

Multivariate Gaussian

Christian Osendorfer

Technische Universität München

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{|2\pi\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\mathbf{x} \in \mathbb{R}^d, \quad \boldsymbol{\mu} \in \mathbb{R}^d, \quad \boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$$

$\boldsymbol{\Sigma}$ is assumed to be symmetric (which means, that $\boldsymbol{\Sigma}^{-1}$ is also symmetric) and positive definite.

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

$$\text{Cov}[\mathbf{X}] = E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{X} - E[\mathbf{X}])^T] = \boldsymbol{\Sigma}$$

Diagonal elements of $\boldsymbol{\Sigma}$?

Expectation is linear. E.g. $E[\mathbf{a}^T \mathbf{X}] = \mathbf{a}^T E[\mathbf{X}]$.

In $2d$, a *bivariate Gaussian* is depicted as an ellipse. Why?

Σ is *real* and *symmetric* and therefore

$$\Sigma = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

\mathbf{U} is orthonormal ($\mathbf{U}\mathbf{U}^T = \mathbf{I}$).

$\mathbf{\Lambda}$ is a diagonal matrix of eigenvalues.

And thus:

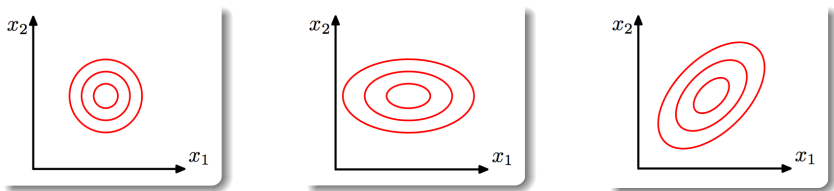
$$\Sigma^{-1} = \mathbf{U}^{-T}\mathbf{\Lambda}^{-1}\mathbf{U}^{-1} = \mathbf{U}\mathbf{\Lambda}^{-1}\mathbf{U}^T = \sum_{i=1}^d \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$$

Full rank is assumed.

The *Mahalanobis distance* can be rewritten:

$$(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T \left(\sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right) (\mathbf{x} - \boldsymbol{\mu}) =$$
$$\sum_{i=1}^D \frac{1}{\lambda_i} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{u}_i \underbrace{\mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})}_{y_i \in \mathbb{R}} = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$$

Mahalanobis distance: The euclidean distance of \mathbf{x} from $\boldsymbol{\mu}$ in a rotated and scaled coordinate system.



Linear transformation of a Gaussian

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

$$\mathbf{Y} = \mathbf{A}\mathbf{X}$$

\mathbf{Y} is a Gaussian with

$$E[\mathbf{Y}] = \mathbf{A}\boldsymbol{\mu}$$

$$\text{Cov}[\mathbf{Y}] = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T$$

MLE for μ

$$\mathbf{x}_i \sim \mathcal{N}(\mu, \Sigma), \mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$$

Likelihood

$$\prod_{i=1}^n \mathcal{N}(\mathbf{x}_i | \mu, \Sigma)$$

and thus the *negative* loglikelihood is

$$\underbrace{\frac{nd}{2} \log 2\pi}_{\text{const.}} + \underbrace{\frac{n}{2} \log |\Sigma|}_{\text{depends on } \Sigma} + \underbrace{\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu)}_{\text{depends on } \mu, \Sigma} := \ell(\mu, \Sigma)$$

We need the derivative w.r.t. μ , and therefore first consider:

$$(\mathbf{x}_i - \mu)^T \Sigma^{-1} (\mathbf{x}_i - \mu) = \mathbf{x}_i^T \Sigma^{-1} \mathbf{x}_i - 2\mathbf{x}_i^T \Sigma^{-1} \mu + \mu^T \Sigma^{-1} \mu$$

The derivative of this term w.r.t. to μ is (helpful: $\frac{\partial \mathbf{a}^T \mathbf{y}}{\partial \mathbf{y}} = \mathbf{a}$, $\frac{\partial \mathbf{y}^T \mathbf{A} \mathbf{y}}{\partial \mathbf{y}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{y}$)

$$2\Sigma^{-1} \mu - 2\Sigma^{-1} \mathbf{x}_i = 2\Sigma^{-1} (\mu - \mathbf{x}_i)$$

So

$$\frac{\partial \ell(\boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} = \frac{1}{2} \sum_{i=1}^n 2\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \mathbf{x}_i) = \boldsymbol{\Sigma}^{-1} \sum_{i=1}^n (\boldsymbol{\mu} - \mathbf{x}_i)$$

Optimum at

$$\boldsymbol{\mu}_{\text{MLE}}^* = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

(second derivative at $\boldsymbol{\mu}_{\text{MLE}}$ is $2\boldsymbol{\Sigma}^{-1}$)

The trace operator: $tr(\mathbf{A}) := \sum_i \mathbf{A}_{ii}$ has a *cyclic* property

$$tr(\mathbf{ABC}) = tr(\mathbf{BCA}) = tr(\mathbf{CAB})$$

that allows the re-casting:

$$\mathbf{x}^T \mathbf{Ax} = tr(\mathbf{x}^T \mathbf{Ax}) = tr(\mathbf{Axx}^T)$$

i.e.

$$tr((\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})) = tr(\boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T)$$

One more fact:

$$|\boldsymbol{\Sigma}| = \frac{1}{|\boldsymbol{\Sigma}^{-1}|}$$

$$\ell(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}) = \text{const.} - \frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}| + \frac{1}{2} \sum_i \text{tr}(\boldsymbol{\Sigma}^{-1}(\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T)$$

$$\left(\frac{\partial \log |\mathbf{A}|}{\partial \mathbf{A}} = \mathbf{A}^{-T}, \frac{\partial \text{tr}(\mathbf{B}\mathbf{A})}{\partial \mathbf{A}} = \mathbf{B}^T \right)$$

$$\frac{\partial \ell(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1})}{\partial \boldsymbol{\Sigma}^{-1}} = -\frac{n}{2} \boldsymbol{\Sigma}^T + \frac{1}{2} \sum_i (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$$

Extremum at

$$\boldsymbol{\Sigma}_{\text{MLE}}^* = \frac{1}{n} \sum_i (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^T$$

Central Limit Theorem

Let $(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n)$ be i.i.d random variables with finite mean $\boldsymbol{\mu}$ and finite covariance $\boldsymbol{\Sigma}$, then

$$\mathbf{S}_n := \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i - \boldsymbol{\mu} \right) \Rightarrow \mathbf{S}_n \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$$

if you average i.i.d. variables, then only mean and covariance are retained (everything else is smoothed away) and a gaussian remains.

Maximum entropy distributions

How much uncertainty is in a distribution?

Differential entropy for continuous distribution P :

$$\mathcal{H}[P] = \int p(\mathbf{x})(-\log p(\mathbf{x}))d\mathbf{x}$$

Given mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ (and nothing else!), what distribution has highest differential entropy?

$$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \arg \max_P (\mathcal{H}[P] | E[\mathbf{X}] = \boldsymbol{\mu}, \text{Cov}[\mathbf{X}] = \boldsymbol{\Sigma})$$

Upper bound on entropy:

$$\frac{1}{2} \log |2\pi e \boldsymbol{\Sigma}|$$

How to sample from $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$\boldsymbol{\Sigma}$ is real, symmetric and positive definite. Thus, *Cholesky decomposition* exists:

$$\mathbf{L}\mathbf{L}^T = \boldsymbol{\Sigma}$$

where \mathbf{L} is lower triangular with strictly positive diagonal entries.

If we can easily sample \mathbf{Z} from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ (we usually can, why?), i.e. $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then

$$\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{Z} \Rightarrow \mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Why?

How to evaluate $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Evaluate *density* at some point \mathbf{x} , e.g. to compute log likelihood.

We need to compute $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$

$$\boldsymbol{\Sigma}^{-1} = \mathbf{L}^{-T} \mathbf{L}^{-1}$$

and compute in a *numerically stable way*

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

Also

$$|\boldsymbol{\Sigma}| = \prod_{j=1}^D L_{jj}^2$$

(prefer to work in the log domain!)

Normalization factor

There is the elementary result for the 1d normal distribution:

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

Remember the eigendecomposition of the covariance matrix:

$$\mathbf{\Sigma} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$

Denote by

$$f(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

the unnormalized probability density function of \mathbf{x} .

We need to compute

$$\int f(\mathbf{x}) d\mathbf{x}$$

Normalization factor

Defining $\mathbf{y} = \mathbf{U}^T(\mathbf{x} - \boldsymbol{\mu})$, the integral changes as follows (*change of variable*):

$$\int f(\mathbf{x})d\mathbf{x} = \int f(\mathbf{x}(\mathbf{y})) \left| \frac{d\mathbf{x}}{d\mathbf{y}} \right| d\mathbf{y} = \int f(\mathbf{x}(\mathbf{y})) |\mathbf{U}| d\mathbf{y} = \int \prod_i e^{-\frac{y_i^2}{2\lambda_i}} d\mathbf{y}$$

Using the elementary result:

$$\int \prod_i e^{-\frac{y_i^2}{2\lambda_i}} d\mathbf{y} = \prod_i \int e^{-\frac{y_i^2}{2\lambda_i}} dy_i = \prod_i \sqrt{2\pi\lambda_i} = \sqrt{|2\pi\boldsymbol{\Sigma}|}$$

Interesting Observation:

$$\text{diagonal } \boldsymbol{\Sigma} \Rightarrow p(\mathbf{x}) = \prod_i p(x_i)$$

In words: For Gaussians, uncorrelated components induce independent components (what is the *general* rule?).

Products of Gaussians

What is the product of two Gaussian pdf?

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \cdot \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

Difficult to answer in *moment parameterization* form:

$$\propto e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

Natural (or: canonical) parameterization:

$$\propto e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A}\mathbf{x} + \mathbf{r}^T \mathbf{x}}, \quad \mathbf{A} = \boldsymbol{\Sigma}^{-1}, \mathbf{r} = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$$

(\propto ignores all constants (i.e. terms without \mathbf{x}). Do the transformation to natural parameters on your own!)

A Gaussian pdf is written in *information form* (or: exponential family form), if natural parameterization is used.

Products of Gaussians

Using natural parameters, we can write:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \cdot \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \propto e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A}_1 \mathbf{x} + \mathbf{r}_1^T \mathbf{x}} \cdot e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A}_2 \mathbf{x} + \mathbf{r}_2^T \mathbf{x}}$$

The result is *again* a Gaussian, because we can write it in information form:

$$\propto e^{-\frac{1}{2}\mathbf{x}^T (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{x} + (\mathbf{r}_1 + \mathbf{r}_2)^T \mathbf{x}}$$

Converting it back into moment parameterization gives a new $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$:

$$\boldsymbol{\Sigma} = (\boldsymbol{\Sigma}_1^{-1} + \boldsymbol{\Sigma}_2^{-1})^{-1}, \quad \boldsymbol{\mu} = \boldsymbol{\Sigma}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_2^{-1} \boldsymbol{\mu}_2)$$

Now we can compute the missing normalization constant for the resulting Gaussian.

(do back-transformations on your own)

Marginalisation

$$I \subset \{1, 2, \dots, d\}, \quad \mathbf{X}_I := (X_i)_{i \in I}$$

What is $p(\mathbf{x}_I)$?

Linear Transformation with a *selection matrix*:

$$\mathbf{X}_I = \mathbf{I}_I \mathbf{X}$$

That is, \mathbf{X}_I is Gaussian:

$$p(\mathbf{x}_I) = \mathcal{N}(\boldsymbol{\mu}_I, \boldsymbol{\Sigma}_I)$$

Conditioning

$$I \subset \{1, 2, \dots, d\}, \quad R = \{1, 2, \dots, d\} \setminus I$$

What is $p(\mathbf{x}_I | \mathbf{x}_R)$?

A Gaussian, but how does it look like?

Rather straightforward when using natural paramterisation by remembering

$$p(\mathbf{x}) = p(\mathbf{x}_I | \mathbf{x}_R) p(\mathbf{x}_R)$$

Basically, read off result after some algebraic reformulations.

Not immediatelly obvious, if we want the results in moment parametrisation (needs *Schur* complement for inverting partitioned matrices). Thus, just the results:

$$\boldsymbol{\mu}_{I|R} = \boldsymbol{\mu}_I + \boldsymbol{\Sigma}_{IR} \boldsymbol{\Sigma}_{RR}^{-1} (\mathbf{x}_R - \boldsymbol{\mu}_R)$$

$$\boldsymbol{\Sigma}_{I|R} = \boldsymbol{\Sigma}_{II} - \boldsymbol{\Sigma}_{IR} \boldsymbol{\Sigma}_{RR}^{-1} \boldsymbol{\Sigma}_{RI}$$

Linear Gaussian systems

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_X)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{A}\mathbf{x} + \mathbf{b}, \boldsymbol{\Sigma}_{Y|X})$$

(\mathbf{x} , \mathbf{y} can have different dimensionalities)

What is $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$?

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_{X|Y}, \boldsymbol{\Sigma}_{X|Y})$$

$$\boldsymbol{\Sigma}_{X|Y} = (\boldsymbol{\Sigma}_X^{-1} + \mathbf{A}^T \boldsymbol{\Sigma}_{Y|X}^{-1} \mathbf{A})^{-1}$$

$$\boldsymbol{\mu}_{X|Y} = \boldsymbol{\Sigma}_{X|Y} \left(\mathbf{A}^T \boldsymbol{\Sigma}_{Y|X}^{-1} (\mathbf{y} - \mathbf{b}) + \boldsymbol{\Sigma}_X^{-1} \boldsymbol{\mu}_X \right)$$

Condition \mathbf{x} on a noisy observation of itself.

Let $\mathbf{Z} = (\mathbf{X}, \mathbf{Y})^T$:

$$p(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}_Z, \boldsymbol{\Sigma}_Z)$$

$$\boldsymbol{\mu}_Z = \begin{pmatrix} \boldsymbol{\mu}_X \\ \mathbf{A}\boldsymbol{\mu}_X + \mathbf{b} \end{pmatrix}$$

$$\boldsymbol{\Sigma}_Z = \begin{pmatrix} \boldsymbol{\Sigma}_X & \boldsymbol{\Sigma}_X \mathbf{A}^T \\ \mathbf{A}\boldsymbol{\Sigma}_X & \boldsymbol{\Sigma}_Y + \mathbf{A}\boldsymbol{\Sigma}_X \mathbf{A}^T \end{pmatrix}$$

And finally, \mathbf{Y} :

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\boldsymbol{\mu}_X + \mathbf{b}, \boldsymbol{\Sigma}_Y + \mathbf{A}\boldsymbol{\Sigma}_X \mathbf{A}^T)$$

Inferring an unknown vector from noisy measurements

Assume, that \mathbf{X} represents the true, but unknown location of some object (e.g. could be 2d/3d position). Model this by

$$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$$

We make noisy observations \mathbf{Y}_i of \mathbf{X} :

$$\mathbf{Y}_i \sim \mathcal{N}(\mathbf{X}, \boldsymbol{\Sigma}_Y)$$

This means, we know in what way our sensor errs.

Compared to the general form above, $\mathbf{A} = \mathbf{I}$ and $\mathbf{b} = \mathbf{0}$.

$$p(\mathbf{x} | \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n) = \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$$

$$\boldsymbol{\Sigma}_n = (\boldsymbol{\Sigma}_0^{-1} + n\boldsymbol{\Sigma}_y^{-1})^{-1}$$

$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_n(\boldsymbol{\Sigma}_y^{-1}(\sum_i \mathbf{y}_i) + \boldsymbol{\Sigma}_0^{-1}\boldsymbol{\mu}_0)$$

You can use the same idea to do sensor fusion (different kinds of sensors with different kinds of measure noise).

Bayes for Gaussian

With the previous formulae, we finally can do a Bayesian approach for Gaussians. To simplify the derivation, we only consider the case $p(\boldsymbol{\mu}|\mathcal{D}, \boldsymbol{\Sigma})$.

That is, we want to determine the posterior distribution for $\boldsymbol{\mu}$ from observations $\mathcal{D} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$, where we assume that the covariance $\boldsymbol{\Sigma}$ of these observations is known.

Gaussian prior for $\boldsymbol{\mu}$:

$$p(\boldsymbol{\mu}) = \mathcal{N}(\boldsymbol{\mu}_0, \mathbf{V}_0)$$

Then

$$\begin{aligned} p(\boldsymbol{\mu}|\mathcal{D}, \boldsymbol{\Sigma}) &= \mathcal{N}(\boldsymbol{\mu}_n, \mathbf{V}_n) \\ \mathbf{V}_n &= (\mathbf{V}_0^{-1} + n\boldsymbol{\Sigma}^{-1})^{-1} \\ \boldsymbol{\mu}_n &= \mathbf{V}_n(\boldsymbol{\Sigma}^{-1}(\sum_i \mathbf{x}_i) + \mathbf{V}_0^{-1}\boldsymbol{\mu}_0) \end{aligned}$$

Bayes for Gaussian

If we don't know anything about the prior (uninformative prior, i.e. $\mathbf{V}_0^{-1} = \mathbf{0I}$), then

$$p(\boldsymbol{\mu}|\mathcal{D}, \boldsymbol{\Sigma}) = \mathcal{N}\left(\frac{1}{n} \sum_i \mathbf{x}_i, \frac{1}{n} \boldsymbol{\Sigma}\right)$$

Remember the MLE of $\boldsymbol{\mu}$?