

# The Gaussian distribution

## Probability distributions

Francesco Corona

# Outline

## The Gaussian distribution

- Conditional Gaussian distributions
- Marginal Gaussian distributions
- Bayes' theorem for Gaussian variables
- Maximum likelihood for the Gaussian
- Bayesian inference for the Gaussian
- Student's t-distribution

## Mixtures of Gaussians

# The Gaussian distribution

## Probability distributions

## The Gaussian distribution

The **Gaussian** or **normal distribution**, is a classic model for the distribution of continuous variables

In the case of a single variable  $x$ , the Gaussian distribution can be written as

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) \quad (1)$$

- ▶  $\mu$  is the mean and  $\sigma^2$  is the variance

In the case of a  $D$ -dimensional variable  $\mathbf{x}$ , the Gaussian distribution is

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

- ▶  $\boldsymbol{\mu}$  is the  $D$ -dimensional mean vector
- ▶  $\boldsymbol{\Sigma}$  is the  $D \times D$  covariance matrix
- ▶  $|\boldsymbol{\Sigma}|$  is the determinant of  $\boldsymbol{\Sigma}$

## The Gaussian distribution (cont.)

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \quad (3)$$

Let us start by considering the geometric form of the Gaussian distribution

- ▶ The Gaussian depends on  $\mathbf{x}$  is through the quadric form

$$\Delta^2 = (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (4)$$

- ▶ Quantity  $\Delta$  is the **Mahalanobis distance** from  $\boldsymbol{\mu}$  to  $\mathbf{x}$
- ▶ It reduces to the Euclidean distance when  $\boldsymbol{\Sigma} = \mathbf{I}$

The Gaussian is constant on surfaces in  $\mathbf{x}$  space for which Eq. 4 is constant

## The Gaussian distribution - Eigenequation

For a square matrix  $\mathbf{A}$  of size  $M \times M$ , the **eigenvector equation** is

$$\mathbf{A}\mathbf{u}_i = \lambda_i\mathbf{u}_i, \quad \text{with } i = 1, \dots, M$$

- ▶  $\mathbf{u}_i$  is an **eigenvector** and  $\lambda_i$  is the corresponding **eigenvalue**
- ▶ The condition for a solution is the characteristic equation

$$|\mathbf{A} - \lambda_i\mathbf{I}| = 0$$

- ▶ Generally, the eigenvalues of a matrix are complex numbers
- ▶ The **rank** of  $\mathbf{A}$  equals the number of its nonzero eigenvalues

$\text{Rank}(\mathbf{A})$  is the number of linearly independent row/columns in  $\mathbf{A}$

- ▶ Given a set of vectors  $\{\mathbf{a}_1, \dots, \mathbf{a}_K\}$ , the set is said to be **linearly independent** if  $\sum_k \alpha_k \mathbf{a}_k = 0$  holds only when all  $\alpha_k = 0$
- ▶ No vector  $\mathbf{a}_k$  can be expressed as linear combination of the others

## The Gaussian distribution - Eigenequation (cont.)

A symmetric matrix  $\mathbf{A}$ , a covariance of the Gaussian, is such that  $A_{ij} = A_{ji}$

$$\mathbf{A}^T = \mathbf{A}$$

The inverse of a symmetric matrix is also symmetric  $(\mathbf{A}^{-1})^T = \mathbf{A}^{-1}$ , with

$$\mathbf{A}\mathbf{A}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbf{I}$$

The eigenvectors  $\mathbf{u}_i$  of a real symmetric matrix can be chosen to be orthonormal

$$\mathbf{u}_i^T \mathbf{u}_j = l_{ij} \quad \text{with } l_{ij} \text{ elements of } \mathbf{I} \text{ such that } l_{ij} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases}$$

The eigenvalues  $\lambda_i$  of a symmetric matrix are real numbers

## The Gaussian distribution - Eigenequation (cont.)

Since the eigenvectors  $\mathbf{u}_i$  can be chosen to be orthogonal and of unit length, we can take them to be the columns of an orthogonal  $M \times M$  matrix  $\mathbf{U}$ ,  $\mathbf{U}\mathbf{U}^T = \mathbf{I}$

- Incidentally, note that also the rows of  $\mathbf{U}$  are orthogonal,  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$

We can use  $\mathbf{U}$  to transform a vector  $\mathbf{x}$  into a new vector  $\tilde{\mathbf{x}}$ , that is  $\tilde{\mathbf{x}} = \mathbf{U}\mathbf{x}$

- The length of the vector is preserved:  $\tilde{\mathbf{x}}^T\tilde{\mathbf{x}} = \mathbf{x}^T\mathbf{U}^T\mathbf{U}\mathbf{x} = \mathbf{x}^T\mathbf{x}$
- The angle between two vectors is preserved:  $\tilde{\mathbf{x}}^T\tilde{\mathbf{y}} = \mathbf{x}^T\mathbf{U}^T\mathbf{U}\mathbf{y} = \mathbf{x}^T\mathbf{y}$

Multiplication by  $\mathbf{U}$  represents a rigid rotation of the coordinate system



## The Gaussian distribution - Eigenequation (cont.)

We can write the eigenequation  $\mathbf{A}\mathbf{u}_i = \lambda_i\mathbf{u}_i$ , with  $i = 1, \dots, N$ , in terms of  $\mathbf{U}$

$$\mathbf{A}\mathbf{U} = \mathbf{U}\mathbf{\Lambda}$$

▶  $\mathbf{\Lambda}$  is a  $M \times M$  diagonal matrix, with diagonal  $\text{diag}(\mathbf{\Lambda}) = (\lambda_1, \dots, \lambda_M)^T$

▶  $\mathbf{U}^T \mathbf{A} \mathbf{U} = \underbrace{\mathbf{U}^T \mathbf{U}}_{\mathbf{I}} \mathbf{\Lambda} \implies \mathbf{U}^T \mathbf{A} \mathbf{U} = \mathbf{\Lambda}$ , matrix  $\mathbf{A}$  is diagonalised by  $\mathbf{U}$

▶  $\underbrace{\mathbf{U}\mathbf{U}^T}_{\mathbf{I}} \mathbf{A} \underbrace{\mathbf{U}\mathbf{U}^T}_{\mathbf{I}} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \implies \mathbf{A} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \implies \mathbf{A} = \sum_{i=1}^M \lambda_i \mathbf{u}_i \mathbf{u}_i^T$

▶  $\mathbf{A}^{-1} = \mathbf{U}\mathbf{\Lambda}^{-1}\mathbf{U}^T \implies \mathbf{A}^{-1} = \sum_{i=1}^M \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$

▶  $|\mathbf{A}| = \prod_{i=1}^M \lambda_i$

▶  $\text{Trace}(\mathbf{A}) = \sum_{i=1}^M \lambda_i$

## The Gaussian distribution (cont.)

Consider the eigenvector equation for a real symmetric covariance matrix  $\Sigma$

$$\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i, \quad \text{with } i = 1 \dots, D \quad (5)$$

Real eigenvalues and its eigenvectors form an orthonormal set

$$\mathbf{u}_i^T \mathbf{u}_j = I_{ij} = \begin{cases} 1, & \text{if } i = j \\ 0, & \text{if } i \neq j \end{cases} \quad (6)$$

$\Sigma$  can be expressed as an expansion of its eigenvectors ( $\star$ )

$$\Sigma = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T \quad (7)$$

The inverse covariance matrix  $\Sigma^{-1}$  can be expressed ( $\star$ ) as

$$\Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (8)$$

## The Gaussian distribution (cont.)

By substituting the inverse covariance matrix  $\Sigma^{-1}$  into the quadratic form  $\Delta^2$

$$\begin{aligned}\Delta^2 &= (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \\ &= (\mathbf{x} - \boldsymbol{\mu})^T \left( \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right) (\mathbf{x} - \boldsymbol{\mu}) \\ &= \sum_{i=1}^D \frac{(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{u}_i \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})}{\lambda_i} \\ &= \sum_{i=1}^D \frac{y_i^2}{\lambda_i}, \quad \text{with } y_i = \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})\end{aligned}\tag{9}$$

## The Gaussian distribution (cont.)

We interpret  $\{y_i\}$  as a new coordinate system defined by orthonormal vectors  $\mathbf{u}_i$ , that are shifted by  $\boldsymbol{\mu}$  and rotated with respect to the original coordinates  $x_i$

- ▶ Forming the vector  $\mathbf{y} = (y_1, \dots, y_D)^T$ , we have

$$\mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu}) \quad (10)$$

$\mathbf{U}$  is an orthogonal matrix whose rows are  $\mathbf{u}_i^T$  (i.e.,  $\mathbf{U}\mathbf{U}^T = \mathbf{I}$  and  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ )

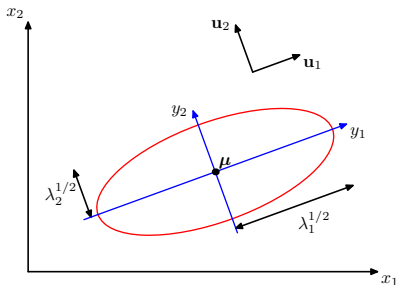
## The Gaussian distribution (cont.)

The quadratic form and thus the Gaussian is constant on surfaces for which  $\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}$  is constant

For positive  $\lambda_i$ , the surfaces are ellipsoids

- ▶ Centred in  $\mu$  and axis oriented along  $u_i$
- ▶ The scaling factors in the directions of the axes are  $\lambda_i^{1/2}$

The red curve is the elliptical surface of a 2D Gaussian in which the density is  $\exp(-1/2)$  of its value at  $\mathbf{x} = \mu$



## The Gaussian distribution (cont.)

For the Gaussian to be well-defined, all of the eigenvalues  $\lambda_i$  of the covariance matrix need be strictly positive<sup>1</sup>, otherwise it cannot be properly normalised

- ▶ A Gaussian for which one or more eigenvalues are zero<sup>2</sup> is singular
- ▶ It is confined to a subspace of lower dimensionality

---

<sup>1</sup>A matrix whose eigenvalues are strictly positive is called **positive definite**

<sup>2</sup>A matrix in which all of the eigenvalues are nonnegative is called **positive semidefinite**

## The Gaussian distribution - Jacobian factor

Under a change of variable, a density does not transform like a regular function

For a change of variables  $x = g(y)$ , a function  $f(x)$  becomes  $\tilde{f}(y) = f(g(y))$

Consider a probability density  $p_x(x)$  that corresponds to a density  $p_y(y)$  wrt a new variable  $y$  and notice that  $p_x(x)$  and  $p_y(y)$  are different densities

- ▶ Observations falling in a range  $(x, x + \delta x)$  have probability  $p_x(x)\delta x$
- ▶ By transforming them, we make them fall in the range  $(y, y + \delta y)$
- ▶ Observations falling in a range  $(y, y + \delta y)$  have probability  $p_y(y)\delta y$

$$p_x(x)\delta x \simeq p_y(y)\delta y$$

$$p_y(y) = p_x(x) \left| \frac{dx}{dy} \right| = p_x(x) \left| \frac{dg(y)}{dy} \right| = p_x(x) |g'(y)|$$

## The Gaussian distribution (cont.)

Consider now the Gaussian in the new coordinate system defined by  $y_i$

In going from  $\mathbf{x}$  to  $\mathbf{y}$ , we have a Jacobian matrix  $\mathbf{J}$

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = U_{ji} \quad (11)$$

Thus, the elements of the  $\mathbf{J}$  are the elements of  $\mathbf{U}^T$

Using the orthonormality property of  $\mathbf{U}$ , the square of the determinant of  $\mathbf{J}$

$$|\mathbf{J}^2| = |\mathbf{U}^T|^2 = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{I}| = 1 \quad \implies |\mathbf{J}| = 1 \quad (12)$$



## The Gaussian distribution (cont.)

We also have that the determinant  $|\Sigma|$  of the covariance matrix can be written as the product of its eigenvalues, and thus we also have

$$|\Sigma| = \prod_{j=1}^D \lambda_j \quad \implies \quad |\Sigma|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2} \quad (13)$$

Hence, the Gaussian distribution in the coordinate system  $y_i$  becomes

$$p(\mathbf{y}) = p(\mathbf{x})|\mathbf{J}| = \prod_{j=1}^D \frac{1}{(2\pi\lambda_j)^{1/2}} \exp\left(-\frac{y_j^2}{2\lambda_j}\right) \quad (14)$$

Here, it is the product of  $D$  independent univariate Gaussian distributions

- ▶ The eigenvectors define a new set of shifted and rotated coordinates
- ▶ Here, the joint distribution factorises into independent distributions

## The Gaussian distribution (cont.)

Using a result derived for the normalisation of the univariate Gaussian,

$$\int p(\mathbf{y}) d\mathbf{y} = \prod_{j=1}^D \int_{-\infty}^{+\infty} \frac{1}{(2\pi\lambda_j)^{1/2}} \exp\left(-\frac{y_j^2}{2\lambda_j}\right) dy_j = 1 \quad (15)$$

which is the integral of the distribution in the  $\mathbf{y}$  coordinate system

## The Gaussian distribution (cont.)

We now look at the moments of the Gaussian distribution (back to the  $\mathbf{x}$ -space)

The expectation of  $\mathbf{x}$  under the Gaussian distribution is given by

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \mathbf{x} d\mathbf{x} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left(-\frac{1}{2}\mathbf{z}^T \Sigma^{-1}\mathbf{z}\right) (\mathbf{z} + \boldsymbol{\mu}) d\mathbf{z}\end{aligned}\quad (16)$$

- We have changed variables using  $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$

The exponent is an even function<sup>3</sup> of the components  $\mathbf{z}$ , which for integrals taken in  $(-\infty, +\infty)$  will make the term in  $\mathbf{z}$  in the factor  $(\mathbf{z} + \boldsymbol{\mu})$  vanish

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} \quad (17)$$

---

<sup>3</sup>A real-valued function is said to be even if  $f(x) = f(-x)$ , or  $f(x) - f(-x) = 0$

## The Gaussian distribution (cont.)

There are  $D^2$  second order moments  $\mathbb{E}[x_i x_j]$  under the Gaussian distribution

$$\begin{aligned}\mathbb{E}[\mathbf{x}\mathbf{x}^T] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \mathbf{x}\mathbf{x}^T d\mathbf{x} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left(-\frac{1}{2}\mathbf{z}^T \Sigma^{-1}\mathbf{z}\right) (\mathbf{z} + \boldsymbol{\mu})(\mathbf{z} + \boldsymbol{\mu})^T d\mathbf{z} \quad (18)\end{aligned}$$

- ▶ The cross-terms involving  $\boldsymbol{\mu}\mathbf{z}^T$  and  $\mathbf{z}\boldsymbol{\mu}^T$  will vanish by symmetry
- ▶ The term  $\boldsymbol{\mu}\boldsymbol{\mu}^T$  is constant and can be taken outside the integral

## The Gaussian distribution (cont.)

As for the term involving  $\mathbf{z}\mathbf{z}^T$ , using the eigenvector expansion of the covariance matrix together with the completeness of the set of eigenvectors

$$\mathbf{z} = \sum_{j=1}^D y_j \mathbf{u}_j, \quad \text{with } y_j = \mathbf{u}_j^T \mathbf{z} \quad (19)$$

$$\begin{aligned} & \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp\left(-\frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z}\right) \mathbf{z} \mathbf{z}^T d\mathbf{z} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \sum_{i=1}^D \sum_{j=1}^D \mathbf{u}_i \mathbf{u}_j^T \int \exp\left(-\sum_{k=1}^D \frac{2y_k^2}{2\lambda_k}\right) y_i y_j d\mathbf{y} \\ &= \sum_{i=1}^D \mathbf{u}_i \mathbf{u}_i^T \lambda_i = \Sigma \end{aligned} \quad (20)$$

## The Gaussian distribution (cont.)

We used the eigenvector equation  $\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i$  and the fact that the integral on the right-hand side of the middle line vanishes by symmetry unless  $i = j$

In the final line we used  $\mathbb{E}[x^2] = \int_{-\infty}^{+\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2$  and  $|\Sigma|^{1/2} = \prod_{j=1}^D \lambda_j^{1/2}$ , together with  $\Sigma = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T$

As a result, we have

$$\mathbb{E}[\mathbf{x} \mathbf{x}^T] = \mu \mu^T + \Sigma \quad (21)$$

For single random variables, we subtracted the mean before taking second moments and define a variance

Similarly, in the multivariate case it is again convenient to subtract off the mean, giving rise to the covariance of a random vector  $\mathbf{x}$  defined by

$$\text{cov}[\mathbf{x}] = \text{cov}[\mathbf{x}, \mathbf{x}] = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x}^T - \mathbb{E}[\mathbf{x}^T])] \quad (22)$$

## The Gaussian distribution (cont.)

For the specific case of a Gaussian distribution, we use of the first-order moment  $\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$  and the second-order moment  $\mathbb{E}[\mathbf{x}\mathbf{x}^T] = \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\Sigma}$

$$\text{cov}[\mathbf{x}] = \boldsymbol{\Sigma} \quad (23)$$

The parameter matrix  $\boldsymbol{\Sigma}$  governs the covariance of  $\mathbf{x}$  under the Gaussian distribution, hence it is called the **covariance matrix**

## The Gaussian distribution (cont.)

The Gaussian distribution is widely used as a density model

- ▶ Though, it suffers from some significant limitations

Consider the number of free parameters  $D(D + 3)/2$  in the distribution (★)

- ▶ A general symmetric covariance matrix  $\Sigma$  has  $D(D + 1)/2$  independent parameters (★) and there are another  $D$  independent parameters in  $\mu$

For large  $D$ , this number grows quadratically with  $D$ , and the computational task of manipulating and inverting large matrices can become prohibitive

A way to *address* this problem is to use covariance matrices of restricted form



## The Gaussian distribution (cont.)

We can consider covariance matrices that are diagonal ( $\Sigma = \text{diag}(\sigma_i^2)$ )

- ▶ A total of  $2D$  independent parameters in the density model
- ▶ Axis-aligned ellipsoids of constant density

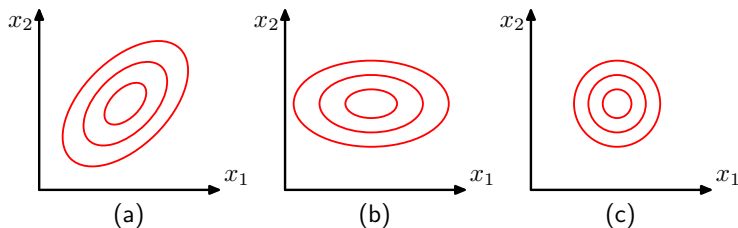
We can restrict covariance matrices to be proportional to the identity matrix ( $\Sigma = \sigma^2 \mathbf{I}$ , or isotropic covariance), we then get  $D + 1$  independent parameters

- ▶ Spherical surfaces of constant density

Such approaches limit the number of degrees of freedom in the distribution and make inversion of the covariance matrix much faster

They also greatly restrict the form of the probability density and limit its ability to capture interesting correlations in the data

## The Gaussian distribution (cont.)



- ▶ (a): general  $\Sigma$
- ▶ (b):  $\Sigma = \text{diag}(\sigma^2)$
- ▶ (c):  $\Sigma = \sigma^2 \mathbf{I}$

# Conditional Gaussian distributions

## The Gaussian distribution

## Conditional Gaussian distributions

Property of the Gaussian: If two sets of variables are jointly Gaussian, then the conditional distribution of one set conditioned on the other is again Gaussian

Suppose  $\mathbf{x}$  is a  $D$ -dimensional vector with Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- ▶ We partition  $\mathbf{x}$  into two disjoint subsets  $\mathbf{x}_a$  and  $\mathbf{x}_b$

Without loss of generality:

- ▶  $\mathbf{x}_a$  comprises the first  $M$  components of  $\mathbf{x}$
- ▶  $\mathbf{x}_b$  comprises the remaining  $D - M$  ones
- ▶

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad (24)$$

## Conditional Gaussian distributions (cont.)

We also define corresponding partitions of

- ▶ the mean vector  $\mu$

$$\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \quad (25)$$

- ▶ the covariance matrix  $\Sigma$

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad (26)$$

Because of the symmetry  $\Sigma^T = \Sigma$ , we have that  $\Sigma_{aa}$  and  $\Sigma_{bb}$  are also symmetric and  $\Sigma_{ba} = \Sigma_{ab}^T$

## Conditional Gaussian distributions (cont.)

In many situations, it is convenient to work with the **precision matrix**

$$\mathbf{\Lambda} = \mathbf{\Sigma}^{-1} \quad (27)$$

The corresponding partitioned precision matrix is given by the form

$$\mathbf{\Lambda} = \begin{pmatrix} \mathbf{\Lambda}_{aa} & \mathbf{\Lambda}_{ab} \\ \mathbf{\Lambda}_{ba} & \mathbf{\Lambda}_{bb} \end{pmatrix} \quad (28)$$

Because the inverse of a symmetric matrix is also symmetric,  $\mathbf{\Lambda}_{aa}$  and  $\mathbf{\Lambda}_{bb}$  are also symmetric and  $\mathbf{\Lambda}_{ba} = \mathbf{\Lambda}_{ab}^T$  ( $\star$ )

## Conditional Gaussian distributions (cont.)

We begin by finding an expression for the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$

From the product rule of probability, this conditional distribution can be evaluated from the joint distribution  $p(\mathbf{x}) = p(\mathbf{x}_a, \mathbf{x}_b)$  by fixing  $\mathbf{x}_b$  and normalising the result to obtain a valid probability distribution over  $\mathbf{x}_a$

Instead of performing this normalisation explicitly, we can obtain the solution more efficiently by using the quadratic form in the exponent of the Gaussian

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

and reinstating the normalisation coefficient at the end of the manipulations

## Conditional Gaussian distributions (cont.)

We use the partitioning:  $\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$   $\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}$   $\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) &= \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

We see that as a function of  $\mathbf{x}_a$ , this is again a quadratic form, and hence the corresponding conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  will be Gaussian

This distribution is completely characterised by its mean and its covariance, and our goal is to identify expressions for the mean and covariance of  $p(\mathbf{x}_a|\mathbf{x}_b)$



## Conditional Gaussian distributions (cont.)

- ▶ We are given a quadratic form defining the exponent terms in a Gaussian, and we need to determine corresponding mean and covariance

The exponent in a general Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  can be written

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const} \quad (29)$$

where 'const' denotes terms which are independent on  $\mathbf{x}$  (i.e.,  $-\frac{1}{2}\boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}$ )

If we take the general quadratic form and express it in the form given by the right-hand side of Eq. 29, then we can equate the matrix of coefficients entering the second-order term in  $\mathbf{x}$  to the inverse covariance matrix  $\boldsymbol{\Sigma}^{-1}$  and the coefficient of the linear term in  $\mathbf{x}$  to  $\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}$ , from which we can obtain  $\boldsymbol{\mu}$

## Conditional Gaussian distributions (cont.)

We apply this procedure to the conditional Gaussian distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$ :

- The quadratic form in the exponent

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) &= \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

We will denote mean and covariance of this distribution by  $\boldsymbol{\mu}_{a|b}$  and  $\boldsymbol{\Sigma}_{a|b}$

Consider the functional dependence on  $\mathbf{x}_a$  in which  $\mathbf{x}_b$  is regarded as a constant

- If we pick out all second-order terms in  $\mathbf{x}_a$ , we have

$$-\frac{1}{2}\mathbf{x}_a^T \boldsymbol{\Lambda}_{aa}\mathbf{x}_a \tag{30}$$

- From which, the covariance of  $p(\mathbf{x}_a|\mathbf{x}_b)$  is given as

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1} \tag{31}$$

## Conditional Gaussian distributions (cont.)

$$\begin{aligned}
 -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) &= \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\
 &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b)
 \end{aligned}$$

- Consider all terms that are linear in  $\mathbf{x}_a$

$$\mathbf{x}_a^T \left( \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \right) \quad (32)$$

From our discussion, the coefficient of  $\mathbf{x}_a$  in this expression must equal  $\boldsymbol{\Sigma}_{a|b}^{-1} \boldsymbol{\mu}_{a|b}$

$$\begin{aligned}
 \boldsymbol{\mu}_{a|b} &= \boldsymbol{\Sigma}_{a|b} \left( \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \right) \\
 &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)
 \end{aligned} \quad (33)$$

## Conditional Gaussian distributions (cont.)

Mean and variance of conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  are expressed in terms of partitioned precision matrix of the original joint distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$

- We can express these results also in terms of partitioned covariance matrix

To do this  $\star$ , we make use of this identity for the inverse of a partitioned matrix

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix} \quad (34)$$

with  $\mathbf{M} = (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}$

$\mathbf{M}^{-1}$  is the **Schur complement** of  $\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1}$  with respect to sub-matrix  $\mathbf{D}$

## Conditional Gaussian distributions (cont.)

Using the definition

$$\begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix} \quad (35)$$

and using

$$\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$$

we have that

$$\Lambda_{aa} = (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1} \quad (36)$$

$$\Lambda_{ab} = -(\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}\Sigma_{ab}\Sigma_{bb}^{-1} \quad (37)$$

## Conditional Gaussian distributions (cont.)

From these we obtain the following expressions for the mean and variance of the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  in terms of the partitioned covariance matrix

$$\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b) \quad (38)$$

$$\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba} \quad (39)$$

Compare  $\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Lambda}_{aa}^{-1}$  with  $\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$

- ▶ The conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  takes a simpler form when expressed in terms of the partitioned precision matrix than when it is expressed in terms of the partitioned covariance matrix

The mean  $\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$  of the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  is a linear function of  $\mathbf{x}_b$

The covariance  $\boldsymbol{\Sigma}_{a|b} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}$  of the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  is independent of  $\mathbf{x}_b$

- ▶ It is an example of **Linear-Gaussian model**

# **Marginal Gaussian distributions**

## **The Gaussian distribution**

## Marginal Gaussian distributions

If a joint distribution  $p(\mathbf{x}_a, \mathbf{x}_b)$  is Gaussian, then  $p(\mathbf{x}_a|\mathbf{x}_b)$ , the conditional distribution, is also Gaussian

We discuss the marginal distribution  $p(\mathbf{x}_a)$

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b \quad (40)$$

and we shall see that this is also Gaussian

The strategy for evaluating this distribution efficiently focuses on the quadratic form in the exponent of the joint distribution

- to identify mean and covariance of the marginal distribution  $p(\mathbf{x}_a)$



## Marginal Gaussian distributions (cont.)

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) &= \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{aa}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

The goal is to integrate out  $\mathbf{x}_b$ , because  $p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$

Picking out those terms that involve  $\mathbf{x}_b$ , we have

$$-\frac{1}{2}\mathbf{x}_b^T \boldsymbol{\Lambda}_{bb}\mathbf{x}_b + \mathbf{x}_b^T \mathbf{m} = -\frac{1}{2}(\mathbf{x}_b - \boldsymbol{\Lambda}_{bb}^{-1}\mathbf{m})^T \boldsymbol{\Lambda}_{bb}(\mathbf{x}_b - \boldsymbol{\Lambda}_{bb}^{-1}\mathbf{m}) + \frac{1}{2}\mathbf{m}^T \boldsymbol{\Lambda}_{bb}^{-1}\mathbf{m} \quad (41)$$

$$\text{with } \mathbf{m} = \boldsymbol{\Lambda}_{bb}\boldsymbol{\mu}_b - \boldsymbol{\Lambda}_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a) \quad (42)$$

The dependence on  $\mathbf{x}_b$  has been cast into the standard quadratic form of a Gaussian distribution corresponding to the first term on the right-hand side above, plus a term that does not depend on  $\mathbf{x}_b$  but that does depend on  $\mathbf{x}_a$

## Marginal Gaussian distributions (cont.)

When we take the exponential of this quadratic form, we see that the integration over  $\mathbf{x}_b$  required by  $p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b$  takes the form

$$\int \exp \left( -\frac{1}{2} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m}) \right) d\mathbf{x}_b \quad (43)$$

It is the integral over an unnormalised Gaussian, and so the result will be the reciprocal of the normalisation coefficient

$$\frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}}$$

which is independent of the mean, and it depends only on the determinant of the covariance matrix

## Marginal Gaussian distributions (cont.)

By completing the square wrt  $\mathbf{x}_b$ , we can integrate out  $\mathbf{x}_b$  and the only term remaining that depends on  $\mathbf{x}_a$  from the contributions on the left-hand side of

$$-\frac{1}{2}\mathbf{x}_b^T \Lambda_{bb} \mathbf{x}_b + \mathbf{x}_b \mathbf{m} = -\frac{1}{2}(\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m})^T \Lambda_{bb} (\mathbf{x}_b - \Lambda_{bb}^{-1} \mathbf{m}) + \frac{1}{2} \mathbf{m}^T \Lambda_{bb}^{-1} \mathbf{m}$$

is the last term on the right-hand side because  $\mathbf{m} = \Lambda_{bb} \boldsymbol{\mu}_b - \Lambda_{ba}(\mathbf{x}_a - \boldsymbol{\mu}_a)$

We then combine this term with the remaining terms that depend on  $\mathbf{x}_a$  from

$$\begin{aligned} -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Lambda (\mathbf{x} - \boldsymbol{\mu}) &= \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \Lambda_{aa} (\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_a - \boldsymbol{\mu}_a)^T \Lambda_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &\quad - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \Lambda_{ba} (\mathbf{x}_a - \boldsymbol{\mu}_a) - \frac{1}{2}(\mathbf{x}_b - \boldsymbol{\mu}_b)^T \Lambda_{bb} (\mathbf{x}_b - \boldsymbol{\mu}_b) \end{aligned}$$

## Marginal Gaussian distributions (cont.)

$$\begin{aligned} & \frac{1}{2} \left( \Lambda_{bb} \mu_b - \Lambda_{ba} (\mathbf{x}_a - \mu_a) \right)^T \Lambda_{bb}^{-1} \left( \Lambda_{bb} \mu_b - \Lambda_{ba} (\mathbf{x}_a - \mu_a) \right) \\ & - \frac{1}{2} \mathbf{x}_a^T \Lambda_{aa} \mathbf{x}_a + \mathbf{x}_a^T (\Lambda_{aa} \mu_a + \Lambda_{ab} \mu_b) + \text{const} \\ & = -\frac{1}{2} \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mathbf{x}_a \\ & \quad + \mathbf{x}_a^T (\Lambda_{aa} - \Lambda_{ab} \Lambda_{bb}^{-1} \Lambda_{ba}) \mu_a + \text{const} \quad (44) \end{aligned}$$

where 'const' denotes quantities independent of  $\mathbf{x}_a$

## Marginal Gaussian distributions (cont.)

Again, by comparison with the general exponent in a Gaussian

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const}$$

- ▶ The covariance of the marginal distribution  $p(\mathbf{x}_a)$  is given by

$$\boldsymbol{\Sigma}_a = (\boldsymbol{\Lambda}_{aa} - \boldsymbol{\Lambda}_{ab}\boldsymbol{\Lambda}_{bb}^{-1}\boldsymbol{\Lambda}_{ba})^{-1} \quad (45)$$

- ▶ The mean of the marginal distribution  $p(\mathbf{x}_a)$  is given by

$$\boldsymbol{\mu}_a = \boldsymbol{\Sigma}_a(\boldsymbol{\Lambda}_{aa} - \boldsymbol{\Lambda}_{ab}\boldsymbol{\Lambda}_{bb}^{-1}\boldsymbol{\Lambda}_{ba})^{-1}\boldsymbol{\mu}_a \quad (46)$$

Covariance is written in terms of partitioned precision matrix  $\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$

We can write this in terms of partitioned covariance matrix  $\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$

## Marginal Gaussian distributions (cont.)

$$\begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix} \quad (47)$$

Using again  $\begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \mathbf{M} & -\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}\mathbf{M} & \mathbf{D}^{-1} + \mathbf{D}^{-1}\mathbf{C}\mathbf{M}\mathbf{B}\mathbf{D}^{-1} \end{pmatrix}$  we have:

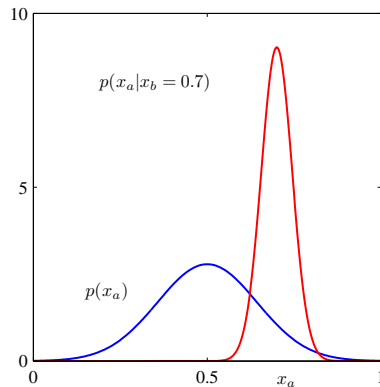
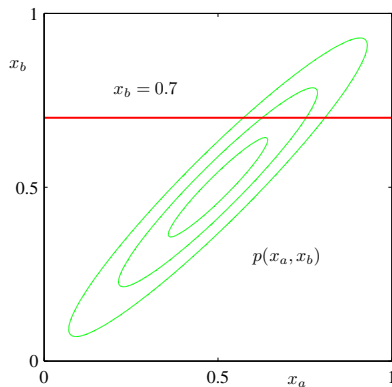
$$\Sigma_{aa} = (\Lambda_{aa} - \Lambda_{ab}\Lambda_{bb}^{-1}\Lambda_{ba}) \quad (48)$$

We obtain an intuitively satisfactory result that the marginal distribution  $p(\mathbf{x}_a)$

$$\mathbb{E}[\mathbf{x}_a] = \boldsymbol{\mu}_a \quad (49)$$

$$\text{cov}[\mathbf{x}_a] = \Sigma_{aa} \quad (50)$$

## Marginal Gaussian distributions (cont.)



# Bayes's theorem for Gaussian variables

## The Gaussian distribution



## Bayes' theorem for Gaussian variables

We considered a Gaussian  $p(\mathbf{x})$  in which we partitioned vector  $\mathbf{x}$  into two subvectors  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$  and then found expressions for

- ▶ the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_{a|b}, \boldsymbol{\Lambda}_{aa}^{-1})$
- ▶ the marginal distribution  $p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$

We also noted that the mean  $\boldsymbol{\mu}_{a|b}$  of the conditional distribution  $p(\mathbf{x}_a|\mathbf{x}_b)$  is a linear function of  $\mathbf{x}_b$ ,  $\boldsymbol{\mu}_{a|b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$

A Gaussian marginal distribution  $p(\mathbf{x})$ , Gaussian conditional distribution  $p(\mathbf{y}|\mathbf{x})$

- ▶  $p(\mathbf{y}|\mathbf{x})$  with mean a linear function of  $\mathbf{x}$  and a covariance independent of  $\mathbf{x}$

We wish to find:

- ▶ the marginal distribution  $p(\mathbf{y})$
- ▶ the conditional distribution  $p(\mathbf{x}|\mathbf{y})$

## Bayes' theorem for Gaussian variables (cont.)

We take the marginal and conditional distributions to be

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \quad (51)$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{Ax} + \mathbf{b}, \mathbf{L}^{-1}) \quad (52)$$

- ▶  $\boldsymbol{\mu}$ ,  $\mathbf{A}$  and  $\mathbf{b}$  are parameters governing the means
- ▶  $\boldsymbol{\Lambda}$  and  $\mathbf{L}$  are precision matrices

If  $\mathbf{x}$  is  $M$ -dimensional and  $\mathbf{y}$  is  $D$ -dimensional, then  $\mathbf{A}$  is  $D \times M$

## Bayes' theorem for Gaussian variables (cont.)

First we find an expression for the joint distribution over  $\mathbf{x}$  and  $\mathbf{y}$ , by defining

$$\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \quad (53)$$

and considering the log of the joint distribution  $p(\mathbf{z}) = p(\mathbf{x}, \mathbf{y}) = p(\mathbf{y}|\mathbf{x})p(\mathbf{x})$

$$\begin{aligned} \ln p(\mathbf{z}) &= \ln p(\mathbf{x}) + \ln p(\mathbf{y}|\mathbf{x}) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) \\ &\quad -\frac{1}{2}(\mathbf{y} - \mathbf{Ax} - \mathbf{b})^T \mathbf{L}(\mathbf{y} - \mathbf{Ax} - \mathbf{b}) \\ &\quad + \text{const} \end{aligned} \quad (54)$$

with 'const' denoting terms independent of  $\mathbf{x}$  and  $\mathbf{y}$

- ▶ It is again a quadratic function of the components of  $\mathbf{z}$
- ▶ Thus  $p(\mathbf{z})$  is again a Gaussian distribution

## Bayes' theorem for Gaussian variables (cont.)

To find the precision of the Gaussian  $p(\mathbf{z})$ , consider the second order terms of

$$\begin{aligned}\ln p(\mathbf{z}) &= \ln p(\mathbf{x}) + \ln p(\mathbf{y}|\mathbf{x}) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{Ax} - \mathbf{b})^T \mathbf{L}(\mathbf{y} - \mathbf{Ax} - \mathbf{b}) + \text{const}\end{aligned}$$

which can be written as

$$\begin{aligned}-\frac{1}{2}\mathbf{x}^T(\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{LA})\mathbf{x} - \frac{1}{2}\mathbf{y}^T \mathbf{Ay} + \frac{1}{2}\mathbf{y}^T \mathbf{LAx} + \frac{1}{2}\mathbf{x}^T \mathbf{A}^T \mathbf{Ly} \\ = \frac{1}{2} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{LA} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{LA} & \mathbf{L} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = -\frac{1}{2}\mathbf{z}^T \mathbf{Rz}\end{aligned}$$

So the Gaussian distribution over  $\mathbf{z}$  has precision (inverse covariance) matrix  $\mathbf{R}$

## Bayes' theorem for Gaussian variables (cont.)

The covariance matrix is found by taking the inverse of the precision matrix

$$\text{cov}[\mathbf{z}] = \mathbf{R}^{-1} = \begin{pmatrix} \mathbf{\Lambda}^{-1} & \mathbf{\Lambda}^{-1} \mathbf{A}^T \\ \mathbf{A} \mathbf{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A} \mathbf{\Lambda}^{-1} \mathbf{A}^T \end{pmatrix} \quad (55)$$

- The matrix inversion formula seen earlier is used ( $\star$ )

## Bayes' theorem for Gaussian variables (cont.)

To find the mean of the Gaussian  $p(\mathbf{z})$ , consider the linear terms of

$$\begin{aligned}\ln p(\mathbf{z}) &= \ln p(\mathbf{x}) + \ln p(\mathbf{y}|\mathbf{x}) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{Ax} - \mathbf{b})^T \mathbf{L}(\mathbf{y} - \mathbf{Ax} - \mathbf{b}) + \text{const}\end{aligned}$$

which can be written as

$$\mathbf{x}^T \boldsymbol{\Lambda} \boldsymbol{\mu} - \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{b} + \mathbf{y}^T \mathbf{L} \mathbf{b} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} \quad (56)$$

## Bayes' theorem for Gaussian variables (cont.)

Using our earlier result in Equation 29<sup>4</sup>, we find that the mean of  $\mathbf{z}$  is

$$\mathbb{E}[\mathbf{z}] = \mathbf{R}^{-1} \begin{pmatrix} \Lambda \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} \quad (57)$$

We can now make use of Equation 55<sup>5</sup> for the covariance of  $\mathbf{z}$  to get

$$\mathbb{E}[\mathbf{z}] = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \end{pmatrix} \quad (58)$$

---

<sup>4</sup>  $-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \text{const}$

<sup>5</sup>  $\text{cov}[\mathbf{z}] = \mathbf{R}^{-1} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1}\mathbf{A}^T \\ -\mathbf{A}\Lambda^{-1} & \mathbf{L}^{-1} + \mathbf{A}\Lambda^{-1}\mathbf{A}^T \end{pmatrix}$

## Bayes' theorem for Gaussian variables (cont.)

We turn now our attention to the marginal distribution  $p(\mathbf{y})$

- ▶ which we have marginalised over  $\mathbf{x}$

The marginal distribution over a subset of components of the random vector takes a simple form when expressed in terms of partitioned covariance matrix

$$\mathbb{E}[\mathbf{x}_a] = \boldsymbol{\mu}_a \quad \text{and} \quad \text{cov}[\mathbf{x}_a] = \boldsymbol{\Sigma}_{aa}$$

Use  $\text{cov}[\mathbf{z}] = \mathbf{R}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}^{-1} & \boldsymbol{\Lambda}^{-1}\mathbf{L} \\ \mathbf{A}\boldsymbol{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T \end{pmatrix}$  and  $\mathbb{E}[\mathbf{z}] = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{pmatrix}$  then

$$\mathbb{E}[\mathbf{y}] = \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \tag{59}$$

$$\text{cov}[\mathbf{y}] = \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T \tag{60}$$

are the searched mean and covariance of the marginal distribution  $p(\mathbf{y})$



## Bayes' theorem for Gaussian variables (cont.)

Finally, we seek an expression for the conditional  $p(\mathbf{x}|\mathbf{y})$

The conditional distribution is easier in terms of partitioned precision matrix

$$\Sigma_{a|b} = \Lambda_{aa}^{-1} \quad \text{and} \quad \mu_{a|b} = \mu_a - \Sigma_{aa}^{-1} \Sigma_{ab} (\mathbf{x}_b - \mu_b)$$

Use  $\text{cov}[\mathbf{z}] = \mathbf{R}^{-1} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} \mathbf{L} \\ \mathbf{A} \Lambda^{-1} & \mathbf{L}^{-1} + \mathbf{A} \Lambda^{-1} \mathbf{A}^T \end{pmatrix}$  and  $\mathbb{E}[\mathbf{z}] = \begin{pmatrix} \mu \\ \mathbf{A} \mu + \mathbf{b} \end{pmatrix}$  then

$$\mathbb{E}[\mathbf{x}|\mathbf{y}] = (\Lambda + \mathbf{A}^T \mathbf{L} \mathbf{A}^{-1}) (\mathbf{A}^T \mathbf{L} (\mathbf{y} - \mathbf{b}) + \Lambda \mu) \quad (61)$$

$$\text{cov}[\mathbf{x}|\mathbf{y}] = (\Lambda + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1} \quad (62)$$

are the searched mean and covariance of the conditional distribution  $p(\mathbf{x}|\mathbf{y})$

## Bayes' theorem for Gaussian variables (cont.)

So what about the Bayes' theorem?

The evaluation of the conditional  $p(\mathbf{x}|\mathbf{y})$  is an example of Bayes' theorem

- ▶ We can interpret the distribution  $p(\mathbf{x})$  as a prior over  $\mathbf{x}$

The conditional distribution  $p(\mathbf{x}|\mathbf{y})$  is the corresponding posterior over  $\mathbf{x}$

- ▶ If  $\mathbf{y}$  is observed

Having found the marginal and conditional distributions, we effectively expressed the joint distribution  $p(\mathbf{z}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x})$  in the form  $p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$

# Maximum likelihood for the Gaussian

## The Gaussian distribution

## Maximum likelihood for the Gaussian

Given a data set  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$  in which the observations  $\{\mathbf{x}_n\}$  are assumed to be drawn independently from a multivariate Gaussian distribution, we can estimate the parameters of the distribution by maximum likelihood

- ▶ The log likelihood function is

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \quad (63)$$

The likelihood function depends on the data set on thru terms

$$\sum_{n=1}^N \mathbf{x}_n \quad \text{and} \quad \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \quad (64)$$

These are the **sufficient statistics** for the Gaussian distribution

## Maximum likelihood for the Gaussian (cont.)

The derivative of log likelihood with respect to  $\mu$  can be written as <sup>6</sup>

$$\frac{\partial}{\partial \mu} \ln p(\mathbf{X}|\mu, \Sigma) = \sum_{n=1}^N \Sigma^{-1}(\mathbf{x}_n - \mu) \quad (65)$$

Set to zero, it gets us the maximum likelihood estimate of the mean

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \quad (66)$$

which is, as expected, the mean of the observed set of data points

---

<sup>6</sup>We used  $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{a}) = \frac{\partial}{\partial \mathbf{x}}(\mathbf{a}^T \mathbf{x}) = \mathbf{a}$

## Maximum likelihood for the Gaussian (cont.)

Maximisation of the log likelihood function  $\ln p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  wrt to  $\boldsymbol{\Sigma}$  is tough (★)  
As expected, the result takes the form

$$\boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \quad (67)$$

It involves  $\boldsymbol{\mu}_{ML}$  as a result of the joint maximisation wrt  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , whereas the solution for  $\boldsymbol{\mu}_{ML}$  does not depend on  $\boldsymbol{\Sigma}_{ML}$  (we first solve for  $\boldsymbol{\mu}_{ML}$  and then  $\boldsymbol{\Sigma}_{ML}$ )

## Maximum likelihood for the Gaussian (cont.)

If we evaluate the expectations of the maximum likelihood solutions under the true distribution, we obtain the following results

$$\mathbb{E}[\mu_{ML}] = \mu \quad (68)$$

$$\mathbb{E}[\Sigma_{ML}] = \frac{N-1}{N}\Sigma \quad (69)$$

- ▶ The expectation of the maximum likelihood estimate for the mean is equal to the true mean, unbiased
- ▶ The maximum likelihood estimate for the covariance underestimates the true covariance, it is biased

We correct the bias, by defining the estimator  $\tilde{\Sigma}$

$$\tilde{\Sigma} = \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \mu_{ML})(\mathbf{x}_n - \mu_{ML})^T \quad (70)$$

which has an expectation that equals the true  $\Sigma$

# Bayesian inference for the Gaussian

## The Gaussian distribution



## Bayesian inference for the Gaussian

The maximum likelihood framework gives point estimates for  $\mu$  and  $\Sigma$

Now we develop a Bayesian treatment by introducing prior distributions

- ▶ over these parameters

We start simple, with a single Gaussian random variable  $x$

- ▶ We will suppose that the variance  $\sigma^2$  is known
- ▶ We consider the task of inferring the mean  $\mu$

We are also given a set of  $N$  observations  $\mathbf{x} = \{x_1, \dots, x_n\}$

## Bayesian inference for the Gaussian (cont.)

The likelihood function can be viewed as a function of  $\mu$  and takes the form

$$p(\mathbf{x}|\mu) = \prod_{n=1}^N p(x_n|\mu) = \frac{1}{(2\pi\sigma^2)^N} \exp\left(-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2\right) \quad (71)$$

It takes the form of the exponential of a quadratic form in  $\mu$  and if we choose a Gaussian prior  $p(\mu)$ , it is a conjugate distribution for the likelihood function

- ▶ The posterior is a product of two exponentials of quadratic functions of  $\mu$
- ▶ Thus, it is also a Gaussian

## Bayesian inference for the Gaussian (cont.)

We take our prior distribution to be

$$p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2) \quad (72)$$

The posterior distribution is given by

$$p(\mu|\mathbf{x}) \propto p(\mathbf{x}|\mu)p(\mu) \quad (73)$$

Some manipulations ( $\star$ ) involving completing the square in the exponent allow to show that the posterior distribution is given by

$$p(\mu|\mathbf{x}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2) \quad (74)$$

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2}\mu_0 + \frac{N\sigma^2}{N\sigma^2 + \sigma^2}\mu_{ML} \quad (75)$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \quad (76)$$

$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$  is the maximum likelihood estimate for  $\mu$ , the sample mean

## Bayesian inference for the Gaussian (cont.)

Note that the mean of the posterior distribution  $p(\mu|\mathbf{x})$  is a compromise between the prior mean  $\mu_0$  and the maximum likelihood solution  $\mu_{ML}$

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma_0^2}\mu_0 + \frac{N\sigma^2}{N\sigma^2 + \sigma^2}\mu_{ML}$$

- ▶ It reduces to the prior mean, if the number of observed data points  $N = 0$
- ▶ For  $N \rightarrow \infty$ , the posterior mean equals the maximum likelihood solution

## Bayesian inference for the Gaussian (cont.)

Consider the result for the variance of the posterior distribution  $p(\mu|\mathbf{x})$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

It is most naturally expressed in terms of inverse variance, or precision

Precisions are additive, so the precision of the posterior is given by

- ▶ the precision of the prior, plus one contribution of the data precision from each of the observed data points

## Bayesian inference for the Gaussian (cont.)

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

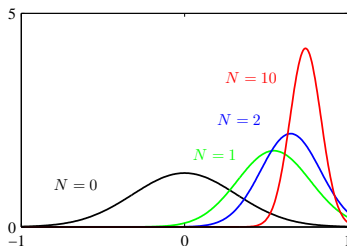
As we increase the number of observed data points, the precision steadily increases, giving a posterior distribution with steadily decreasing variance

- ▶ With no observed data points, we have the prior variance
- ▶ If  $N \rightarrow \infty$ , variance  $\sigma_N^2 \rightarrow 0$  and the posterior distribution becomes infinitely peaked around the ML solution  $\mu_{ML}$

## Bayesian inference for the Gaussian (cont.)

Bayesian inference for the mean  $\mu$  of a Gaussian with known variance  $\sigma$

- ▶ Data from a Gaussian of mean 0.8 and variance 0.1



The prior is chosen to have mean 0

In both prior and likelihood function, the variance is set to the true value

- ▶ the prior distribution over  $\mu$ , the curve  $N = 0$ , itself Gaussian
- ▶ the posterior distribution for increasing numbers  $N$  of points

We see that the maximum likelihood result of a point estimate for  $\mu$  is recovered precisely from the Bayesian formalism in the limit  $N \rightarrow \infty$

For finite  $N$ , in the limit  $\sigma_0^2 \rightarrow \infty$  in which the prior has infinite variance then the posterior mean is the ML result, while the posterior variance is  $\sigma_N^2 = \sigma^2/N$

## Bayesian inference for the Gaussian (cont.)

The analysis of Bayesian inference for the mean of a  $D$ -dimensional Gaussian random variable  $\mathbf{x}$  with known covariance and unknown mean is straightforward

(★)



## Bayesian inference for the Gaussian (cont.)

So far, we have assumed that the variance of the Gaussian distribution over the data is known and our goal is to infer the mean

- ▶ Let us suppose that the mean is known and we wish to infer the variance

Calculations are simpler, if we choose a conjugate form for the prior distribution

- ▶ The likelihood function for  $\lambda = 1/\sigma^2$  is

$$p(\mathbf{x}|\lambda) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \lambda^{-1}) \propto \lambda^{N/2} \exp\left(-\frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2\right) \quad (77)$$

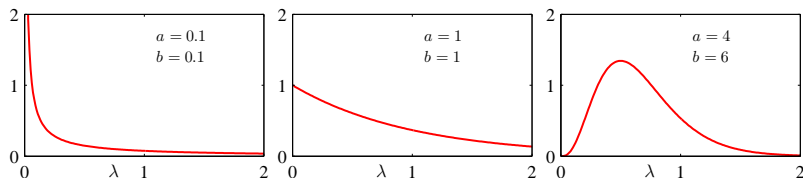
The corresponding conjugate prior should therefore be proportional to the product of a power of  $\lambda$  and the exponential of a linear function of  $\lambda$

## Bayesian inference for the Gaussian (cont.)

This corresponds to the **gamma distribution**, which is defined as

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) \quad (78)$$

The gamma function  $\Gamma(\cdot)$  assures a correct normalisation



The integral of the gamma distribution is finite if  $a > 0$

The distribution itself is finite for  $a \geq 1$

The mean and variance of the gamma distribution are  $\begin{cases} \mathbb{E}[\lambda] = a/b \\ \text{var}[\lambda] = a/b^2 \end{cases}$

## Bayesian inference for the Gaussian (cont.)

Consider a prior distribution  $\text{Gam}(\lambda|a_0, b_0)$ , multiply by the likelihood function

- Obtain a posterior distribution

$$p(\lambda|\mathbf{x}) \propto \lambda^{a_0-1} \lambda^{N/2} \exp\left(-b_0\lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2\right) \quad (79)$$

- Recognise the Gamma distribution  $\text{Gam}(\lambda|a_N, b_N)$ , with

$$a_N = a_0 + \frac{N}{2} \quad (80)$$

$$b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{ML}^2 \quad (81)$$

where  $\sigma_{ML}^2$  is the maximum likelihood estimator of the variance

## Bayesian inference for the Gaussian (cont.)

Note that there is no need to keep track of the normalisation constants in the prior and in the likelihood function in

$$p(\lambda|\mathbf{x}) \propto \lambda^{a_0-1} \lambda^{N/2} \exp\left(-b_0\lambda - \frac{\lambda}{2} \sum_{n=1}^N (x_n - \mu)^2\right)$$

Because, if required, the correct coefficient can be found at the end using the normalised form of the gamma distro

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

## Bayesian inference for the Gaussian (cont.)

The effect of observing  $N$  points is to increase the value of coefficient  $a$  by  $N/2$

$$a_N = a_0 + \frac{N}{2}$$

We can interpret  $a_0$  in the prior in terms of  $2a_0$  'effective' prior observations

Similarly, the  $N$  points contribute  $N\sigma_{ML}^2/2$  to parameter  $b$ ,  $\sigma_{ML}^2$  is the variance

$$b_N = b_0 + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2 = b_0 + \frac{N}{2} \sigma_{ML}^2$$

We interpret  $b_0$  in the prior as arising from the  $2a_0$  'effective' prior observations having variance  $2b_0/(2a_0) = b_0/a_0$

## Bayesian inference for the Gaussian (cont.)

Suppose now that both mean and variance are unknown

To find a conjugate prior, consider the dependence of the likelihood on  $\mu$  and  $\lambda$

$$\begin{aligned} p(\mathbf{x}|\mu, \lambda) &= \prod_{n=1}^N \left( \frac{\lambda}{2\pi} \right) \exp \left( -\frac{\lambda}{2} (x_n - \mu)^2 \right) \\ &\propto \left( \lambda^{1/2} \exp \left( -\frac{\lambda \mu^2}{2} \right) \right)^N \exp \left( \lambda \mu \sum_{n=1}^N x_n - \frac{\lambda}{2} \sum_{n=1}^N x_n^2 \right) \quad (82) \end{aligned}$$

We wish to identify the prior distribution  $p(\mu, \lambda)$  with the same functional dependence on  $\mu$  and  $\lambda$  as the likelihood function above

## Bayesian inference for the Gaussian (cont.)

$$\begin{aligned} p(\mu, \lambda) &\propto \left( \lambda^{1/2} \exp \left( -\frac{\lambda \mu^2}{2} \right) \right)^\beta \exp(c\lambda\mu - d\lambda) \\ &= \exp \left( -\frac{\beta\lambda}{2} (\mu - c/\beta)^2 \right) \lambda^{\beta/2} \exp \left( -\left( d - \frac{c^2}{2\beta} \right) \lambda \right) \quad (83) \end{aligned}$$

where  $c$ ,  $d$  and  $\beta$  are constant

Because  $p(\mu, \lambda) = p(\mu|\lambda)p(\lambda)$ , we can find  $p(\mu|\lambda)$  and  $p(\lambda)$  by inspection

- ▶  $p(\mu|\lambda)$  is a Gaussian whose precision is a linear function of  $\lambda$
- ▶  $p(\lambda)$  is a gamma distribution

## Bayesian inference for the Gaussian (cont.)

The normalised prior takes the form

$$p(\mu, \lambda) = \mathcal{N}(\mu | \mu_0, (\beta\lambda)^{-1}) \text{Gam}(\lambda | a, b) \quad (84)$$

where we defined the new constants

- ▶  $\mu_0 = c/\beta$
- ▶  $a = (1 + \beta)/2$
- ▶  $b = d - c^2/2\beta$

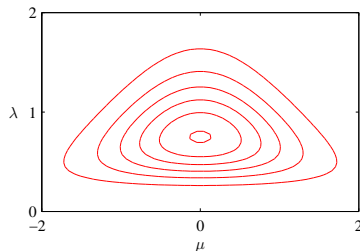
This distribution is called **normal-gamma** or **Gaussian-gamma distribution**

- ▶ It is not simply the product of an independent Gaussian prior over  $\mu$  and gamma prior over  $\lambda$ , because the precision of  $\mu$  is a linear function of  $\lambda$
- ▶ Even if we choose a prior in which  $\mu$  and  $\lambda$  are independent, the posterior distribution would exhibit a coupling between precision of  $\mu$  and value of  $\lambda$



## Bayesian inference for the Gaussian (cont.)

Contour plot of the normal-gamma distribution



Parameter values  $\mu_0 = 0$ ,  $\beta = 2$ ,  $a = 5$  and  $b = 6$

## Bayesian inference for the Gaussian (cont.)

In the case of a  $D$ -variate Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$  over a variable  $\mathbf{x}$

- ▶ the conjugate prior distribution for the mean  $\boldsymbol{\mu}$  assuming a known precision is again a Gaussian distribution
- ▶ the conjugate prior distribution for the precision matrix  $\boldsymbol{\Lambda}$  assuming a known mean is a **Wishart distribution**

$$\mathbf{W}(\boldsymbol{\Lambda}|\mathbf{W}, \nu) = B|\boldsymbol{\Lambda}|^{(\nu-D-1)/2} \exp\left(-\frac{1}{2}\text{Tr}(\mathbf{W}^{-1}\boldsymbol{\Lambda})\right) \quad (85)$$

where  $\nu$  is the **number of degrees of freedom** of the distribution,  $\mathbf{W}$  is a  $D \times D$  scale matrix and  $\text{Tr}(\cdot)$  denotes the trace of a matrix

The normalisation constant  $B$  is given by

$$B(\mathbf{W}, \nu) = |\mathbf{W}|^{-\nu/2} \left(2^{\nu D/2} \pi^{D(D-1)/4} \prod_{i=1}^D \Gamma\left(\frac{\nu+1-i}{2}\right)\right)^{-1} \quad (86)$$

## Bayesian inference for the Gaussian (cont.)

- ▶ the conjugate prior distribution assuming both mean  $\mu$  and precision matrix  $\Lambda$  unknown is a **normal-Wishart** or **Gaussian-Wishart distribution**

$$p(\mu, \Lambda | \mu_0, \beta, \mathbf{W}, \nu) = \mathcal{N}(\mu | \mu_0, (\beta \Lambda^{-1})) \mathcal{W}(\Lambda | \mathbf{W}, \nu) \quad (87)$$

# **Student's t-distribution**

## **The Gaussian distribution**

## Student's t-distribution

The conjugate prior for the precision of a Gaussian is a gamma distribution

If we have a univariate Gaussian  $\mathcal{N}(x|\mu, \tau^{-1})$  together with a gamma prior  $\text{Gam}(\tau|a, b)$  and we integrate out the precision, we obtain the marginal distribution of  $x$  in the form

$$\begin{aligned} p(x|\mu, a, b) &= \int_0^{+\infty} \mathcal{N}(x|\mu, \tau^{-1}) \text{Gam}(\tau|a, b) d\tau \\ &= \int_0^{+\infty} \frac{b^a e^{-b\tau} \tau^{a-1}}{\Gamma(a)} \left(\frac{\tau}{2\pi}\right)^{1/2} \exp\left(-\frac{\tau}{2}(x-\mu)^2\right) d\tau \\ &= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \left(b + \frac{(x-\mu)^2}{2}\right)^{-a-1/2} \Gamma(a+1/2) \quad (88) \end{aligned}$$

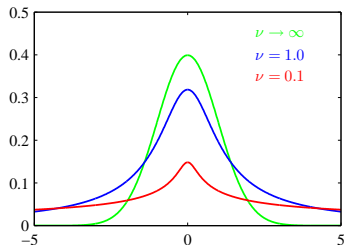
where we can make the change of variables  $z = \tau(b + (x - \mu)^2/2)$

## Student's t-distribution (cont.)

By convention, we define new parameters  $\nu = 2a$  and  $\lambda = a/b$  to get

$$p(x|\mu, a, b) = \text{Stu}(x|\mu, \lambda, \nu) = \frac{\Gamma(\frac{\nu}{2} + \frac{1}{2})}{\Gamma(\frac{\nu}{2})} \left(\frac{\lambda}{\pi\nu}\right)^{1/2} \left(1 + \frac{\lambda(x - \mu)^2}{\nu}\right)^{-\nu/2 - 1/2} \quad (89)$$

which is known as **Student's t-distribution** with parameter  $\lambda$  being the **precision** and parameter  $\nu$  being called the **degree of freedom**



Student's t-distribution for  $\mu = 0$ ,  $\lambda = 1$  and for various values of  $\nu$

The limit  $\nu \rightarrow \infty$  equals a Gaussian distribution with mean  $\mu$  and precision  $\lambda$

For  $\nu = 1$ , the t-distribution reduces to the Cauchy distribution

## Student's t-distribution (cont.)

$$\int_0^{+\infty} \mathcal{N}(x|\mu, \tau^{-1}) \text{Gam}(\tau|a, b) d\tau$$

The Student's t-distribution is obtained by adding up an infinite number of Gaussian distributions having the same mean but different precisions

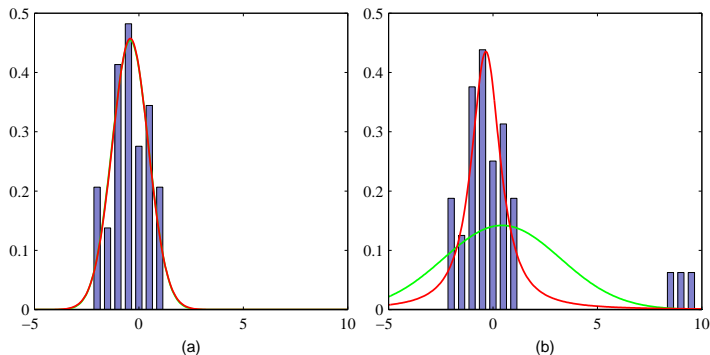
- ▶ This can be interpreted as an infinite mixture of Gaussians

The result is a distribution that in general has longer 'tails' than a Gaussian

- ▶ the t-distribution has an important property called robustness
- ▶ it is less sensitive than the Gaussian to the presence of outliers

## Student's t-distribution (cont.)

Maximum likelihood fits for a Gaussian (green) and a t-distribution (red)



The Gaussian is strongly distorted by the presence of outlying points



## Student's t-distribution (cont.)

Substituting parameters  $\nu = 2a$ ,  $\lambda = 1/b$  and  $\eta = \tau b/a$ , the t-distribution is

$$\text{St}(x|\mu, \lambda, \nu) = \int_0^{+\infty} \mathcal{N}(x|\mu, (\eta\lambda)^{-1}) \text{Gam}(\frac{\eta/\nu}{2}, \nu/2) d\eta \quad (90)$$

$$= \frac{\Gamma(\frac{\nu}{2} + \frac{1}{2})}{\Gamma(\frac{\nu}{2})} \left(\frac{\lambda}{\pi\nu}\right)^{1/2} \left(1 + \frac{\lambda(x - \mu)^2}{\nu}\right)^{-\nu/2 - 1/2} \quad (91)$$

We can also generalise this to a  $D$ -dimensional Gaussian  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda})$  to get

$$\text{St}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}, \nu) = \int_0^{+\infty} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, (\eta\boldsymbol{\Lambda})^{-1}) \text{Gam}(\frac{\eta/\nu}{2}, \nu/2) d\eta \quad (92)$$

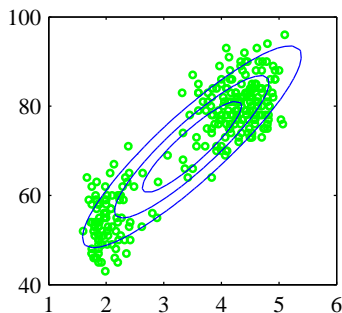
$$= \frac{\Gamma(\frac{D}{2} + \frac{\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\boldsymbol{\Lambda}|^{1/2}}{(\pi\nu)^{D/2}} \left(1 + \frac{(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda} (\mathbf{x} - \boldsymbol{\mu})}{\nu}\right)^{-D/2 - \nu/2} \quad (93)$$

# Mixtures of Gaussians

## The Gaussian distribution

## Mixtures of Gaussians

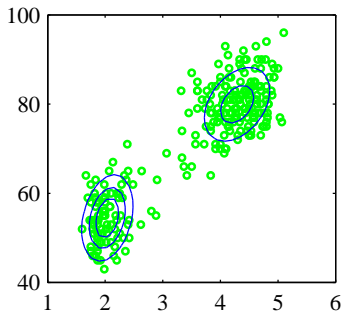
While the Gaussian distribution has some important analytical properties, it suffers from significant limitations when it comes to modelling real data sets



A Gaussian distribution fitted to the data using maximum likelihood

This distribution fails to capture the two clumps in the data and places much of its probability mass in the central region between the clumps where the data are relatively sparse

## Mixtures of Gaussians (cont.)



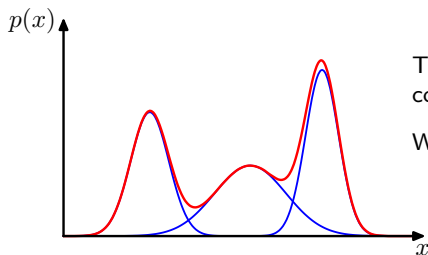
A linear combination of two Gaussians fitted using maximum likelihood

A linear superposition gives a better characterisation of the data set

Superpositions, formed by taking linear combinations of basic distributions such as Gaussians, can be formulated as probabilistic models (**Mixture distributions**)

## Mixtures of Gaussians (cont.)

A Gaussian mixture distribution in one dimension



Three Gaussians (blue, scaled by a coefficient) and their sum (red)

We can get very complex densities

By using a sufficient number of Gaussians, and by adjusting their means and covariances as well as the coefficients in the linear combination, almost any continuous density can be approximated to arbitrary accuracy

## Mixtures of Gaussians (cont.)

We consider a linear superposition of  $K$  Gaussian densities of the form

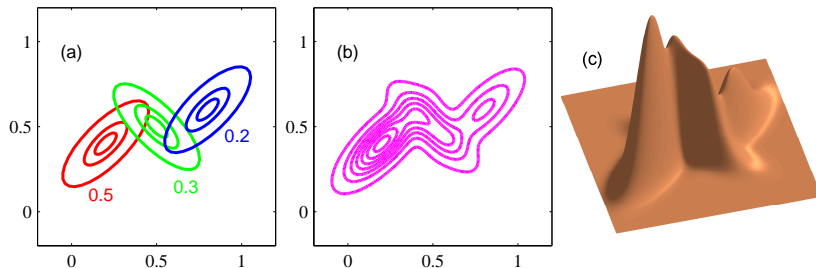
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (94)$$

Such a model is a **mixture of Gaussians**, each Gaussian density  $\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  is a **component** of the mixture, with its own mean  $\boldsymbol{\mu}_k$  and covariance  $\boldsymbol{\Sigma}_k$

The parameters  $\pi_k$  in  $p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  are **mixing coefficients**

## Mixtures of Gaussians (cont.)

A Gaussian mixture distribution in two dimensions



- ▶ The contours at constant density for each of the mixture components
- ▶ The contours at constant density of the mixture distribution  $p(\mathbf{x})$
- ▶ The surface plot of the mixture distribution  $p(\mathbf{x})$

The numbers in the first plot are the mixing coefficients

## Mixtures of Gaussians (cont.)

If we integrate both sides of  $p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  with respect to  $\mathbf{x}$  and note that both  $p(\mathbf{x})$  and the Gaussian components are normalised, we get that

$$\sum_{k=1}^K \pi_k = 1 \quad (95)$$

The requirements that  $p(\mathbf{x}) \geq 0$  and  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \geq 0$  imply that  $0 \leq \pi_k \leq 1$



## Mixtures of Gaussians (cont.)

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$\pi_k \in [0, 1]$  and  $\sum \pi_k = 1$ : The mixing coefficients can be seen as probabilities

From sum and probabilities rules, the marginal density is given by

$$p(\mathbf{x}) = \sum_{k=1}^K p(k) p(\mathbf{x} | k) \quad (96)$$

By  $p(k) = \pi_k$ , we view coefficients as prior probabilities (of picking the  $k$ -th component) and by  $p(\mathbf{x} | k) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  the probability of  $\mathbf{x}$  conditioned on  $k$

## Mixtures of Gaussians (cont.)

An important role is played by posterior probabilities (or **responsibilities**)  $p(k|\mathbf{x})$

By the Bayes' theorem, the posterior probabilities can be written as

$$\begin{aligned}\gamma_k(\mathbf{x}) &\equiv p(k|\mathbf{x}) \\ &= \frac{p(k)p(\mathbf{x}|k)}{\sum_l p(l)p(\mathbf{x}|l)} \\ &= \frac{\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_l \pi_l \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}\end{aligned}\tag{97}$$

## Mixtures of Gaussians (cont.)

The Gaussian mixture distribution is governed by the parameters  $\pi$ ,  $\mu$  and  $\Sigma$

$$\pi \equiv \{\pi_1, \dots, \pi_K\}$$

$$\mu \equiv \{\mu_1, \dots, \mu_K\}$$

$$\Sigma \equiv \{\Sigma_1, \dots, \Sigma_K\}$$

One way to set the values of the parameters is to use (log) maximum likelihood

$$\ln p(\mathbf{X}|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left( \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right) \quad (98)$$

which is bad, very bad, because of the summation over  $k$  inside the logarithm

- ▶ No longer closed-form solution
- ▶ Iterative numerical optimisation
- ▶ Expectation maximisation (ok!)