

Sparse Cholesky Factorization by Greedy Conditional Selection

Stephen Huan

Theory Club

February 28, 2022

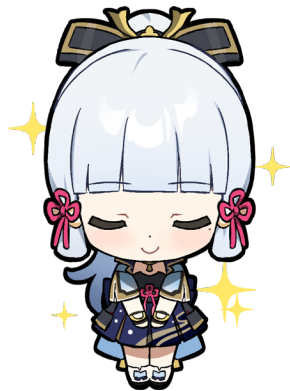


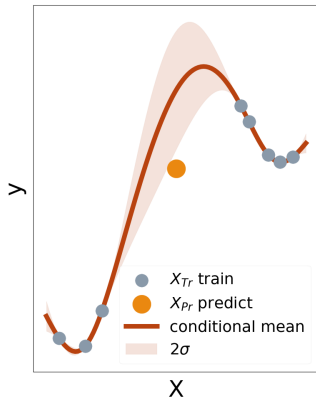
Table of Contents

1. High-level Summary
2. Cholesky Factorization
3. Schur Complement
4. Multivariate Gaussians
5. Gaussian Process Regression
6. Sparse Cholesky Factorization
7. References



The Problem: Gaussian Process Regression

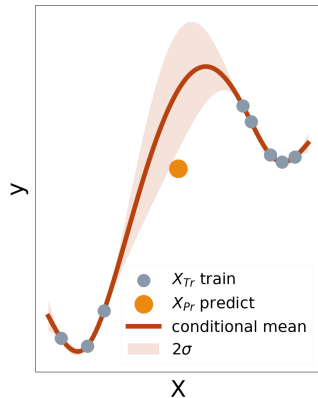
Measurements \mathbf{y}_{Tr} at N points X_{Tr}



The Problem: Gaussian Process Regression

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

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

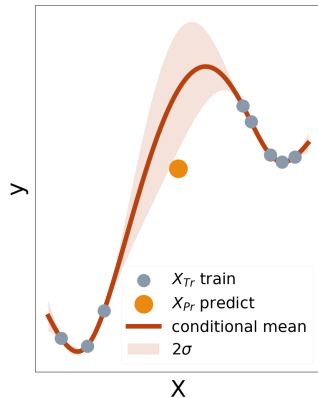


The Problem: Gaussian Process Regression

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}



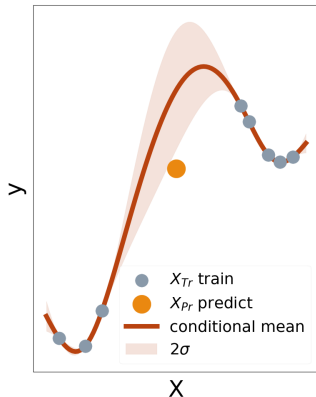
The Problem: Gaussian Process Regression

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}

Computational cost scales as N^3



The Problem: Gaussian Process Regression

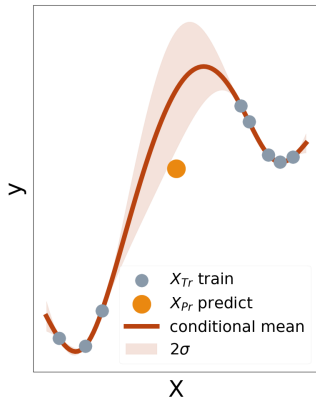
Measurements \mathbf{y}_{Tr} at N points X_{Tr}

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

Model as Gaussian process
→ condition on \mathbf{y}_{Tr}

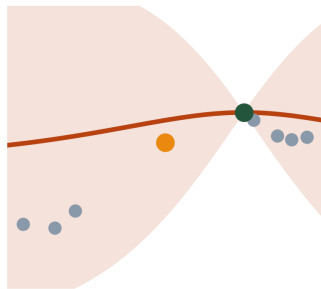
Computational cost scales as N^3

Choose k most informative points!



Conditional k -th Nearest Neighbors

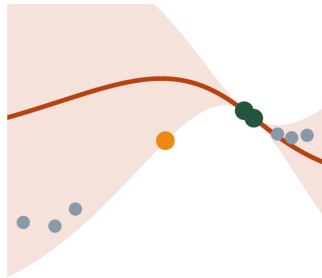
Naive: select k closest points



Conditional k -th Nearest Neighbors

Naive: select k closest points

Chooses redundant information

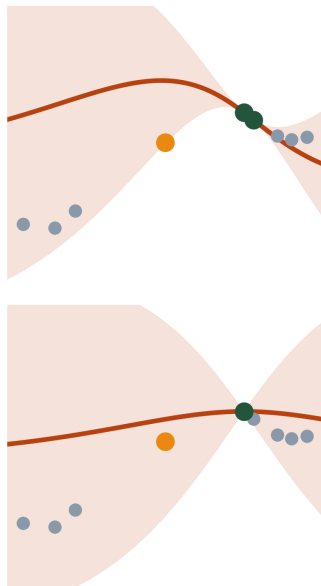


Conditional k -th Nearest Neighbors

Naive: select k closest points

Chooses redundant information

Maximize *mutual information*!

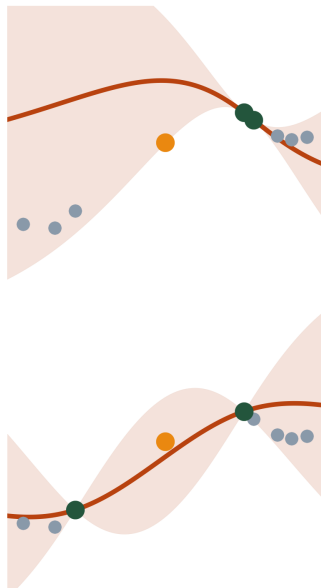


Conditional k -th Nearest Neighbors

Naive: select k closest points

Chooses redundant information

Maximize *mutual information*!



Conditional k -th Nearest Neighbors

Naive: select k closest points

Chooses redundant information

Maximize *mutual information*!

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



Conditional k -th Nearest Neighbors

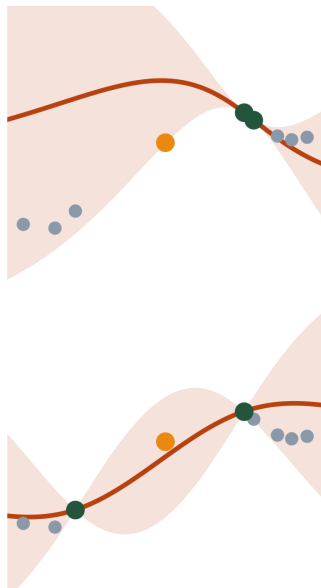
Naive: select k closest points

Chooses redundant information

Maximize *mutual information*!

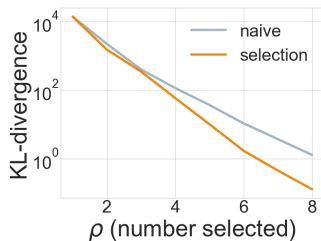
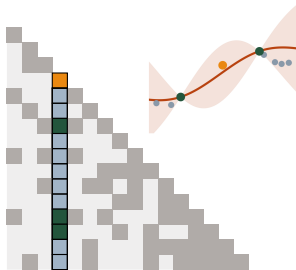
Direct computation: $\mathcal{O}(Nk^4)$

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



Cholesky Factorization by Selection

Apply column-wise
→ sparse approx. of GP



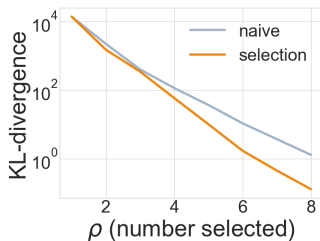
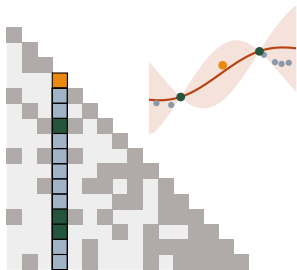
Cholesky Factorization by Selection

Apply column-wise

→ sparse approx. of GP

Maximum mutual information

→ minimum KL divergence



Cholesky Factorization by Selection

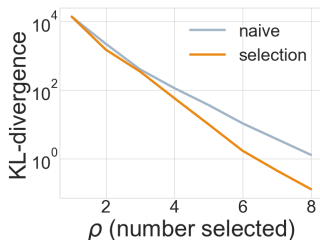
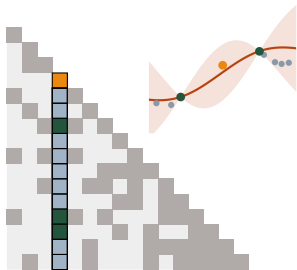
Apply column-wise

→ sparse approx. of GP

Maximum mutual information

→ minimum KL divergence

Improves approx. algorithm of ¹



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

Table of Contents

1. High-level Summary
2. Cholesky Factorization
3. Schur Complement
4. Multivariate Gaussians
5. Gaussian Process Regression
6. Sparse Cholesky Factorization
7. References



LU Decomposition

... and its symmetric counterpart

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

LU Decomposition

... and its symmetric counterpart

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

Not always possible, need PLU in general!



LU Decomposition

... and its symmetric counterpart

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 $M = M^T$ 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 .



LU Decomposition

... and its symmetric counterpart

LU where L is lower triangular and U is upper triangular

... always possible, need PLU in general!

Special case for (square) symmetric matrices:

Theorem

If $M = M^T$ 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 $M = M^T \implies (LDK) = (LDK)^T$.

From matrix multiplication, able to see $K = L^T$. □

Cholesky Factorization

Let M be (symmetric) *positive definite*.

Cholesky Factorization



be (symmetric) *positive definite*.

Then $M = LDL^T$ becomes LL^T :

$$\begin{aligned}M &= LDL^T \\&= LD^{\frac{1}{2}}D^{\frac{1}{2}}L^T \\&= LD^{\frac{1}{2}}(LD^{\frac{1}{2}})^T \\&= L'L'^T\end{aligned}$$

Cholesky Factorization



M be (symmetric) *positive definite*.

Then $M = LDL^T$ becomes LL^T :

$$\begin{aligned}M &= LDL^T \\ &= LD^{\frac{1}{2}}D^{\frac{1}{2}}L^T \\ &= LD^{\frac{1}{2}}(LD^{\frac{1}{2}})^T \\ &= L'L'^T\end{aligned}$$

This is the Cholesky factorization!

Why Do We Care?

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

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

Matrix-vector product $\Theta\mathbf{v} \rightarrow L(L^\top\mathbf{v})$

$N^2 \rightarrow Ns$

Solving linear system $\Theta^{-1}\mathbf{v} \rightarrow L^{-\top}(L^{-1}\mathbf{v})$

$N^3 \rightarrow Ns$

Log determinant $\log\det \Theta \rightarrow 2 \log\det L = 2 \sum_{i=1}^N \log L_{ii}$

$N^3 \rightarrow N$

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

??? $\rightarrow Ns$

Why Do We Care?

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

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

Matrix-vector product $\Theta v \rightarrow L(L^\top v)$

$N^2 \rightarrow Ns$

Solving linear system $\Theta^{-1}v \rightarrow L^{-\top}(L^{-1}v)$

$N^3 \rightarrow Ns$

Log determinant $\log \det \Theta \rightarrow 2 \log \det L = 2 \sum_{i=1}^N \log L_{ii}$

$N^3 \rightarrow N$

Sampling from $x \sim \mathcal{N}(\mu, \Theta) \rightarrow z \sim \mathcal{N}(\mathbf{0}, I), x = Lz + \mu$

??? $\rightarrow Ns$



Computing the Cholesky Factorization

Down-looking

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
```

Computing the Cholesky Factorization

Down-looking

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
```

Computing the Cholesky Factorization

Up-looking

Let L' be blocked according to:

$$\begin{aligned}L' &= \begin{pmatrix} L & \mathbf{0} \\ \mathbf{r}^\top & d \end{pmatrix} \\L'L'^\top &= \begin{pmatrix} L & \mathbf{0} \\ \mathbf{r}^\top & d \end{pmatrix} \begin{pmatrix} L^\top & \mathbf{r} \\ \mathbf{0}^\top & d \end{pmatrix} \\&= \begin{pmatrix} LL^\top & L\mathbf{r} \\ \mathbf{r}^\top L^\top & \mathbf{r}^\top \mathbf{r} + d^2 \end{pmatrix}\end{aligned}$$

So if we have a Cholesky factor for a principle submatrix of Θ , we can extend it inductively by reading off the appropriate data!

$$\begin{aligned}\begin{pmatrix} LL^\top & L\mathbf{r} \\ \mathbf{r}^\top L^\top & \mathbf{r}^\top \mathbf{r} + d^2 \end{pmatrix} &= \begin{pmatrix} \Theta & \mathbf{c} \\ \mathbf{c}^\top & C \end{pmatrix} \\ \mathbf{r} &= L^{-1}\mathbf{c} \\ d &= \sqrt{C - \mathbf{r}^\top \mathbf{r}}\end{aligned}$$

Computing the Cholesky Factorization

Up-looking

Let L' be blocked according to:



$$\begin{aligned}L' &= \begin{pmatrix} L & \mathbf{0} \\ \mathbf{r}^\top & d \end{pmatrix} \\L'L'^\top &= \begin{pmatrix} L & \mathbf{0} \\ \mathbf{r}^\top & d \end{pmatrix} \begin{pmatrix} L^\top & \mathbf{r} \\ \mathbf{0}^\top & d \end{pmatrix} \\ &= \begin{pmatrix} LL^\top & L\mathbf{r} \\ \mathbf{r}^\top L^\top & \mathbf{r}^\top \mathbf{r} + d^2 \end{pmatrix}\end{aligned}$$

So if we have a Cholesky factor for a principle submatrix of Θ , we can extend it inductively by reading off the appropriate data!

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

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

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

Computing the Cholesky Factorization


Up-looking

```
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
```

Computing the Cholesky Factorization

Up-looking

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



Computing the Cholesky Factorization

Right-looking

$$\begin{aligned} L &= (\mathbf{l}_1 \quad \mathbf{l}_2 \quad \cdots \quad \mathbf{l}_N) \\ LL^\top &= (\mathbf{l}_1 \quad \mathbf{l}_2 \quad \cdots \quad \mathbf{l}_N) \begin{pmatrix} \mathbf{l}_1^\top \\ \mathbf{l}_2^\top \\ \vdots \\ \mathbf{l}_N^\top \end{pmatrix} \\ &= \mathbf{l}_1\mathbf{l}_1^\top + \mathbf{l}_2\mathbf{l}_2^\top + \cdots + \mathbf{l}_N\mathbf{l}_N^\top = \Theta \end{aligned}$$

From lower triangularity, nested submatrices!

Computing the Cholesky Factorization

Right-looking



$$L = (\mathbf{l}_1 \quad \mathbf{l}_2 \quad \cdots \quad \mathbf{l}_N)$$

$$LL^T = (\mathbf{l}_1 \quad \mathbf{l}_2 \quad \cdots \quad \mathbf{l}_N) \begin{pmatrix} \mathbf{l}_1^T \\ \mathbf{l}_2^T \\ \vdots \\ \mathbf{l}_N^T \end{pmatrix}$$

$$= \mathbf{l}_1\mathbf{l}_1^T + \mathbf{l}_2\mathbf{l}_2^T + \cdots + \mathbf{l}_N\mathbf{l}_N^T = \Theta$$

From lower triangularity, nested submatrices!

Computing the Cholesky Factorization

Right-looking

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

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

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

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

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

$$\begin{aligned} \mathbf{l}_2 \mathbf{l}_2^\top + \cdots + \mathbf{l}_N \mathbf{l}_N^\top &= \Theta - \left(\frac{\Theta_1}{\sqrt{\Theta_{11}}} \right) \left(\frac{\Theta_1}{\sqrt{\Theta_{11}}} \right)^\top \\ &= \Theta - \frac{\Theta_1 \Theta_1^\top}{\Theta_{11}} \end{aligned}$$

Proceed inductively on rank-one update

Computing the Cholesky Factorization

Right-looking

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



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

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

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

$$l_1 = \frac{\Theta_1}{l_1} = \frac{\Theta_1}{\sqrt{\Theta_{11}}}$$

$$\begin{aligned} \mathbf{l}_2 \mathbf{l}_2^\top + \cdots + \mathbf{l}_N \mathbf{l}_N^\top &= \Theta - \left(\frac{\Theta_1}{\sqrt{\Theta_{11}}} \right) \left(\frac{\Theta_1}{\sqrt{\Theta_{11}}} \right)^\top \\ &= \Theta - \frac{\Theta_1 \Theta_1^\top}{\Theta_{11}} \end{aligned}$$

Proceed inductively on rank-one update

Computing the Cholesky Factorization

Right-looking

```
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
```

Computing the Cholesky Factorization

Left-looking

Recall:

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

Look at \mathbf{l}_i :

$$\begin{aligned}\mathbf{l}_i\mathbf{l}_i^\top &= \left(\Theta - (\mathbf{l}_1\mathbf{l}_1^\top + \mathbf{l}_2\mathbf{l}_2^\top + \cdots + \mathbf{l}_{i-1}\mathbf{l}_{i-1}^\top) \right)_i \\ &= \Theta_i - (l_{1i}\mathbf{l}_1 + l_{2i}\mathbf{l}_2 + \cdots + l_{i-1,i}\mathbf{l}_{i-1}) \\ &= \Theta_i - (\mathbf{l}_1 \quad \mathbf{l}_2 \quad \cdots \quad \mathbf{l}_{i-1}) \begin{pmatrix} l_{1i} \\ l_{2i} \\ \vdots \\ l_{i,i-1} \end{pmatrix} \\ &= \Theta_i - L_{:,i}L_{i,:i}\end{aligned}$$

Don't need to store modified Θ in memory!

Computing the Cholesky Factorization

Left-looking

Recall:

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

Look at \mathbf{l}_i :

$$\begin{aligned} \mathbf{l}_i \mathbf{l}_i^\top &= \left(\Theta - (\mathbf{l}_1 \mathbf{l}_1^\top + \mathbf{l}_2 \mathbf{l}_2^\top + \cdots + \mathbf{l}_{i-1} \mathbf{l}_{i-1}^\top) \right)_i \\ &= \Theta_i - (l_{1i} \mathbf{l}_1 + l_{2i} \mathbf{l}_2 + \cdots + l_{i-1,i} \mathbf{l}_{i-1}) \\ &= \Theta_i - \begin{pmatrix} \mathbf{l}_1 & \mathbf{l}_2 & \cdots & \mathbf{l}_{i-1} \end{pmatrix} \begin{pmatrix} l_{1i} \\ l_{2i} \\ \vdots \\ l_{i,i-1} \end{pmatrix} \\ &= \Theta_i - L_{:,i} L_{i,:i} \end{aligned}$$



Don't need to store modified Θ in memory!

Computing the Cholesky Factorization

Left-looking

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

Computing the Cholesky Factorization

Left-looking

```
1 def left_cholesky(theta: np.ndarray) -> np.ndarray:
2     n = len(theta)
3     L = np.zeros((n, n))
4     for i in range(n):
5         L[:, i] = theta[:, i] - L[:, :i]@L[i, :i]
6         L[i, i] = np.sqrt(L[i, i])
```



Table of Contents

1. High-level Summary
2. Cholesky Factorization
3. **Schur Complement**
4. Multivariate Gaussians
5. Gaussian Process Regression
6. Sparse Cholesky Factorization
7. References



Schur Complement

or recursive Cholesky factorization

Block Θ 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 Θ 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 Θ on Θ_{11}

Schur Complement

or recursive Cholesky factorization

Block Θ 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 Θ 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 Θ on Θ_{11}



Proper Determinant of Block Matrix

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

$$\det(\Theta) = ?$$

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

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

Schur complement gives proper answer:

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

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

Proper Determinant of Block Matrix

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

$$\det(\Theta) = ?$$

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

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

Schur complement gives proper answer:

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

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



Proper Submatrix of Inverse

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

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

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

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,

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

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

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

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,

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

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

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



A Few Important Questions...

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

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

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?

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 Θ blocked as

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

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 Θ 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 Θ complemented on Θ_{11} and then on Θ_{22} the same as

Θ complemented on $\begin{pmatrix} \Theta_{11} & \Theta_{12} \\ \Theta_{21} & \Theta_{22} \end{pmatrix}$?

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 Θ 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 Θ complemented on Θ_{11} and then on Θ_{22} the same as

Θ complemented on $\begin{pmatrix} \Theta_{11} & \Theta_{12} \\ \Theta_{21} & \Theta_{22} \end{pmatrix}$?

Intuitively, it should be, but tedious to prove

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 Θ 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 Θ complemented on Θ_{11} and then on Θ_{22} the same as

Θ 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

1. High-level Summary
2. Cholesky Factorization
3. Schur Complement
4. **Multivariate Gaussians**
5. Gaussian Process Regression
6. Sparse Cholesky Factorization
7. References



The Multivariate Gaussian

Recall: Gaussian (or normal) distribution:

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

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.

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

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

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

Defining Everything

Naturally,

$$\mu_i = \mathbf{E}[x_i]$$

$$\boldsymbol{\mu} = \mathbf{E}[\mathbf{x}]$$

$$\Sigma_{ij} = \text{Cov}[x_i, x_j]$$

$$= \mathbf{E}[(x_i - \mathbf{E}[x_i])(x_j - \mathbf{E}[x_j])]$$

$$= \mathbf{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]$$

Defining Everything

Naturally,

$$\mu_i = \mathbb{E}[x_i]$$

$$\boldsymbol{\mu} = \mathbb{E}[\boldsymbol{x}]$$

$$\Sigma_{ij} = \text{Cov}[x_i, x_j]$$

$$= \mathbb{E}[(x_i - \mathbb{E}[x_i])(x_j - \mathbb{E}[x_j])]$$

$$= \mathbb{E}[(\boldsymbol{x} - \boldsymbol{\mu})(\boldsymbol{x} - \boldsymbol{\mu})^\top]$$

Two natural (and fundamental) questions from here:

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

Defining Everything

Naturally,

$$\mu_i = \mathbb{E}[x_i]$$

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}]$$

$$\Sigma_{ij} = \text{Cov}[x_i, x_j]$$

$$= \mathbb{E}[(x_i - \mathbb{E}[x_i])(x_j - \mathbb{E}[x_j])]$$

$$= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]$$

Two natural (and fundamental) questions from here:

1. What is the probability density function $f(\mathbf{x})$?
2. How can we sample from $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$?

Surprisingly enough, Cholesky factorization answers both!

Defining Everything

Naturally,



$$\mu_i = \mathbb{E}[x_i]$$

$$\boldsymbol{\mu} = \mathbb{E}[\mathbf{x}]$$

$$\begin{aligned}\Sigma_{ij} &= \text{Cov}[x_i, x_j] \\ &= \mathbb{E}[(x_i - \mathbb{E}[x_i])(x_j - \mathbb{E}[x_j])] \\ &= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top]\end{aligned}$$

Two natural (and fundamental) questions from here:

1. What is the probability density function $f(\mathbf{x})$?
2. How can we sample from $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$?

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!

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 μ and Σ , distribution is determined.

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 μ and Σ , 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 Σ 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!

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 μ and Σ , 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 Σ 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!



Completely Independent Variables

Well, if Σ has particular structure, it's actually trivial:

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

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

$$\begin{aligned} f(\mathbf{z}) &= \prod_{i=1}^N f(z_i) \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_i^2} \\ &= \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}(z_1^2 + z_2^2 + \dots + z_N^2)} \\ &= \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}\mathbf{z}^\top \mathbf{z}} \end{aligned}$$

Completely Independent Variables

Well, if Σ has particular structure, it's actually trivial:

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

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

$$\begin{aligned} f(\mathbf{z}) &= \prod_{i=1}^N f(z_i) \\ &= \prod_{i=1}^N \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z_i^2} \\ &= \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}(z_1^2 + z_2^2 + \dots + z_N^2)} \\ &= \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}\mathbf{z}^\top \mathbf{z}} \end{aligned}$$



Moment Matching

How can we generalize to arbitrary Σ ?

Moment match!

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

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

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[L\mathbf{z} + \boldsymbol{\mu}] = L\mathbb{E}[\mathbf{z}] + \boldsymbol{\mu} = \boldsymbol{\mu}$$

$$\text{Cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top]$$

$$= \mathbb{E}[L\mathbf{z}(L\mathbf{z})^\top]$$

$$= \mathbb{E}[L\mathbf{z}\mathbf{z}^\top L^\top]$$

$$= L\mathbb{E}[\mathbf{z}\mathbf{z}^\top]L^\top$$

$$= LL^\top$$

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

Moment Matching

How can we generalize to arbitrary Σ ?

Moment match!

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

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

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[L\mathbf{z} + \boldsymbol{\mu}] = L\mathbb{E}[\mathbf{z}] + \boldsymbol{\mu} = \boldsymbol{\mu}$$

$$\text{Cov}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top]$$

$$= \mathbb{E}[L\mathbf{z}(L\mathbf{z})^\top]$$

$$= \mathbb{E}[L\mathbf{z}\mathbf{z}^\top L^\top]$$

$$= L\mathbb{E}[\mathbf{z}\mathbf{z}^\top]L^\top$$

$$= LL^\top$$

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

Sampling with Cholesky Factorization

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

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

Sampling with Cholesky Factorization

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

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

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

$$\begin{aligned}\Sigma &= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \\ \mathbf{y}^\top \Sigma \mathbf{y} &= \mathbf{y}^\top \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \mathbf{y} \\ &= \mathbb{E}[\mathbf{y}^\top (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{y}] \\ &= \mathbb{E}[\left((\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{y}\right)^\top (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{y}] \\ &= \mathbb{E}[\|(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{y}\|^2] \geq 0\end{aligned}$$

Sampling with Cholesky Factorization

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

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

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

$$\begin{aligned}\Sigma &= \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \\ \mathbf{y}^\top \Sigma \mathbf{y} &= \mathbf{y}^\top \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] \mathbf{y} \\ &= \mathbb{E}[\mathbf{y}^\top (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{y}] \\ &= \mathbb{E}[\|(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{y}\|^2] \geq 0\end{aligned}$$



Probability Density Function from Sampling

What's the probability density function $f(\boldsymbol{x})$?

Probability Density Function from Sampling

What's the probability density function $f(\boldsymbol{x})$?

Idea: view \boldsymbol{x} resulting from an invertible transformation from \boldsymbol{z} .

Probability Density Function from Sampling

What's the probability density function $f(\boldsymbol{x})$?

Idea: view \boldsymbol{x} resulting from an invertible transformation from \boldsymbol{z} .

We know $f(\boldsymbol{z})$, so $f(\boldsymbol{x})$ should be similar!

Probability Density Function from Sampling

What's the probability density function $f(\mathbf{x})$?

Idea: view \mathbf{x} resulting from an invertible transformation from \mathbf{z} .

We know $f(\mathbf{z})$, so $f(\mathbf{x})$ should be similar!

In scalars:

$$z \sim \mathcal{N}(0, 1)$$

$$x = \sigma z + \mu$$

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

$$z = \frac{x - \mu}{\sigma}$$

PDF from Sampling — Scalar Edition

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

$$1 = \int_{-\infty}^{\infty} f(z) dz = \int_{-\infty}^{\infty} f(z) \frac{dz}{dx} dx$$

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

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

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

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

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

PDF from Sampling — Vector Edition

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

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

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

$$\begin{aligned} 1 &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\mathbf{z}) \, d\mathbf{z} \\ &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\mathbf{z}) \frac{d\mathbf{z}}{d\mathbf{x}} \, d\mathbf{x} \end{aligned} \quad \text{(informal)}$$

$$= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(\mathbf{z}) |\det(J_{\mathbf{z}})| \, d\mathbf{x} \quad \text{(formal)}$$

$$= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} \underbrace{f(L^{-1}(\mathbf{x} - \boldsymbol{\mu})) \det(L^{-1})}_{\text{PDF of } \mathbf{x}} \, d\mathbf{x}$$

PDF from Sampling — Vector Edition

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

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

$$\begin{aligned} &= \frac{1}{\det(L)} f(L^{-1}(\mathbf{x} - \boldsymbol{\mu})) \\ &= \frac{1}{\det(L)} \frac{1}{\sqrt{(2\pi)^N}} e^{-\frac{1}{2}(L^{-1}(\mathbf{x} - \boldsymbol{\mu}))^\top (L^{-1}(\mathbf{x} - \boldsymbol{\mu}))} \end{aligned}$$

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

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

PDF from Sampling — Vector Edition

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

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

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

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

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

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

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



Summary

Compare PDFs of multivariate normal and scalar normal:

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

Compare to scalar:

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

Summary

Compare PDFs of multivariate normal and scalar normal:

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

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!



Cholesky Factorization for Gaussians

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

Density computation:

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

Back-substitution, $\mathcal{O}(N_s)$

Cholesky Factorization for Gaussians

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

Density computation:

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

Back-substitution, $\mathcal{O}(N_s)$



Closure of Multivariate Gaussians

Many statistical operations preserve distribution

Closure of Multivariate Gaussians

Many statistical operations preserve distribution

Affine transformation

Closure of Multivariate Gaussians

Many statistical operations preserve distribution

Affine transformation

Joint distribution & marginalization:

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_{11})$$

$$\mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_{22})$$

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

Closure of Multivariate Gaussians

Many statistical operations preserve distribution

Affine transformation

Joint distribution & marginalization:

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_{11})$$

$$\mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_{22})$$

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

Conditioning

Closure of Multivariate Gaussians

Many statistical operations preserve distribution

Affine transformation

Joint distribution & marginalization:

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_{11})$$

$$\mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_{22})$$

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

Conditioning



Conditioning

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

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

$$\pi(\mathbf{x}_2 | \mathbf{x}_1) = \frac{\pi(\mathbf{x}_1 | \mathbf{x}_2)\pi(\mathbf{x}_2)}{\pi(\mathbf{x}_1)} = \frac{\pi(\mathbf{x}_1, \mathbf{x}_2)}{\pi(\mathbf{x}_1)}$$

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

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

$$\mathbf{x}_2 | \mathbf{x}_1 \sim \mathcal{N}(-Q_{22}^{-1}Q_{21}\mathbf{x}_1, Q_{22}^{-1})$$

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

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

Conditioning with Schur Complements

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

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

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

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

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

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

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

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

Conditioning with Schur Complements

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

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

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

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

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

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

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

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



Statistical Interpretation

From conditioning,

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

Statistical Interpretation

From conditioning,

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

Schur complement \iff conditional covariance!

Statistical Interpretation

From conditioning,

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

Schur complement \iff conditional covariance!

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

Statistical Interpretation

From conditioning,

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

Schur complement \iff conditional covariance!

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

Quotient rule statistically trivial:

$$\pi((x_1 \mid x_2) \mid x_3) = \pi(x_1 \mid x_2, x_3)$$

Statistical Interpretation

From conditioning,

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

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

Statistical Interpretation

From conditioning,

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

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

1. High-level Summary
2. Cholesky Factorization
3. Schur Complement
4. Multivariate Gaussians
5. **Gaussian Process Regression**
6. Sparse Cholesky Factorization
7. References



Gaussian Processes

Probability distribution over *vectors*

Gaussian Processes

Probability distribution over *vectors*

Extend to distribution over *functions*?

Gaussian Processes

Probability distribution over *vectors*

Extend to distribution over *functions*?

Idea: for finite set of points, function simply vector

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

$$\mathbf{y} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$$

Gaussian Processes

Probability distribution over *vectors*

Extend to distribution over *functions*?

Idea: for finite set of points, function simply vector

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

$$\mathbf{y} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$$

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

Gaussian Processes

Probability distribution over *vectors*

Extend to distribution over *functions*?

Idea: for finite set of points, function simply vector

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

$$\mathbf{y} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$$

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

How to assign mean and covariance in a sensible way?

Gaussian Processes

Probability distribution over *vectors*

Extend to distribution over *functions*?

Idea: for finite set of points, function simply vector

$$\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

$$\mathbf{y} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$$

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

How to assign mean and covariance in a sensible way?



Gaussian Process Definition

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

We say

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

If for all point sets X ,

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

$$\mathbf{y} = \{f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)\}$$

$$\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \Theta)$$

where

$$\boldsymbol{\mu}_i = \mu(\mathbf{x}_i)$$

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

Regression with Gaussian Processes

Simply condition prediction points on training points:

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

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

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

Regression with Gaussian Processes

Simply condition prediction points on training points:

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

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

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

Nonparametric! No training! Uncertainty quantification!



Regression with Gaussian Processes

Simply condition prediction points on training points:

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

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

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

Nonparametric! No training! Uncertainty quantification!

... $\mathcal{O}(N^3)$ to compute $\Theta_{\text{Tr},\text{Tr}}^{-1}$



Regression with Gaussian Processes

Simply condition prediction points on training points:

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

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

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

Nonparametric! No training! Uncertainty quantification!

... $\mathcal{O}(N^3)$ to compute $\Theta_{\text{Tr},\text{Tr}}^{-1}$

And we're back to the starting problem



Screening Effect

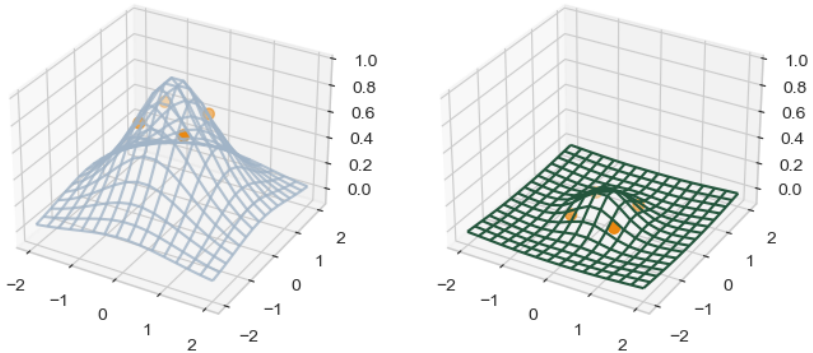


Figure: Conditional on nearby points, far away points have less covariance

Table of Contents

1. High-level Summary
2. Cholesky Factorization
3. Schur Complement
4. Multivariate Gaussians
5. Gaussian Process Regression
6. Sparse Cholesky Factorization
7. References



Cholesky Factorization by KL Minimization

Measure approximation error by KL divergence:

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

Cholesky Factorization by KL Minimization

Measure approximation error by KL divergence:

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

Re-write KL divergence:

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

where Θ_1 and Θ_2 are both of size $N \times N$

Cholesky Factorization by KL Minimization

Measure approximation error by KL divergence:

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

Re-write KL divergence:

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

where Θ_1 and Θ_2 are both of size $N \times N$



Cholesky Factorization as GP Regression

Theorem

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

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

Cholesky Factorization as GP Regression

Theorem

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

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

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

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

Cholesky Factorization as GP Regression

Theorem

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

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

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

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

But marginalization in covariance is conditioning in precision!

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

Cholesky Factorization as GP Regression

Theorem

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

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

Plugging the optimal L back into the KL divergence

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

But marginalization in covariance is conditioning in precision!

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

This is precisely sparse Gaussian process regression!



Table of Contents

1. High-level Summary
2. Cholesky Factorization
3. Schur Complement
4. Multivariate Gaussians
5. Gaussian Process Regression
6. Sparse Cholesky Factorization
7. References



References

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

Thank You!

Thank You!

