

# Fast Gaussian process regression by Greedy Conditional Selection

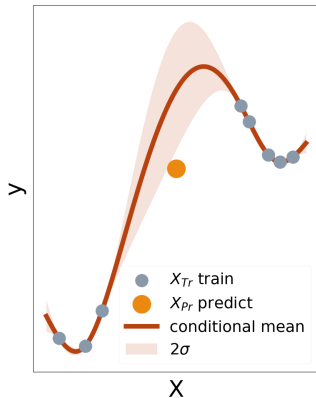
Stephen Huan    Florian Schäfer

Short & Sweet seminar

April 22, 2022

# The problem: Gaussian process regression

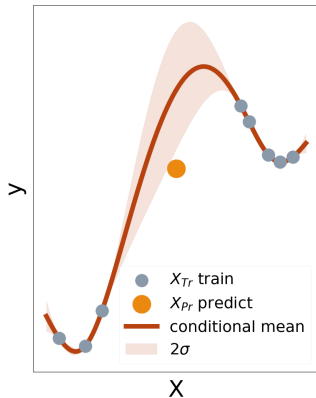
Measurements  $y_{Tr}$  at  $N$  points  $X_{Tr}$



# The problem: Gaussian process regression

Measurements  $\mathbf{y}_{Tr}$  at  $N$  points  $X_{Tr}$

Estimate unseen data  $\mathbf{y}_{Pr}$  at  $X_{Pr}$



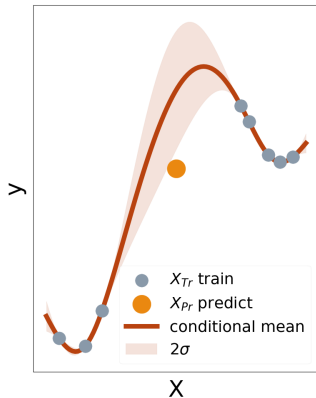
# The problem: Gaussian process regression

Measurements  $\mathbf{y}_{Tr}$  at  $N$  points  $X_{Tr}$

Estimate unseen data  $\mathbf{y}_{Pr}$  at  $X_{Pr}$

Model as Gaussian process

→ condition on  $\mathbf{y}_{Tr}$



## Cubic bottleneck

Closed-form conditional distribution:

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

$$\Theta_{Pr,Pr|Tr} := \text{COV}[\mathbf{y}_{Pr} \mid \mathbf{y}_{Tr}] = \Theta_{Pr,Pr} - \Theta_{Pr,Tr} \Theta_{Tr,Tr}^{-1} \Theta_{Tr,Pr}$$

## Cubic bottleneck

Closed-form conditional distribution:

$$E[\mathbf{y}_{Pr} | \mathbf{y}_{Tr}] = \boldsymbol{\mu}_{Pr} + \Theta_{Pr,Tr} \Theta_{Tr,Tr}^{-1} (\mathbf{y}_{Tr} - \boldsymbol{\mu}_{Tr})$$

$$\Theta_{Pr,Pr|Tr} := \text{COV}[\mathbf{y}_{Pr} | \mathbf{y}_{Tr}] = \Theta_{Pr,Pr} - \Theta_{Pr,Tr} \Theta_{Tr,Tr}^{-1} \Theta_{Tr,Pr}$$

Kernel function  $K(\cdot, \cdot)$ :  $\Theta_{i,j} := K(\mathbf{x}_i, \mathbf{x}_j)$

## Cubic bottleneck

Closed-form conditional distribution:

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

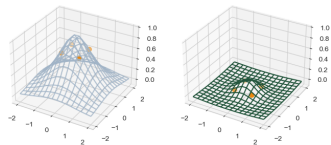
$$\Theta_{Pr,Pr|Tr} := \text{COV}[\mathbf{y}_{Pr} | \mathbf{y}_{Tr}] = \Theta_{Pr,Pr} - \Theta_{Pr,Tr} \Theta_{Tr,Tr}^{-1} \Theta_{Tr,Pr}$$

Kernel function  $K(\cdot, \cdot)$ :  $\Theta_{i,j} := K(\mathbf{x}_i, \mathbf{x}_j)$

Computational cost scales as  $N^3$

# Screening effect

“Screening effect”

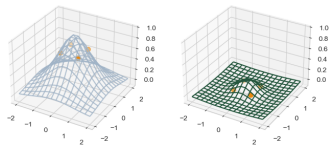




# Screening effect

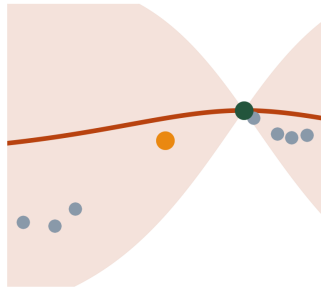
“Screening effect”

Choose  $k$  most informative points



## Conditional $k$ -th nearest neighbors

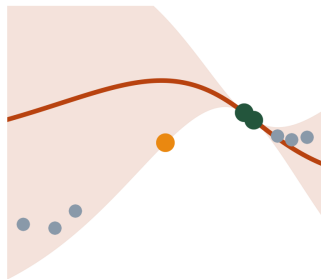
Naive: select  $k$  closest points



## Conditional $k$ -th nearest neighbors

Naive: select  $k$  closest points

Chooses redundant information

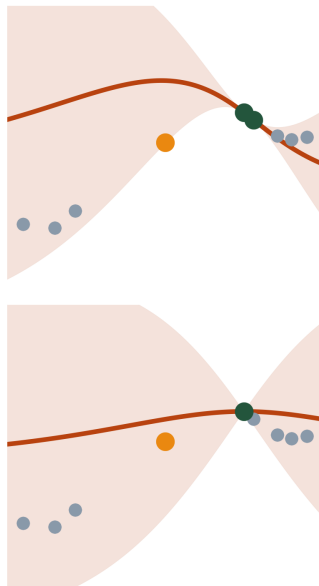


## Conditional $k$ -th nearest neighbors

Naive: select  $k$  closest points

Chooses redundant information

Maximize *mutual information!*

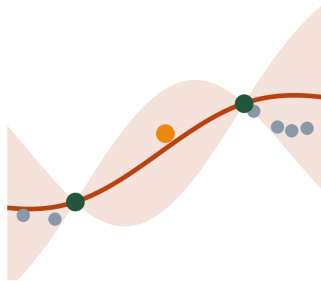
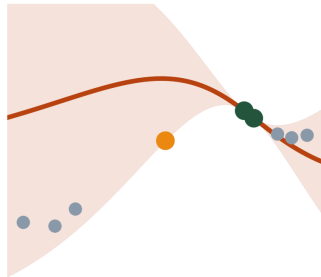


## Conditional $k$ -th nearest neighbors

Naive: select  $k$  closest points

Chooses redundant information

Maximize *mutual information!*



# Mutual information

Mutual information or information gain:

$$I[\mathbf{y}_{Pr}; \mathbf{y}_{Tr}] = H[\mathbf{y}_{Pr}] - H[\mathbf{y}_{Pr} | \mathbf{y}_{Tr}]$$

## Mutual information

Mutual information or information gain:

$$I[\mathbf{y}_{Pr}; \mathbf{y}_{Tr}] = H[\mathbf{y}_{Pr}] - H[\mathbf{y}_{Pr} | \mathbf{y}_{Tr}]$$

Entropy increases with log determinant of covariance

## Mutual information

Mutual information or information gain:

$$I[\mathbf{y}_{Pr}; \mathbf{y}_{Tr}] = H[\mathbf{y}_{Pr}] - H[\mathbf{y}_{Pr} | \mathbf{y}_{Tr}]$$

Entropy increases with log determinant of covariance

Information-theoretic EV-VE identity:

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



## Greedy mutual information maximization

Greedy selection, maintain selected indices  $I$

## Greedy mutual information maximization

Greedy selection, maintain selected indices  $I$

Criterion simplifies to:

$$\operatorname{argmax}_{j \notin I} \frac{\Theta_{j, \Pr|I}^2}{\Theta_{j,j|I}}$$

## Greedy mutual information maximization

Greedy selection, maintain selected indices  $I$

Criterion simplifies to:

$$\operatorname{argmax}_{j \notin I} \frac{\Theta_{j, \Pr|I}^2}{\Theta_{j,j|I}}$$

Direct computation:  $\mathcal{O}(Nk^4)$

## Greedy mutual information maximization

Greedy selection, maintain selected indices  $I$

Criterion simplifies to:

$$\operatorname{argmax}_{j \notin I} \frac{\Theta_{j, \text{Pr}|I}^2}{\Theta_{j,j|I}}$$

Direct computation:  $\mathcal{O}(Nk^4)$

Storing a partial Cholesky factor:  $\mathcal{O}(Nk^2)$

## Conditioning as rank-one update

Key idea: assume we have  $\Theta_{|I}$ , rank-one update to  $\Theta_{|I \cup \{k\}}$

$$\Theta_{:, : | I \cup \{k\}} = \Theta_{:, : | I} - \Theta_{:, k | I} \Theta_{k, k | I}^{-1} \Theta_{k, : | I}$$

$$\mathbf{u} = \frac{\Theta_{:, k | I}}{\sqrt{\Theta_{k, k | I}}}$$

$$\Theta_{|I \cup \{k\}} = \Theta_{|I} - \mathbf{u} \mathbf{u}^\top$$

## Efficient computation from Cholesky factor

Statistical interpretation of Cholesky factorization:

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

## Efficient computation from Cholesky factor

Statistical interpretation of Cholesky factorization:

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

Store partial Cholesky factor  $L$

$$L_i = \frac{\Theta_{:,k|i}}{\sqrt{\Theta_{kk|i}}}$$

## Algorithm

Indices  $I$ , select  $k$ ,  $i$ th iteration, have:

Conditional covariances  $\Theta_{:,Pr|I}$

Conditional variances  $\text{diag}(\Theta_{:, :|I})$

First  $i - 1$  columns of  $L$



## Algorithm

Indices  $I$ , select  $k$ ,  $i$ th iteration, have:

Conditional covariances  $\Theta_{:,Pr|I}$

Conditional variances  $\text{diag}(\Theta_{:, :|I})$

First  $i - 1$  columns of  $L$

Update  $L$

$$L_{:,i} \leftarrow \Theta_{:,k} - L_{:,i-1} L_{k,:i-1}^\top$$

## Algorithm

Indices  $I$ , select  $k$ ,  $i$ th iteration, have:

Conditional covariances  $\Theta_{:,Pr|I}$

Conditional variances  $\text{diag}(\Theta_{:, :|I})$

First  $i - 1$  columns of  $L$

Update  $L$

$$L_{:,i} \leftarrow \Theta_{:,k} - L_{:, :i-1} L_{k, :i-1}^\top$$

Update conditional values for candidate  $j$

$$\begin{aligned}\Theta_{jj|I \cup \{k\}} &\leftarrow \Theta_{jj|I} - L_{j,i}^2 \\ \Theta_{j,Pr|I \cup \{k\}} &\leftarrow \Theta_{j,Pr|I} - L_{j,i} L_{Pr,i}\end{aligned}$$

## Extending to multiple prediction points

Objective conditional log determinant of prediction points

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

## Extending to multiple prediction points

Objective conditional log determinant of prediction points

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

By the matrix determinant lemma,

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

## Extending to multiple prediction points

Objective conditional log determinant of prediction points

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

By the matrix determinant lemma,

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

By the quotient rule, we combine the conditioning:

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

## Algorithm for multiple prediction points

Final objective simplifies to:

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

## Algorithm for multiple prediction points

Final objective simplifies to:

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

Store *two* factors (one for numerator, one for denominator)

## Algorithm for multiple prediction points

Final objective simplifies to:

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

Store *two* factors (one for numerator, one for denominator)

“Pre-condition” numerator factor on prediction points

$$\Theta_{kk | I, Pr} = \Theta_{kk | Pr, I}$$



## Algorithm for multiple prediction points

Final objective simplifies to:

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

Store *two* factors (one for numerator, one for denominator)

“Pre-condition” numerator factor on prediction points

$$\Theta_{kk | I, Pr} = \Theta_{kk | Pr, I}$$

Complexity of  $\mathcal{O}(Nk^2 + Nm^2 + m^3)$  for  $m$  prediction points

Global approximation by KL-minimization

Approximate GP by sparse Cholesky factor of its precision

## Global approximation by KL-minimization

Approximate GP by sparse Cholesky factor of its precision

Measure resulting approximation accuracy by KL divergence:

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

## Global approximation by KL-minimization

Approximate GP by sparse Cholesky factor of its precision

Measure resulting approximation accuracy by KL divergence:

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

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

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

## Global approximation by KL-minimization

Approximate GP by sparse Cholesky factor of its precision

Measure resulting approximation accuracy by KL divergence:

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

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

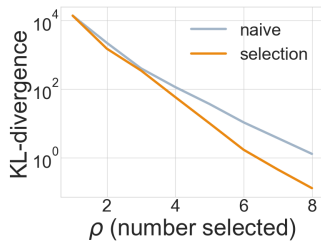
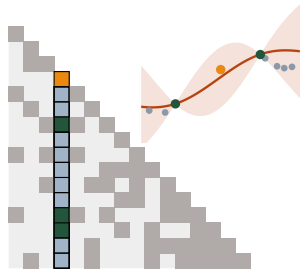
$$L_{s_i, i} = \frac{\Theta_{s_i, s_i}^{-1} \mathbf{e}_1}{\sqrt{\mathbf{e}_1^\top \Theta_{s_i, s_i}^{-1} \mathbf{e}_1}}$$

Minimize variance of  $i$ th point, conditional on selected!

$$2\mathbb{D}_{\text{KL}} = \sum_{i=1}^N \left[ \log(\Theta_{ii|s_i - \{i\}}) - \log(\Theta_{ii|i+1:}) \right]$$

# Cholesky factorization by selection

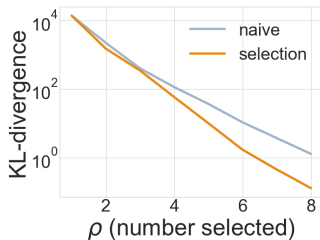
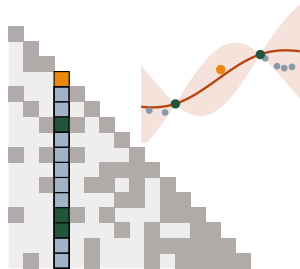
Apply column-wise directly



# Cholesky factorization by selection

Apply column-wise directly

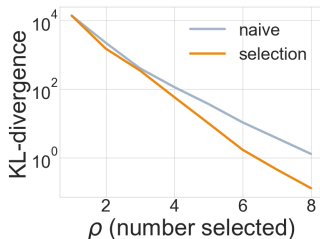
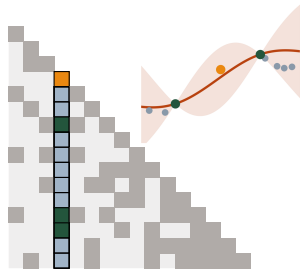
Improves approx. algorithm of <sup>1</sup>



# Cholesky factorization by selection

Apply column-wise directly

Improves approx. algorithm of <sup>1</sup>



---

<sup>1</sup>F. Schäfer, M. Katzfuss, and H. Owhadi, "Sparse Cholesky factorization by Kullback-Leibler minimization," *arXiv preprint arXiv:2004.14455*, 2020



## GP regression by Cholesky factorization

Lower triangular factor for precision:  $LL^T = \Theta^{-1}$

## GP regression by Cholesky factorization

Lower triangular factor for precision:  $LL^{\top} = \Theta^{-1}$

Upper triangular factor for covariance:  $L^{-\top}L^{-1} = \Theta$

## GP regression by Cholesky factorization

Lower triangular factor for precision:  $LL^\top = \Theta^{-1}$

Upper triangular factor for covariance:  $L^{-\top}L^{-1} = \Theta$

$U = L^{-\top}$ , look at submatrices:

$$\begin{aligned}\Theta &= UU^\top = \begin{pmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{pmatrix} \begin{pmatrix} U_{11}^\top & 0 \\ U_{12}^\top & U_{22}^\top \end{pmatrix} \\ &= \begin{pmatrix} U_{11}U_{11}^\top + U_{12}U_{12}^\top & U_{12}U_{22}^\top \\ U_{22}U_{12}^\top & U_{22}U_{22}^\top \end{pmatrix}\end{aligned}$$

$$U_{22} = \text{chol}(\Theta_{22})$$

$$U_{12} = \Theta_{12}U_{22}^{-\top}$$

$$U_{11} = \text{chol}(\Theta_{11|2})$$

## GP regression by Cholesky factorization

Write conditional terms:

$$\begin{aligned} \mathbb{E}[\mathbf{y}_1 \mid \mathbf{y}_2] &= \Theta_{12}\Theta_{22}^{-1}\mathbf{y}_2 \\ &= U_{12}U_{22}^{-1}\mathbf{y}_2 \end{aligned}$$

$$\begin{aligned} \text{Cov}[\mathbf{y}_1 \mid \mathbf{y}_2] &= \Theta_{11} - \Theta_{12}\Theta_{22}^{-1}\Theta_{21} \\ &= U_{11}U_{11}^\top \end{aligned}$$

## GP regression by Cholesky factorization

Write conditional terms:

$$\begin{aligned} \mathbb{E}[\mathbf{y}_1 \mid \mathbf{y}_2] &= \Theta_{12}\Theta_{22}^{-1}\mathbf{y}_2 \\ &= U_{12}U_{22}^{-1}\mathbf{y}_2 \\ \text{Cov}[\mathbf{y}_1 \mid \mathbf{y}_2] &= \Theta_{11} - \Theta_{12}\Theta_{22}^{-1}\Theta_{21} \\ &= U_{11}U_{11}^\top \end{aligned}$$

Recall:  $U = L^{-\top}$  so  $UL^\top = L^\top U = I$

$$\begin{aligned} \begin{pmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{pmatrix} \begin{pmatrix} L_{11}^\top & L_{21}^\top \\ 0 & L_{22}^\top \end{pmatrix} &= \begin{pmatrix} L_{11}^\top & L_{21}^\top \\ 0 & L_{22}^\top \end{pmatrix} \begin{pmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{pmatrix} = I \\ \begin{pmatrix} U_{11}L_{11}^\top & U_{11}L_{21}^\top + U_{12}L_{22}^\top \\ 0 & U_{22}L_{22}^\top \end{pmatrix} &= \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} = I \\ \begin{pmatrix} L_{11}^\top U_{11} & L_{11}^\top U_{12} + L_{21}^\top U_{22} \\ 0 & L_{22}^\top U_{22} \end{pmatrix} &= \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix} = I \end{aligned}$$

## GP regression by Cholesky factorization

Reading from submatrices,

$$U_{11} = L_{11}^{-\top}$$

$$U_{22} = L_{22}^{-\top}$$

$$U_{12} = -L_{11}^{-\top} L_{21}^{\top} L_{22}^{-\top}$$

## GP regression by Cholesky factorization

Reading from submatrices,

$$U_{11} = L_{11}^{-\top}$$

$$U_{22} = L_{22}^{-\top}$$

$$U_{12} = -L_{11}^{-\top} L_{21}^{\top} L_{22}^{-\top}$$

Re-write conditional terms:

$$\begin{aligned} \mathbb{E}[\mathbf{y}_1 \mid \mathbf{y}_2] &= U_{12} U_{22}^{-1} \mathbf{y}_2 \\ &= (-L_{11}^{-\top} L_{21}^{\top} L_{22}^{-\top}) L_{22}^{\top} \mathbf{y}_2 \\ &= -L_{11}^{-\top} L_{21}^{\top} \mathbf{y}_2 \end{aligned}$$

$$\begin{aligned} \text{Cov}[\mathbf{y}_1 \mid \mathbf{y}_2] &= U_{11} U_{11}^{\top} \\ &= L_{11}^{-\top} L_{11}^{-1} \end{aligned}$$

$$\mathbf{e}_i^{\top} \text{Cov}[\mathbf{y}_1 \mid \mathbf{y}_2] \mathbf{e}_j = (L_{11}^{-1} \mathbf{e}_i)^{\top} (L_{11}^{-1} \mathbf{e}_j)$$

## Summary

Selection algorithm for Gaussian process regression

Drop-in replacement for  $k$ -th nearest neighbors

Leverage GP regression for sparse Cholesky factorization

Leverage Cholesky factorization for GP regression

Thank you!