

# Probability Distributions

Sadegh Eskandari

Department of Computer Science, University of Guilan

# Probability Distributions: Introduction

---

- **Remember:** Probability theory provides a consistent framework for the quantification and manipulation of uncertainty in data
- **Density Estimation:** given a finite set  $\mathbf{x}_1, \dots, \mathbf{x}_N$  of observations for a random variable  $\mathbf{x}$ , the goal is to model the probability distribution  $p(\mathbf{x})$ .
- We will assume that the data points are independent and identically distributed (iid).

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^N p(\mathbf{x}_n)$$

## Density Estimation

- **Parametric**
  - Selecting a common distribution and estimating the parameters for the density function from the data
  - binomial and multinomial distributions for discrete random variables
  - Gaussian distribution for continuous random variables.
  - Parameter estimation procedure: maximum likelihood, Bayesian method
- **Non-Parametric**
  - Histograms, Nearest-Neighbours, Kernels

# Binary Variables

---

## Bernoulli Distribution

- Consider a single binary random variable  $x \in \{0,1\}$
- For example,  $x$  might describe the outcome of flipping a coin, with  $x = 1$  representing ‘heads’, and  $x = 0$  representing ‘tails’.
- The probability of  $x = 1$  will be denoted by the parameter  $0 \leq \mu \leq 1$  so that:

$$p(x = 1|\mu) = \mu$$

- And hence:

$$p(x = 0|\mu) = 1 - \mu$$

- Therefore, the probability distribution over  $x$  can be written in the form:

$$\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x}$$

# Binary Variables

---

## Bernoulli Distribution

$$\begin{aligned}\mathbb{E}[x] &= \sum_x xp(x) \\ &= \sum_x x \text{Bern}(x|\mu) = 0 \times \text{Bern}(x=0|\mu) + 1 \times \text{Bern}(x=1|\mu) = \mu\end{aligned}$$

$$\begin{aligned}\text{var}(x) &= \mathbb{E}[(x - \mathbb{E}[x])^2] = \mathbb{E}[(x - \mu)^2] = \sum_x (x - \mu)^2 p(x) \\ &= (0 - \mu)^2 \text{Bern}(x=0|\mu) + (1 - \mu)^2 \text{Bern}(x=1|\mu) \\ &= \mu^2(1 - \mu) + (1 - \mu)^2 \mu = \mu(1 - \mu)\end{aligned}$$

# Binary Variables

---

## Bernoulli Distribution

- Now suppose we have a data set  $\mathcal{D} = \{x_1, \dots, x_N\}$  of observed values of  $x$
- We know that  $\mathcal{D}$  is derived from a Bernoulli distribution
- But, we do not know the parameter  $\mu$
- So, we want to estimate  $\mu$  using  $\mathcal{D}$
- The **maximum likelihood** approach:

$$\mu_{ML} = \arg \max_{\mu} p(\mathcal{D}|\mu) = \arg \max_{\mu} \prod_{n=1}^N p(x_n|\mu) = \arg \max_{\mu} \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n}$$

$$= \arg \max_{\mu} \ln \left( \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n} \right)$$

$$= \arg \max_{\mu} \underbrace{\left( \sum_{n=1}^N \{x_n \ln \mu + (1 - x_n) \ln(1 - \mu)\} \right)}_{f(\mu)}$$

$$\frac{\partial f}{\partial \mu} = 0$$

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n = \frac{m}{N}$$

the number of observations of  $x = 1$



# Binary Variables

---

## Bernoulli Distribution

- Now suppose we flip a coin, say, 3 times and happen to observe 3 heads.
- Then  $N = m = 3$  and  $\mu_{ML} = 1$
- Then the maximum likelihood result would predict that all future observations should give heads!
- In fact this is an extreme example of the **over-fitting** associated with maximum likelihood.
- **Solution:** Bayesian Approach

$$p(\mu|\mathcal{D}) = \frac{p(\mathcal{D}|\mu) \times p(\mu)}{p(\mathcal{D})}$$

Diagram illustrating the Bayesian approach equation:

- Likelihood** (yellow text) points to  $p(\mathcal{D}|\mu)$
- Prior Probability** (yellow text) points to  $p(\mu)$
- Posterior Probability** (yellow text) points to  $p(\mu|\mathcal{D})$

# Binary Variables

## Bernoulli Distribution

- Step1: Likelihood function

$$p(D|\mu) = \prod_{n=1}^N p(x_n|\mu) = \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n} = \mu^{\sum_{n=1}^N x_n} (1-\mu)^{N-\sum_{n=1}^N x_n} = \mu^m (1-\mu)^{N-m}$$

- Step2: Prior Probability

□ In this step we need to introduce a prior distribution  $p(\mu)$  over the parameter  $\mu$ . But how?

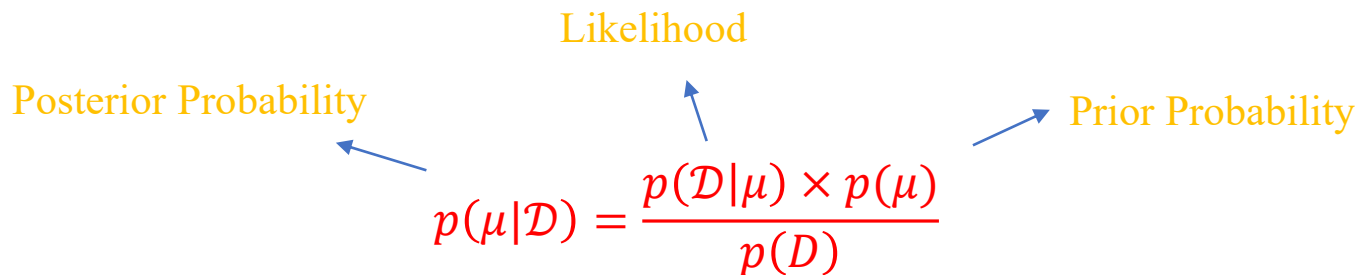
□ **Conjugacy**: As the likelihood function takes the form of powers of  $\mu$  and  $1-\mu$ , then, if we choose a prior to be proportional to powers of  $\mu$  and  $1-\mu$ , then the posterior distribution will have the same functional form as the prior.

□ Here we choose a prior called Beta distribution:

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

$$\mathbb{E}[\mu] = \frac{a}{a+b}$$

$$\text{var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$



$$\Gamma(x) \equiv \int_0^{\infty} u^{x-1} e^{-u} du$$

$$\Gamma(x+1) = x\Gamma(x)$$

$$\Gamma(1) = 1$$

**Proof: Homework**

$$\Gamma(x) = x!, \forall x \in \mathbb{Z}$$

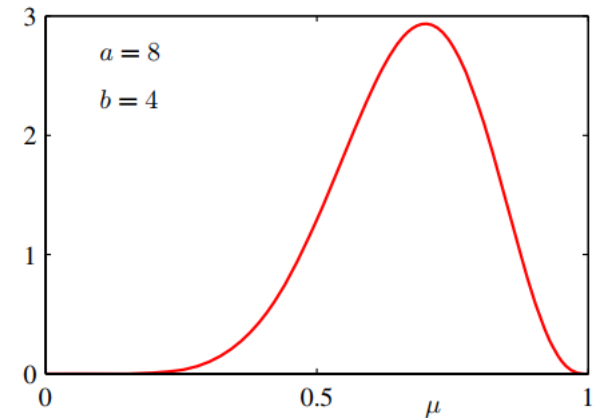
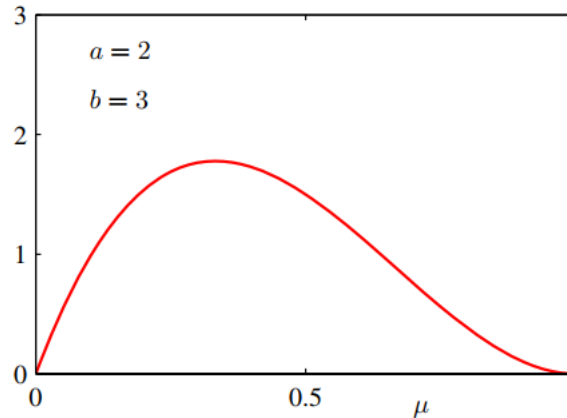
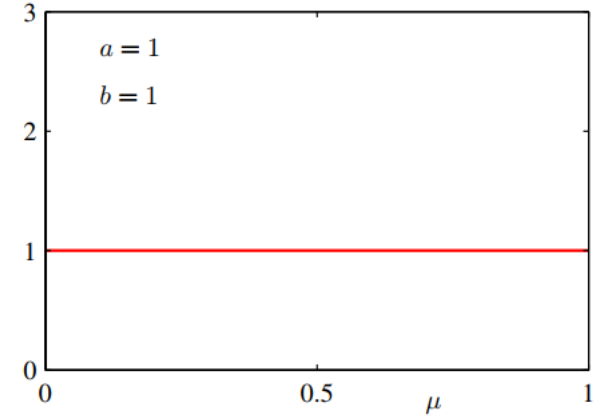
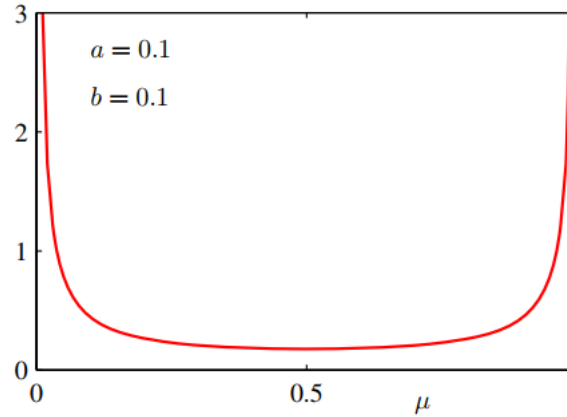
# Binary Variables

## Bernoulli Distribution

- Plots of  $Beta(\mu|a, b)$  given by for various values of the hyperparameters  $a$  and  $b$ .

### Step3: Posterior Probability

$$p(\mu|m, l, a, b) = \frac{\mu^m(1-\mu)^l \times Beta(\mu|a, b)}{\int_0^1 \mu^m(1-\mu)^l \times Beta(\mu|a, b)d\mu}$$
$$= \frac{\Gamma(a+m+b+l)}{\Gamma(a+m)\Gamma(b+l)} \mu^{a+m-1} (1-\mu)^{b+l-1}$$



- Sequential Learning:** The posterior distribution can act as the prior if we subsequently observe additional data (applicable for big data).



# Binary Variables

---

## Bernoulli Distribution

- **Question:** How to predict the outcome of the next trial of  $x$ , given the observed data set  $\mathcal{D}$ ?

$$\begin{aligned} & p(x = 1|\mathcal{D}) \\ &= \int_0^1 p(x = 1, \mu|\mathcal{D}) d\mu \\ &= \int_0^1 p(x = 1 | \mu, \mathcal{D}) p(\mu|\mathcal{D}) d\mu \\ &= \int_0^1 p(x = 1 | \mu) p(\mu|\mathcal{D}) d\mu \\ &= \int_0^1 \mu p(\mu|\mathcal{D}) d\mu \\ &= \mathbb{E}[\mu|\mathcal{D}] = \frac{m + a}{m + a + l + b} \end{aligned}$$

# Binary Variables

---

## Binomial Distribution

- The number  $m$  of observations of  $x = 1$ , given that the data set has size  $N$

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

- Where

$$\binom{N}{m} \equiv \frac{N!}{(N-m)! m!}$$

$$\mathbb{E}[m] = \sum_{m=0}^N m \text{Bin}(m|N, \mu) = N\mu$$

$$\text{var}[m] = N\mu(1 - \mu)$$

Proof: Homework

# Binary Variables

---

## Generalized Bernoulli Distribution

- Often, we encounter discrete variables that can take on one of  $K$  possible mutually exclusive states.

$$x \in \{s_1, s_2, \dots, s_K\} \quad \xrightarrow{\text{1-of-K encoding}} \quad \mathbf{x} \in \left\{ \begin{array}{l} (1, 0, 0, \dots, 0)^T \\ (0, 1, 0, \dots, 0)^T \\ \vdots \\ (0, 0, 0, \dots, 1)^T \end{array} \right\}$$

- The random vector  $\mathbf{x}$  can be described by  $K$  binary variables  $x_1, x_2, \dots, x_K$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_K \end{bmatrix} \quad \text{Such that} \quad \sum_{k=1}^K x_k = 1 \quad \xrightarrow{\begin{array}{l} p(x_k = 1 | \mu_k) = \mu_k \\ \mu_k \geq 0, \quad \sum_{k=1}^K \mu_k = 1 \end{array}} \quad p(\mathbf{x} | \boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad \boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_K \end{bmatrix}$$

# Binary Variables

## Generalized Bernoulli Distribution

- Now consider a data set  $\mathcal{D}$  of  $N$  independent observations  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N$ .
- So, we want to estimate the vector  $\boldsymbol{\mu}$  using  $\mathcal{D}$
- The **Maximum Likelihood** Approach

$$\boldsymbol{\mu}_{ML} = \arg \max_{\boldsymbol{\mu}} p(\mathcal{D}|\boldsymbol{\mu}) = \arg \max_{\boldsymbol{\mu}} \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\mu}) = \arg \max_{\mu_1, \dots, \mu_K} \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \arg \max_{\mu_1, \dots, \mu_K} \prod_{k=1}^K \mu_k^{\sum_n x_{nk}} = \arg \max_{\mu_1, \dots, \mu_K} \prod_{k=1}^K \mu_k^{m_k}$$

$$m_k = \sum_{n=1}^N x_{nk}$$

$$= \arg \max_{\mu_1, \dots, \mu_K} \ln \left( \prod_{k=1}^K \mu_k^{m_k} \right) = \arg \max_{\mu_1, \dots, \mu_K} \sum_{k=1}^K m_k \ln \mu_k$$

$$\text{s.t.} \quad \sum_{k=1}^K \mu_k = 1$$

Lagrange

$$\mu_k^{ML} = \frac{m_k}{N}$$

# Binary Variables

## Generalized Bernoulli Distribution

- The Bayesian Approach

Remember

$$p(\boldsymbol{\mu}|\mathcal{D}) = \frac{p(\mathcal{D}|\boldsymbol{\mu}) \times p(\boldsymbol{\mu})}{p(\mathcal{D})}$$

$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{m_k}$$

Conjugacy  
(Dirichlet Distribution)

$$p(\boldsymbol{\mu}) = \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

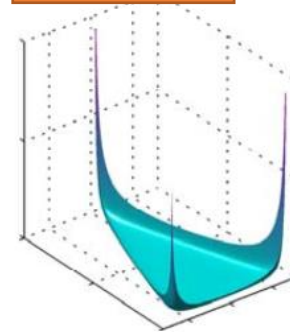
$$\alpha_0 = \sum_{k=1}^K \alpha_k$$

$$p(\boldsymbol{\mu}|\mathcal{D}) = \text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha} + \mathbf{m})$$

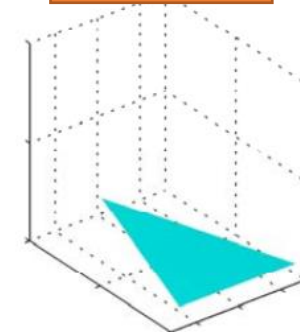
$$= \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_1 + m_1) \dots \Gamma(\alpha_K + m_N)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$

$$\mathbf{m} = (m_1, \dots, m_K)^T$$

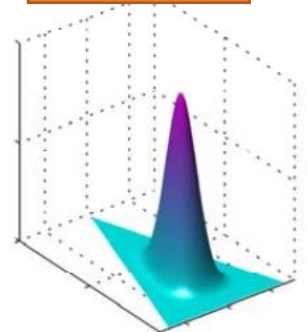
$\alpha_k = 0.1$



$\alpha_k = 1$



$\alpha_k = 10$



Dirichlet Distribution for  $K = 3$

# Binary Variables

---

## Multinomial Distribution

- The joint distribution of the quantities  $m_1, \dots, m_K$ , conditioned on the parameters  $\boldsymbol{\mu}$  and on the total number  $N$  of observations:

$$\text{Mult}(m_1, m_2, \dots, m_K | \boldsymbol{\mu}) = \binom{N}{m_1 m_2 \dots m_K} \prod_{k=1}^K \mu_k^{m_k}$$

- Where

$$\binom{N}{m_1 m_2 \dots m_K} \equiv \frac{N!}{m_1! m_2! \dots m_K!}$$

$$\sum_{k=1}^K m_k = N$$

# **Gaussian Distribution**

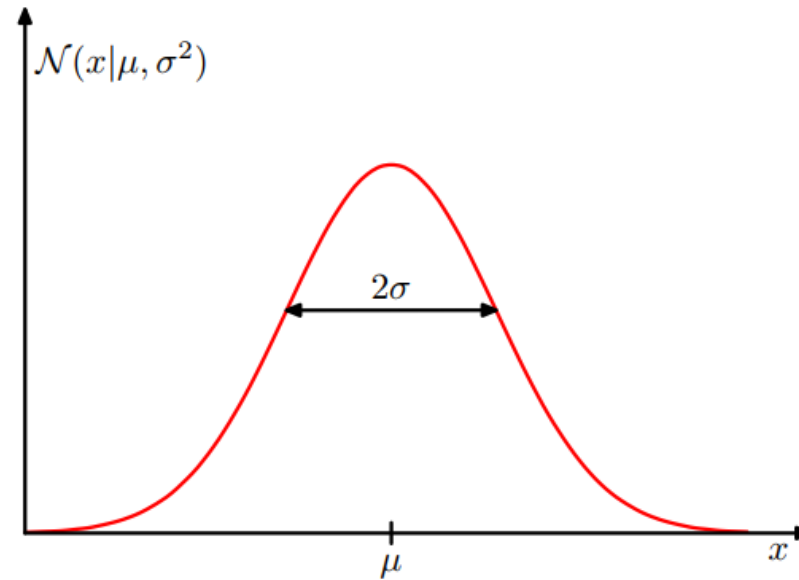
# The Gaussian distribution

---

- The **Gaussian or normal distribution** is the most important distribution for **continuous variables**.
- For the case of a single real-valued variable  $x$ , the Gaussian distribution is defined by

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

- $\mu$ : mean
  - $\sigma^2$ : variance
  - $\sigma$ : standard deviation
  - $\frac{1}{\sigma^2}$ : precision
- $\mathcal{N}(x|\mu, \sigma^2) \geq 0$
  - $\int_{-\infty}^{+\infty} \mathcal{N}(x|\mu, \sigma^2) = 1$



$$p(\mu - \sigma \leq x \leq \mu + \sigma) \approx 0.68$$

$$p(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 0.95$$

$$p(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 0.99$$



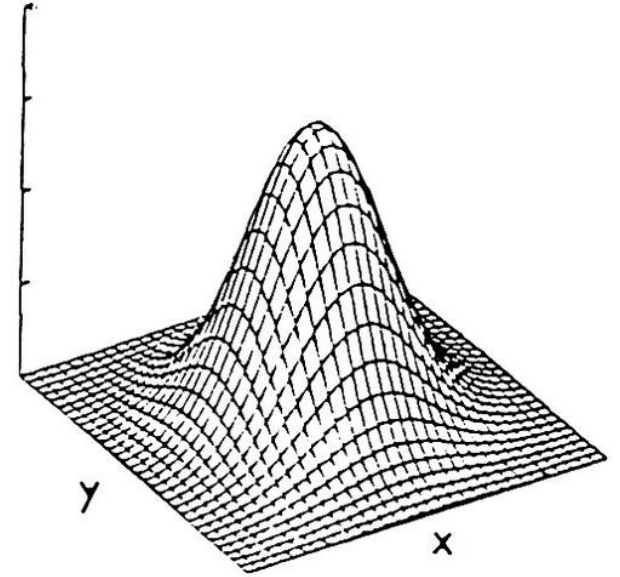
# The Gaussian distribution

---

- Gaussian distribution over a  $D$ -dimensional vector  $\mathbf{x}$  of continuous variables

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

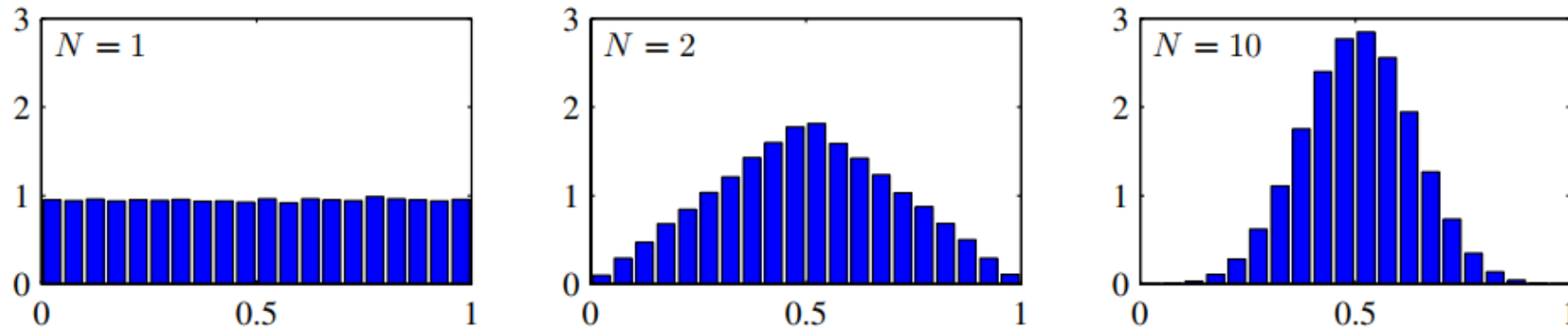
- $\boldsymbol{\mu}$ : a  $D \times 1$  mean vector
- $\boldsymbol{\Sigma}$ : a  $D \times D$  covariance matrix
- $|\boldsymbol{\Sigma}|$ : The determinant of  $\boldsymbol{\Sigma}$



# The Central Limit Theorem

---

- Mean of a set of random variables, which is of course itself a random variable, has a distribution that becomes increasingly Gaussian as the number of terms in the sum increases



```
import numpy as np
from matplotlib import pyplot as plt
N = int(input())
means = []
for i in range(1,100000)
    means.append(np.mean(np.random.random(size=(N,))))

plt.hist(means,bins = 100, range=(0,1))
```

# Prerequisites ...

## Eigenvectors and Eigenvalues

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

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

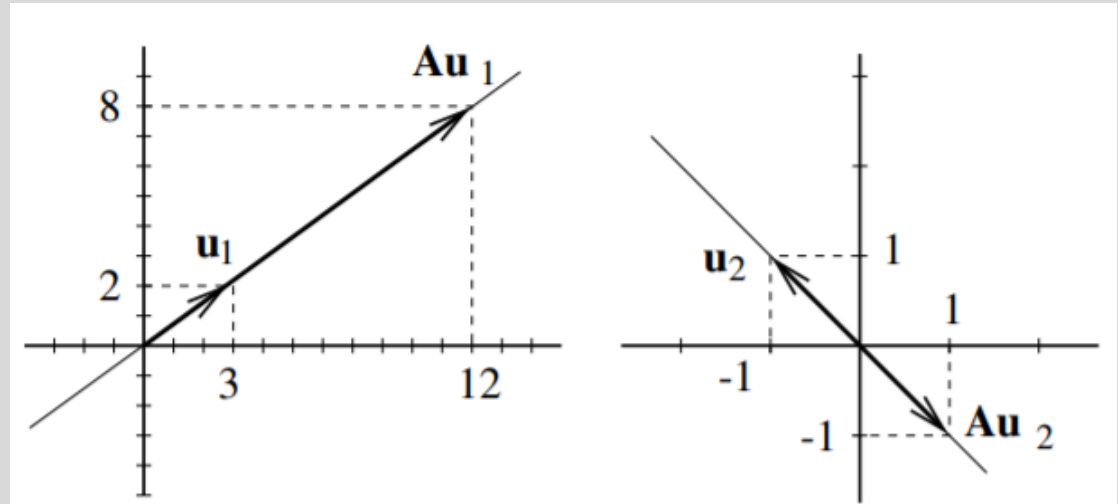
$\mathbf{u}_i$ : Eigenvector       $\lambda_i$ : Eigenvalue

- **Characteristic Equation:**

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

- Example:

$$\mathbf{A} = \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix} \left\{ \begin{array}{l} \mathbf{u}_1 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}, \lambda_1 = 4 \\ \mathbf{u}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \lambda_2 = -1 \end{array} \right.$$



# Prerequisites ...

---

## Eigenvectors and Eigenvalues

- For most applications we normalize the eigenvectors (i.e., transform them such that their length is equal to one)

$$\mathbf{u}_i \mathbf{u}_i^T = 1$$

- To normalize, we simply divide  $\mathbf{u}_i$  by its length  $|\mathbf{u}_i|$

- Example:

$$\mathbf{A} = \begin{bmatrix} 2 & 3 \\ 2 & 1 \end{bmatrix}$$

$$\left\{ \begin{array}{l} \mathbf{u}_1 = \begin{bmatrix} 3 \\ 2 \end{bmatrix}, \lambda_1 = 4 \\ \mathbf{u}_2 = \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \lambda_2 = -1 \end{array} \right.$$

$$|\mathbf{u}_1| = \sqrt{3^2 + 2^2} = \sqrt{13}$$

$$|\mathbf{u}_2| = \sqrt{-1^2 + 1^2} = \sqrt{2}$$

Normalized eigenvectors

$$\mathbf{u}_1 = \begin{bmatrix} 3/\sqrt{13} \\ 2/\sqrt{13} \end{bmatrix} = \begin{bmatrix} 0.8331 \\ 0.5547 \end{bmatrix}$$

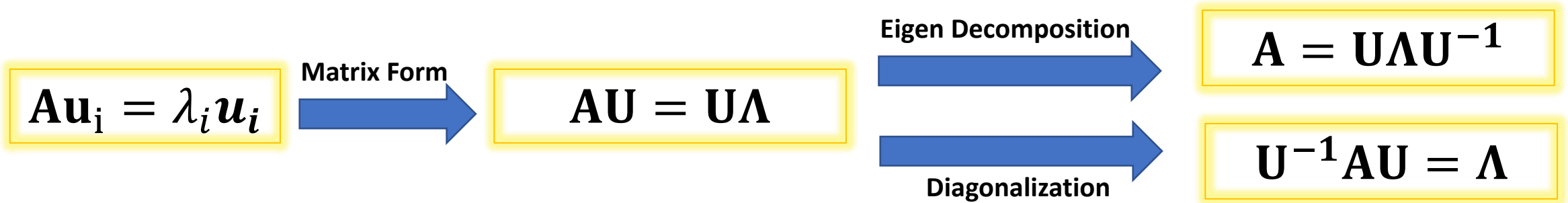
$$\mathbf{u}_2 = \begin{bmatrix} -1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix} = \begin{bmatrix} -0.7071 \\ 0.7071 \end{bmatrix}$$

# Prerequisites ...

---

## Eigenvectors and Eigenvalues

- We can re-write the eigenvector equation in matrix form:



$$\mathbf{U} = \begin{pmatrix} \mathbf{u}_1 & \dots & \mathbf{u}_M \end{pmatrix} \quad \mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_M \end{pmatrix}$$

$$\mathbf{A}\mathbf{u}_i = \lambda_i\mathbf{u}_i \Rightarrow \mathbf{A}^{-1}\mathbf{A}\mathbf{u}_i = \mathbf{A}^{-1}\lambda_i\mathbf{u}_i \Rightarrow \mathbf{u}_i = \lambda_i\mathbf{A}^{-1}\mathbf{u}_i \Rightarrow \frac{1}{\lambda_i}\mathbf{u}_i = \mathbf{A}^{-1}\mathbf{u}_i$$

# Prerequisites ...

---

## Eigenvectors and Eigenvalues

- If  $\mathbf{A}$  is a **real symmetric matrix**, then its eigenvalues are **real** and can be chosen to form **orthonormal set**, so that

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1 & , \text{if } i = j \\ 0 & , \text{otherwise} \end{cases}$$

Proof: Homework

- Or

$$\mathbf{U}^T \mathbf{U} = \mathbf{I} \quad \Rightarrow \quad \mathbf{U}^T \mathbf{U} \mathbf{U}^{-1} = \mathbf{U}^{-1} = \mathbf{U}^T$$

- Then

$$\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^{-1} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T = \sum_{i=1}^M \lambda_i \mathbf{u}_i \mathbf{u}_i^T \quad \text{A very nice property ☺}$$

$$\mathbf{A}^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^{-1} = \mathbf{U} \mathbf{\Lambda}^{-1} \mathbf{U}^T = \sum_{i=1}^M \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad \text{Another nice property ☺}$$

# Prerequisites ...

---

## Eigenvectors and Eigenvalues

- The **rank of matrix  $\mathbf{A}$**  is equal to the number of nonzero eigenvalues.
- A matrix  $\mathbf{A}$  is called **positive definite** if its eigenvalues are strictly positive.
- A matrix  $\mathbf{A}$  is called **positive semidefinite** if its eigenvalues are nonnegative.
- The product of the eigenvalues of  $\mathbf{A}$  is the same as  $|\mathbf{A}|$ . Therefore,  $\mathbf{A}$  is invertible if and only if it **does not have** a zero eigenvalue (its rank equals  $M$ )
- Generally the **covariance matrix for the Gaussian distribution ( $\Sigma$ )** is symmetric and positive definite.

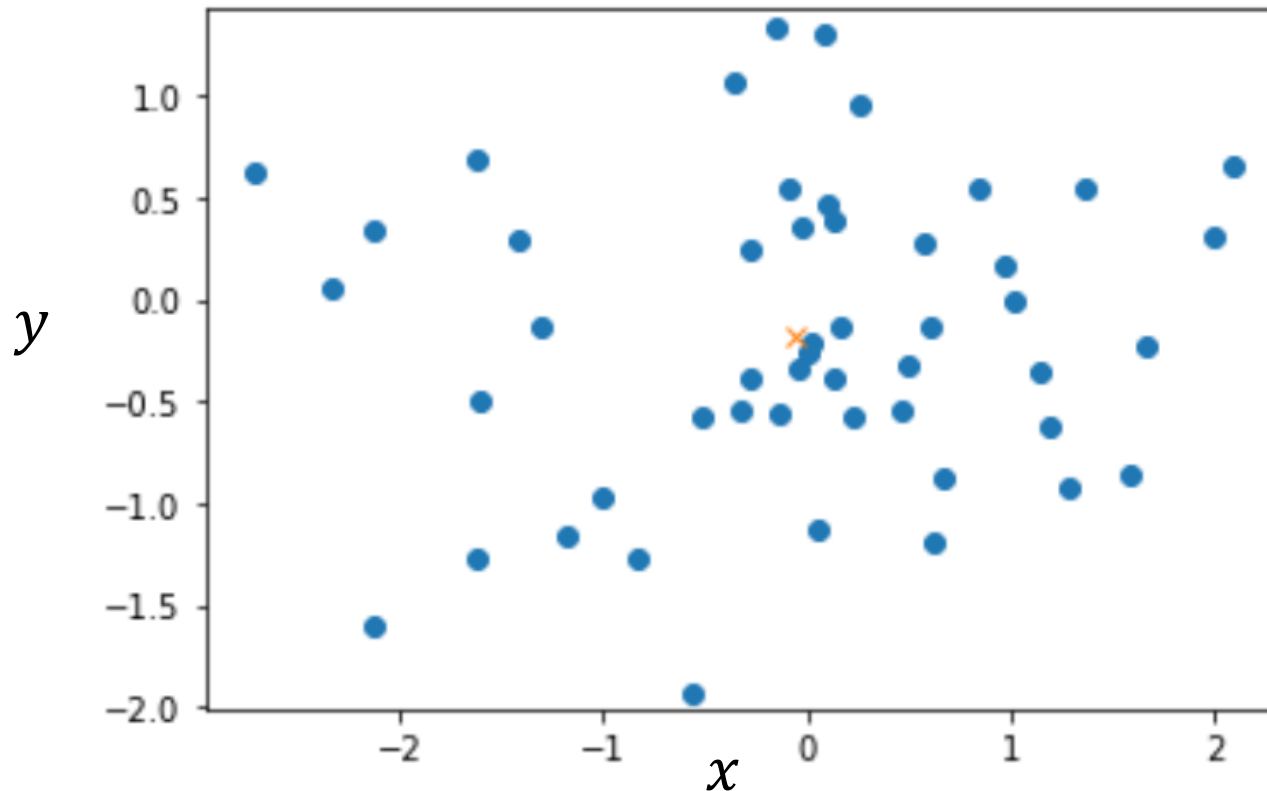
# Prerequisites ...

---

## Mahalanobis Distance

- The **Euclidean distance** of a point from the mean (example for a 2D variable):

$$\sqrt{(x - \bar{x})^2 + (y - \bar{y})^2}$$



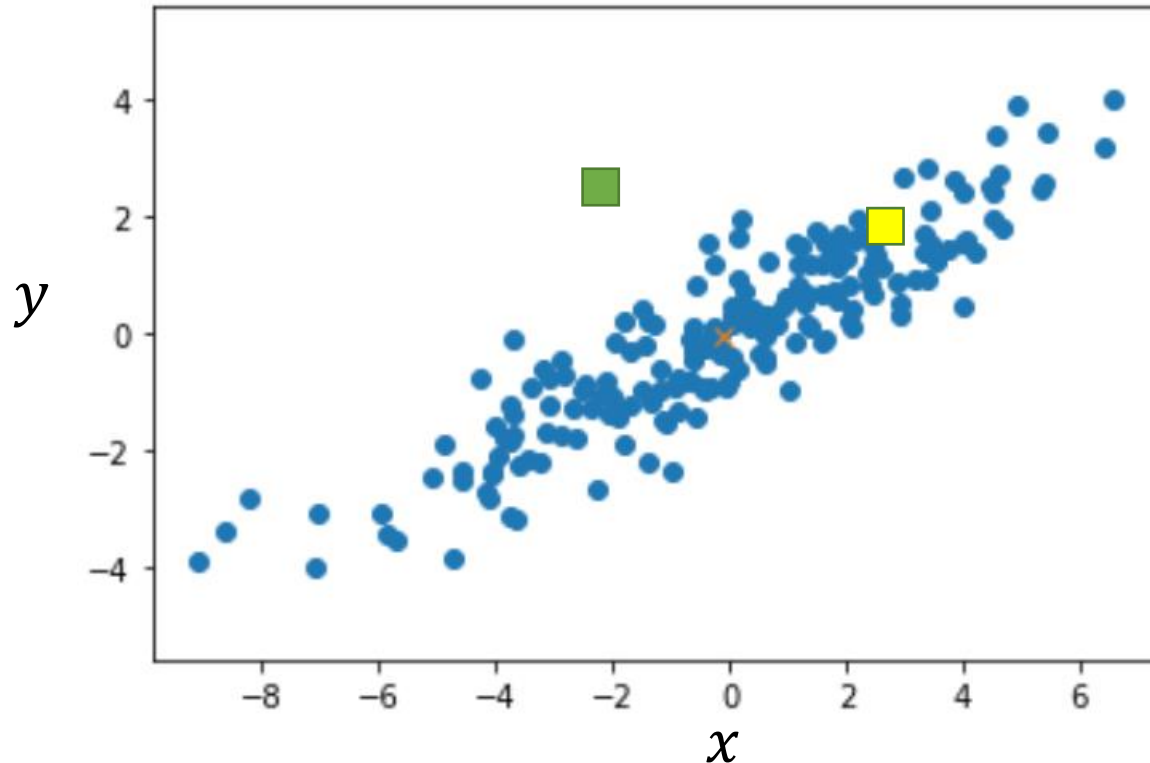


# Prerequisites ...

---

## Mahalanobis Distance

- However, Euclidean distance has limitations in real datasets, which often have some degree of covariance



# Prerequisites ...

---

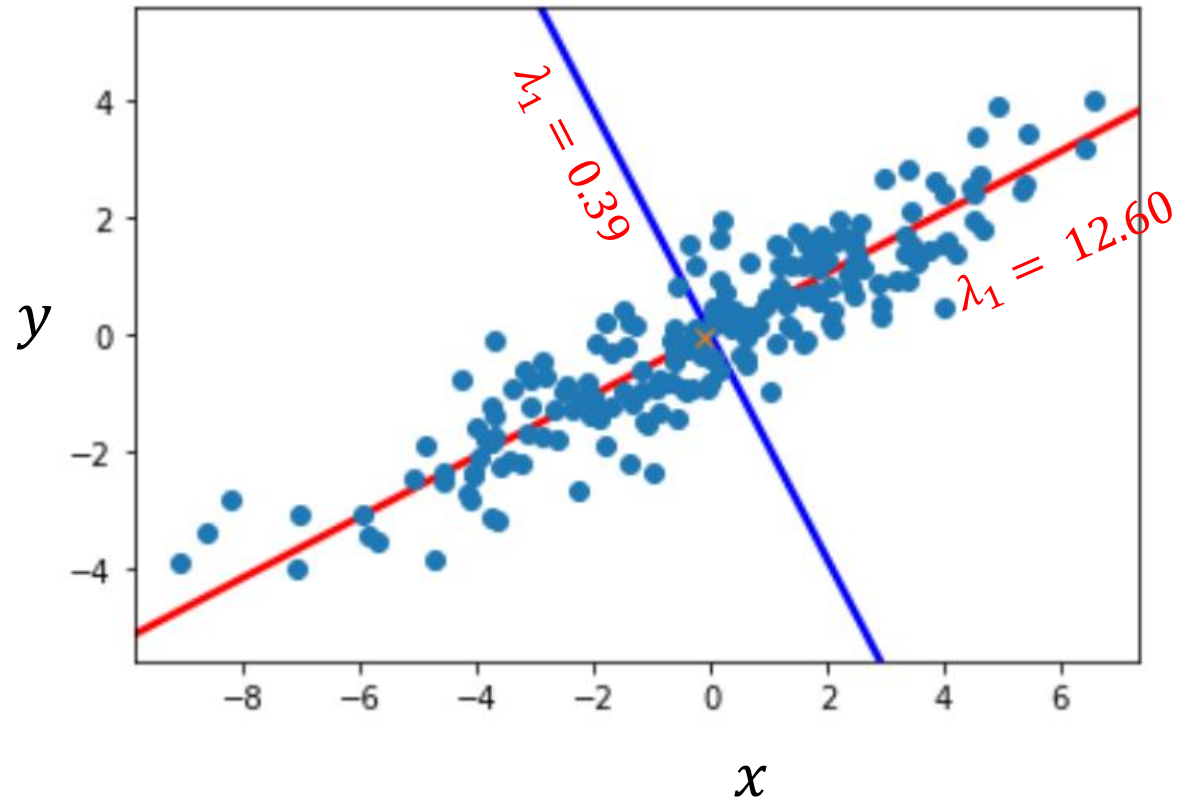
## Mahalanobis Distance

- The idea of **Mahalanobis distance** is to remove the covariance by treating **each eigenvector as a new axis**, shrink the axis by  $\sqrt{\lambda_i}$ , then calculate distance between points

$$D^2 = (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{S}^{-1} (\mathbf{x} - \bar{\mathbf{x}})$$

Data mean vector

Data covariance matrix



# Prerequisites ...

---

## Jacobian Factor

- Under a nonlinear change of variable, a probability density transforms differently from a simple function, due to the **Jacobian factor**.
- For instance, if we consider a change of variables  $x = g(y)$ , then a function  $f(x)$  becomes  $h(y) = f(g(y))$
- Now consider a probability density  $p_x(x)$

- Observations falling in the 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 the range  $(y, y + \delta y)$  have probability  $p_y(y)\delta y$

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

- In the case of multivariate probabilities, in going from  $\mathbf{x}$  to  $\mathbf{y}$  coordinate system, we have:

$$p(\mathbf{y}) = p(\mathbf{x}) |J|$$

Where

$$J_{ij} = \frac{\partial \mathbf{x}_i}{\partial \mathbf{y}_j}$$

# The Gaussian distribution

---

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

- The geometric form of the Gaussian distribution

- The Gaussian distribution depends on  $\mathbf{x}$  is through the quadratic form

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

- The quantity  $\Delta$  is the [Mahalanobis distance](#) from  $\boldsymbol{\mu}$  to  $\mathbf{x}$
- This quantity reduces to the Euclidean distance when  $\boldsymbol{\Sigma} = \mathbf{I}$
- The Gaussian distribution will be constant on surfaces in  $\mathbf{x}$ -space for which  $\Delta^2$  is constant.

# The Gaussian distribution

---

- Consider the eigenvector equation for  $\Sigma$  (this matrix is real symmetric)

$$\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i \quad , i = 1, \dots, D$$

- Eigenvectors form an orthonormal set

$$\mathbf{u}_i^T \mathbf{u}_j = \begin{cases} 1 & , \text{if } i = j \\ 0 & , \text{otherwise} \end{cases}$$

- $\Sigma$  can be expressed as an expansion of its eigenvectors

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

- The inverse covariance matrix can be expressed as

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

# The Gaussian distribution

---

- By substituting the inverse covariance matrix into the quadratic form  $\Delta^2$

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

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

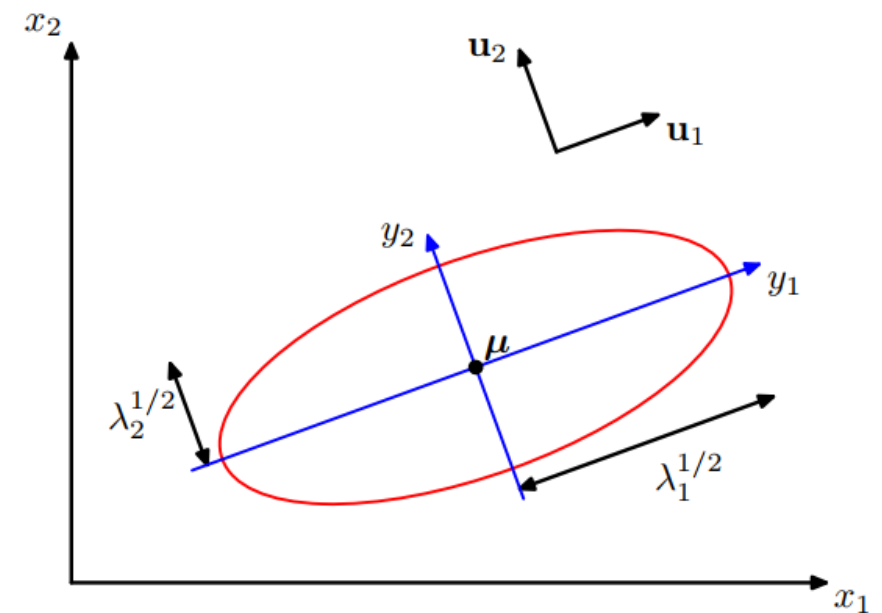
$$\mathbf{y} = \mathbf{U}(\mathbf{x} - \boldsymbol{\mu})$$

- $\mathbf{U}$  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

---

- We can interpret  $\{y_i\}$  as a new coordinate system defined by the orthonormal vectors  $\mathbf{u}_i$  that are shifted and rotated with respect to the original  $x_i$  coordinates.
- The quadratic form, and hence the Gaussian density, will be 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
  - Centered in  $\boldsymbol{\mu}$  and axis oriented along  $\mathbf{u}_i$ .
  - The scaling factor in the directions of the axis are  $\lambda_i^{\frac{1}{2}}$



# The Gaussian distribution

---

- Now consider the form of the Gaussian distribution in the new coordinate system defined by the  $y_i$ .
- In going from the  $\mathbf{x}$  to the  $\mathbf{y}$  coordinate system, we have a Jacobian matrix  $\mathbf{J}$

$$J_{ij} = \frac{\partial x_i}{\partial y_j} = \mathbf{U}_{ij}^T$$

$$\begin{aligned} \mathbf{y} &= \mathbf{U}(\mathbf{x} - \boldsymbol{\mu}) \\ \Rightarrow \mathbf{x} &= \mathbf{U}^{-1}\mathbf{y} + \boldsymbol{\mu} = \mathbf{U}^T\mathbf{y} + \boldsymbol{\mu} \end{aligned}$$

- Using the orthonormality property of the matrix  $\mathbf{U}$ :

$$|\mathbf{J}|^2 = |\mathbf{U}^T|^2 = |\mathbf{U}^T| |\mathbf{U}^T| = |\mathbf{U}^T| |\mathbf{U}| = |\mathbf{U}^T \mathbf{U}| = |\mathbf{I}| = 1 \Rightarrow |\mathbf{J}| = 1$$

- Moreover

$$|\Sigma| = \prod_{j=1}^D \lambda_j \Rightarrow |\Sigma|^{\frac{1}{2}} = \prod_{j=1}^D \lambda_j^{\frac{1}{2}}$$



# The Gaussian distribution

---

- Thus in the  $y_j$  coordinate system, the Gaussian distribution takes the form

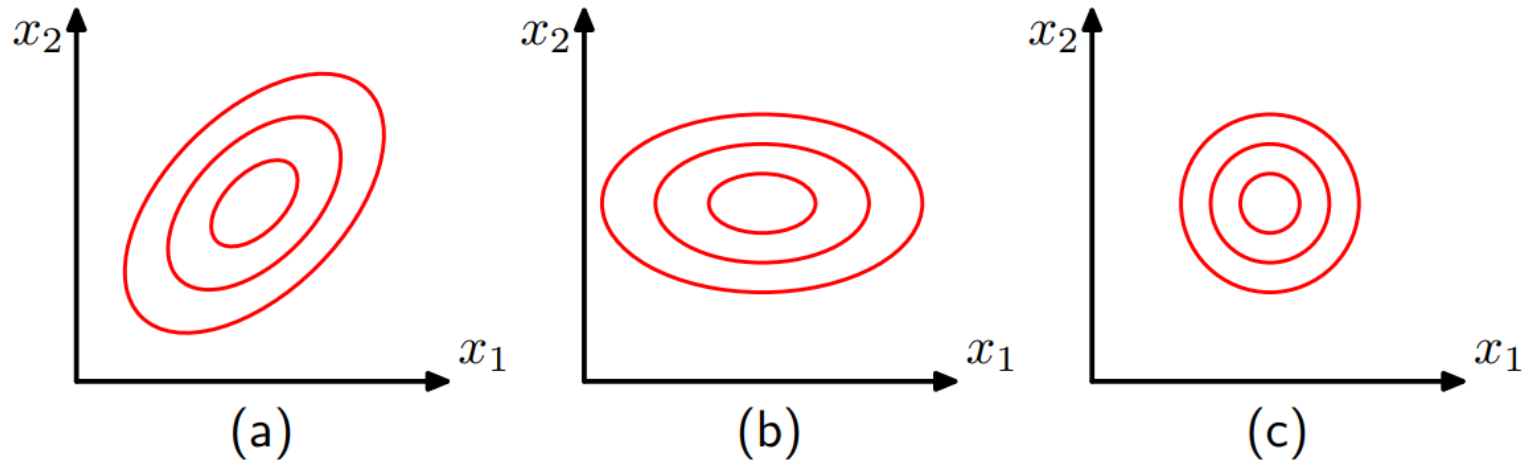
$$p(\mathbf{y}) = p(\mathbf{x}) |\mathbf{J}| = p(\mathbf{x})$$

$$\begin{aligned} \Rightarrow p(\mathbf{y}) &= \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}\Delta^2\right\} \xrightarrow{\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \text{ and } |\boldsymbol{\Sigma}|^{\frac{1}{2}} = \prod_{j=1}^D \lambda_j^{\frac{1}{2}}} p(\mathbf{y}) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{\prod_{j=1}^D \lambda_j^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} \sum_{i=1}^D \frac{y_i^2}{\lambda_i}\right\} \\ &= \frac{1}{\prod_{j=1}^D (2\pi)^{\frac{D}{2}} \lambda_j^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} \sum_{i=1}^D \frac{y_i^2}{\lambda_i}\right\} = \prod_{j=1}^D \frac{1}{(2\pi\lambda_j)^{D/2}} \exp\left\{-\frac{y_i^2}{2\lambda_i}\right\} \end{aligned}$$

- Therefore  $p(\mathbf{y})$  is the product of  $D$  independent univariate Gaussian distributions

# The Gaussian distribution

---



□ a) General  $\Sigma$   $\longrightarrow$   $\frac{D(D+3)}{2}$  parameters

□ b) Diagonal  $\Sigma$   $\longrightarrow$   $2D$  parameters

□ c)  $\Sigma = \sigma^2 \mathbf{I}$   $\longrightarrow$   $D + 1$  parameters

# Conditional Gaussian distribution

**Property 1 of the Gaussian Distribution:** If two sets of variables are jointly Gaussian, then the conditional distribution of one set conditioned on the other is again Gaussian.

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N} \left( (\mathbf{x}, \mathbf{y}) \mid \boldsymbol{\mu}_{(\mathbf{x}, \mathbf{y})}, \boldsymbol{\Sigma}_{(\mathbf{x}, \mathbf{y})} \right) \implies p(\mathbf{x} \mid \mathbf{y}) = \mathcal{N} \left( (\mathbf{x} \mid \mathbf{y}) \mid \boldsymbol{\mu}_{\mathbf{x} \mid \mathbf{y}}, \boldsymbol{\Sigma}_{\mathbf{x} \mid \mathbf{y}} \right)$$

- Suppose  $\mathbf{x}$  is a  $D$ -dimensional vector with Gaussian distribution  $\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma})$ 
  - we partition  $\mathbf{x}$  into two disjoint subsets  $\mathbf{x}_a$  and  $\mathbf{x}_b$ .
  - Without loss of generality, we can take  $\mathbf{x}_a$  to form the first  $M$  components of  $\mathbf{x}$ , with  $\mathbf{x}_b$  comprising the remaining  $D - M$  components,
  - Then

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$$

# Conditional Gaussian distribution

---

○ We also define corresponding partitions of

□ the mean vector  $\boldsymbol{\mu}$  given by  $\boldsymbol{\mu} = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix}$       □ The covariance matrix  $\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$

□ Because  $\boldsymbol{\Sigma}$  is symmetric ( $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^T$ ) then  $\boldsymbol{\Sigma}_{aa}$  and  $\boldsymbol{\Sigma}_{bb}$  are also symmetric and  $\boldsymbol{\Sigma}_{ab} = \boldsymbol{\Sigma}_{ba}^T$

○ In many situations, it is convenient to work with the **precision matrix**:  $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$

□ The corresponding partition for  $\boldsymbol{\Lambda}$ :  $\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$

□ Because the inverse of a symmetric matrix is also symmetric then  $\boldsymbol{\Lambda}_{aa}$  and  $\boldsymbol{\Lambda}_{bb}$  are also symmetric and  $\boldsymbol{\Lambda}_{ab} = \boldsymbol{\Lambda}_{ba}^T$

□ It should be stressed that, for instance,  $\boldsymbol{\Lambda}_{aa}$  is not simply given by the inverse of  $\boldsymbol{\Sigma}_{aa}$ .

# Conditional Gaussian distribution

---

- We have

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

$$\Rightarrow p(\mathbf{x}_a|\mathbf{x}_b) = \frac{p(\mathbf{x}_a, \mathbf{x}_b)}{p(\mathbf{x}_b)} = \frac{p(\mathbf{x}_a, \mathbf{x}_b)}{\int_{-\infty}^{+\infty} p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_a}$$

Not an easy approach



## Better approach (Analytical Method):

- Remember

- The Gaussian distribution depends on  $\mathbf{x}$  is through the quadratic form

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

- Therefore, to show that  $p(\mathbf{x}_a|\mathbf{x}_b)$  is Gaussian, we need to proof that  $p(\mathbf{x}_a|\mathbf{x}_b)$  has a similar quadratic form with respect to  $\mathbf{x}_a$ .

# Conditional Gaussian distribution

---

- We have

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

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

# Conditional Gaussian distribution

---

○ **Question:** How to find  $\boldsymbol{\mu}_{\mathbf{x}_a|\mathbf{x}_b}$  and  $\boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b}$  for  $p(\mathbf{x}_a|\mathbf{x}_b)$ ?

○ **Answer:** Using an approach called **Completing the Square**

□ For a general Gaussian distribution  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , the exponent can be written as:

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

□ For  $p(\mathbf{x}_a|\mathbf{x}_b)$ , we have:

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

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

$$-\frac{1}{2} \mathbf{x}_a^T \boldsymbol{\Lambda}_{aa} \mathbf{x}_a \quad \longrightarrow \quad \boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b}^{-1} = \boldsymbol{\Lambda}_{aa} \implies \boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b} = \boldsymbol{\Lambda}_{aa}^{-1}$$

# Conditional Gaussian distribution

---

□ For  $p(\mathbf{x}_a|\mathbf{x}_b)$ , we have:

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

□ If we pick out all terms that are linear in  $\mathbf{x}_a$ , we have

$$\begin{aligned}\mathbf{x}_a^T \{\boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)\} = \mathbf{x}_a^T \boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b}^{-1} \boldsymbol{\mu}_{\mathbf{x}_a|\mathbf{x}_b} & \xrightarrow{\boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b} = \boldsymbol{\Lambda}_{aa}^{-1}} \boldsymbol{\mu}_{\mathbf{x}_a|\mathbf{x}_b} = \boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b} \{\boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)\} \\ & = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)\end{aligned}$$



# Conditional Gaussian distribution

○ **Summary:**  $\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$      $\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}$      $\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$      $\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$

$$\boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b} = \boldsymbol{\Lambda}_{aa}^{-1} \quad \boldsymbol{\mu}_{\mathbf{x}_a|\mathbf{x}_b} = \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b)$$

○ **Question:** How to find  $\boldsymbol{\mu}_{\mathbf{x}_a|\mathbf{x}_b}$  and  $\boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b}$  in terms of  $\boldsymbol{\Sigma}$  (not  $\boldsymbol{\Lambda}$ )

□ We can use the following identity:

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

□ 
$$\begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix} \quad \longrightarrow \quad \begin{aligned} \boldsymbol{\Lambda}_{aa} &= (\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1} \\ \boldsymbol{\Lambda}_{ab} &= -(\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})^{-1}\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1} \end{aligned}$$

□ 
$$\boldsymbol{\Sigma}_{\mathbf{x}_a|\mathbf{x}_b} = \boldsymbol{\Lambda}_{aa}^{-1} = (\boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba})$$
    
$$\boldsymbol{\mu}_{\mathbf{x}_a|\mathbf{x}_b} = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b)$$

# Marginal Gaussian distribution

**Property 2 of the Gaussian Distribution:** If two sets of variables are jointly Gaussian, then the marginal distributions is again Gaussian.

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left((\mathbf{x}, \mathbf{y}) \mid \boldsymbol{\mu}_{(\mathbf{x}, \mathbf{y})}, \boldsymbol{\Sigma}_{(\mathbf{x}, \mathbf{y})}\right) \implies p(\mathbf{x}) = \mathcal{N}\left(\mathbf{x} \mid \boldsymbol{\mu}_{\mathbf{x}}, \boldsymbol{\Sigma}_{\mathbf{x}}\right)$$

○

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix} \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

○ Similar to conditional probability, we can prove that

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \longrightarrow \quad p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a \mid \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

# Bayes' Theorem for Gaussian Variables

**Property 3 of the Gaussian Distribution:** Given a marginal Gaussian distribution for  $\mathbf{x}$  and a conditional Gaussian distribution for  $\mathbf{y}$  given  $\mathbf{x}$  in the form

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

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

The joint distribution of  $\mathbf{x}$  and  $\mathbf{y}$  is given by

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \mid \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A}\boldsymbol{\mu} + \mathbf{b} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Lambda}^{-1} & \boldsymbol{\Lambda}^{-1}\mathbf{A}^T \\ \mathbf{A}\boldsymbol{\Lambda}^{-1} & \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T \end{pmatrix}\right)$$

○ We find an expression for the joint distribution  $p(\mathbf{x}, \mathbf{y})$ .

□ To do this, we define

$$\mathbf{z} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}$$

# Bayes' Theorem for Gaussian Variables

---

□ Then we have

$$p(\mathbf{z}) = p(\mathbf{x})p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \times \mathcal{N}(\mathbf{y}|\mathbf{Ax} + \mathbf{b}, \mathbf{L}^{-1})$$

□ Considering the log of the joint distribution

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

□ As before, we see that this is a quadratic function of the components of  $\mathbf{z}$ , and hence  $p(\mathbf{z})$  is Gaussian distribution.

□ To find the precision of this Gaussian, we consider the second order terms

$$\begin{aligned} &-\frac{1}{2}\mathbf{x}^T(\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})\mathbf{x} - \frac{1}{2}\mathbf{y}^T\mathbf{L}\mathbf{y} + \frac{1}{2}\mathbf{y}^T\mathbf{L}\mathbf{A}\mathbf{x} + \frac{1}{2}\mathbf{x}^T\mathbf{A}^T\mathbf{L}\mathbf{y} \\ &= -\frac{1}{2}\begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A} & -\mathbf{A}^T\mathbf{L} \\ -\mathbf{L}\mathbf{A} & \mathbf{L} \end{pmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} = -\frac{1}{2}\mathbf{z}^T\boldsymbol{\Sigma}_z^{-1}\mathbf{z} \end{aligned}$$

**Remember:** Completing the Squares

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

$$\boldsymbol{\Sigma}_z^{-1} = \begin{pmatrix} \boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A} & -\mathbf{A}^T\mathbf{L} \\ -\mathbf{L}\mathbf{A} & \mathbf{L} \end{pmatrix}$$

# Bayes' Theorem for Gaussian Variables

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

□ Then

$$\Sigma_z = \begin{pmatrix} \Lambda + \mathbf{A}^T \mathbf{L} \mathbf{A} & -\mathbf{A}^T \mathbf{L} \\ -\mathbf{L} \mathbf{A} & \mathbf{L} \end{pmatrix}^{-1} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} \mathbf{A}^T \\ \mathbf{A} \Lambda^{-1} & \mathbf{L}^{-1} + \mathbf{A} \Lambda^{-1} \mathbf{A}^T \end{pmatrix}$$

□ Similarly, we can find the mean of the Gaussian distribution over  $z$  by identifying the linear terms

Remember: Completing the Squares

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

$$\mathbf{x}^T \Lambda \boldsymbol{\mu} - \mathbf{x}^T \mathbf{A}^T \mathbf{L} \mathbf{b} + \mathbf{y}^T \mathbf{L} \mathbf{b} = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^T \begin{pmatrix} \Lambda \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} = \mathbf{z}^T \Sigma_z^{-1} \boldsymbol{\mu}_z \quad \rightarrow$$

$$\boldsymbol{\mu}_z = \Sigma_z \begin{pmatrix} \Lambda \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} = \begin{pmatrix} \Lambda^{-1} & \Lambda^{-1} \mathbf{A}^T \\ \mathbf{A} \Lambda^{-1} & \mathbf{L}^{-1} + \mathbf{A} \Lambda^{-1} \mathbf{A}^T \end{pmatrix} \begin{pmatrix} \Lambda \boldsymbol{\mu} - \mathbf{A}^T \mathbf{L} \mathbf{b} \\ \mathbf{L} \mathbf{b} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu} \\ \mathbf{A} \boldsymbol{\mu} + \mathbf{b} \end{pmatrix}$$

# Bayes' Theorem for Gaussian Variables

**Property 4 of the Gaussian Distribution:** Given a marginal Gaussian distribution for  $\mathbf{x}$  and a conditional Gaussian distribution for  $\mathbf{y}$  given  $\mathbf{x}$  in the form

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

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

The marginal distribution of  $\mathbf{y}$  is given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T)$$

- It is obvious from properties 2 and 3.

**Remember**

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$$

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix}$$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix}$$

$$\boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \implies p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

# Bayes' Theorem for Gaussian Variables

**Property 5 of the Gaussian Distribution:** Given a marginal Gaussian distribution for  $\mathbf{x}$  and a conditional Gaussian distribution for  $\mathbf{y}$  given  $\mathbf{x}$  in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

The conditional distribution of  $\mathbf{x}$  given  $\mathbf{y}$  is

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \quad \text{Where} \quad \boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1}$$

- It is obvious from properties 1 and 3.

**Remember**

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{pmatrix} \quad \boldsymbol{\Lambda} = \begin{pmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{pmatrix}$$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \implies p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b - \boldsymbol{\mu}_b), \boldsymbol{\Lambda}_{aa}^{-1})$$

# 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**.

$$\begin{aligned} \boldsymbol{\mu}_{ML}, \boldsymbol{\Sigma}_{ML} &= \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} p(\mathbf{X} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \\ &= \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} \ln \prod_{n=1}^N \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \\ &= \arg \max_{\boldsymbol{\mu}, \boldsymbol{\Sigma}} -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\boldsymbol{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \end{aligned}$$
$$\left. \begin{array}{l} \boldsymbol{\mu}_{ML} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \\ \boldsymbol{\Sigma}_{ML} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu}_{ML})(\mathbf{x}_n - \boldsymbol{\mu}_{ML})^T \end{array} \right\}$$

Proof: Homework



# Maximum Likelihood for the Gaussian

---

- In the case of  $D = 1$  (univariate Gaussian distribution):

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \qquad \sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

- **Question:** Are  $\boldsymbol{\mu}_{ML}$  and  $\boldsymbol{\Sigma}_{ML}$  good estimations for  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ ?

- A good estimation  $\hat{\alpha}$  for parameter  $\alpha$  should be **unbiased to the data set**

$$\mathbb{E}[\hat{\alpha}] = \alpha$$

- For the cases of  $\boldsymbol{\mu}_{ML}$  and  $\boldsymbol{\Sigma}_{ML}$  we have:

$$\mathbb{E}[\boldsymbol{\mu}_{ML}] = \boldsymbol{\mu} \qquad \mathbb{E}[\boldsymbol{\Sigma}_{ML}] = \left(\frac{N-1}{N}\right) \boldsymbol{\Sigma}$$

- A better (unbiased) estimation for  $\sigma^2$

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

# Bayesian Inference for the Gaussian

---

- The maximum likelihood framework gave point estimates for the parameters  $\mu$  and  $\Sigma$ . Now we develop a **Bayesian treatment** by introducing **prior distributions** over these parameters.
- We will consider the following cases
  - The variance is known, and we consider the task of inferring the mean
  - The mean is known, and we consider the task of inferring the variance

# Bayesian Inference for the Gaussian

---

- **Case 1:** The variance is known, and we consider the task of inferring the mean
  - Let us begin with a simple example in which we consider a single Gaussian random variable  $x$  ( $D = 1$ ).
  - We shall suppose that the variance  $\sigma^2$  is known, and we consider the task of inferring the mean  $\mu$  given a set of  $N$  observations  $\mathbf{X} = \{x_1, \dots, x_N\}$ .

- The likelihood function:

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

**Remember**  $p(\mu|\mathcal{D}) = \frac{p(\mathcal{D}|\mu) \times p(\mu)}{p(\mathcal{D})}$

- The likelihood function takes the form of the exponential of a quadratic form in  $\mu$ . If we choose a Gaussian prior  $p(\mu)$ , it will be a conjugate distribution for the likelihood function.
- The posterior is a product of two exponentials of quadratic functions of  $\mu$  and hence will also be Gaussian.

# Bayesian Inference for the Gaussian

---

- We take our prior distribution to be

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

- The posterior distribution is given by

$$p(\mu | \mathbf{X}) \propto p(\mathbf{X} | \mu) p(\mu)$$

- Some manipulations involving completing the square in the exponent allow to show that the posterior distribution is given by

$$p(\mu | \mathbf{X}) = \mathcal{N}(\mu | \mu_N, \sigma_N^2)$$
$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML}$$

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

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$$

# Bayesian Inference for the Gaussian

$$\mu_N = \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_{ML} \qquad \frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

□ Note that  $\mu_N$  (the mean of the posterior distribution) is a compromise between the prior mean ( $\mu_0$ ) and the maximum likelihood solution  $\mu_{ML}$ .

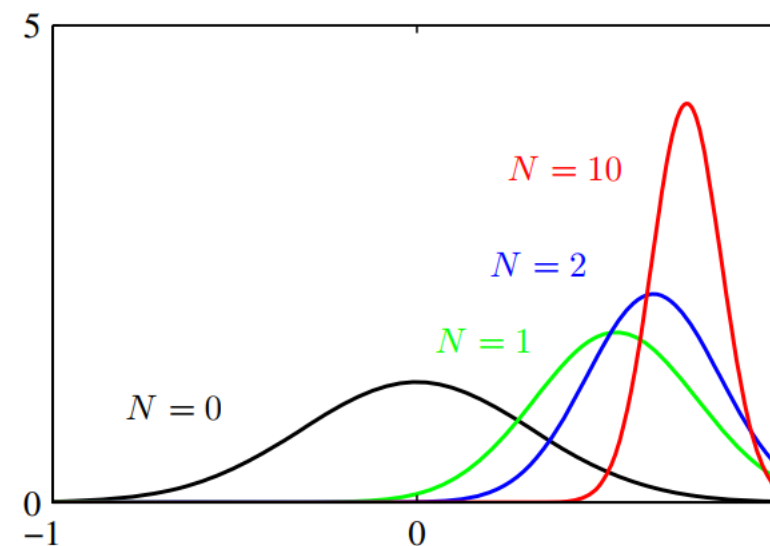
➤ If  $N = 0$ ,  $\mu_N$  reduces to the prior mean.

➤ For  $N \rightarrow \infty$ , the posterior mean equals the maximum likelihood solution.

□ As we increase the number of observed data points, the precision steadily increases, corresponding to a posterior distribution with steadily decreasing variance.

➤ With no observed data points, we have the prior variance

➤ If  $N \rightarrow \infty$ , the variance  $\sigma_N^2 \rightarrow 0$  and the posterior distribution becomes infinitely peaked around  $\mu_{ML}$



**Figure.** The data points are generated from a Gaussian of mean 0.8 and variance 0.1, and the prior is chosen to have mean 0.

# Bayesian Inference for the Gaussian

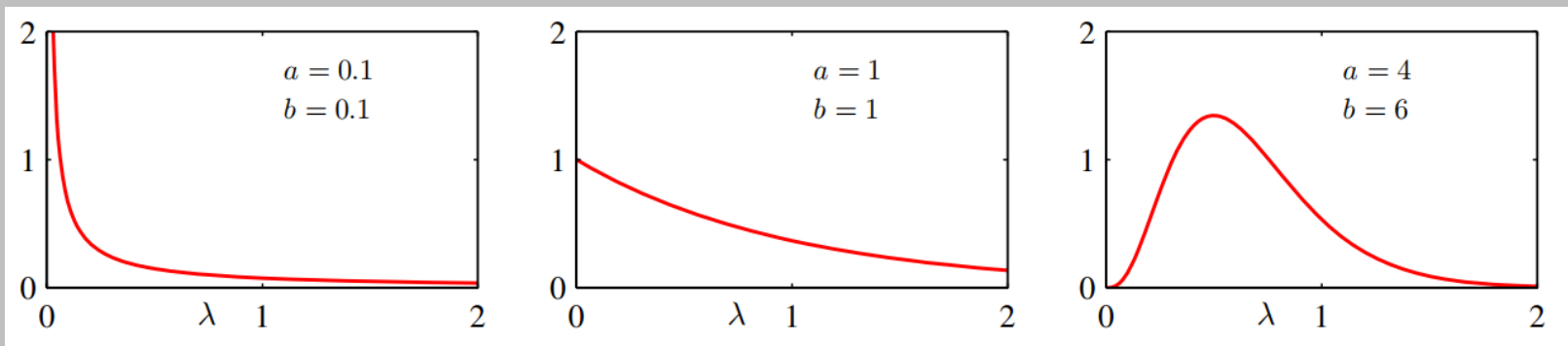
- **Case 2:** The mean is known, and we consider the task of inferring the variance

- The likelihood function (It turns out to be most convenient to work with the precision  $\lambda \equiv \frac{1}{\sigma^2}$ ):

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

- The corresponding conjugate prior should therefore be proportional to the product of a power of  $\lambda$  and the exponential of a linear function of  $\lambda$ .
- This corresponds to the gamma distribution which is defined by

$$\text{Gam}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda) \quad \mathbb{E}[\lambda] = \frac{a}{b}, \text{var}[\lambda] = \frac{a}{b^2}$$



# Bayesian Inference for the Gaussian

---

- Consider a prior distribution  $\text{Gam}(\lambda|a_0, b_0)$ . Multiplying by the likelihood function, we obtain a the following posterior distribution

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

- which we recognize as a gamma distribution of the form  $\text{Gam}(\lambda|a_N, b_N)$  where

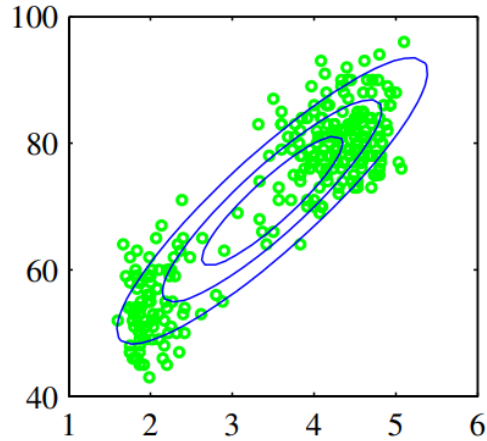
$$a_N = a_0 + \frac{N}{2}$$

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

# Mixtures of Gaussians

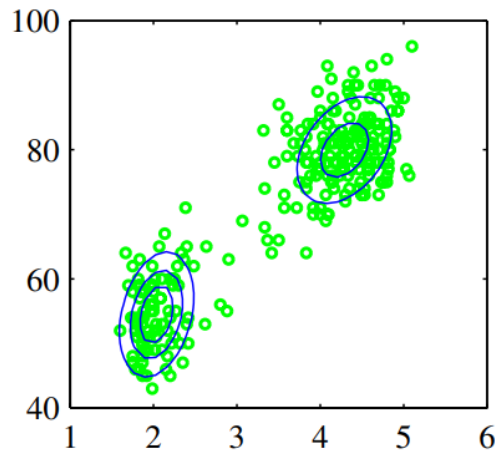
---

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



**Figure:** A single Gaussian distribution fitted to the data using maximum likelihood.

- Note that this distribution fails to capture the two clumps in the data and indeed places much of its probability mass in the central region between the clumps where the data are relatively sparse.



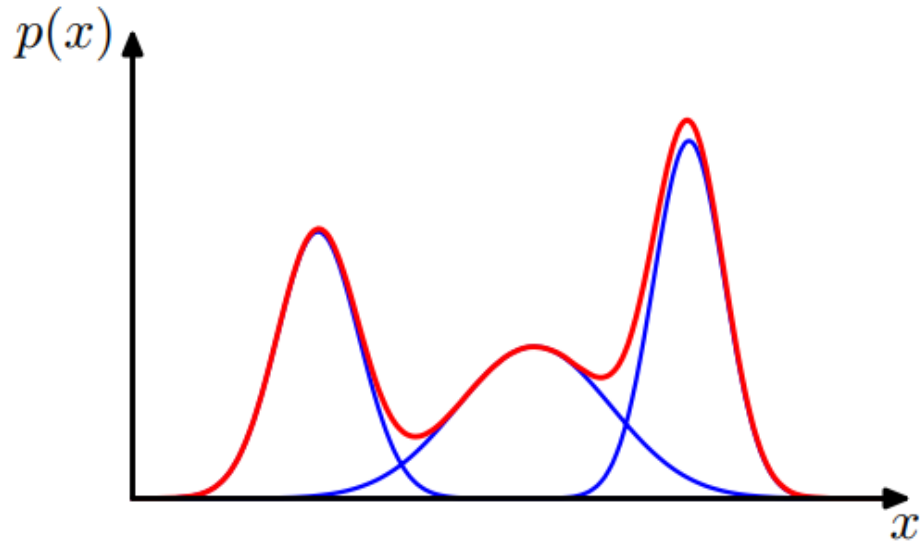
**Figure:** A linear combination of two Gaussians fitted using maximum likelihood

- Such superpositions, formed by taking linear combinations of more basic distributions such as Gaussians, can be formulated as probabilistic models known as **mixture distributions**



# Mixtures of Gaussians

---



**Figure:** Example of a Gaussian mixture distribution in one dimension

- Three Gaussians (blue) 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.
- We consider a mixture of  $K$  Gaussian densities of the form

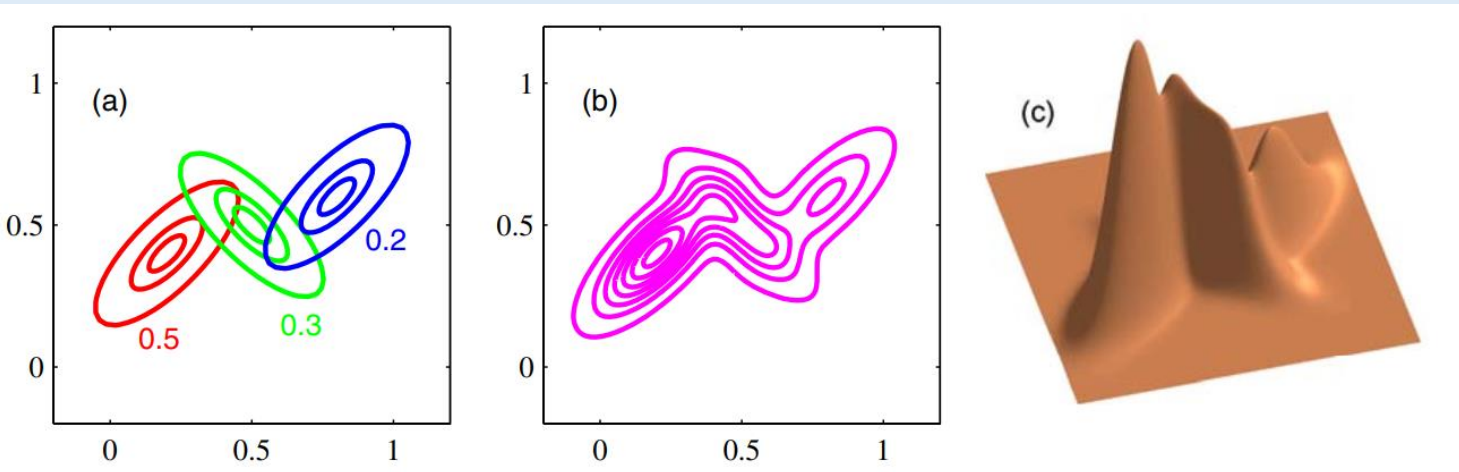
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

# Mixtures of Gaussians

- We consider a mixture of  $K$  Gaussian densities of the form

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

- The parameters  $\pi_k$  are called **mixing coefficients**.  $0 \leq \pi_k \leq 1$  and  $\sum_{k=1}^K \pi_k = 1$



**Figure:** A mixture of 3 Gaussians in a two-dimensional space

- (a) Contours of constant density for each of the mixture components (b) Contours of the mixture distribution  $p(\mathbf{x})$ . (c) A surface plot of the distribution  $p(\mathbf{x})$ .