

Multivariate Gaussian Cheat Sheet

GEOFF PLEISS

Definition (Multivariate Gaussian). Let \mathbf{y} be a d -dimensional vector-valued random variable. \mathbf{y} is multivariate Gaussian if and only if all linear combination of its entries are univariate Gaussian; i.e. for all $\mathbf{c} \in \mathbb{R}^d$, we have that $p(\mathbf{c}^\top \mathbf{y} = a) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{1}{2\sigma^2} (a - \mu)^2\right)$ for some $\mu, \sigma \in \mathbb{R}$.

1) Multivariate Gaussian Density

Let \mathbf{y} be a multivariate Gaussian random variable with mean $\mathbb{E}[\mathbf{y}] = \boldsymbol{\mu}$ and covariance $\mathbb{E}[(\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})^\top] = \mathbf{K}$. The probability density of \mathbf{y} is given by:

$$p(\mathbf{y} = \mathbf{a}) = \mathcal{N}(\mathbf{a}; \boldsymbol{\mu}, \mathbf{K}) := |2\pi\mathbf{K}|^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{a} - \boldsymbol{\mu})^\top \mathbf{K}^{-1}(\mathbf{a} - \boldsymbol{\mu})\right). \quad (1)$$

We will use the notation $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$, $p(\mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$, and $p(\mathbf{y} = \mathbf{a}) = \mathcal{N}(\mathbf{a}; \boldsymbol{\mu}, \mathbf{K})$ interchangeably. All should be read as “ \mathbf{y} is a multivariate Gaussian random variable with mean $\boldsymbol{\mu}$ and covariance \mathbf{K} .”

2) Important Multivariate Gaussian Closures

Many important operations on multivariate Gaussians preserve Gaussianity.

1. **Closure under affine transformation.** Let $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$. Given matrix \mathbf{A} and vector \mathbf{b} , we have:

$$(\mathbf{A}\mathbf{y} + \mathbf{b}) \sim \mathcal{N}\left(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\mathbf{K}\mathbf{A}^\top\right). \quad (2)$$

(For Gaussian processes with appropriate regularity conditions, this property can be generalized to *closure under arbitrary linear operations*.)

2. **Closure under linear combination.** Now let $\mathbf{y}' \sim \mathcal{N}(\boldsymbol{\mu}', \mathbf{K}'')$. If $\mathbf{y} \perp \mathbf{y}'$ (read: \mathbf{y} and \mathbf{y}' are independent random variables), then

$$(\mathbf{y} + \mathbf{y}') \sim \mathcal{N}(\boldsymbol{\mu} + \boldsymbol{\mu}', \mathbf{K} + \mathbf{K}''). \quad (3)$$

(This property is analogously extended to Gaussian processes.)

3. **Closure under marginalization.** Given the following block multivariate Gaussian random variable, we have:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}' \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}' \\ \mathbf{K}'^\top & \mathbf{K}'' \end{bmatrix}\right) \implies p(\mathbf{y}) = \int p\left(\begin{bmatrix} \mathbf{y} \\ \mathbf{y}' \end{bmatrix}\right) d\mathbf{y}' = \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}). \quad (4)$$

(This property naturally applies to Gaussian processes by definition: any finite subset of Gaussian process evaluations are multivariate Gaussian distributed.)

4. **Closure under conditioning.** With the same block random variable, we have:

$$(\mathbf{y}' \mid \mathbf{y} = \mathbf{a}) = \mathcal{N}\left(\mathbf{K}'^\top \mathbf{K}^{-1} \mathbf{a}, \mathbf{K}'' - \mathbf{K}'^\top \mathbf{K}^{-1} \mathbf{K}'\right). \quad (5)$$

We will often drop the $= \mathbf{a}$ and simply write $\mathbf{y}' \mid \mathbf{y}$.

3) Other Useful Properties

1. **Uncorrelation implies independence.** Two random variables \mathbf{y} and \mathbf{y}' are uncorrelated if $\mathbb{E}[(\mathbf{y} - \mathbb{E}[\mathbf{y}])(\mathbf{y}' - \mathbb{E}[\mathbf{y}'])^\top] = \mathbf{0}$. For arbitrary random variables, uncorrelation does not imply independence. However, for multivariate Gaussians:

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{y}' \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{K} & \mathbf{K}' \\ \mathbf{K}'^\top & \mathbf{K}'' \end{bmatrix} \right) : \quad \mathbf{K}' = \mathbf{0} \quad \Leftrightarrow \quad \mathbf{y} \perp \mathbf{y}' \quad (6)$$

2. **Conditional variances from the Cholesky factorization.** The *Cholesky factor* of the (positive-definite) covariance matrix \mathbf{K} is the unique lower triangular matrix \mathbf{L} such that 1) $\mathbf{L}\mathbf{L}^\top = \mathbf{K}$ and 2) $L_{ii} > 0$ for all i . Consider the N -dimensional multivariate Gaussian random variable $[y_1 \ \cdots \ y_N] \mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$. We can write the density of \mathbf{y} in *autoregressive form*:

$$p(\mathbf{y}) = p(\underbrace{y_1}_{:=z_1}) \times p(\underbrace{y_2 | y_1}_{:=z_2}) \times p(\underbrace{y_3 | y_1, y_2}_{:=z_3}) \times \cdots \times p(\underbrace{y_N | y_1, \dots, y_{N-1}}_{:=z_N}). \quad (7)$$

We can interpret Eq. (7) as sequential Bayesian inference. We first consider the random variable y_1 . After observing y_1 , we then consider y_2 (conditioned on our observation). After observing y_2 , we then consider y_3 (conditioning on our observations). And so on.

By Eq. (4), z_1, \dots, z_N are univariate Gaussians. The Cholesky factor gives us a convenient way to automatically compute the variance of these conditional random variables:

$$\mathbb{V}[z_i] = \mathbb{V}[y_i | y_1, \dots, y_{i-1}] = L_{ii}^2. \quad (8)$$

3. **Sampling.** To draw a sample from $\mathcal{N}(\boldsymbol{\mu}, \mathbf{K})$, we often use the following computational routine:

- Use e.g. `np.random.randn` to draw a sample $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Compute $\mathbf{L}\boldsymbol{\epsilon} + \boldsymbol{\mu}$, where \mathbf{L} is the Cholesky factor of \mathbf{K} .

We do not specifically need to use the matrix \mathbf{L} in the second step; we could instead use any matrix \mathbf{A} such that $\mathbf{A}\mathbf{A}^\top = \mathbf{K}$. (Note that all matrices \mathbf{A} that satisfy this condition are equivalent to \mathbf{L} up to an orthogonal rotation. In other words, there exists some orthogonal matrix \mathbf{Q} so that $\mathbf{A} = \mathbf{Q}\mathbf{L}$.)

4. **Sequences.** Consider a sequence of multivariate Gaussian variables $\{\mathbf{y}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \mathbf{K}_i)\}$, where $\{\boldsymbol{\mu}_i\}$ and $\{\mathbf{K}_i\}$ represent a sequence of means and covariances, respectively. (Weak) convergence of \mathbf{y}_i is uniquely determined by convergence of the mean/covariance sequences:

$$\{\boldsymbol{\mu}_i\} \rightarrow \boldsymbol{\mu}, \{\mathbf{K}_i\} \rightarrow \mathbf{K} \quad \Longrightarrow \quad \{p(\mathbf{y}_i)\} \xrightarrow{\text{dist.}} \mathcal{N}(\boldsymbol{\mu}, \mathbf{K}) \quad (9)$$