# Sparse Cholesky Factorization by Greedy Conditional Selection

Stephen Huan

Theory Club

February 28, 2022

# Table of Contents

# The Problem: Gaussian Process Regression

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

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

Model as Gaussian process
$\rightarrow$ condition on $\boldsymbol{y}_{\mathsf{Tr}}$

Computational cost scales as $N^3$

Choose $k$ most informative points!

# Conditional $k$-th Nearest Neighbors



Naive: select $k$ closest points

Chooses redundant information

Maximize *mutual information*!

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

Store Cholesky factor $\rightarrow \mathcal{O}(Nk^2)$!
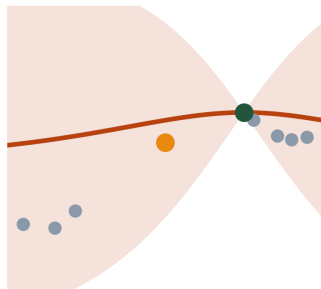
# Conditional $k$-th Nearest Neighbors



Naive: select $k$ closest points

Chooses redundant information

Maximize *mutual information*!

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

Store Cholesky factor $\rightarrow \mathcal{O}(Nk^2)$!

# Conditional $k$-th Nearest Neighbors

Naive: select $k$ closest points

Chooses redundant information

Maximize *mutual information*!

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

Store Cholesky factor $\rightarrow \mathcal{O}(Nk^2)$!
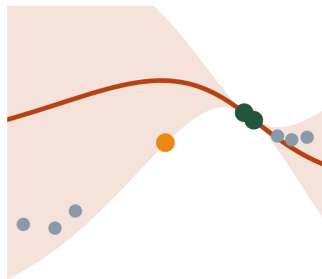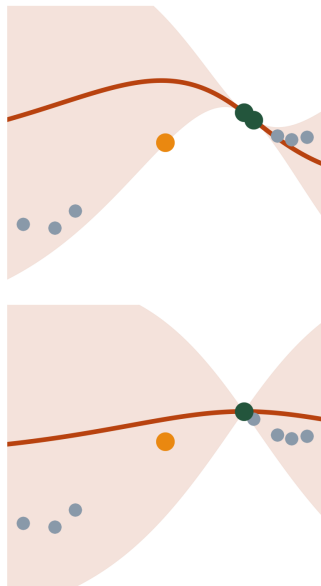
# Conditional $k$-th Nearest Neighbors



Naive: select $k$ closest points

Chooses redundant information

Maximize *mutual information*!

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

Store Cholesky factor $\rightarrow \mathcal{O}(Nk^2)$!

# Cholesky Factorization by Selection
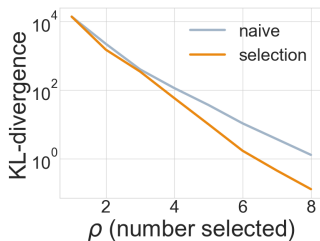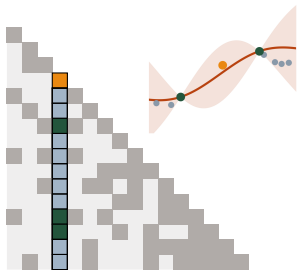


Apply column-wise
$\rightarrow$ sparse approx. of GP

Maximum mutual information
$\rightarrow$ minimum KL-divergence

Improves approx. algorithm of [1]



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

# Table of Contents

$M = LU$ where $L$ is lower triangular and $U$ is upper triangular

Not always possible, need $PLU$ in general!

Special case for (square) symmetric matrices:

## Theorem
*If $M = M^\top$ and $\det(M) \neq 0$, then $M = LDL^T$ where $L$ is from the LU decomposition of $M$ and $D$ is the diagonal of $U$.*

## Proof sketch.
(MATH3406 Fall 2021, Prof. Wing Li) Let $M = LDK$. Just do matrix multiplication on $M = M^\top \implies (LDK) = (LDK)^T$. From matrix multiplication, able to see $K = L^\top$. $\qquad\square$

Let $M$ be (symmetric) *positive definite*.

Then $M = LDL^\top$ becomes $LL^\top$:

$$M = LDL^\top$$
$$= LD^{\frac{1}{2}}D^{\frac{1}{2}}L^\top$$
$$= LD^{\frac{1}{2}}(LD^{\frac{1}{2}})^\top$$
$$= L'L'^\top$$

This is the Cholesky factorization!

$\Theta = LL^\top$, $L$ has $N$ columns, $s$ non-zero entries per column

$L\boldsymbol{v}$ and $L^{-1}\boldsymbol{v}$ both cost $\mathcal{O}(Ns)$

Matrix-vector product $\Theta\boldsymbol{v} \to L(L^\top \boldsymbol{v})$
$$N^2 \to Ns$$

Solving linear system $\Theta^{-1}\boldsymbol{v} \to L^{-\top}(L^{-1}\boldsymbol{v})$
$$N^3 \to Ns$$

Log determinant $\mathrm{logdet}\,\Theta \to 2\,\mathrm{logdet}\,L = 2\sum_{i=1}^{N} \log L_{ii}$
$$N^3 \to N$$

Sampling from $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Theta) \to \boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, I), \boldsymbol{x} = L\boldsymbol{z} + \boldsymbol{\mu}$
$$??? \to Ns$$

Like LU

Gaussian elimination downwards

```
1   def down_cholesky(theta: np.ndarray) -> np.ndarray:
2       M, n = np.copy(theta), len(theta)
3       L = np.identity(n)
4       for i in range(n):
5           for j in range(i + 1, n):
6               L[j, i] = M[j, i]/M[i, i]
7               # zero out everything below
8               M[j] -= L[j, i]*M[i]
9           # update L
10          L[:, i] *= np.sqrt(M[i, i])
11      return L
```

Let $L'$ be blocked according to:

$$L' = \begin{pmatrix} L & \mathbf{0} \\ \mathbf{r}^\top & d \end{pmatrix}$$

$$L'L'^\top = \begin{pmatrix} L & \mathbf{0} \\ \mathbf{r}^\top & d \end{pmatrix} \begin{pmatrix} L^\top & \mathbf{r} \\ \mathbf{0}^\top & d \end{pmatrix}$$

$$= \begin{pmatrix} LL^\top & L\mathbf{r} \\ \mathbf{r}^\top L^\top & \mathbf{r}^\top \mathbf{r} + d^2 \end{pmatrix}$$

So if we have a Cholesky factor for a principle submatrix of $\Theta$, we can extend it inductively by reading off the appropiate data!

$$\begin{pmatrix} LL^\top & L\mathbf{r} \\ \mathbf{r}^\top L^\top & \mathbf{r}^\top \mathbf{r} + d^2 \end{pmatrix} = \begin{pmatrix} \Theta & \mathbf{c} \\ \mathbf{c}^\top & C \end{pmatrix}$$

$$\mathbf{r} = L^{-1}\mathbf{c}$$

$$d = \sqrt{C - \mathbf{r}^\top \mathbf{r}}$$

```python
1   def Lsolve(L: np.ndarray, y: np.ndarray) -> np.ndarray:
2       """ Solves Lx = y for lower triangular L. """
3       n = len(y)
4       x = np.zeros(n)
5       for i in range(n):
6           x[i] = (y[i] - L[i, :i].dot(x[:i]))/L[i, i]
7       return x
8
9   def up_cholesky(theta: np.ndarray) -> np.ndarray:
10      n = len(theta)
11      L = np.zeros((n, n))
12      for i in range(n):
13          row = Lsolve(L, theta[:i, i])
14          L[i, :i] = row
15          L[i, i] = np.sqrt(theta[i, i] - row.dot(row))
16      return L
```

$$L = \begin{pmatrix} \boldsymbol{l}_1 & \boldsymbol{l}_2 & \cdots & \boldsymbol{l}_N \end{pmatrix}$$

$$LL^\top = \begin{pmatrix} \boldsymbol{l}_1 & \boldsymbol{l}_2 & \cdots & \boldsymbol{l}_N \end{pmatrix} \begin{pmatrix} \boldsymbol{l}_1^\top \\ \boldsymbol{l}_2^\top \\ \vdots \\ \boldsymbol{l}_N^\top \end{pmatrix}$$

$$= \boldsymbol{l}_1\boldsymbol{l}_1^\top + \boldsymbol{l}_2\boldsymbol{l}_2^\top + \cdots + \boldsymbol{l}_N\boldsymbol{l}_N^\top = \Theta$$

From lower triangularity, nested submatrices!

$$\boldsymbol{l}_1\boldsymbol{l}_1^\top + \boldsymbol{l}_2\boldsymbol{l}_2^\top + \cdots + \boldsymbol{l}_N\boldsymbol{l}_N^\top = \Theta$$

$$\boldsymbol{l}_1\boldsymbol{l}_1^\top = \Theta_1$$

$$l_1^2 = \Theta_{11}$$

$$l_1 = \sqrt{\Theta_{11}}$$

$$\boldsymbol{l}_1 = \frac{\Theta_1}{l_1} = \frac{\Theta_1}{\sqrt{\Theta_{11}}}$$

$$\boldsymbol{l}_2\boldsymbol{l}_2^\top + \cdots + \boldsymbol{l}_N\boldsymbol{l}_N^\top = \Theta - \left(\frac{\Theta_1}{\sqrt{\Theta_{11}}}\right)\left(\frac{\Theta_1}{\sqrt{\Theta_{11}}}\right)^\top$$

$$= \Theta - \frac{\Theta_1\Theta_1^\top}{\Theta_{11}}$$

Proceed inductively on rank-one update

# Computing the Cholesky Factorization

```python
1  def right_cholesky(theta: np.ndarray) -> np.ndarray:
2      M, n = np.copy(theta), len(theta)
3      L = np.zeros((n, n))
4      for i in range(n):
5          L[:, i] = M[:, i]/np.sqrt(M[i, i])
6          M -= np.outer(L[:, i], L[:, i])
7      return L
```

Recall:
$$\boldsymbol{l}_1\boldsymbol{l}_1^\top + \boldsymbol{l}_2\boldsymbol{l}_2^\top + \cdots + \boldsymbol{l}_N\boldsymbol{l}_N^\top = \Theta$$

Look at $\boldsymbol{l}_i$:

$$\boldsymbol{l}_i\boldsymbol{l}_i^\top = \left(\Theta - (\boldsymbol{l}_1\boldsymbol{l}_1^\top + \boldsymbol{l}_2\boldsymbol{l}_2^\top + \cdots + \boldsymbol{l}_{i-1}\boldsymbol{l}_{i-1}^\top)\right)_i$$

$$= \Theta_i - (l_{1i}\boldsymbol{l}_1 + l_{2i}\boldsymbol{l}_2 + \cdots + l_{i-1,i}\boldsymbol{l}_{i-1})$$

$$= \Theta_i - \begin{pmatrix} \boldsymbol{l}_1 & \boldsymbol{l}_2 & \cdots & \boldsymbol{l}_{i-1} \end{pmatrix} \begin{pmatrix} l_{1i} \\ l_{2i} \\ \vdots \\ l_{i,i-1} \end{pmatrix}$$

$$= \Theta_i - L_{:,:i}L_{i,:i}$$

Don't need to store modified $\Theta$ in memory!

# Computing the Cholesky Factorization

```python
def left_cholesky(theta: np.ndarray) -> np.ndarray:
    n = len(theta)
    L = np.zeros((n, n))
    for i in range(n):
        L[:, i] = theta[:, i] - L[:, :i]@L[i, :i]
        L[:, i] /= np.sqrt(L[i, i])
    return L
```

# Table of Contents

# Schur Complement
## or recursive Cholesky factorization

Block $\Theta$ as follows:

$$\Theta = \begin{pmatrix} \Theta_{11} & \Theta_{12} \\ \Theta_{21} & \Theta_{22} \end{pmatrix}$$

Then proceed by one step of Gaussian elimination:

$$\begin{pmatrix} \Theta_{11} & \Theta_{12} \\ \mathbf{0} & \Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12} \end{pmatrix}$$

Thus,

$$= \begin{pmatrix} I & 0 \\ \Theta_{21}\Theta_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} \Theta_{11} & 0 \\ 0 & \Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12} \end{pmatrix} \begin{pmatrix} I & \Theta_{11}^{-1}\Theta_{12} \\ 0 & I \end{pmatrix}$$

so we see the Cholesky factorization of $\Theta$ is

$$\begin{pmatrix} I & 0 \\ \Theta_{21}\Theta_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} \text{chol}(\Theta_{11}) & 0 \\ 0 & \text{chol}(\Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12}) \end{pmatrix}$$

The term in blue is the *Schur complement* of $\Theta$ on $\Theta_{11}$

## Proper Determinant of Block Matrix

$$\Theta = \begin{pmatrix} \Theta_{11} & \Theta_{12} \\ \Theta_{21} & \Theta_{22} \end{pmatrix}$$

$\det(\Theta) = ?$

$\quad = \det(\Theta_{11})\det(\Theta_{22}) - \det(\Theta_{21})\det(\Theta_{12})?$      wrong!

$\quad = \det(\Theta_{11}\Theta_{22} - \Theta_{21}\Theta_{12})?$                wrong!

Schur complement gives proper answer:

$$\Theta = \begin{pmatrix} I & 0 \\ \Theta_{21}\Theta_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} \Theta_{11} & 0 \\ 0 & \Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12} \end{pmatrix} \begin{pmatrix} I & \Theta_{11}^{-1}\Theta_{12} \\ 0 & I \end{pmatrix}$$

$$\det(\Theta) = \det(\Theta_{11})\det(\Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12})$$

$$\Theta = \begin{pmatrix} \Theta_{11} & \Theta_{12} \\ \Theta_{21} & \Theta_{22} \end{pmatrix}$$

$$(\Theta^{-1})_{22} = ?$$

$$= (\Theta_{22})^{-1}? \qquad \text{wrong!}$$

Schur complement to the rescue again!

## Proper Submatrix of Inverse

$$\Theta = \begin{pmatrix} I & 0 \\ \Theta_{21}\Theta_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} \Theta_{11} & 0 \\ 0 & \Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12} \end{pmatrix} \begin{pmatrix} I & \Theta_{11}^{-1}\Theta_{12} \\ 0 & I \end{pmatrix}$$

For notational convenience, we denote the Schur complement $\Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12}$ as $\Theta_{22|1}$. Inverting both sides of the equation,

$$\Theta^{-1} = \begin{pmatrix} I & -\Theta_{11}^{-1}\Theta_{12} \\ 0 & I \end{pmatrix} \begin{pmatrix} \Theta_{11}^{-1} & 0 \\ 0 & \Theta_{22|1}^{-1} \end{pmatrix} \begin{pmatrix} I & 0 \\ -\Theta_{21}\Theta_{11}^{-1} & I \end{pmatrix}$$

$$= \begin{pmatrix} \Theta_{11}^{-1} + (\Theta_{11}^{-1}\Theta_{12})\Theta_{22|1}^{-1}(\Theta_{21}\Theta_{11}^{-1}) & -(\Theta_{11}^{-1}\Theta_{12})\Theta_{22|1}^{-1} \\ -\Theta_{22|1}^{-1}(\Theta_{21}\Theta_{11}^{-1}) & \Theta_{22|1}^{-1} \end{pmatrix}$$

So $(\Theta^{-1})_{22}$ can be read off as $\Theta_{22|1}^{-1}$,

$$= \left(\Theta_{22} - \Theta_{21}\Theta_{11}^{-1}\Theta_{12}\right)^{-1}$$

## A Few Important Questions...

Is the Schur complement symmetric positive definite (s.p.d.)?

    If it isn't, we're kinda screwed — have been assuming so

Is Schur complementing transitive?

    i.e. suppose we have $\Theta$ blocked as

$$\Theta = \begin{pmatrix} \Theta_{11} & \Theta_{12} & \Theta_{13} \\ \Theta_{21} & \Theta_{22} & \Theta_{23} \\ \Theta_{31} & \Theta_{32} & \Theta_{33} \end{pmatrix}$$

    Is $\Theta$ complemented on $\Theta_{11}$ and then on $\Theta_{22}$ the same as
    $\Theta$ complemented on $\begin{pmatrix} \Theta_{11} & \Theta_{12} \\ \Theta_{21} & \Theta_{22} \end{pmatrix}$?

    Intuitively, it should be, but tedious to prove

New perspective which changes everything!

# Table of Contents

## The Multivariate Gaussian

Recall: Gaussian (or normal) distribution:

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$$

Important (defining?) property: completely determined by mean and variance, all higher-order cumulants zero.

We're going to extend this to higher dimensions. Consider

$$\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

where $\boldsymbol{x}$ ("variables") is a $N \times 1$ vector, $\boldsymbol{\mu}$ ("mean vector") is a $N \times 1$ vector, and $\Sigma$ ("covariance matrix") is a $N \times N$ matrix

Naturally,

$$\mu_i = \mathrm{E}[x_i]$$
$$\boldsymbol{\mu} = \mathrm{E}[\boldsymbol{x}]$$
$$\Sigma_{ij} = \mathrm{Cov}[x_i, x_j]$$
$$= \mathrm{E}[(x_i - \mathrm{E}[x_i])(x_j - \mathrm{E}[x_j])]$$
$$= \mathrm{E}[(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^\top]$$

Two natural (and fundamental) questions from here:
1. What is the probability density function $f(\boldsymbol{x})$?
2. How can we sample from $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$?

Surprisingly enough, Cholesky factorization answers both!

## Independent Variables

Gaussian has the (unique?) property if $\Sigma_{ij} = 0$, then $x_i$ and $x_j$ are statistically independent. This is not true in general!

Key property we will make heavy use of: moment matching. If we know $\mu$ and $\Sigma$, distribution is determined.

Consider: if $x_i$ and $x_j$ *were* independent, then $\Sigma_{ij} = 0$. So suppose $x_i$ and $x_j$ are not independent but $\Sigma_{ij} = 0$. It's the same $\Sigma$ as when they were independent. So $x_i$ and $x_j$ must be distributed like they're independent. By contradiction, they must have been independent in the first place!

Well, if $\Sigma$ has particular structure, it's actually trivial:

$$z \sim \mathcal{N}(\mathbf{0}, I_N)$$

$$z_i \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$$

$$f(z) = \prod_{i=1}^{N} f(z_i)$$

$$= \prod_{i=1}^{N} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} z_i^2}$$

$$= \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}(z_1^2 + z_2^2 + \cdots + z_N^2)}$$

$$= \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2} z^\top z}$$

How can we generalize to arbitrary $\Sigma$?

Moment match!

$$\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, I_N)$$
$$\boldsymbol{x} = L\boldsymbol{z} + \boldsymbol{\mu}$$
$$\mathrm{E}[\boldsymbol{x}] = \mathrm{E}[L\boldsymbol{z} + \boldsymbol{\mu}] = L\,\mathrm{E}[\boldsymbol{z}] + \boldsymbol{\mu} = \boldsymbol{\mu}$$
$$\mathrm{Cov}[\boldsymbol{x}] = \mathrm{E}[(\boldsymbol{x} - \mathrm{E}[\boldsymbol{x}])(\boldsymbol{x} - \mathrm{E}[\boldsymbol{x}])^\top]$$
$$= \mathrm{E}[L\boldsymbol{z}(L\boldsymbol{z})^\top]$$
$$= \mathrm{E}[L\boldsymbol{z}\boldsymbol{z}^\top L^\top]$$
$$= L\,\mathrm{E}[\boldsymbol{z}\boldsymbol{z}^\top]L^\top$$
$$= LL^\top$$

so $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, LL^\top)$. We want $\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$, so $\Sigma = LL^\top$

## Sampling with Cholesky Factorization

As we just saw, we can sample $x \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$ by instead sampling $z \sim \mathcal{N}(\mathbf{0}, I_N)$ and computing $x = Lz + \boldsymbol{\mu}$.

Since $LL^\top = \Sigma$, a natural pick is $L = \text{chol}(\Sigma)$.

Why is $\Sigma$ s.p.d.? Because it's a covariance/Gram matrix!

$$\Sigma = \text{E}[(x - \boldsymbol{\mu})(x - \boldsymbol{\mu})^\top]$$
$$y^\top \Sigma y = y^\top \text{E}[(x - \boldsymbol{\mu})(x - \boldsymbol{\mu})^\top] y$$
$$= \text{E}[y^\top (x - \boldsymbol{\mu})(x - \boldsymbol{\mu})^\top y]$$
$$= \text{E}[((x - \boldsymbol{\mu})^\top y)^\top (x - \boldsymbol{\mu})^\top y]$$
$$= \text{E}[\|(x - \boldsymbol{\mu})^\top y\|^2] \geq 0$$

## Probability Density Function from Sampling

What's the probability density function $f(x)$?

Idea: view $x$ resulting from a invertible transformation from $z$.

We know $f(z)$, so $f(x)$ should be similar!

In scalars:

$$z \sim \mathcal{N}(0, 1)$$
$$x = \sigma z + \mu$$
$$x \sim \mathcal{N}(\mu, \sigma^2)$$
$$z = \frac{x - \mu}{\sigma}$$

Since $f(z)$ is a valid probability density function,

$$1 = \int_{-\infty}^{\infty} f(z)\, \mathrm{d}z = \int_{-\infty}^{\infty} f(z) \frac{\mathrm{d}z}{\mathrm{d}x}\, \mathrm{d}x$$

We now perform the change of variables $z = \frac{x-\mu}{\sigma}$

$$= \underbrace{\int_{-\infty}^{\infty} f\left(\frac{x-\mu}{\sigma}\right) \frac{1}{\sigma}\, \mathrm{d}x}_{\text{PDF of } x}$$

$$f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}$$

$$\frac{1}{\sigma} f\left(\frac{x-\mu}{\sigma}\right) = \frac{1}{\sigma} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$$

## PDF from Sampling — Vector Edition

$$\boldsymbol{x} = L\boldsymbol{z} + \boldsymbol{\mu}$$

$$\boldsymbol{z} = L^{-1}(\boldsymbol{x} - \boldsymbol{\mu})$$

Since $f(\boldsymbol{z})$ is a valid probability density function,

$$
\begin{aligned}
1 &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\boldsymbol{z}) \, \mathrm{d}\boldsymbol{z} \\
&= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\boldsymbol{z}) \frac{\mathrm{d}\boldsymbol{z}}{\mathrm{d}\boldsymbol{x}} \, \mathrm{d}\boldsymbol{x} && \text{(informal)} \\
&= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\boldsymbol{z}) |\det(J_{\boldsymbol{z}})| \, \mathrm{d}\boldsymbol{x} && \text{(formal)} \\
&= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \underbrace{f(L^{-1}(\boldsymbol{x} - \boldsymbol{\mu})) \det(L^{-1})}_{\text{PDF of } \boldsymbol{x}} \, \mathrm{d}\boldsymbol{x}
\end{aligned}
$$

## PDF from Sampling — Vector Edition

$$f(\boldsymbol{z}) = \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}\boldsymbol{z}^\top \boldsymbol{z}}$$

Expanding $\det(L^{-1}) f(L^{-1}(\boldsymbol{x} - \boldsymbol{\mu}))$,

$$= \frac{1}{\det(L)} f(L^{-1}(\boldsymbol{x} - \boldsymbol{\mu}))$$

$$= \frac{1}{\det(L)} \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}(L^{-1}(\boldsymbol{x}-\boldsymbol{\mu}))^\top (L^{-1}(\boldsymbol{x}-\boldsymbol{\mu}))}$$

Since $LL^\top = \Sigma$, $\det(\Sigma) = \det(L)^2$

$$= \frac{1}{\sqrt{(2\pi)^N \det(\Sigma)}} e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top L^{-T} L^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}$$

$$= \frac{1}{\sqrt{(2\pi)^N \det(\Sigma)}} e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}$$

Compare PDFs of multivariate normal and scalar normal:

$$\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

$$f(\boldsymbol{x}) = \frac{1}{\sqrt{(2\pi)^N \det(\Sigma)}} e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}$$

Compare to scalar:

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}}$$

Remarkable similarity!

Sampling: $\boldsymbol{x} = L\boldsymbol{z} + \mu$, matrix-vector product, $\mathcal{O}(Ns)$

Density computation:

$$
\begin{aligned}
(\boldsymbol{x} - \boldsymbol{\mu})^{\top} \Sigma^{-1} (\boldsymbol{x} - \boldsymbol{\mu}) &= (\boldsymbol{x} - \boldsymbol{\mu})^{\top} L^{-\top} L^{-1} (\boldsymbol{x} - \boldsymbol{\mu}) \\
&= (L^{-1}(\boldsymbol{x} - \boldsymbol{\mu}))^{\top} L^{-1}(\boldsymbol{x} - \boldsymbol{\mu}) \\
&= \| L^{-1}(\boldsymbol{x} - \boldsymbol{\mu}) \|^2
\end{aligned}
$$

Back-substitution, $O(Ns)$

Many statistical operations preserve distribution

Affine transformation

Joint distribution & marginalization:

$$\boldsymbol{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_{11})$$
$$\boldsymbol{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_{22})$$
$$\begin{pmatrix} \boldsymbol{x}_1 \\ \boldsymbol{x}_2 \end{pmatrix} \sim \mathcal{N}\left( \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right)$$

Conditioning

Assume $\boldsymbol{\mu} = \mathbf{0}$ and use precision instead of covariance!

$$Q = \Sigma^{-1} = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix}$$

$$\pi(\boldsymbol{x}_2 \mid \boldsymbol{x}_1) = \frac{\pi(\boldsymbol{x}_1 \mid \boldsymbol{x}_2)\pi(\boldsymbol{x}_2)}{\pi(\boldsymbol{x}_1)} = \frac{\pi(\boldsymbol{x}_1, \boldsymbol{x}_2)}{\pi(\boldsymbol{x}_1)}$$

$$\propto \pi(\boldsymbol{x}_1, \boldsymbol{x}_2)$$

$$\propto e^{-\frac{1}{2}\boldsymbol{x}_2^\top Q_{22}\boldsymbol{x}_2 - (Q_{21}\boldsymbol{x}_1)^\top \boldsymbol{x}_2}$$

$$\boldsymbol{x}_2 \mid \boldsymbol{x}_1 \sim \mathcal{N}\left(-Q_{22}^{-1}Q_{21}\boldsymbol{x}_1, Q_{22}^{-1}\right)$$

If $\boldsymbol{\mu} \neq \mathbf{0}$, shift $\boldsymbol{x}^* = \boldsymbol{x} - \boldsymbol{\mu}$, $\mathrm{E}[\boldsymbol{x}^*] = \mathbf{0}$

$$\boldsymbol{x}_2 \mid \boldsymbol{x}_1 \sim \mathcal{N}\left(\boldsymbol{\mu}_2 - Q_{22}^{-1}Q_{21}(\boldsymbol{x}_1 - \boldsymbol{\mu}_1), Q_{22}^{-1}\right)$$

## Conditioning with Schur Complements

$$\boldsymbol{x}_2 \mid \boldsymbol{x}_1 \sim \mathcal{N}\left(\boldsymbol{\mu}_2 - Q_{22}^{-1}Q_{21}(\boldsymbol{x}_1 - \boldsymbol{\mu}_1), Q_{22}^{-1}\right)$$

$$Q = \Sigma^{-1} = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix}$$

$$= \begin{pmatrix} \Sigma_{11}^{-1} + \left(\Sigma_{11}^{-1}\Sigma_{12}\right)\Sigma_{22|1}^{-1}\left(\Sigma_{21}\Sigma_{11}^{-1}\right) & -\left(\Sigma_{11}^{-1}\Sigma_{12}\right)\Sigma_{22|1}^{-1} \\ -\Sigma_{22|1}^{-1}\left(\Sigma_{21}\Sigma_{11}^{-1}\right) & \Sigma_{22|1}^{-1} \end{pmatrix}$$

$$Q_{22}^{-1} = (\Sigma_{22|1}^{-1})^{-1} = \Sigma_{22|1}$$

$$= \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$$

$$Q_{22}^{-1}Q_{21} = -\Sigma_{22|1}(\Sigma_{22|1}^{-1}\Sigma_{21}\Sigma_{11}^{-1})$$

$$= -\Sigma_{21}\Sigma_{11}^{-1}$$

$$\boldsymbol{x}_2 \mid \boldsymbol{x}_1 \sim \mathcal{N}\left(\boldsymbol{\mu}_2 + \Sigma_{21}\Sigma_{11}^{-1}(\boldsymbol{x}_1 - \boldsymbol{\mu}_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right)$$

# Statistical Interpretation

From conditioning,

$$\boldsymbol{x}_2 \mid \boldsymbol{x}_1 \sim \mathcal{N}\left(\boldsymbol{\mu}_2 + \Sigma_{21}\Sigma_{11}^{-1}(\boldsymbol{x}_1 - \boldsymbol{\mu}_1), \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\right)$$

Schur complement $\iff$ conditional covariance!

s.p.d. because covariance matrices s.p.d.

Quotient rule statistically trivial:
$\pi((x_1 \mid x_2) \mid x_3) = \pi(x_1 \mid x_2, x_3)$

Conditioning in covariance $\iff$ marginalization in precision

# Table of Contents

Probability distribution over *vectors*

Extend to distribution over *functions*?

Idea: for finite set of points, function simply vector

$$X = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N\}$$
$$\boldsymbol{y} = \{f(\boldsymbol{x}_1), f(\boldsymbol{x}_2), \ldots, f(\boldsymbol{x}_N)\}$$

Idea: for points we're not given, marginalization is trivial

How to assign mean and covariance in a sensible way?

## Gaussian Process Definition

Let $\mu(\boldsymbol{x})$ be the *mean function* and
$K(\boldsymbol{x}, \boldsymbol{x}')$ be the *covariance function* or *kernel function*

We say

$$f(\boldsymbol{x}) \sim \mathcal{GP}(\mu(\boldsymbol{x}), K(\boldsymbol{x}, \boldsymbol{x}'))$$

If for all point sets $X$,

$$X = \{\boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_N\}$$
$$\boldsymbol{y} = \{f(\boldsymbol{x}_1), f(\boldsymbol{x}_2), \ldots, f(\boldsymbol{x}_N)\}$$
$$\boldsymbol{y} \sim \mathcal{N}(\boldsymbol{\mu}, \Theta)$$

where

$$\boldsymbol{\mu}_i = \mu(\boldsymbol{x}_i)$$
$$\Theta_{ij} = K(\boldsymbol{x}_i, \boldsymbol{x}_j)$$

## Regression with Gaussian Processes

Simply condition prediction points on training points:

$$\Theta = \begin{pmatrix} \Theta_{\mathsf{Tr,Tr}} & \Theta_{\mathsf{Tr,Pr}} \\ \Theta_{\mathsf{Pr,Tr}} & \Theta_{\mathsf{Pr,Pr}} \end{pmatrix}$$

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

$$\mathrm{Cov}[\boldsymbol{y}_{\mathsf{Pr}} \mid \boldsymbol{y}_{\mathsf{Tr}}] = \Theta_{\mathsf{Pr,Pr}} - \Theta_{\mathsf{Pr,Tr}} \Theta_{\mathsf{Tr,Tr}}^{-1} \Theta_{\mathsf{Tr,Pr}}$$

Nonparametric! No training! Uncertainty quantification!

... $\mathcal{O}(N^3)$ to compute $\Theta_{\mathsf{Tr,Tr}}^{-1}$
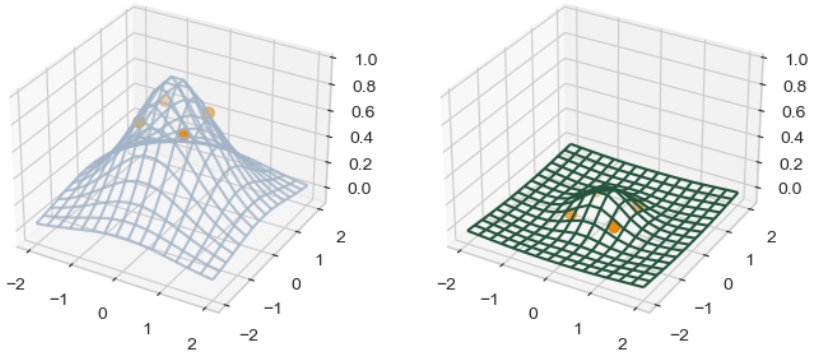
And we're back to the starting problem

# Screening Effect



Figure: Conditional on nearby points, far away points have less covariance

# Table of Contents

## Cholesky Factorization by KL Minimization

Measure approximation error by KL-divergence:

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

Re-write KL-divergence:

$$2\mathbb{D}_{\mathsf{KL}} \left( \mathcal{N}(\mathbf{0}, \Theta_1) \, \middle\| \, \mathcal{N}(\mathbf{0}, \Theta_2) \right) =$$
$$\operatorname{trace}(\Theta_2^{-1}\Theta_1) + \operatorname{logdet}(\Theta_2) - \operatorname{logdet}(\Theta_1) - N$$

where $\Theta_1$ and $\Theta_2$ are both of size $N \times N$

# Cholesky Factorization as GP Regression

*[1]. The non-zero entries of the $i$th column of $L$ are:*

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

Plugging the optimal $L$ back into the KL-divergence, we obtain:

$$\sum_{i=1}^{N} \left[ \log \left( (\boldsymbol{e}_1^\top \Theta_{s_i,s_i}^{-1} \boldsymbol{e}_1)^{-1} \right) \right] - \text{logdet}(\Theta)$$

But marginalization in covariance is conditioning in precision!

$$(\boldsymbol{e}_1^\top \Theta_{s_i,s_i}^{-1} \boldsymbol{e}_1)^{-1} = \Theta_{ii|s_i-\{i\}}$$

This is precisely sparse Gaussian process regression!

# Table of Contents

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

# Thank You!