

Gaussian processes, part 1: multivariate Gaussians

Stephen Huan

<https://cgdct.moe>

Theory club 2023-09-22

Overview

Cholesky factorization

Characterizations

References

Symmetric LU decomposition

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

Not always possible, need PLU in general!

Special case for 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 $M = M^T \Rightarrow (LDK) = (LDK)^T$. From matrix multiplication, able to see $K = L^T$. □

Cholesky factorization

Let M be (symmetric) *positive definite*, that is,

$$\langle \mathbf{x}, M\mathbf{x} \rangle = \langle M\mathbf{x}, \mathbf{x} \rangle \geq 0 \text{ for all } \mathbf{x} \in \mathbb{R}^n.$$

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

$$M = LDL^\top = L(\sqrt{D}\sqrt{D})L^\top = (L\sqrt{D})(L\sqrt{D})^\top = L'L'^\top$$

for $L' := L\sqrt{D}$ and $(\sqrt{D})_{i,j} := \sqrt{D_{i,j}}$.

This is the Cholesky factorization!

Why should 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 $\mathcal{O}(Ns)$ floating-point operations

- 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 $\log\det \Theta \rightarrow 2\log\det L = 2\sum_{i=1}^N \log L_{i,i}$
 - $N^3 \rightarrow N$
- Sampling from $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Theta) \rightarrow \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \text{Id}), \mathbf{x} = L\mathbf{z} + \boldsymbol{\mu}$
 - ??? $\rightarrow Ns$

Here, we care about its statistical interpretation

Schur complement

Block $\Theta = \begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \Theta_{2,1} & \Theta_{2,2} \end{pmatrix}$ then perform a step of Gaussian elimination

$$\begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \mathbf{0} & \Theta_{2,2} - \Theta_{2,1}\Theta_{1,1}^{-1}\Theta_{1,2} \end{pmatrix}.$$

Denote the term in blue, the *Schur complement* of Θ on $\Theta_{1,1}$, as

$$\Theta = \begin{pmatrix} \text{Id} & 0 \\ \Theta_{2,1}\Theta_{1,1}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \Theta_{1,1} & 0 \\ 0 & \Theta_{2,2|1} \end{pmatrix} \begin{pmatrix} \text{Id} & \Theta_{1,1}^{-1}\Theta_{1,2} \\ 0 & \text{Id} \end{pmatrix},$$

so we see the Cholesky factorization of Θ is

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

Recurring on both diagonal blocks finishes the construction.

Determinant of a block matrix

$$\Theta = \begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \Theta_{2,1} & \Theta_{2,2} \end{pmatrix}$$

$$\begin{aligned} \det(\Theta) &= \det(\Theta_{1,1}) \det(\Theta_{2,2}) - \det(\Theta_{2,1}) \det(\Theta_{1,2})? \\ &= \det(\Theta_{1,1}\Theta_{2,2} - \Theta_{2,1}\Theta_{1,2})? \end{aligned}$$

wrong!

wrong!

Schur complement gives proper answer

$$\Theta = \begin{pmatrix} \text{Id} & 0 \\ \Theta_{2,1}\Theta_{1,1}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \Theta_{1,1} & 0 \\ 0 & \Theta_{2,2|1} \end{pmatrix} \begin{pmatrix} \text{Id} & \Theta_{1,1}^{-1}\Theta_{1,2} \\ 0 & \text{Id} \end{pmatrix}$$

$$\det(\Theta) = \det(\Theta_{1,1}) \det(\Theta_{2,2|1})$$

Submatrix of inverse

$$\Theta = \begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \Theta_{2,1} & \Theta_{2,2} \end{pmatrix}$$

$$(\Theta^{-1})_{2,2} = ?$$

$$= (\Theta_{2,2})^{-1}?$$

wrong!

Schur complement to the rescue again!

Proper Submatrix of Inverse

Recall the factorization of Θ as

$$\begin{pmatrix} \text{Id} & 0 \\ \Theta_{2,1}\Theta_{1,1}^{-1} & \text{Id} \end{pmatrix} \begin{pmatrix} \Theta_{1,1} & 0 \\ 0 & \Theta_{2,2|1} \end{pmatrix} \begin{pmatrix} \text{Id} & \Theta_{1,1}^{-1}\Theta_{1,2} \\ 0 & \text{Id} \end{pmatrix}.$$

Inverting both sides of the equation, Θ^{-1} is

$$\begin{aligned} & \begin{pmatrix} \text{Id} & -\Theta_{1,1}^{-1}\Theta_{1,2} \\ 0 & \text{Id} \end{pmatrix} \begin{pmatrix} \Theta_{1,1}^{-1} & 0 \\ 0 & \Theta_{2,2|1}^{-1} \end{pmatrix} \begin{pmatrix} \text{Id} & 0 \\ -\Theta_{2,1}\Theta_{1,1}^{-1} & \text{Id} \end{pmatrix} \\ &= \begin{pmatrix} \Theta_{1,1}^{-1} + \left(\Theta_{1,1}^{-1}\Theta_{1,2}\right)\Theta_{2,2|1}^{-1}\left(\Theta_{2,1}\Theta_{1,1}^{-1}\right) & -\left(\Theta_{1,1}^{-1}\Theta_{1,2}\right)\Theta_{2,2|1}^{-1} \\ -\Theta_{2,2|1}^{-1}\left(\Theta_{2,1}\Theta_{1,1}^{-1}\right) & \Theta_{2,2|1}^{-1} \end{pmatrix} \end{aligned}$$

So $(\Theta^{-1})_{2,2}$ can be read off as $\Theta_{2,2|1}^{-1} := (\Theta_{2,2|1})^{-1}$.

Leaving linear algebra...

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

- Have been assuming so throughout

Is Schur complementing transitive?

- i.e. suppose we have Θ blocked as

$$\Theta = \begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} & \Theta_{1,3} \\ \Theta_{2,1} & \Theta_{2,2} & \Theta_{2,3} \\ \Theta_{3,1} & \Theta_{3,2} & \Theta_{3,3} \end{pmatrix}$$

- Is Θ complemented on $\Theta_{1,1}$ and then on $\Theta_{2,2}$ the same as Θ complemented on $\begin{pmatrix} \Theta_{1,1} & \Theta_{1,2} \\ \Theta_{2,1} & \Theta_{2,2} \end{pmatrix}$? Intuitively yes, but tedious

Intuitive and economic proofs from statistical perspective

Multivariate Gaussian

Recall: Gaussian (or normal) distribution:

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

(alternative) 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,

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

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

$$\Sigma_{i,j} := \text{Cov}[x_i, x_j]$$

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

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

Two natural questions from here.

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

A technical path

Recall *moment/cumulant generating functions*

$$M_{\mathbf{X}}(\boldsymbol{\xi}) := \mathbb{E}[\exp(\boldsymbol{\xi}^\top \mathbf{x})]$$

$$K_{\mathbf{X}}(\boldsymbol{\xi}) := \log M_{\mathbf{X}}(\boldsymbol{\xi})$$

One characterization: multivariate Gaussian has

$$M_{\mathbf{x}}(\boldsymbol{\xi}) = \exp(\boldsymbol{\xi}^\top \boldsymbol{\mu} + \boldsymbol{\xi}^\top \boldsymbol{\Theta}^{-1} \boldsymbol{\xi})$$

$$K_{\mathbf{x}}(\boldsymbol{\xi}) = \boldsymbol{\xi}^\top \boldsymbol{\mu} + \boldsymbol{\xi}^\top \boldsymbol{\Theta}^{-1} \boldsymbol{\xi}$$

Invert the integral transformation to recover PDF

Instead: exploit moment matching and independence

Identity covariance

If Σ has the simplest structure possible, it's trivial:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \text{Id}_N)$$

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

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

Moment matching

Generalize to arbitrary Σ by moment matching

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \text{Id}_N)$$

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

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

$$\begin{aligned}\text{Cov}[\mathbf{x}] &= \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^\top] \\ &= \mathbb{E}[L\mathbf{z}(L\mathbf{z})^\top] \\ &= \mathbb{E}[L\mathbf{z}\mathbf{z}^\top L^\top] \\ &= L\mathbb{E}[\mathbf{z}\mathbf{z}^\top]L^\top \\ &= LL^\top\end{aligned}$$

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 $x = Lz + \boldsymbol{\mu}$.

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.

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

Pushforward

Integrate against test function ψ

$$\mathbb{E}_\eta[\psi \circ T] := \int [\psi \circ T](\mathbf{z}) \eta(\mathbf{z}) d\mathbf{z} = \int \psi(\mathbf{x}) T_\# \eta(\mathbf{x}) d\mathbf{x}$$

with change of variables $\mathbf{x} = T(\mathbf{z})$ and pushforward

$$T_\# \eta := (\eta \circ T^{-1}) |\det \nabla T^{-1}|$$

Therefore we have

$$\mathbb{E}_\eta[\psi \circ T] = \int \psi(\mathbf{x}) T_\# \eta(\mathbf{x}) d\mathbf{x} = \int \psi(\mathbf{x}) \pi(\mathbf{x}) d\mathbf{x} = \mathbb{E}_\pi[\psi]$$

if and only if $T_\# \eta = \pi$.

PDF from transport

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

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

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

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, $\mathcal{O}(N_s)$

(Log)-likelihood 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

Affine transformation

Joint distribution & marginalization:

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_{1,1})$$

$$\mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_{2,2})$$

$$\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}, \begin{pmatrix} \Sigma_{1,1} & \Sigma_{1,2} \\ \Sigma_{2,1} & \Sigma_{2,2} \end{pmatrix} \right)$$

Conditioning

Conditioning

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

$$Q = \Sigma^{-1} = \begin{pmatrix} Q_{1,1} & Q_{1,2} \\ Q_{2,1} & Q_{2,2} \end{pmatrix}$$

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

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

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

$$\mathbf{x}_2 | \mathbf{x}_1 \sim \mathcal{N}\left(-Q_{2,2}^{-1}Q_{2,1}\mathbf{x}_1, Q_{2,2}^{-1}\right)$$

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

$$\mathbf{x}_2 | \mathbf{x}_1 \sim \mathcal{N}\left(\boldsymbol{\mu}_2 - Q_{2,2}^{-1}Q_{2,1}(\mathbf{x}_1 - \boldsymbol{\mu}_1), Q_{2,2}^{-1}\right)$$

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,

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

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

References

-  Bishop, Christopher M. (2006). *Pattern Recognition and Machine Learning*. Information Science and Statistics. New York: Springer. ISBN: 978-0-387-31073-2.
-  Kallenberg, Olav (2021). *Foundations of Modern Probability*. Vol. 99. Probability Theory and Stochastic Modelling. Cham: Springer International Publishing. ISBN: 978-3-030-61870-4 978-3-030-61871-1. DOI: 10.1007/978-3-030-61871-1. (Visited on 09/15/2023).
-  McCullagh, P (2018). *Tensor Methods in Statistics: Monographs on Statistics and Applied Probability*. Boca Raton, FL: CRC Press. ISBN: 978-1-351-08556-4 978-1-351-07711-8 978-1-351-09401-6 978-1-315-89801-8. (Visited on 04/17/2022).
-  Rue, Havard and Leonhard Held (Feb. 2005). *Gaussian Markov Random Fields: Theory and Applications*. New York: Chapman and Hall/CRC. ISBN: 978-0-429-20882-9. DOI: 10.1201/9780203492024.