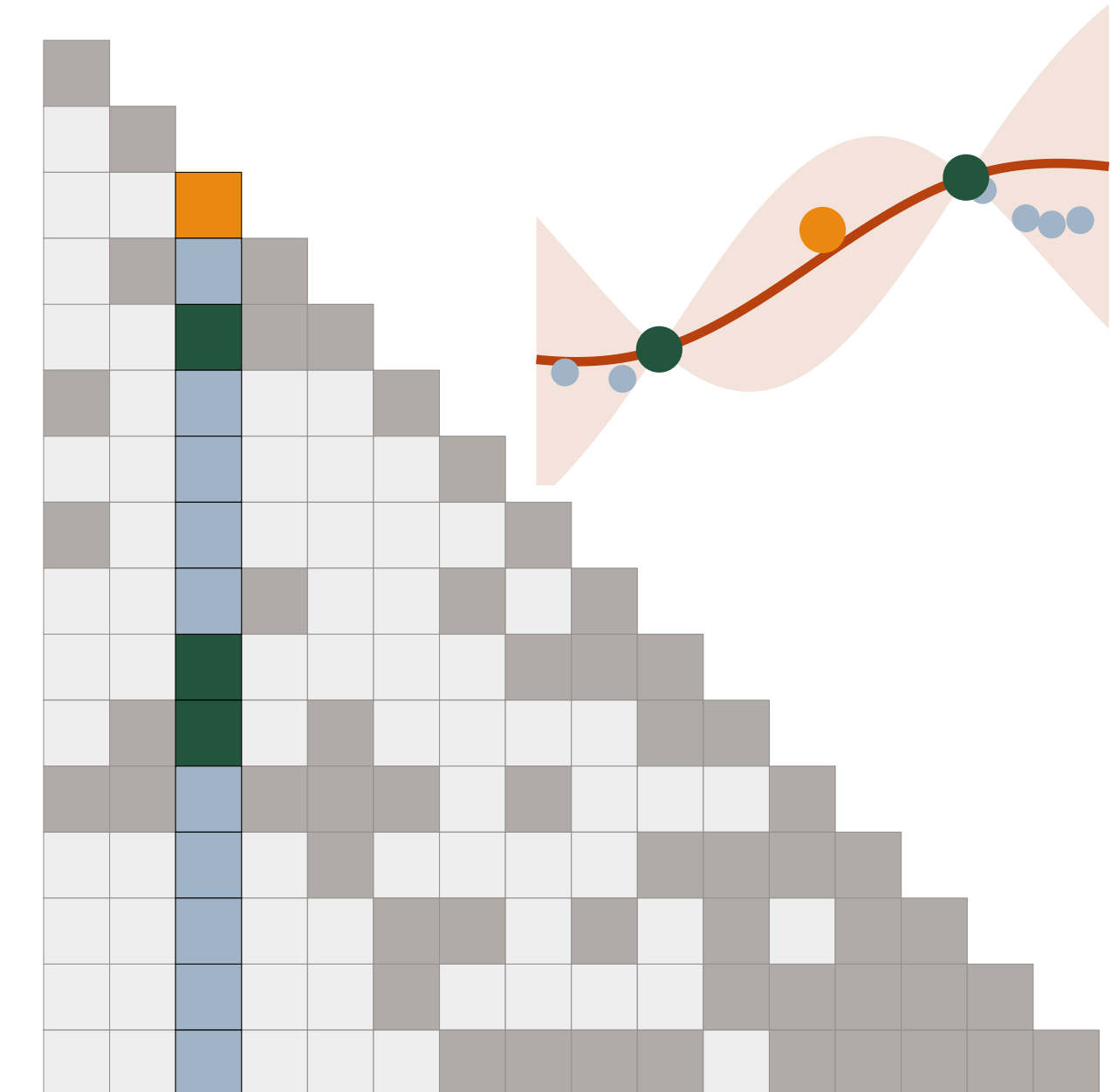
Objective becomes minimize variance of $i$th point, conditional on selected!

$$\mathbb{D}_{\mathrm{KL}}\big(\Theta \; \| \; (LL^{\top})^{-1}\big) \propto \sum_{i=1}^{N} \big[\log(\Theta_{i,i|s_i\setminus\{i\}}) - \log(\Theta_{i,i|i+1:})\big]$$
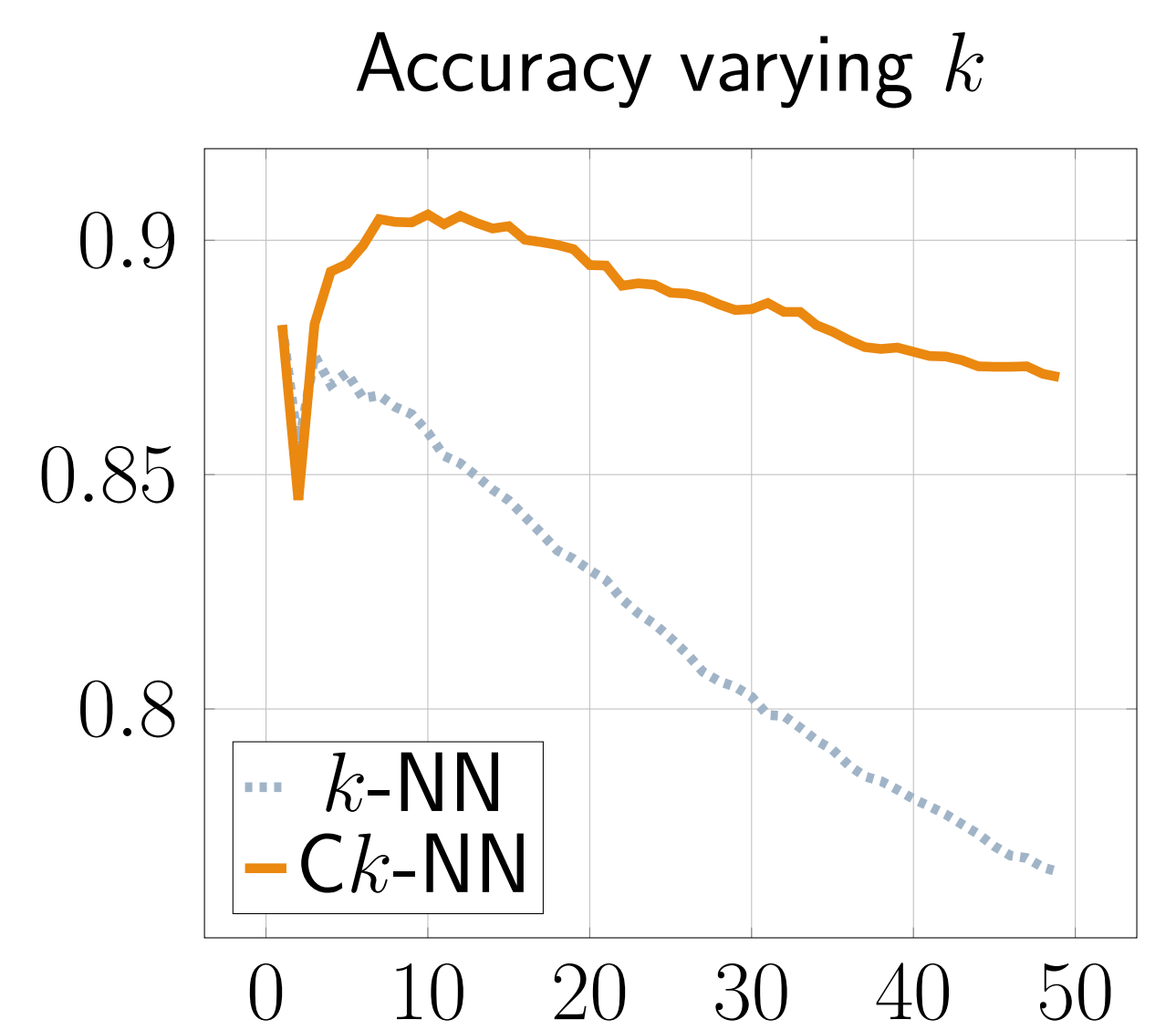
## Applying selection to Cholesky factorization

For a column in isolation, unknown point is the diagonal entry, below it are candidates, and add selected entries to the sparsity pattern $s_i$.

However, for aggregated columns (supernodes), a candidate can be added between prediction points. By careful application of rank-one downdating, this structure can be preserved at no additional cost.
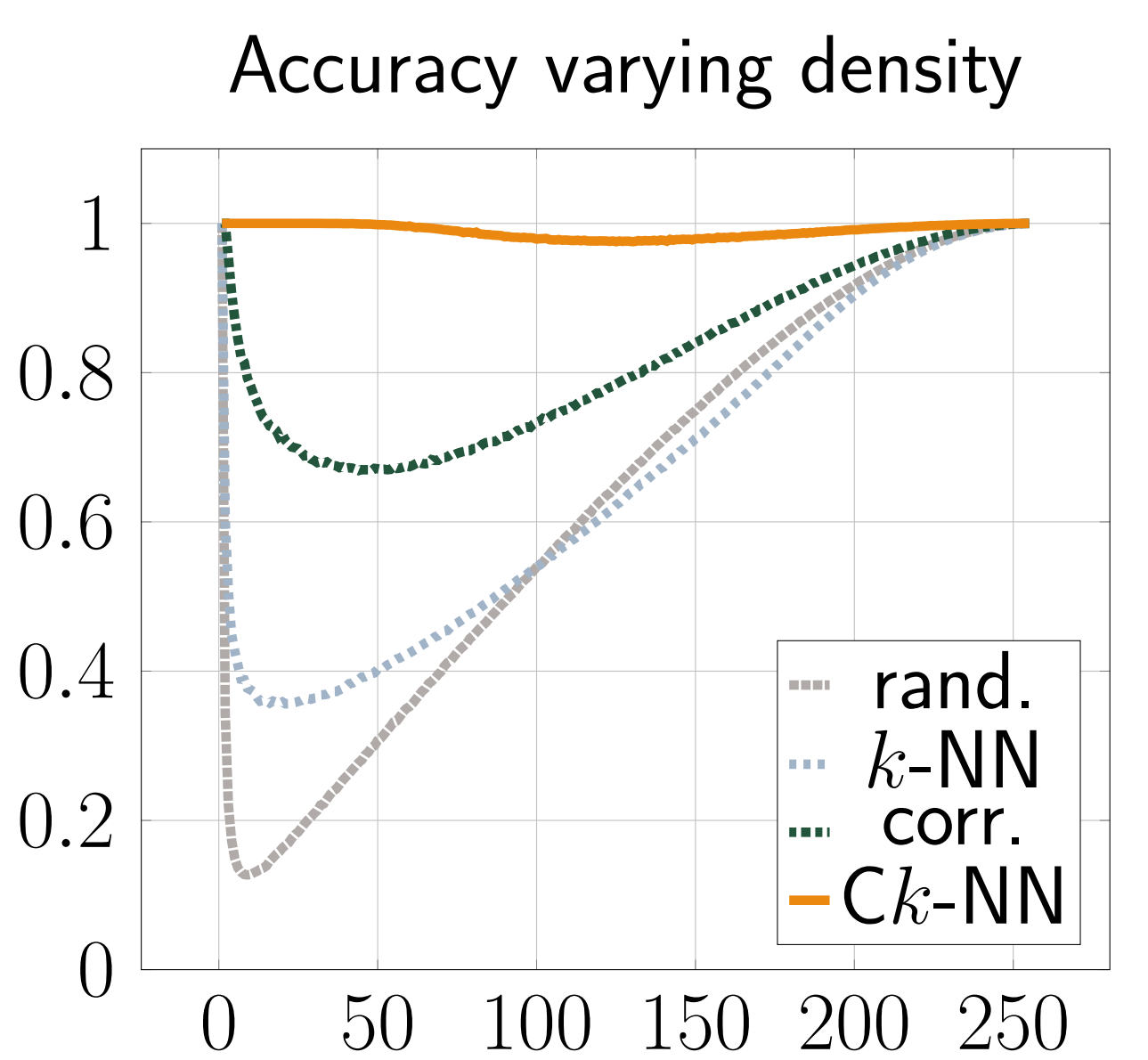
## Drop-in replacement for $k$-NN on MNIST

We classify an image by taking the mode label in $k$ selected images. C$k$-NN gives better accuracies on the MNIST dataset for every $k > 2$.
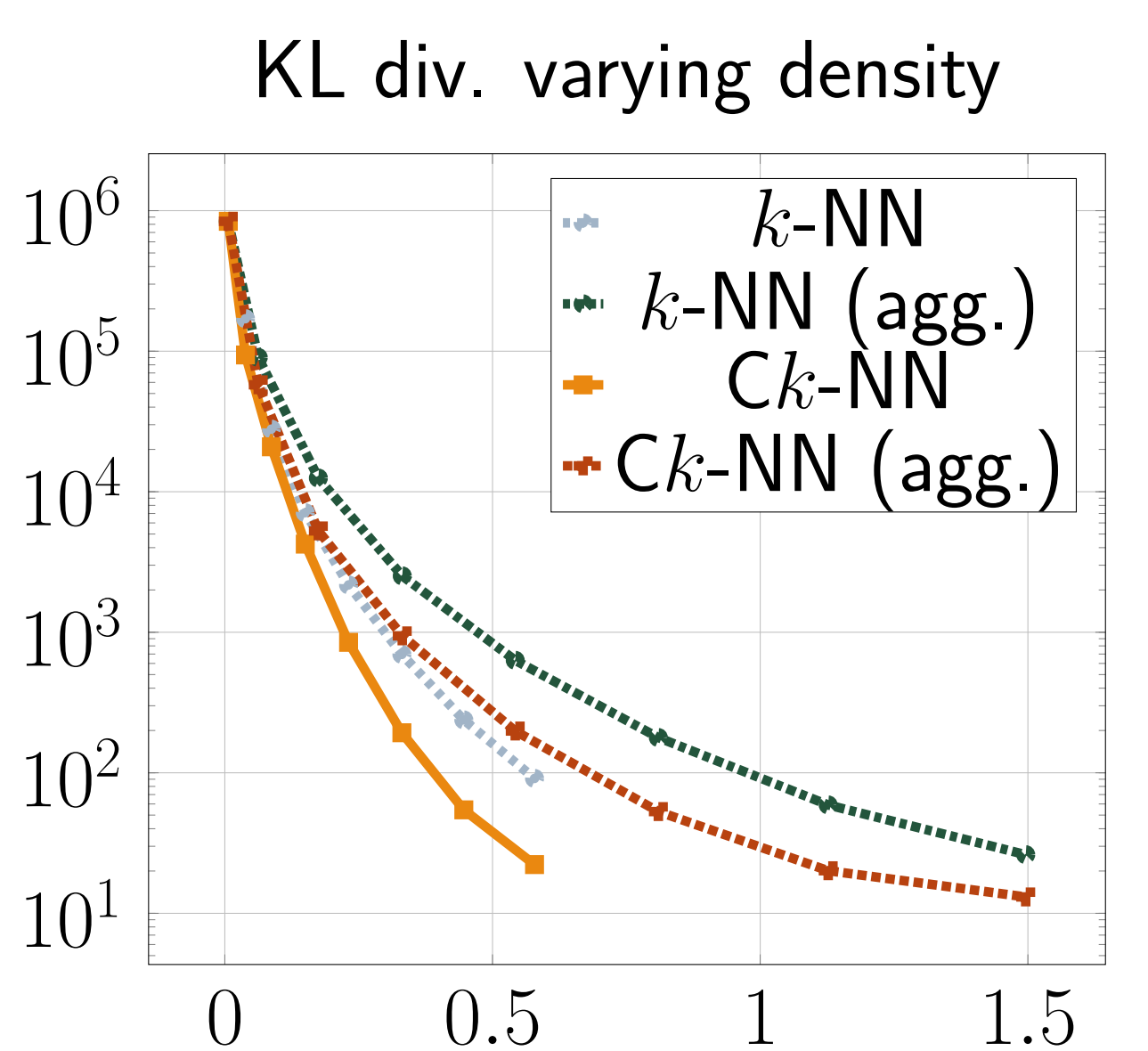

Accuracy varying $k$

## Recovery of sparse factors

Motivated by compressive sensing, we generate sparse factors $L$ to be recovered from measurements $LL^{\top}$. C$k$-NN recovers $L$ with near perfect accuracy over varying densities.


Accuracy varying density

## Better KL divergence with sparser factors

Plugging the selection algorithm into Cholesky factorization leads to better KL divergence for the same number of nonzero entries as $k$-NN.


KL div. varying density

## Preconditioning with conjugate gradient

Because minimizing KL divergence minimizes the Kaporin condition number, our method needs fewer iterations of the conjugate gradient to solve linear systems $\Theta\boldsymbol{x} = \boldsymbol{y}$.


Residual over iterations