# (Multivariate) Gaussian (Normal) Probability Densities

Carl Edward Rasmussen, José Miguel Hernández-Lobato & Richard Turner

August 5th, 2020

# Gaussian Density

The probability density of a D-dimensional Gaussian with mean vector $\boldsymbol{\mu}$ and covariance matrix $\Sigma$ is given by

$$p(\mathbf{x}|\boldsymbol{\mu}, \Sigma) \;=\; \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma) \;=\; \frac{1}{(2\pi)^{D/2}|\Sigma|^{1/2}} \exp\big(-\tfrac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})\big),$$

and we also write

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

The covariance matrix $\Sigma$ must be symmetric and positive definite.

In the special (scalar) case where $D = 1$ we have

$$p(x|\mu, \sigma^2) \;=\; \frac{1}{\sqrt{2\pi\sigma^2}} \exp\big(-\tfrac{1}{2}(x-\mu)^2/\sigma^2\big),$$

where $\sigma^2$ is the variance and $\sigma$ is the standard deviation.

The *standard* Gaussian has $\boldsymbol{\mu} = \mathbf{0}$ and $\Sigma = I$ (the unit matrix), shorthand

$$\mathcal{N}(\mathbf{x}) \;=\; \mathcal{N}(\mathbf{x}|\boldsymbol{\mu} = \mathbf{0}, \Sigma = I).$$

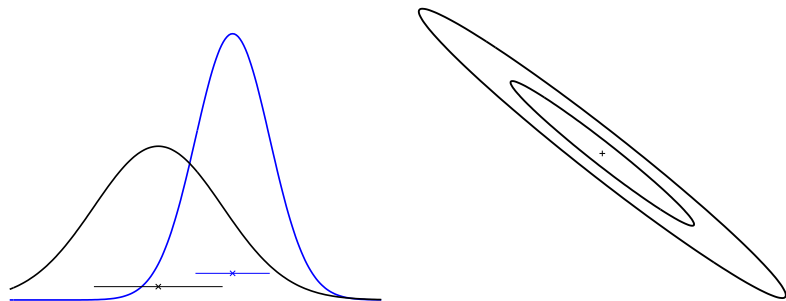# Parametrisation

There are two commonly used parametrisations of Gaussians

- *standard* parametrisation:
    - *mean* $\mu$ and
    - *covariance* $\Sigma$
- *natural* parametrisation:
    - *natural mean* $\nu = \Sigma^{-1}\mu$ and
    - *precision* matrix $R = \Sigma^{-1}$.

Different operations are more convenient in either parametrisation.
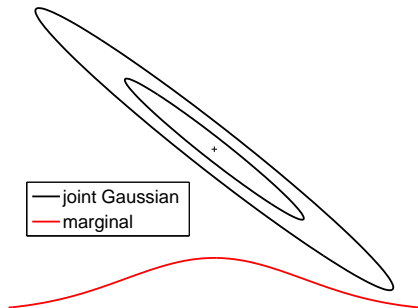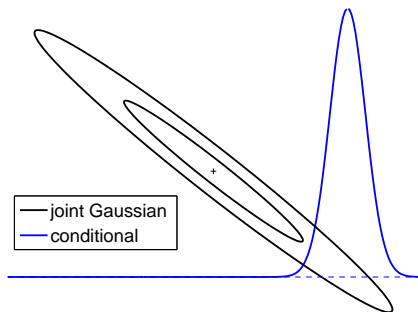
# Gaussian Pictures



The mean corresponds to the location or center of the distribution.

In one dimension, the square root of the variance corresponds to the *width* of the distribution.

In multiple dimensions, the eigen-vectors of the covariance matrix give the principal axis of the elliptical equi-probability contours of the distribution, and the square root of the eigenvalues the width of the distribution in the corresponding directions.

# Conditionals and Marginals of a Gaussian, pictorial



Both the conditionals $p(x|y)$ and the marginals $p(x)$ of a joint Gaussian $p(x, y)$ are again Gaussian.

# Conditionals and Marginals of a Gaussian, algebra

If $\mathbf{x}$ and $\mathbf{y}$ are jointly Gaussian

$$p(\mathbf{x}, \mathbf{y}) = p\left(\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}\right),$$

we get the marginal distribution of $\mathbf{x}$, $p(\mathbf{x})$ by

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}\right) \implies p(\mathbf{x}) = \mathcal{N}(\mathbf{a}, A),$$

and the conditional distribution of $\mathbf{x}$ given $\mathbf{y}$ by

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}\right) \implies p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{a} + BC^{-1}(\mathbf{y} - \mathbf{b}), \ A - BC^{-1}B^\top),$$

where $\mathbf{x}$ and $\mathbf{y}$ can be scalars or vectors.

# Kullback-Leibler Divergence (Relative Entropy)

The Kullback-Leibler (KL) divergence between continuous distributions is

$$\mathcal{KL}(q(x)\|p(x)) = \int q(x) \log \frac{q(x)}{p(x)} dx.$$

The KL divergence is an asymmetric measure of distance between distributions.
The KL divergence between two Gaussians is

$$\mathcal{KL}(\mathcal{N}_0\|\mathcal{N}_1) = \tfrac{1}{2}\log|\Sigma_1\Sigma_0^{-1}| + \tfrac{1}{2}\operatorname{tr}\left(\Sigma_1^{-1}\big((\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^\top + \Sigma_0 - \Sigma_1\big)\right).$$

# KL matching constrained Gaussians

It is often convenient to approximate one distribution with another, simpler one, by finding the *closest match* within a constrained family.

Minimizing KL divergence between a general Gaussian $\mathcal{N}_g$ and a factorized Gaussian $\mathcal{N}_f$ will match the means $\mu_f = \mu_g$ and for the covariances either:

$$\frac{\partial \mathcal{KL}(\mathcal{N}_f \| \mathcal{N}_g)}{\partial \Sigma_f} = -\tfrac{1}{2}\Sigma_f^{-1} + \tfrac{1}{2}\Sigma_g^{-1} = 0 \;\Rightarrow\; (\Sigma_f)_{ii} = 1/(\Sigma_g^{-1})_{ii},$$
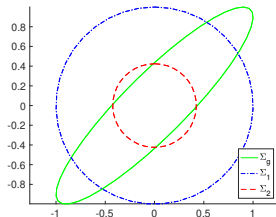
or

$$\frac{\partial \mathcal{KL}(\mathcal{N}_g \| \mathcal{N}_f)}{\partial \Sigma_f} = \tfrac{1}{2}\Sigma_f^{-1} - \tfrac{1}{2}\Sigma_f^{-1}\Sigma_g\Sigma_f^{-1} = 0 \;\Rightarrow\; (\Sigma_f)_{ii} = (\Sigma_g)_{ii}.$$

Interpretation:



- averaging wrt the *factorized* Gaussian, the fitted variance equals the *conditional* variance of $\Sigma_g$,

- averaging wrt the *general* Gaussian, the fitted variance equals the *marginal* variance of $\Sigma_g$,

with straight forward generalization to block diagonal Gaussians.

# Incomplete (truncated) scalar Gaussian integrals

Let $\Phi(z)$ be the standard cumulative Gaussian

$$\Phi(z) = \int_{-\infty}^{z} \mathcal{N}(x)dx = \int_{-\infty}^{z} \frac{1}{\sqrt{2\pi}}\exp(-\tfrac{1}{2}x^2)dx.$$

We then have the following incomplete Gaussian integrals

$$\int_{a}^{b} \mathcal{N}(x|\mu, \sigma^2)dx = \Phi(\beta) - \Phi(\alpha), \text{ where } \alpha = \frac{a - \mu}{\sigma} \text{ and } \beta = \frac{b - \mu}{\sigma}.$$

Further

$$\int_{a}^{b} \frac{x - \mu}{\sigma}\mathcal{N}(x|\mu, \sigma^2)dx = \mathcal{N}(\alpha) - \mathcal{N}(\beta),$$

and

$$\int_{a}^{b} \left(\frac{x - \mu}{\sigma}\right)^2 \mathcal{N}(x|\mu, \sigma^2)dx = \alpha\mathcal{N}(\alpha) - \beta\mathcal{N}(\beta) + \Phi(\beta) - \Phi(\alpha),$$

which can both be shown by integration by parts. Both expressions have the expected behaviour as $a \to -\infty$ and/or $b \to \infty$ (one sided Gaussians).

# Appendix: Some useful Gaussian identities

If x is multivariate Gaussian with mean $\mu$ and covariance matrix $\Sigma$

$$p(\mathbf{x}; \mu, \Sigma) = (2\pi|\Sigma|)^{-D/2} \exp\left(-(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)/2\right),$$

then

$$\mathbb{E}[\mathbf{x}] = \mu,$$
$$\mathbb{V}[\mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])^2] = \Sigma.$$

For any matrix $A$, if $\mathbf{z} = A\mathbf{x} + \mathbf{b}$ then $\mathbf{z}$ is Gaussian and

$$\mathbb{E}[\mathbf{z}] = A\mu + \mathbf{b},$$
$$\mathbb{V}[\mathbf{z}] = A\Sigma A^\top.$$

# Matrix and Gaussian identities cheat sheet

Matrix identities

- Matrix inversion lemma (Woodbury, Sherman & Morrison formula)

$$(Z + UWV^\top)^{-1} = Z^{-1} - Z^{-1}U(W^{-1} + V^\top Z^{-1}U)^{-1}V^\top Z^{-1}$$

- A similar equation exists for determinants

$$|Z + UWV^\top| = |Z|\,|W|\,|W^{-1} + V^\top Z^{-1}U|$$

The product of two Gaussian density functions

$$\mathcal{N}(\mathbf{x}|\mathbf{a}, A)\,\mathcal{N}(P^\top \mathbf{x}|\mathbf{b}, B) = z_c\,\mathcal{N}(\mathbf{x}|\mathbf{c}, C)$$

- is proportional to a Gaussian density function with covariance and mean

$$C = \left(A^{-1} + P\,B^{-1}P^\top\right)^{-1} \qquad \mathbf{c} = C\,\left(A^{-1}\mathbf{a} + P\,B^{-1}\,\mathbf{b}\right)$$

- and has a normalizing constant $z_c$ that is Gaussian both in $\mathbf{a}$ and in $\mathbf{b}$

$$z_c = (2\,\pi)^{-\frac{m}{2}}|B + P^\top A\,P|^{-\frac{1}{2}}\exp\left(-\frac{1}{2}(\mathbf{b} - P^\top\,\mathbf{a})^\top\left(B + P^\top A\,P\right)^{-1}(\mathbf{b} - P^\top\,\mathbf{a})\right)$$