

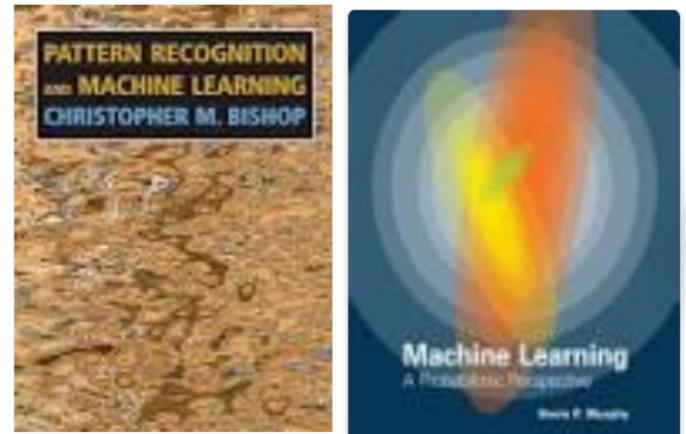
# Uncertainty Quantification

Sep 2022

# Uncertainty Quantification

Recommended for reading:

- Bishop, Pattern Recognition and Machine Learning
- Murphy, Machine Learning A probabilistic approach
- These notes (from various sources)



Prerequisites:

Probability

Numerical methods

Knowledge of deep learning

Participation 20%

Final project 80%

# Uncertainty Quantification

- Introduction
  - Uncertainty and predictive modeling
  - RV and distribution
  - MLE and statistical decision
  - Bayesian perspective
- Simulation
  - RV generation and CLT
  - Monte Carlo estimation
  - Laplace approximation and adjoint method

# Uncertainty Quantification

- Regression
  - Linear and logistic regression
  - Generalized linear model
  - Bayesian linear regression
  - Gaussian prior and Gaussian regression
- Bayesian modeling
  - Bayesian learning
  - Prediction with posterior distribution
  - Bayesian neural networks

# Uncertainty Quantification

- Probabilistic graphical models
  - Bayesian network
  - Directed and undirected graphical models
  - Hidden Markov model and EM
  - Exponential families
- Further topics
  - MCMC and sequential Monte Carlo
  - Variational inference

# Uncertainty Quantification

- Probabilistic graphical models
  - Bayesian network
  - Directed and undirected graphical models
  - Hidden Markov model and EM
  - Exponential families
- Further topics
  - MCMC and sequential Monte Carlo
  - Variational inference

# Probability

**定义 4.8** 设  $\mathcal{F}$  是由  $\Omega$  的一些子集组成的集合(这种由集合组成的集合一般叫做集合系),  $P=P(\cdot)$  是  $\mathcal{F}$  上有定义的实值函数. 若定义域  $\mathcal{F}$  和函数  $P$  满足下列条件:

$$(1) \Omega \in \mathcal{F}; \quad (4.1)$$

$$(2) \text{若 } A \in \mathcal{F}, \text{ 则 } A^c = \Omega - A \in \mathcal{F}; \quad (4.2)$$

$$(3) \text{若 } A_n \in \mathcal{F} (\text{一切 } n \geq 1), \text{ 则 } \bigcup_{n=1}^{\infty} A_n \in \mathcal{F}; \quad (4.3)$$

$$(4) P(A) \geq 0 (\text{一切 } A \in \mathcal{F}); \quad (4.4)$$

$$(5) P(\Omega) = 1; \quad (4.5)$$

(6) 若  $A_n \in \mathcal{F}$  (一切  $n \geq 1$ ), 且两两不相交, 就有

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} P(A_n), \quad (4.6)$$

则称  $P$  是  $\mathcal{F}$  上的概率测度(简称概率),  $P(A)$  为  $A$  的概率(也称  $A$  发生的概率). (见注)

# Conditional probability

Two children, the older child is a girl.

What is the probability that both children are girls?

Two children, at least one of them is a boy.

What is the probability that both children are boys?

# Conditional probability

Two children, the older child is a girl.

What is the probability that both children are girls?

Two children, at least one of them is a boy.

What is the probability that both children are boys?

$$\frac{\#\{girl-girl\}}{\#\{boy-girl, girl-girl\}} = \frac{1}{2}$$

# Conditional probability

Two children, the older child is a girl.

What is the probability that both children are girls?

Two children, at least one of them is a boy.

What is the probability that both children are boys?

$$\frac{\#\{girl-girl\}}{\#\{boy-girl, girl-girl\}} = \frac{1}{2}$$

$$\frac{\#\{boy-boy\}}{\#\{boy-girl, girl-boy, boy-boy\}} = \frac{1}{3}$$

# Conditional probability

**定义 5.1** 设  $A$  和  $B$  都是条件  $S$  下的事件, 则称条件  $S$  已经实现且  $B$  已发生的情形下  $A$  发生的概率为  $B$  发生的条件下  $A$  的条件概率, 并记为  $P(A|B)$ .

一般情形下, 怎样计算条件概率呢? 有下列计算公式:

$$P(A|B) = \frac{P(AB)}{P(B)} \quad (\text{当 } P(B) \neq 0).$$

$$P(AB) = P(B) \cdot P(A|B).$$

# Product formula

**定理 5.1(一般乘法公式)** 设  $A_1, A_2, \dots, A_n$  是  $n$  个事件 ( $n \geq 2$ ), 满足  $P(A_1A_2\cdots A_{n-1}) \neq 0$ , 则

$$P(A_1 \cdots A_n) = P(A_1)P(A_2 | A_1)P(A_3 | A_1A_2) \cdots \\ \cdot P(A_n | A_1 \cdots A_{n-1}). \quad (5.3)$$

**定义 5.2** 设  $A$  和  $B$  都是条件  $S$  下的随机事件, 若满足

$$P(AB) = P(A)P(B),$$

则称  $A$  与  $B$  相互独立. (“相互独立”简称“独立”.)

# Law of total probability

**定理 6.1(全概公式)<sup>①</sup>** 如果事件组  $B_1, B_2, \dots, B_n (n \geq 2)$  满足下列条件：

- (1)  $B_1, B_2, \dots, B_n$  两两不相容且  $P(B_i) > 0 (i = 1, 2, \dots, n)$ ；
- (2)  $B_1 \cup B_2 \cup \dots \cup B_n$  是必然事件，

则对任何事件  $A$  皆有

$$P(A) = \sum_{i=1}^n P(B_i)P(A|B_i). \quad (6.2)$$

# Bayesian formula

**定理 6.2(逆概公式)** 如果事件组  $B_1, \dots, B_n (n \geq 2)$  满足下列条件：

(1)  $B_1, \dots, B_n$  两两不相容且  $P(B_i) > 0 (i = 1, 2, \dots, n)$ ；

(2)  $\bigcup_{i=1}^n B_i$  是必然事件，

则对任一事件  $A$ , 只要  $P(A) > 0$ , 就有

$$P(B_k | A) = \frac{P(B_k)P(A | B_k)}{\sum_{i=1}^n P(B_i)P(A | B_i)}. \quad (6.6)$$

# Bayesian formula

The "Posterior"

$$P(H|E)$$

The "Likelihood"

$$\frac{P(E|H)P(H)}{P(E)}$$

The "Prior"

**H** represents our **hypothesis** we wish to examine.

**E** represents **evidence** for that hypothesis (presumably evidence we obtain through experiment and/or observation).

**P(H/E)** is called the **posterior**. This value gives us the probability of our hypothesis given our obtained evidence. In a sense, it shows the degree that your hypothesis is confirmed from your evidence. This is the unknown end value that we would wish to obtain.

# Bayesian formula

The "Posterior"

$$P(H|E)$$

The "Likelihood"

$$\frac{P(E|H)P(H)}{P(E)}$$

The "Prior"

The "Marginal"

$P(H)$  is called the **prior**. This is how probable our hypothesis was before any evidence is taken into account. The prior is infamous for being hard to define properly. How could you know any prior information about a hypothesis you haven't even tested? Without going into too much detail (because one could go on indefinitely about Bayesian priors), just know that this value has to be assumed. There are methods that can pick out better prior probabilities amongst worse ones, such as examining the results of previously conducted experiments. The prior that is developed from past experiments, however, will still be arbitrary to a degree.

$P(E)$  is called the marginal likelihood/probability or often just the **marginal**. This value represents the probability of our evidence given all possible hypotheses. Most of the time, we do not actually need a direct value of the marginal, as it can often be reduced out of the equation.

# 随机硬币

两枚硬币，一枚均匀，另一枚正面概率为 $3/4$ 。现随机取一枚硬币，连续投掷3次都是正面！取到的是均匀硬币的概率多大？

# 随机硬币

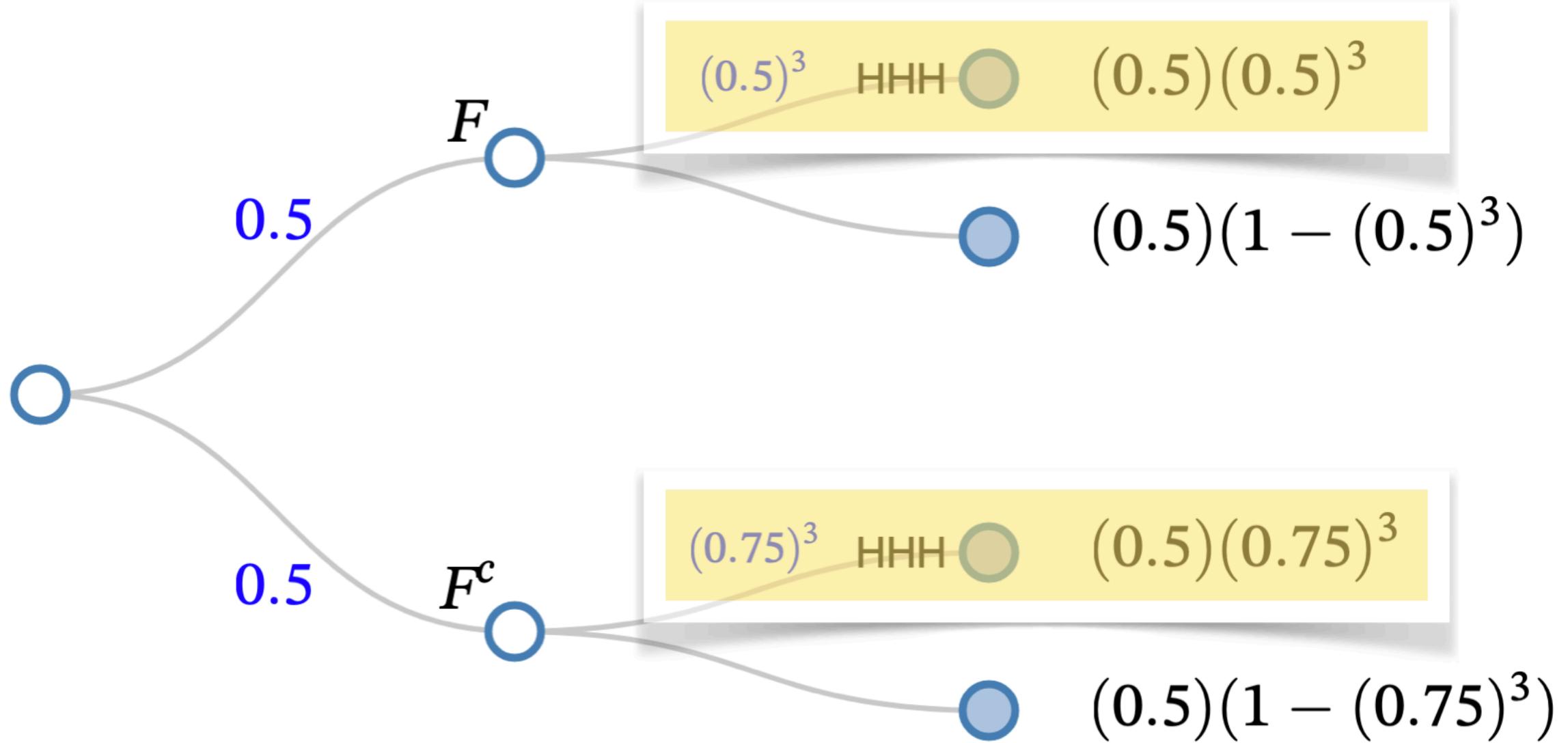
两枚硬币，一枚均匀，另一枚正面概率为 $3/4$ 。现随机取一枚硬币，连续投掷3次都是正面！取到的是均匀硬币的概率多大？

$$A = \{\text{出现三次正面}\}$$

$$F = \{\text{取到均匀硬币}\}$$

$$P(F | A) = ?$$

$$P(F | A) = \frac{P(A | F)P(F)}{P(A | F)P(F) + P(A | F^c)P(F^c)}$$



$$P(F \mid A) = \frac{(0.5)(0.5)^3}{(0.5)(0.5)^3 + (0.5)(0.75)^3}$$

$$P(F \mid A) = \frac{P(A \mid F)P(F)}{P(A \mid F)P(F) + P(A \mid F^c)P(F^c)}$$

## Example: medical diagnosis

As an example of how to use this rule, consider the following medical diagnosis problem. Suppose you are a woman in your 40s, and you decide to have a medical test for breast cancer called a **mammogram**. If the test is positive, what is the probability you have cancer? That obviously depends on how reliable the test is. Suppose you are told the test has a **sensitivity** of 80%, which means, if you have cancer, the test will be positive with probability 0.8. In other words,

$$p(x = 1|y = 1) = 0.8 \tag{2.8}$$

where  $x = 1$  is the event the mammogram is positive, and  $y = 1$  is the event you have breast cancer. Many people conclude they are therefore 80% likely to have cancer. But this is false! It ignores the prior probability of having breast cancer, which fortunately is quite low:

$$p(y = 1) = 0.004 \tag{2.9}$$

Ignoring this prior is called the **base rate fallacy**. We also need to take into account the fact that the test may be a **false positive** or **false alarm**. Unfortunately, such false positives are quite likely (with current screening technology):

$$p(x = 1|y = 0) = 0.1 \tag{2.10}$$

Combining these three terms using Bayes rule, we can compute the correct answer as follows:

$$p(y = 1|x = 1) = \frac{p(x = 1|y = 1)p(y = 1)}{p(x = 1|y = 1)p(y = 1) + p(x = 1|y = 0)p(y = 0)} \tag{2.11}$$

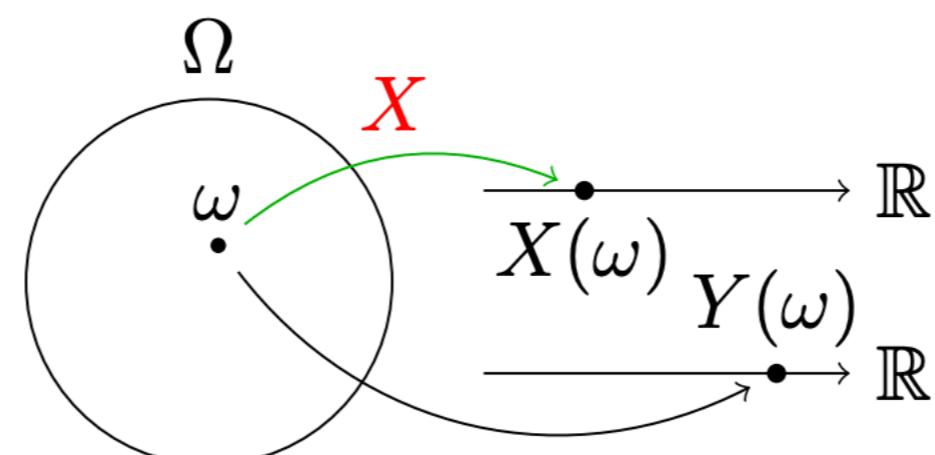
$$= \frac{0.8 \times 0.004}{0.8 \times 0.004 + 0.1 \times 0.996} = 0.031 \tag{2.12}$$

where  $p(y = 0) = 1 - p(y = 1) = 0.996$ . In other words, if you test positive, you only have about a 3% chance of actually having breast cancer!<sup>3</sup>

# Random Variables

**定义 1.1(随机变量的直观描述)** 如果条件  $S$  实现下的情况可以用一个数量  $X$  来描述,  $X$  究竟等于多少不能预先确定, 而随着条件  $S$  下的结果不同而可能变化, 但对任何实数  $c$ , 事件“ $X$  取值不超过  $c$ ”是有概率的, 则把这样一种变量  $X$  叫做随机变量.

**定义 1.1'(随机变量的数学描述)** 如果条件  $S$  下所有可能的结果组成集合  $\Omega = \{\omega\}$ ,  $X = X(\omega)$  是  $\Omega$  上有定义的实值函数, 而且对任何实数  $c$ , 事件  $\{\omega: X(\omega) \leq c\}$  是有概率的, 则称  $X$  是随机变量. (更完全的定义见本章 § 2.4.)



# Conditional Independence

We say  $X$  and  $Y$  are **unconditionally independent** or **marginally independent**, denoted  $X \perp Y$ , if we can represent the joint as the product of the two marginals (see Figure 2.2), i.e.,

$$X \perp Y \iff p(X, Y) = p(X)p(Y) \quad (2.14)$$

Unfortunately, unconditional independence is rare, because most variables can influence most other variables. However, usually this influence is mediated via other variables rather than being direct. We therefore say  $X$  and  $Y$  are **conditionally independent** (CI) given  $Z$  iff the conditional joint can be written as a product of conditional marginals:

$$X \perp Y|Z \iff p(X, Y|Z) = p(X|Z)p(Y|Z) \quad (2.15)$$

# Conditional Independence

Another characterization of CI is this:

**Theorem 2.2.1.**  $X \perp Y|Z$  iff there exist function  $g$  and  $h$  such that

$$p(x, y|z) = g(x, z)h(y, z)$$

for all  $x, y, z$  such that  $p(z) > 0$ .

prove that  $p(x|y, z) = p(x|z)$  when  $x \perp y|z$

# Mass and density function

$$P(x, y) = P(x|y)P(y) : \text{ product rule}$$

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} : \text{ Bayes theorem}$$

$$P(x) = \sum_{y \in \mathcal{Y}} P(x, y) = \sum_{y \in \mathcal{Y}} P(x|y)P(y)$$

# Expectation and variance

$$\mathbb{E}[x] := \int_{-\infty}^{+\infty} xp(x) dx : \text{ mean value,}$$

$$\sigma_x^2 := \int_{-\infty}^{+\infty} (x - \mathbb{E}[x])^2 p(x) dx : \text{ variance}$$

$$\mathbb{E}[f(x)] := \int_{-\infty}^{+\infty} f(x) p(x) dx.$$

$$\mathbb{E}_{x,y}[f(x, y)] = \mathbb{E}_x \left[ \mathbb{E}_{y|x}[f(x, y)] \right]$$

Given two random variables  $x, y$ , their *covariance* is defined as

$$\text{cov}(x, y) := \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])],$$

and their *correlation* as

$$r_{xy} := \mathbb{E}[xy] = \text{cov}(x, y) + \mathbb{E}[x]\mathbb{E}[y].$$

# Expectation and variance

Given two random variables  $x, y$ , their *covariance* is defined as

$$\text{cov}(x, y) := \mathbb{E}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])], \quad (2.28)$$

and their *correlation* as

$$r_{xy} := \mathbb{E}[xy] = \text{cov}(x, y) + \mathbb{E}[x]\mathbb{E}[y]. \quad (2.29)$$

A *random vector* is a collection of random variables,  $\mathbf{x} = [x_1, \dots, x_l]^T$ , and  $p(\mathbf{x})$  is the joint PDF (probability mass for discrete variables),

$$p(\mathbf{x}) = p(x_1, \dots, x_l). \quad (2.30)$$

The *covariance matrix* of a random vector  $\mathbf{x}$  is defined as

$$\boxed{\text{Cov}(\mathbf{x}) := \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}])(\mathbf{x} - \mathbb{E}[\mathbf{x}])^T] : \text{ covariance matrix}}, \quad (2.31)$$

or

$$\text{Cov}(\mathbf{x}) = \begin{bmatrix} \text{cov}(x_1, x_1) & \dots & \text{cov}(x_1, x_l) \\ \vdots & \ddots & \vdots \\ \text{cov}(x_l, x_1) & \dots & \text{cov}(x_l, x_l) \end{bmatrix}. \quad (2.32)$$

# Expectation and variance

Another symbol that will be used to denote the covariance matrix is  $\Sigma_x$ . Similarly, the *correlation matrix* of a random vector  $\mathbf{x}$  is defined as

$$R_x := \mathbb{E} [\mathbf{x}\mathbf{x}^T] : \text{ correlation matrix,} \quad (2.33)$$

or

$$\begin{aligned} R_x &= \begin{bmatrix} \mathbb{E}[x_1, x_1] & \dots & \mathbb{E}[x_1, x_l] \\ \vdots & \ddots & \vdots \\ \mathbb{E}[x_l, x_1] & \dots & \mathbb{E}[x_l, x_l] \end{bmatrix} \\ &= \text{Cov}(\mathbf{x}) + \mathbb{E}[\mathbf{x}] \mathbb{E}[\mathbf{x}^T]. \end{aligned} \quad (2.34)$$

$$\mathbb{E} [xx^T] = \Sigma + \mu\mu^T \quad \text{Cov} [\mathbf{A}\mathbf{x} + \mathbf{b}] = \mathbf{A}\text{Cov} [\mathbf{x}] \mathbf{A}^T$$

# Gaussian distribution

write  $x \sim \mathcal{N}(\mu, \sigma^2)$  or  $\mathcal{N}(x|\mu, \sigma^2)$ , if

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

It can be shown that the corresponding mean and variance are

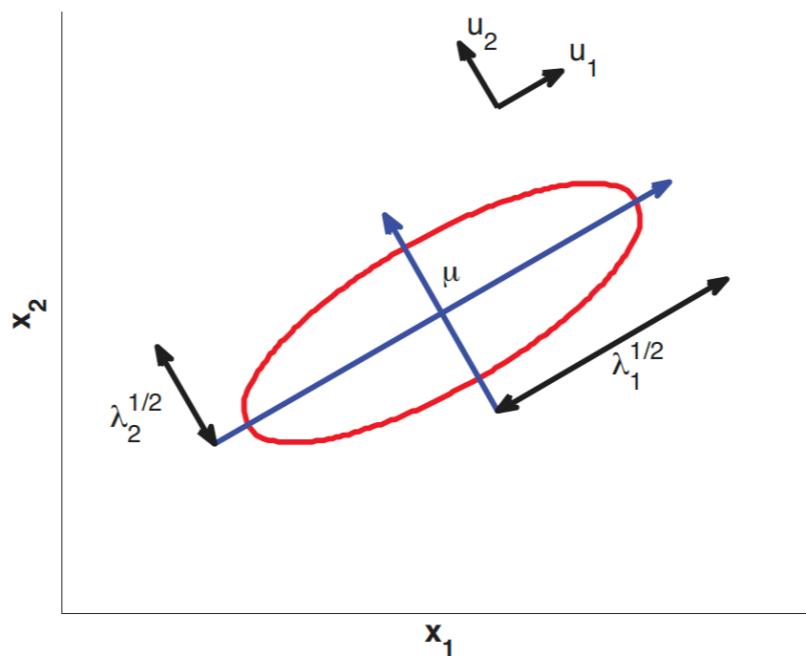
$$\mathbb{E}[x] = \mu \quad \text{and} \quad \sigma_x^2 = \sigma^2.$$

$$\int_{-\infty}^{+\infty} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) dx = \sqrt{2\pi}\sigma.$$

# Gaussian distribution

*Gaussian* or *normal* distribution,  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with parameters  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ , which is defined as

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{l/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) : \text{ Gaussian}$$



Marginal independent if and only if marginals are uncorrelated  
Not true in general !

# Gaussian distribution

**eigendecomposition** of  $\Sigma$ . That is, we write  $\Sigma = \mathbf{U}\Lambda\mathbf{U}^T$ , where  $\mathbf{U}$  is an orthonormal matrix of eigenvectors satisfying  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$ , and  $\Lambda$  is a diagonal matrix of eigenvalues.

Using the eigendecomposition, we have that

$$\Sigma^{-1} = \mathbf{U}^{-T}\Lambda^{-1}\mathbf{U}^{-1} = \mathbf{U}\Lambda^{-1}\mathbf{U}^T = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \quad (4.2)$$

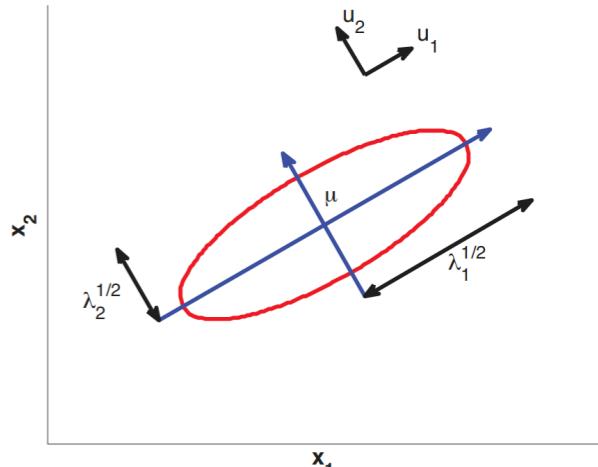
where  $\mathbf{u}_i$  is the  $i$ 'th column of  $\mathbf{U}$ , containing the  $i$ 'th eigenvector. Hence we can rewrite the Mahalanobis distance as follows:

$$(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^T \left( \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T \right) (\mathbf{x} - \boldsymbol{\mu}) \quad (4.3)$$

$$= \sum_{i=1}^D \frac{1}{\lambda_i} (\mathbf{x} - \boldsymbol{\mu})^T \mathbf{u}_i \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu}) = \sum_{i=1}^D \frac{y_i^2}{\lambda_i} \quad (4.4)$$

where  $y_i \triangleq \mathbf{u}_i^T (\mathbf{x} - \boldsymbol{\mu})$ . Recall that the equation for an ellipse in 2d is

$$\frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} = 1$$



# Parametric modeling - An overview

a *nonnegative* (loss) function,

$$\mathcal{L}(\cdot, \cdot) : \mathbb{R} \times \mathbb{R} \longmapsto [0, \infty),$$

and compute  $\theta_*$  so as to minimize the total loss, or as we say the *cost*, over all the data points, or

$$f(\cdot) := f_{\theta_*}(\cdot) : \theta_* = \arg \min_{\theta \in \mathcal{A}} J(\theta),$$

where

$$J(\theta) := \sum_{n=1}^N \mathcal{L}(y_n, f_\theta(x_n)),$$

In linear regression

$$\hat{y} = \hat{\theta}_0 + \hat{\theta}_1 x_1 + \cdots + \hat{\theta}_l x_l := \hat{\theta}^T x$$

$$J(\theta) = \sum_{n=1}^N (y_n - \theta^T x_n)^2$$